#### **RESEARCH ARTICLE**



# Fast and stable nonconvex constrained distributed optimization: the ELLADA algorithm

Wentao Tang<sup>1,2</sup> · Prodromos Daoutidis <sup>1</sup>

Received: 4 May 2020 / Revised: 29 November 2020 / Accepted: 29 November 2020 © The Author(s), under exclusive licence to Springer Science+Business Media, LLC part of Springer Nature 2021

#### **Abstract**

Distributed optimization using multiple computing agents in a localized and coordinated manner is a promising approach for solving large-scale optimization problems, e.g., those arising in model predictive control (MPC) of large-scale plants. However, a distributed optimization algorithm that is computationally efficient, globally convergent, amenable to nonconvex constraints remains an open problem. In this paper, we combine three important modifications to the classical alternating direction method of multipliers for distributed optimization. Specifically, (1) an extralayer architecture is adopted to accommodate nonconvexity and handle inequality constraints, (2) equality-constrained nonlinear programming (NLP) problems are allowed to be solved approximately, and (3) a modified Anderson acceleration is employed for reducing the number of iterations. Theoretical convergence of the proposed algorithm, named ELLADA, is established and its numerical performance is demonstrated on a large-scale NLP benchmark problem. Its application to distributed nonlinear MPC is also described and illustrated through a benchmark process system.

**Keywords** Distributed optimization  $\cdot$  Nonconvex optimization  $\cdot$  Model predictive control  $\cdot$  Acceleration

#### 1 Introduction

Distributed optimization (Boyd et al. 2011) refers to methods of performing optimization using a distributed architecture—the monolithic problem is first decomposed into several subproblems, each handled by a corresponding solver (agent), and

Published online: 03 January 2021

Present Address: Surface Operations, Projects and Technology, Shell Global Solutions (U.S.) Inc., Houston, TX 77082, USA



Prodromos Daoutidis daou001@umn.edu

Department of Chemical Engineering and Materials Science, University of Minnesota, Minneapolis, MN 55455, USA

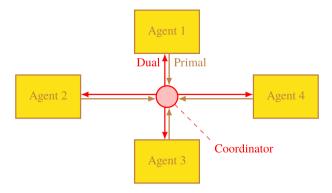


Fig. 1 Primal-dual distributed optimization

necessary information among the agents is communicated to coordinate the distributed computation. The alternating direction method of multipliers (ADMM), as the most classical and representative algorithm for distributed optimization, was proposed in the 1970s (Glowinski and Marroco 1975; Gabay and Mercier 1976) and has underdone significant development in the last decade. ADMM is an iterative primal—dual algorithm that is closely related to the method of multipliers (MM). In each iteration, the blocks of primal variables are alternately optimized, and the dual variables are updated according to the updated values of the primal variables. The underlying architecture of the algorithm is illustrated in Fig. 1, with two distinct groups of components—the distributed agents and the coordinator(s) collecting information about the primal and dual values from each other to solve their own subproblems.

Convergence is the most basic requirement of distributed optimization and also the central theme of a large amount of theoretical research on ADMM. For convex problems and nonconvex problems whose nonconvexity resides in the objective function, by connecting the ADMM algorithm with the monotone operator theory and Douglas-Rachford splitting, the Fejér monotonicity and monotonicity of the augmented Lagrangian become standard arguments of proving theoretical convergence guarantees (He and Yuan 2012; Nishihara et al. 2015; Hong et al. 2016; Hong and Luo 2017; Wang et al. 2019; Themelis and Patrinos 2020). However, nonconvex constraints appear much more difficult to handle. To guarantee convergence, Hours and Jones (2015) suggested dualizing and penalizing all nonconvex constraints, making them undifferentiated and tractable by ADMM; however, this alteration of the problem structure eliminates the option for distributed agents to use any subroutine other than MM. Houska et al. (2016) used a quadratic programming problem to decide the dual variables in the augmented Lagrangian as well as an extrapolation of primal updates; this algorithm, however, involves a central agent that extracts Hessian and gradient information of the subsystem models from the distributed agents in every iteration, and is thus essentially semi-centralized. Scutari et al. (2016) adopted feasibility-preserving convex approximations to approach the solution, which is applicable to problems without nonconvex equality constraints. We note that several



recent papers (Sun and Sun 2019; Jiang et al. 2019; Yang et al. 2020) proposed the idea of placing slack variables corresponding to the inter-subsystem constraints and forcing the decay to zero by tightening the penalty parameters of slack variables. This modification to the ADMM with slack variables and their penalties leads to a globally convergent *extra-layer* augmented Lagrangian-based algorithm with preserved agent-coordinator problem architecture.

Computational efficiency is also of critical importance for distributed optimization. The slothfulness of primal-dual algorithms typically arises from two issues. First, the subgradient (first-order) update of dual variables restricts the number of iterations to be of linear complexity (Hong and Luo 2017). For convex problems, momentum methods (Goldstein et al. 2014; Ouyang et al. 2015) or Krylov subspace methods (Zhang and White 2018) can be adopted for acceleration, which, however, do not directly extend to nonconvex problems. Under nonconvexity, it was only very recently realized that Anderson acceleration, a multi-secant technique for fixedpoint problems, can be generally used to accelerate the dual variables (Zhang et al. 2018, 2019; Fu et al. 2019). The second cause for the high computational cost of distributed optimization is the instruction on the distributed agents to fully solve their subproblems to high precision in each iteration. Such exhaustive efforts may be unnecessary since the dual information to be received from the coordinator will keep changing. For convex problems, it is possible to linearize the augmented Lagrangian and replace the distributed subproblems with Arrow-Hurwicz-Uzawa gradient flows (Dhingra et al. 2019). In the presence of nonconvexity of the objective functions, a dual perturbation technique to restore the convergence of the augmented Lagrangian was proposed (Hajinezhad and Hong 2019). It is yet unknown how to accommodate such gradient flows to nonconvex constraints. A different approach is to allow inexact solution of the subproblems with adaptively tightening tolerances (Eckstein and Yao 2017, 2018). Such an approximate ADMM algorithm allows a better balance between the primal and dual updates, and avoids wasteful computational steps inside the subroutines.

An important engineering application that demands both a convergent and efficient distributed optimization algorithm is model predictive control (MPC), which is an advanced control strategy widely adopted in process industries. In MPC, control decisions are made through solving an optimal control problem minimizing the cost associated with the predicted trajectory in a future horizon subject to the system dynamics and operational constraints (Rawlings et al. 2017). For large-scale systems, it is desirable to seek a decomposition, e.g., using community detection or network block structures (Daoutidis et al. 2018, 2019; Tang et al. 2018), and deploy distributed MPC strategies (Scattolini 2009; Christofides et al. 2013), which promises better performance than fully decentralized MPC by enabling coordination, while avoiding assembling and computing on a monolithic model as in centralized MPC. In general, distributed nonlinear MPC with subsystem interactions should be considered as a distributed optimization problem under nonconvex constraints. Whether the distributed optimization algorithm converges affects the quality of control decisions and hence the closed-loop control performance, whereas computational efficiency determines whether a distributed nonlinear MPC strategy can be practically implemented online.



Although ADMM for distributed MPC has been discussed (Farokhi et al. 2014; Mota et al. 2014), due to the lack of convergent and efficient nonconvex constrained distributed optimization algorithms, its application is so far limited mainly to linear systems (Giselsson et al. 2013; Wang and Ong 2017). The most common approach of distributed nonlinear MPC is to iterate the control inputs among the subsystems (in sequence or in parallel) (Stewart et al. 2010; Liu et al. 2010; Chen et al. 2012). The input iteration routine is typically either semi-decentralized by implicitly assuming that the subsystems interact only through inputs and considering state coupling as disturbances, or semi-centralized by using moving-horizon predictions based on the entire system, which, however, contradicts the fact that the subsystem models should be usually packaged inside the local agents rather than shared over the entire system. In a different vein, efficient centralized MPC algorithms have been proposed which exploit the underlying sparse patterns in the variable-constraint structure of the optimal control problems and decompose the linear algebraic operations involved in the optimization solver (Wang and Boyd 2009; Patterson and Rao 2014; Biegler and Thierry 2018). While these efficient centralized MPC algorithms improve computational efficiency without compromising the solution quality, distributed MPC is still advantageous in its capability of solving very large MPC problems on a localized subsystem basis without the need of engineering the solver details. Moreover, such structure exploitation can also be utilized by the solvers employed in distributed MPC. The crucial computational issue faced by distributed MPC is its iterative complexity and this will be addressed in this paper.

The purpose of this work is to develop a convergent and computationally efficient algorithm for distributed optimization under nonconvex constraints. Although the algorithm is in principle not restricted to any specific problem, we consider distributed nonlinear MPC as an application of special interest. Based on the above discussion, we identify the following modifications to the classical ADMM algorithm as the key to mitigating the challenges in convergence and computational complexity: (1) additional slack variables are placed on the constraints relating the distributed agents and the coordinator, (2) approximate optimization is performed in the distributed agents, and (3) the Anderson acceleration technique is adopted by the coordinator. We therefore combine and extend as appropriate these techniques into a new algorithm with a two-layer augmented Lagrangian-based architecture, in which the outer layer handles the slack variables as well as inequality constraints by using a barrier technique, and the inner layer performs approximate ADMM under an acceleration scheme. With guaranteed stability and elevated speed, to the best knowledge of the authors, the proposed algorithm is the first practical and generic algorithm of its kind for distributed nonlinear MPC with truly localized model information. We name this algorithm as ELLADA (standing for extra-layer augmented Lagrangianbased accelerated distributed approximate optimization).

The paper discusses the motivation, develops the ELLADA algorithm and establishes its theoretical properties, and illustrates its application through case studies. The remainder of this paper is organized as follows. In Sect. 2, we first review the classical ADMM and its modified versions. Then we derive our ELLADA algorithm in Sect. 3 with a trilogy pattern. First, a basic two-layer augmented Lagrangian-based algorithm (ELL) is introduced and its convergence is discussed. Then the approximate solution of



equality-constrained NLP problems and the Anderson acceleration scheme are incorporated to form the ELLA and ELLADA algorithms. In Sect. 4, a large-scale nonlinear optimization benchmark problem is used to demonstrate the numerical performance of the proposed ELLADA algorithm. The implementation of ELLADA on the distributed optimization problem involved in distributed nonlinear MPC is shown in Sect. 5, and a case study on a benchmark process is examined in Sect. 6. Conclusions and discussions are given in Sect. 7.

## 2 ADMM and its modifications

#### **2.1 ADMM**

The alterating direction method of multipliers is the most commonly used algorithm for distributed optimization under linear equality constraints (Boyd et al. 2011). Specifically, consider the following problem

$$\min_{x,\bar{x}} f(x) + g(\bar{x}) \quad \text{s.t.} \quad Ax + B\bar{x} = 0$$
 (1)

with two blocks of variables x and  $\bar{x}$ , where f and g are usually assumed to be convex. [The symbols in (1) are not related to the ones in Sect. 5.] The augmented Lagrangian for such a constrained optimization problem is

$$L(x, \bar{x}; y) = f(x) + g(\bar{x}) + y^{\mathsf{T}} (Ax + B\bar{x}) + \frac{\rho}{2} ||Ax + B\bar{x}||^2,$$
 (2)

in which y stands for the vector of dual variables (Lagrangian multipliers) and  $\rho > 0$  is called the penalty parameter. According to the duality theory, the optimal solution should be determined by a saddle point of the augmented Lagrangian:

$$\sup_{y} \min_{x,\bar{x}} L(x,\bar{x};y). \tag{3}$$

The classical method of multipliers (MM) deals with this saddle point problem with an iterative procedure, where the primal variables are optimized first and then the dual variables are updated with a subgradient ascent (Bertsekas 2016, Chapter 6):

$$(x^{k+1}, \bar{x}^{k+1}) = \arg\min_{x, \bar{x}} L(x, \bar{x}; y^k),$$
  
$$y^{k+1} = y^k + \rho (Ax^{k+1} + B\bar{x}^{k+1}),$$
 (4)

in which the superscript stands for the count of iterations. In a distributed context, x and  $\bar{x}$  usually can not be optimized simultaneously. ADMM is thus an approximation of MM that allows the optimization of x and  $\bar{x}$  to be performed separately, i.e.,



$$x^{k+1} = \arg\min_{x} L(x, \bar{x}^{k}; y^{k}),$$

$$\bar{x}^{k+1} = \arg\min_{\bar{x}} L(x^{k}, \bar{x}; y^{k}),$$

$$y^{k+1} = y^{k} + \rho (Ax^{k+1} + B\bar{x}^{k+1}).$$
(5)

Since the appearance of ADMM in 1970s (Glowinski and Marroco 1975; Gabay and Mercier 1976), there have been many works regarding its theoretical properties, extensions and applications. As we have mentioned in the Introduction, ADMM is known to have a linear convergence rate for convex problems. This does not change when the variables are constrained in convex sets. For example, if  $x \in \mathcal{X}$ , it suffices to modify the corresponding term f(x) in the objective function by adding an indicator function  $\mathbb{I}_{\mathcal{X}}(x)$  (equal to 0 if  $x \in \mathcal{X}$  and  $+\infty$  otherwise), which is still a convex function.

## 2.2 ADMM with approximate updates

Unless the objective terms f(x) and  $g(\bar{x})$  are of simple forms such as quadratic functions, the optimization of x and  $\bar{x}$  in (5) does not have an exact solution. Usually, iterative algorithms for nonlinear programming need to be called for the first two lines of (5), and always searching for a highly accurate solution in each ADMM iteration will result in an excessive computational cost. It is thus desirable to solve the optimization subproblems in ADMM inexactly when the dual variables are yet far from the optimum, i.e., to allow  $x^{k+1}$  and  $\bar{x}^{k+1}$  to be chosen such that

$$d_{v}^{k+1} \in \partial_{x} L(x^{k+1}, \bar{x}^{k}; y^{k}), \quad d_{\bar{v}}^{k+1} \in \partial_{\bar{v}} L(x^{k+1}, \bar{x}^{k+1}; y^{k}), \tag{6}$$

where  $\partial_x$  and  $\partial_{\bar{x}}$  represent the subgradients with respect to x and  $\bar{x}$ , respectively, and  $d_x$  and  $d_{\bar{x}}$  are not exactly 0 but only converging to 0 asymptotically. For example, one can assign externally a shrinking and summable sequence of absolute errors (Eckstein and Bertsekas 1992):

$$||d_x^k|| \le \epsilon_x^k, ||d_{\bar{x}}^k|| \le \epsilon_{\bar{x}}^k, \sum_{k=1}^{\infty} \epsilon_x^k < \infty, \sum_{k=1}^{\infty} \epsilon_{\bar{x}}^k < \infty,$$
 (7)

or a sequence of relative errors to the errors proportional to other variations in the algorithm (Eckstein and Yao 2017; Xie et al. 2017).

It was shown in Eckstein and Yao (2017) that a relative error criterion for terminating the iterations in subproblems, compared to other approximation criteria such as a summable absolute error sequence, better reduces the total number of subroutine iterations throughout the ADMM algorithm. Such a relative error criterion is a *constructive* one, rendered to guarantee the decrease of a quadratic distance between the intermediate solutions  $(x^k, \bar{x}^k, y^k)$  and the optimum  $(x^*, \bar{x}^*, y^*)$ . In the context of distributed optimization problems under nonconvex constraints, since the convergence proof is established on a different basis from the quadratic distance, the construction of such a criterion must be reconsidered. We will address this issue in Sect. 3.2.



## 2.3 Anderson acceleration

Linear convergence of the classical ADMM is essentially the result of subgradient dual update, which uses the information of only the first-order derivatives with respect to the dual variables:  $\partial_y L = Ax + B\bar{x}$ . The idea of creating a *quadratically convergent* algorithm using only first-order derivatives originates back from Nesterov's approach of solving convex optimization problems, which performs iterations based on a linear extrapolation of the previous two iterations instead of the current solution alone (Nesterov 1983). Such a *momentum* method can be used to accelerate the ADMM algorithm, which can be seen as iterations over the second block of primal variables  $\bar{x}$  and the dual variables y (Goldstein et al. 2014). However, such a momentum is inappropriate for nonconvex problems, since the behavior of the extrapolated point can not be well controlled by a bound on the curvature of the objective function.

Therefore, we resort to a different type of technique—Anderson acceleration, which was proposed in Anderson (1965) first and later "rediscovered" in the field of chemical physics (Pulay 1980). Generally speaking, Anderson acceleration is used to solve the fixed-point iteration problem

$$w = h_0(w) \tag{8}$$

for some vector w and non-expansive mapping  $h_0$  (satisfying  $\|h_0(w) - h_0(w')\| \le \|w - w'\|$  for any w and w'). Different from the simple Krasnoselskii–Mann iteration  $w^{k+1} = \kappa w^k + (1-\kappa)h_0(w^k)$  ( $\kappa \in (0,1)$ ), Anderson acceleration takes a quasi-Newton approach, which aims at a nearly quadratic convergence rate (Fang and Saad 2009). Specifically, in each iteration k, the results from the previous m iterations are recalled from memory to form the matrix of secants in w and  $h(w) = w - h_0(w)$ :

$$\Delta_{w}^{k} = \left[\delta_{w}^{k-m} \dots \delta_{w}^{k-1}\right], \delta_{w}^{k'} = w^{k'+1} - w^{k'}, k' = k - m, \dots, k - 1;$$

$$\Delta_{b}^{k} = \left[\delta_{b}^{k-m} \dots \delta_{b}^{k-1}\right], \delta_{b}^{k'} = h(w^{k'+1}) - h(w^{k'}), k' = k - m, \dots, k - 1.$$
(9)

An estimated Jacobian is given by

$$H_k = I + \left(\Delta_h^k - \Delta_w^k\right) \left(\Delta_w^{k\top} \Delta_w^k\right)^{-1} \Delta_w^{k\top},\tag{10}$$

or

$$H_k^{-1} = I + \left(\Delta_w^k - \Delta_h^k\right) \left(\Delta_w^{k\mathsf{T}} \Delta_h^k\right)^{-1} \Delta_w^{k\mathsf{T}},\tag{11}$$

which minimizes the Frobenius norm of  $B_k - I$  subject to  $B_k \Delta_w^k = \Delta_h^k$ . Then the quasi-Newton iteration  $w^{k+1} = w^k - H_k^{-1} h^k$  leads to a weighted sum of the previous m function values:

<sup>&</sup>lt;sup>1</sup> There are two different types of Anderson acceleration. Here we focus on Type I, which was found to have better performance (Fang and Saad 2009) and was improved in Zhang et al. (2018).



$$w^{k+1} = \sum_{m'=0}^{m} \alpha_{m'}^{k} h_0(x^{k-m+m'})$$
 (12)

where the weights  $\{\alpha_{m'}^k\}_{m'=0}^m$  are specified by

$$\alpha_{m'}^{k} = \begin{cases} s_{0}^{k}, & m' = 0\\ s_{m'}^{k} - s_{m'-1}^{k}, & m' = 1, \dots, m-1\\ 1 - s_{m-1}^{k}, & m' = m \end{cases}$$
 (13)

with  $s_{m'}^k$  being the m'th component  $s^k$ :

$$s^k = (\Delta_w^{k\top} \Delta_h^k)^{-1} \Delta_w^{k\top} h^k. \tag{14}$$

Anderson acceleration (12) may not always be convergent, although local convergence was studied in some special cases (Toth and Kelley 2015). Recently, a globally convergent modification of Anderson acceleration was proposed in Zhang et al. (2018), where regularization, restarting, and safeguarding measures are taken to ensure the well-conditioning of the  $\Delta_w^k$  matrix, boundedness of the inverse Jacobian estimate (11), and acceleration only in a safety region, respectively.

The relevance of Anderson acceleration to ADMM lies in that the ADMM algorithm (5) can be seen as fixed-point iterations  $(\bar{x}^k, y^k) \to (\bar{x}^{k+1}, y^{k+1}), k = 0, 1, 2, ...$  (Zhang et al. 2019), which is the same idea underlying the ADMM with Nesterov acceleration. For problems with nonconvex constraints, the iteration mapping h is not necessarily non-expansive, and hence one can not directly establish the convergence of Anderson acceleration with the original techniques used in Zhang et al. (2018). We will address this issue in Sect. 3.3.

#### 2.4 ADMM under nonconvex constraints

The presence of nonconvexity largely increases the difficulty of distributed optimization. Most of the work in nonconvex ADMM considers problems with nonconvex objective function with bounded Hessian eigenvalues or the Kurdyka–Łojasiewicz property assumptions, under which a convergence rate of  $\mathcal{O}\left(1/\sqrt{k}\right)$  (slower than that of convex ADMM,  $\mathcal{O}(1/k)$ ) was established (Li and Pong 2015; Hong et al. 2016; Wang et al. 2019). However, for many distributed optimization problems, e.g., the distributed MPC of nonlinear processes, there exist nonconvex constraints on the variables, which is intrinsically non-equivalent to the problems with nonconvex objective functions. For our problem of interest, the relevant works are scarce.

Here we introduce the algorithm of Sun and Sun (2019) for (1) under nonconvex constraints  $x \in \mathcal{X}$  and  $\bar{x} \in \bar{\mathcal{X}}$ , reformulated with slack variables z:

$$\min_{\substack{x \ \bar{x} \ z}} f(x) + g(\bar{x}) \quad \text{s.t. } Ax + B\bar{x} + z = 0, z = 0, \ x \in \mathcal{X}, \ \bar{x} \in \bar{\mathcal{X}}.$$
 (15)

The augmented Lagrangian is now written as



$$L(x, \bar{x}, z; y, \lambda, \rho, \beta) = f(x) + g(\bar{x}) + \mathbb{I}_{\chi}(x) + \mathbb{I}_{\bar{\chi}}(\bar{x}) + y^{\mathsf{T}}(Ax + B\bar{x} + z) + \frac{\rho}{2} ||Ax + B\bar{x} + z||^{2} + \lambda^{\mathsf{T}}z + \frac{\beta}{2} ||z||^{2}.$$
(16)

The algorithm is a two-layer one, where each outer iteration (indexed by k) contains a series of inner iterations (indexed by r). In the inner iterations, the classical ADMM algorithm is used to update x,  $\bar{x}$ , z and y in sequence, while keeping  $\lambda$  and  $\beta$  unchanged:

$$x^{k,r+1} = \arg\min_{x} L(x, \bar{x}^{k,r}, z^{k,r}; y^{k,r}, \lambda^{k}, \rho^{k}, \beta^{k})$$

$$= \arg\min_{x \in \mathcal{X}} f(x) + \frac{\rho^{k}}{2} \left\| Ax + B\bar{x}^{k,r} + z^{k,r} + \frac{y^{k,r}}{\rho^{k}} \right\|^{2}$$

$$\bar{x}^{k,r+1} = \arg\min_{\bar{x}} L(x^{k,r+1}, \bar{x}, z^{k,r}; y^{k,r}, \lambda^{k}, \rho^{k}, \beta^{k})$$

$$= \arg\min_{\bar{x} \in \bar{\mathcal{X}}} g(\bar{x}) + \frac{\rho^{k}}{2} \left\| Ax^{k,r+1} + B\bar{x} + z^{k,r} + \frac{y^{k,r}}{\rho^{k}} \right\|^{2}$$

$$z^{k,r+1} = \arg\min_{z} L(x^{k,r+1}, \bar{x}^{k,r+1}, z; y^{k,r}, \lambda^{k}, \rho^{k}, \beta^{k})$$

$$= -\frac{\rho^{k}}{\rho^{k} + \beta^{k}} \left( Ax^{k,r+1} + B\bar{x}^{k,r+1} + \frac{y^{k,r}}{\rho^{k}} \right) - \frac{1}{\rho^{k} + \beta^{k}} \lambda^{k}$$

$$y^{k,r+1} = y^{k,r} + \rho^{k} (Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r+1})$$

$$(17)$$

Under mild assumptions, in the presence of slack variables z, it was proved (Sun and Sun 2019) that if one chooses  $\rho^k = 2\beta^k$ , then the inner iterations converge to the set of stationary points  $(x^k, \bar{x}^k, z^k, y^k)$  of the relaxed problem

$$\min_{x,\bar{x},z} f(x) + g(\bar{x}) + \lambda^{k\top} z + \frac{\beta^k}{2} ||z||^2$$
s.t.  $Ax + B\bar{x} + z = 0, x \in \mathcal{X}, \bar{x} \in \bar{\mathcal{X}}.$  (18)

Then in the outer iterations, the dual variables  $\lambda^k$  are updated. To enforce the convergence of the slack variables to zero, the corresponding penalty  $\beta^k$  is amplified by a ratio  $\gamma > 1$  if the returned  $z^k$  from the inner iterations does not decay enough from the previous outer iteration  $z^{k-1}(||z^k|| > \omega||z^{k-1}||, \omega \in (0,1))$ . The outer iteration is written as

$$\lambda^{k+1} = \prod_{\left[\underline{\lambda},\overline{\lambda}\right]} \left(\lambda^k + \beta^k z^k\right), \quad \beta^{k+1} = \begin{cases} \gamma \beta^k, & \|z^k\| > \omega \|z^{k-1}\| \\ \beta^k, & \|z^k\| \le \omega \|z^{k-1}\| \end{cases}$$
(19)

in which the projection  $\Pi$  onto a predefined compact hypercube  $\left[\underline{\lambda},\overline{\lambda}\right]$  is used to guarantee the boundedness of the dual variables and hence the augmented Lagrangian L. If the augmented Lagrangian L remains bounded despite the increase of the penalty parameters  $\rho^k$  and  $\beta^k$ , the algorithm converges to a stationary point of the original problem (1). The iterative complexity of such an algorithm to reach an  $\epsilon$ -approximate stationary point is  $\mathcal{O}(\epsilon^{-4} \ln(\epsilon^{-1}))$ .



In the next section, building on the algorithm of Sun and Sun (2019) that guarantees the convergence of distributed optimization under nonconvex constraints, we propose a new algorithm that integrates into it the ideas of approximate ADMM and Anderson acceleration, aiming at improving the computational efficiency.

## 3 Proposed algorithm

## 3.1 Basic algorithm and its convergence

Consider an optimization problem in the following form:

$$\min_{x,\bar{x}} \quad f(x) + g(\bar{x})$$
s.t.  $Ax + B\bar{x} = 0$ , (20)
$$x \in \mathcal{X} = \{x | \phi(x) < 0, \psi(x) = 0\}, \bar{x} \in \bar{\mathcal{X}}$$

or equivalently with slack variables

$$\min_{x,\bar{x},z} f(x) + g(\bar{x})$$
s.t.  $Ax + B\bar{x} + z = 0, z = 0,$  (21)
$$x \in \mathcal{X} = \{x | \phi(x) \le 0, \psi(x) = 0\}, \bar{x} \in \bar{\mathcal{X}}.$$

We make the following assumptions.

**Assumption 1** Assume that f is lower bounded, i.e.,  $\exists \underline{f}$  such that  $\forall x \in \mathcal{X}, f(x) \geq \underline{f}$ .

**Assumption 2** Function *g* is lower bounded.

Our basic algorithm (Algorithm 1) for (21) is slightly modified from the procedure of Sun and Sun (2019), which considered the case where g(x) = 0 and  $\bar{\mathcal{X}}$  is a hypercube. The algorithm uses an inner loop of ADMM iterations and an outer loop of MM with possibly amplifying penalty parameters. The inner iterations are terminated when the following criterion is met

$$\epsilon_{1}^{k} \geq \epsilon_{1}^{k,r} := \left\| \rho^{k} A^{T} (B \bar{x}^{k,r+1} + z^{k,r+1} - B \bar{x}^{k,r} - z^{k,r}) \right\|, 
\epsilon_{2}^{k} \geq \epsilon_{2}^{k,r} := \left\| \rho^{k} B^{T} (z^{k,r+1} - z^{k,r}) \right\|, 
\epsilon_{3}^{k} \geq \epsilon_{3}^{k,r} := \left\| A x^{k,r+1} + B \bar{x}^{k,r+1} + z^{k,r+1} \right\|.$$
(22)



```
Set: Bound of dual variables \left|\underline{\lambda}, \overline{\lambda}\right|, shrinking ratio of slack variables \omega \in [0, 1), amplifying ratio
       of penalty parameter \gamma > 1, diminishing outer iteration tolerances \{\epsilon_1^k, \epsilon_2^k, \epsilon_3^k\}_{k=1}^{\infty} \downarrow 0,
       terminating tolerances \epsilon_1, \epsilon_2, \epsilon_3 > 0;
2 Initialization: Starting points x^0, \bar{x}^0, z^0, duals and bounds \lambda^1 \in [\underline{\lambda}, \overline{\lambda}], penalty \beta^1 > 0;
3 Outer iteration count k \leftarrow 0;
    while stationarity criterion (25) is not met do
            \rho^k = 2\beta^k:
            Inner iteration count r \leftarrow 0;
            Initialization: x^{k,0}, \bar{x}^{k,0}, z^{k,0}, y^{k,0} satisfying \lambda^k + \beta^k z^{k,0} + y^{k,0} = 0;
 7
            while stopping criterion (22) is not met do
 8
                    x^{k,r+1} = \arg\min_{x \in X} f(x) + \frac{\rho^k}{2} \left\| Ax + B\bar{x}^{k,r} + z^{k,r} + \frac{y^{k,r}}{2} \right\|^2;
 9
                   \bar{x}^{k,r+1} = \arg\min_{\bar{x} \in \bar{X}} g(\bar{x}) + \frac{\rho^k}{2} \left\| Ax^{k,r+1} + B\bar{x} + z^{k,r} + \frac{y^{k,r}}{\rho^k} \right\|^2;
10
                   z^{k,r+1} = -\frac{\rho^k}{\rho^k + \beta^k} \left( A x^{k,r+1} + B \bar{x}^{k,r+1} + \frac{y^{k,r}}{\rho^k} \right) - \frac{1}{\rho^k + \beta^k} \lambda^k;
11
                   y^{k,r+1} = y^{k,r} + \rho^k \left( Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r+1} \right);
12
13
            (x^{k+1}, \bar{x}^{k+1}, z^{k+1}, y^{k+1}) \leftarrow (x^{k,r}, \bar{x}^{k,r}, z^{k,r}, y^{k,r});
            Update \lambda^{k+1} and \beta^{k+1} by (19);
16
17 end
```

**Lemma 1** (Convergence of inner iterations) Suppose that Assumptions 1 and 2 hold. When  $\rho^k = 2\beta^k$ , the inner iterations are terminated at a finite r when (22) is met and the point  $(x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r+1})$  satisfies the following conditions

$$\begin{aligned} d_1^k &\in \partial f \left( x^{k,r+1} \right) + \mathcal{N}_{\mathcal{X}} \left( x^{k,r+1} \right) + A^{\mathsf{T}} y^{k,r+1} \\ d_2^k &\in \partial g \left( \bar{x}^{k,r+1} \right) + \mathcal{N}_{\bar{\mathcal{X}}} \left( \bar{x}^{k,r+1} \right) + B^{\mathsf{T}} y^{k,r+1} \\ 0 &= \lambda^k + \beta^k z^{k,r+1} + y^{k,r+1}, \quad d_3^k = A x^{k,r+1} + B \bar{x}^{k,r+1} + z^{k,r+1} \end{aligned} \tag{23}$$

for some  $d_1^k$ ,  $d_2^k$  and  $d_3^k$  satisfying  $\|d_1^k\| \le \epsilon_1^k$ ,  $\|d_2^k\| \le \epsilon_2^k$  and  $\|d_3^k\| \le \epsilon_3^k$ , respectively.  $\mathcal{N}_{\mathcal{X}}(x)$  ( $\mathcal{N}_{\bar{\mathcal{X}}}(\bar{x})$ ) refers to the normal cone to the set  $\mathcal{X}$  ( $\bar{\mathcal{X}}$ ) at point x ( $\bar{x}$ ):

$$\mathcal{N}_{\mathcal{X}}(x) = \left\{ v \middle| v^{\mathsf{T}} \left( x' - x \right) \le 0, \forall x' \in \mathcal{X} \right\}. \tag{24}$$

The proof of the above lemma is given in "Appendix 1" using the augmented Lagrangian (16) as a decreasing Lyapunov function throughout the inner iterations (Hong et al. 2016; Hong and Luo 2017), which gives the convergence of the inner iterations. It is apparent that if  $e_1^k$ ,  $e_2^k$ ,  $e_3^k$  are all equal to 0, (23) is the Karush–Kuhn–Tucker optimality condition of the relaxed problem (18) (Rockafellar and Wets 1998). To establish the convergence of outer iterations, we need to make the following assumption to restrict the upper level of the augmented Lagrangian.



**Assumption 3** The augmented Lagrangians are uniformly upper bounded at initialization of all inner iterations, i.e., there exists  $\overline{L} \ge L(x^{k,0}, \overline{x}^{k,0}, z^{k,0}, y^{k,0}, \lambda^k, \rho^k, \beta^k)$  for all k.

The above assumption is actually a "warm start" requirement. Suppose that we have a feasible solution  $(x^0, \bar{x}^0)$  to the original problem (20), then we can always choose  $x^{k,0} = x^0$ ,  $\bar{x}^{k,0} = \bar{x}^0$ ,  $z^{k,0} = 0$ ,  $y^{k,0} = -\lambda^k$  to guarantee an  $\bar{L} = f(x^0) + g(\bar{x}^0)$ .

**Lemma 2** (Convergence of outer iterations) Suppose that Assumptions 1, 2 and 3 hold. Then for any  $\epsilon_1$ ,  $\epsilon_2$ , and  $\epsilon_3 > 0$ , within a finite number of outer iterations k, Algorithm 1 finds an approximate stationary point  $(x^{k+1}, \bar{x}^{k+1}, z^{k+1}, y^{k+1})$  of (20), satisfying

$$d_{1} \in \partial f(x^{k+1}) + \mathcal{N}_{\mathcal{X}}(x^{k+1}) + A^{\mathsf{T}} y^{k+1}$$

$$d_{2} \in \partial g(\bar{x}^{k+1}) + \mathcal{N}_{\bar{\mathcal{X}}}(\bar{x}^{k+1}) + B^{\mathsf{T}} y^{k+1}$$

$$d_{3} = A x^{k+1} + B \bar{x}^{k+1}$$
(25)

for some  $d_1, d_2, d_3$  satisfying  $||d_j|| \le \epsilon_j$ , j = 1, 2, 3.

See "Appendix 2" for a proof. In addition to the convergence, we can also establish a theoretical complexity. Previously in Sun and Sun (2019), it was shown that to reach an  $\epsilon$ -approximate stationary point satisfying (25) with  $\epsilon_1, \epsilon_2, \epsilon_3 = \epsilon > 0$ , the total number of inner iterations needed is of the order  $\mathcal{O}(\epsilon^{-4} \ln{(1/\epsilon)})$ . Here, we show that by appropriately choosing the way that the tolerances  $(\epsilon_1^k, \epsilon_2^k, \epsilon_3^k)$  shrink, the iteration complexity can be provably reduced anywhere in  $(\mathcal{O}(\epsilon^{-2}), \mathcal{O}(\epsilon^{-4})]$ , for which a proof is given in "Appendix 3".

**Lemma 3** (Complexity of the basic algorithm) Suppose that Assumptions 1, 2 and 3 hold. For some constant  $\vartheta \in (0, \omega]$ , choose  $\varepsilon_1^k \sim \mathcal{O}(\vartheta^k)$ ,  $\varepsilon_2^k \sim \mathcal{O}(\vartheta^k)$ , and  $\varepsilon_3^k \sim \mathcal{O}((\vartheta/\beta)^k)$ . Then each outer iteration k requires  $R^k \sim \mathcal{O}((\vartheta\omega)^{-2k})$  inner iterations. Hence, for the Algorithm 1 to reach an  $\epsilon$ -approximate stationary point, the total iterations needed is  $R \sim \mathcal{O}(\varepsilon^{-2(1+\varsigma)})$ , where  $\varsigma = \log_\vartheta \omega \in (0, 1]$ .

#### 3.2 Approximate algorithm

We note that the basic algorithm requires undesirable complete minimization of x and  $\bar{x}$  in each inner iteration (Lines 9–10, Algorithm 1). For simplicity, we assume that such a minimization oracle<sup>2</sup> exists for  $\bar{x}$ .

We use the word "oracle" with its typical meaning in mathematics and computer science. An oracle refers to an ad hoc numerical or computational procedure, regarded as a black box mechanism, to generate the needed results as its outputs based on some input information.



**Assumption 4** The minimization of the augmented Lagrangian with respect to  $\bar{x}$  (Line 10, Algorithm 1) admits a unique explicit solution

$$\bar{x}^{k,r+1} = G(B, Ax^{k,r+1} + z^{k,r} + y^{k,r}/\rho^k, \rho^k). \tag{26}$$

Let us also assume that the problem has a smoothness property as follows.

**Assumption 5** Functions f,  $\phi$  and  $\psi$  are continuously differentiable, and  $\mathcal{X}$  has a nonempty interior.

Under this smoothness assumption, the KKT condition for *x*-minimization is written as the following equalities with  $\mu \ge 0$  and  $\nu$  representing the Lagrangian dual variables corresponding to the inequalities  $\phi(x) \le 0$  and  $\psi(x) = 0$ , respectively

$$0 = \nabla f(x^{k,r+1}) + \rho^{k} A^{\top} (Ax^{k,r+1} + B\bar{x}^{k,r} + z^{k,r} + y^{k,r}/\rho^{k})$$

$$+ \sum_{c=1}^{C_{\phi}} \mu_{c} \nabla \phi_{c}(x^{k,r+1}) + \sum_{c=1}^{C_{\psi}} \nu_{c} \nabla \psi_{c}(x^{k,r+1})$$

$$0 = \mu_{c} \phi_{c}(x^{k,r+1}), c = 1, \dots, C_{\phi}, \quad 0 = \psi_{c}(x^{k,r+1}), c = 1, \dots, C_{\psi}.$$

$$(27)$$

Line 9 of Algorithm 1 is thus to solve the above equations for  $x^{k,r+1}$ . This can be achieved through an interior point algorithm, which employs double-layer iterations to find the solution. In the outer iteration, a barrier technique is used to convert the inequality constraints into an additional term in the objective; the optima (or stationary points) of the resulting barrier problems converge to true optima (stationary points) as the barrier parameter converges to 0. In the inner iteration, a proper search method is used to obtain the optimum of the barrier problem. Since both the interior point algorithm and the ELL algorithm 1 have a double-layer structure, using an interior point algorithm as the NLP solver in ELL results in four layers of iterations. For simplification, we consider *matching the outer layers and inner layers of the NLP algorithm and ELL, respectively. Specifically, this implies the following two modifications of the ELL algorithm*.

To incorporate the outer iterations of interior point optimization into the outer iterations of ELL, the function f(x) is appended with a barrier term  $-b^k \sum_{c=1}^{C_{\phi}} \ln \left( -\phi_c(x) \right)$ , where the  $b^k$  is the barrier parameter varying with the outer iteration k and decaying to 0 as k increases. Hence a "barrier augmented Lagrangian" can be defined as

$$L_b = L - b \sum_{c=1}^{C_{\phi}} \ln \left( -\phi_c(x) \right). \tag{28}$$

Based on the arguments in the previous subsection, if the *x*-optimization step returns a  $x^{k,r+1}$  minimizing  $L_b$  with respect to x, then the inner iterations result in the descent of  $L_{b^k}$ , which implies the satisfaction of conditions (23), with f modified by the barrier function. Obviously, if Assumption 3 holds for L, then it also holds for  $L_{b^k}$  when  $\mathcal{X}$  has a nonempty interior. It follows that the outer iterations can find



an approximate stationary point of the original problem with the decay of barrier parameters  $b^k$ .

With the inequality constraints handled by the barrier term, the inner ADMM iterations only deals with equality-constrained NLP problems. To merge the iterative routine to solve equality-constrained NLP problems (Wächter and Biegler 2005) into the ADMM iterations, we allow to perform only *a proper amount of searching steps* instead of the entire equality-constrained NLP in each inner iteration, so that the solution to the equality-constrained NLP problem can be approached asymptotically throughout the inner iterations. For this purpose, we assume that we have at hand a solver that can find an approximate solution of the equality-constrained NLP with any prespecified tolerances of violations to the KKT conditions.

#### **Assumption 6** Assume that for any equality-constrained smooth NLP problem

$$\min_{x} \chi(x) \quad \text{s.t.} \quad \psi(x) = 0 \tag{29}$$

a solver that guarantees the convergence to any approximate stationary point of the above problem with a lower objective function is available. That is, starting from any initial point  $x^0$ , for any tolerances  $\epsilon_4$ ,  $\epsilon_5 > 0$ , within a finite number of searches the solver finds a point (x, v) satisfying

$$d_4 = \nabla \chi(x) + \sum_{c=1}^{C_{\psi}} v_c \nabla \psi_c(x), \quad d_{5c} = \psi_c(x), c = 1, \dots, C_{\psi}.$$
 (30)

for some  $||d_4|| \le \epsilon_4$ ,  $||d_5|| \le \epsilon_5$ , and  $f(x) \le f(x^0)$ . Such an approximate solution is denoted as  $F(x^0; \chi, \psi, \epsilon_4, \epsilon_5)$ .

The above approximate NLP solution oracle is realizable by NLP solvers where the tolerances of the KKT conditions are allowed to be specified by the user, e.g., the IPOPT solver (Wächter and Biegler 2006). Under Assumption 6, the *x*-update step on Line 9 of Algorithm 1 is replaced by an approximate NLP solution

$$x^{k,r+1} = F(x^{k,r}; \chi^{k,r}, \psi, \epsilon_4^{k,r}, \epsilon_5^{k,r}), \tag{31}$$

where the objective function in the current iteration is the part of barrier augmented Lagrangian  $L_{bk}$  that is related to x with the indicator function  $\mathbb{I}_{x}(x)$  excluded:

$$\chi^{k,r}(x) = f(x) - b_k \sum_{c=1}^{C_{\phi}} \ln\left(-\phi_c(x)\right) + \frac{\rho^k}{2} \left\| Ax + B\bar{x}^{k,r} + z^{k,r} + \frac{y^{k,r}}{\rho^k} \right\|^2.$$
 (32)

This approximate algorithm with inexact x-minimization is summarized as Algorithm 2. The inner iterations are performed until  $\epsilon_4^{k,r}$  and  $\epsilon_5^{k,r}$  are lower than  $\epsilon_4^k$  and  $\epsilon_5^k$ , respectively, and (22) holds. The outer iterations are terminated when  $\epsilon_4^k \leq \epsilon_4$ ,  $\epsilon_5^k \leq \epsilon_5$ , the barrier parameter is sufficiently small  $b^k \leq \epsilon_6$ , and (25) holds.



```
Set: Bound of dual variables |\underline{\lambda}, \overline{\lambda}|, shrinking ratio of slack variables \omega \in [0, 1), amplifying ratio
       of penalty parameter \gamma > 1, diminishing outer iteration tolerances \left\{ \epsilon_{1,2,3,4,5,6}^k \right\}_{k=1}^{\infty} \downarrow 0,
       diminishing barrier parameters \{b^k\}_{k=1}^{\infty} \downarrow 0, terminating tolerances \epsilon_{1,2,3,4,5,6} > 0;
 2 Initialization: Starting points x^0, \bar{x}^0, z^0, duals and bounds \lambda^1 \in [\underline{\lambda}, \overline{\lambda}], penalty \beta^1 > 0;
    Outer iteration count k \leftarrow 0;
     while \epsilon_4^k \ge \epsilon_4 or \epsilon_5^k \ge \epsilon_5 or b^k \ge \epsilon_6 or stationarity criterion (25) is not met do
            Set: Diminishing tolerances \left\{ \epsilon_4^{k,r}, \epsilon_5^{k,r} \right\}_{r=1}^{\infty} \downarrow 0;
            Let \rho^k = 2\beta^k;
 6
            Inner iteration count r \leftarrow 0;
            Initialization: x^{k,0}, \bar{x}^{k,0}, z^{k,0}, y^{k,0} satisfying \lambda^k + \beta^k z^{k,0} + y^{k,0} = 0;
            while \epsilon_4^{k,r} \ge \epsilon_4^k or \epsilon_5^{k,r} \ge \epsilon_5^k or stopping criterion (22) is not met do
                   x^{k,r+1} = F\left(x^{k,r}; \chi^{k,r}, \psi, \epsilon_4^{k,r}, \epsilon_5^{k,r}\right), where \chi^{k,r} is given by (32);
10
                   \bar{x}^{r+1} = G\left(B, Ax^{r+1} + z^{k,r} + y^{k,r}/\rho^k, \rho^k\right), where G is given by (26);
11
                   z^{k,r+1} = -\frac{\rho^k}{\rho^k + \beta^k} \left( A x^{k,r+1} + B \bar{x}^{k,r+1} + \frac{y^{k,r}}{\rho^k} \right) - \frac{1}{\rho^k + \beta^k} \lambda^k;
                  y^{k,r+1} = y^{k,r} + \rho^k \left( Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r+1} \right);
13
            (x^{k+1}, \bar{x}^{k+1}, z^{k+1}, y^{k+1}) \leftarrow (x^{k,r}, \bar{x}^{k,r}, z^{k,r}, y^{k,r});
            Update \lambda^{k+1} and \beta^{k+1} by (19);
18 end
```

**Lemma 4** (Convergence of the approximate algorithm) Suppose that Assumptions 1–6 hold. For any outer iteration k, given any positive tolerances  $\{\epsilon_1^k, \ldots, \epsilon_5^k\}$ , within a finite number of inner iterations r, the obtained solution satisfies

$$d_{1}^{k} + d_{4}^{k} = \nabla f(x^{k,r+1}) + \sum_{c=1}^{C_{\phi}} \mu_{c}^{k,r+1} \nabla \phi_{c}(x^{k,r+1})$$

$$+ \sum_{c=1}^{C_{\psi}} v_{c}^{k,r+1} \nabla \psi_{c}(x^{k,r+1}) + A^{\mathsf{T}} y^{k,r+1}$$

$$d_{2}^{k} \in \partial g(\bar{x}^{k,r+1}) + \mathcal{N}_{\bar{x}}(\bar{x}^{k,r+1}) + B^{\mathsf{T}} y^{k,r+1},$$

$$0 = \lambda^{k} + \beta^{k} z^{k,r+1} + y^{k,r+1}$$

$$d_{3}^{k} = A x^{k,r+1} + B \bar{x}^{k,r+1} + z^{k,r+1}, \quad d_{5}^{k} = \psi(x^{k,r+1}),$$

$$-b^{k} = \mu_{c}^{k,r+1} \phi_{c}(x^{k,r+1}), c = 1, \dots, C_{\phi}.$$

$$(33)$$

for some  $d_1^k, \ldots, d_5^k$  with  $\|d_1^k\| \le \epsilon_1^k, \ldots, \|d_5^k\| \le \epsilon_5^k$ . Then, suppose that the outer iteration tolerances  $\{\epsilon_1^k, \ldots, \epsilon_5^k\}$  and barrier parameters  $b^k$  are diminishing with increasing k, given any terminating tolerances  $\epsilon_1, \ldots, \epsilon_6 > 0$ , within a finite number



of outer iterations, Algorithm 2 finds a point  $(x^{k+1}, \bar{x}^{k+1}, z^{k+1}, y^{k+1}, \mu^{k+1}, \nu^{k+1})$  satisfying

$$d_{1} + d_{4} = \nabla f(x^{k+1}) + \sum_{c=1}^{C_{\phi}} \mu_{c}^{k+1} \nabla \phi_{c}(x^{k+1}) + \sum_{c=1}^{C_{\psi}} v_{c}^{k+1} \nabla \psi_{c}(x^{k+1}) + A^{\top} y^{k+1}$$

$$d_{2} \in \partial g(\bar{x}^{k+1}) + \mathcal{N}_{\bar{x}}(\bar{x}^{k+1}) + B^{\top} y^{k+1}, \quad 0 = \lambda^{k} + \beta^{k} z^{k+1} + y^{k+1}$$

$$d_{3} = A x^{k+1} + B \bar{x}^{k+1}, \quad d_{5} = \psi(x^{k+1}),$$

$$-d_{6} = \mu_{c}^{k+1} \phi_{c}(x^{k+1}), c = 1, \dots, C_{\phi}.$$

$$(34)$$

for some  $d_1, \ldots, d_6$  with  $\|d_j\| \le \epsilon_j, j = 1, \ldots, 5, d_6 \in (0, \epsilon_6]$ .

## 3.3 Accelerated algorithm

In the ELLA algorithm, let us make the following assumption regarding our choice of tolerances  $e_4^{k,r}$  and  $e_5^{k,r}$ .

**Assumption 7** Suppose that we choose a continuous and strictly monotonically increasing function  $\pi\colon [0,\infty)\to [0,\infty)$  with  $\pi(0)=0$  such that  $\epsilon_4^{k,r}=\pi\left(\epsilon_4^{k,r}\right)$ , and choose  $\epsilon_4^{k,r+1}$  with an analogous function of  $\left\|\rho^kA^\top\left(B\bar{\chi}^{k,r+1}-B\bar{\chi}^{k,r}+z^{k,r+1}-z^{k,r}\right)\right\|$  when the resulting value is strictly smaller than the previous tolerance  $\epsilon_4^{k,r}$  but not smaller the ultimate one  $\epsilon_4^k$ .

The choice of function  $\pi$  to relate the stationarity tolerance and equality tolerance in NLP subroutine is aimed at balancing the effort to reduce both errors. The choice of  $\epsilon_4^{k,r+1}$  is based on the rationale that after the rth inner iteration, the obtained solution  $x^{k,r+1}$  satisfies the stationarity condition within a tolerance of  $\epsilon_4^{k,r}$ , and after the update of  $\bar{x}$ , z and y variables, the violation to the stationarity condition is bounded by  $\epsilon_4^{k,r} + \|\rho^k A^\top (B\bar{x}^{k,r+1} - B\bar{x}^{k,r} + z^{k,r+1} - z^{k,r})\|$ . Therefore,  $\epsilon_4^{k,r}$  should be balanced with the second term, which, however, is realizable only after the  $\bar{x}$ - and z-updates after the x-update and hence assigned to  $\epsilon_4^{k,r+1}$ .

We note from Algorithm 2 that under Assumption 7, each inner iteration r is a mapping from  $(x^{k,r}, \bar{x}^{k,r}, z^{k,r+1}, y^{k,r+1}, e_4^{k,r})$  to  $(x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r+1}, y^{k,r+1}, e_4^{k,r+1})$ . In fact, despite the dependence of the latter variables on  $x^{k,r}$  and  $e_4^{k,r}$ , such dependence can be ignored in the sense that the descent of the barrier augmented Lagrangian  $L_{b^k}$  will always guide the sequence of intermediate solutions towards the set of  $e_4^k$  -approximate stationary points of the relaxed barrier problem. It follows that under the approximate algorithm, the sequence  $\{(\bar{x}^{k,r}, z^{k,r})\}_{r=1}^{\infty}$  will converge to a fixed point, and the convergence of  $\{y^{k,r}\}$  accompanies the convergence of  $\{z^{k,r}\}$ . It is thus clear that we may resort to Anderson acceleration introduced in Sect. 2.3 by denoting  $w = (\bar{x}, z)$ , the iteration as a mapping  $h_0$ , and  $h(w) = w - h_0(w)$ , and collecting at the rth inner iteration the following multi-secant information about the previous m inner iterations:



$$\Delta_w^{k,r} = \left[\delta_w^{k,r-m} \dots \delta_w^{k,r-1}\right], \Delta_h^{k,r} = \left[\delta_h^{k,r-m} \dots \delta_h^{k,r-1}\right],\tag{35}$$

where 
$$\delta_h^{r-m'} = w^{k,r-m'+1} - w^{k,r-m'}$$
 and  $\delta_h^{r-m'} = h((w^{k,r-m'+1}) - h(w^{k,r-m'}), m' = m-1, \dots, 0.$ 

However, the possibility that  $\Delta_w^k$  may not be of full rank and  $H_{k,r}$  may be singular requires certain modifications to the original acceleration scheme. The following regulation technique was used in Zhang et al. (2018). To ensure the invertibility of  $H_{k,r}$ , the  $\delta_h$  vector is perturbed to

$$\tilde{\delta}_{h}^{k,r-m+m'} = \left(1 - \theta_{k,r}^{m'}\right) \delta_{h}^{k,r-m+m'} + \theta_{k,r}^{m'} \delta_{w}^{k,r-m+m'}, \tag{36}$$

where the perturbation magnitude  $\theta_{k,r}^{m'}$  is determined by

$$\theta_{k,r}^{m'} = \varphi \left( \frac{\left(\hat{\delta}_{w}^{k,r-m+m'}\right)^{\mathsf{T}} \left(H_{k,r}^{m'}\right)^{-1} \delta_{h}^{k,r-m+m'}}{\left\|\hat{\delta}_{w}^{k,r-m+m'}\right\|^{2}}; \eta_{\theta} \right). \tag{37}$$

with regularization hyperparameter  $\eta_{\theta} \in (0, 1)$ . The function  $\varphi(\theta; \eta)$  is defined by

$$\varphi(\theta;\eta) = \begin{cases} (\eta \operatorname{sign}\theta - \theta)/(1 - \theta), & |\theta| \le \eta \\ 0, & |\theta| > \eta \end{cases}$$
 (38)

With regulation,  $H_{k,r}^{-1}$  is induced from  $\left(H_{k,r}^{0}\right)^{-1}=I$  according to

$$\left(H_{k,r}^{m'+1}\right)^{-1} = \left(H_{k,r}^{m'}\right)^{-1}$$

$$+ \frac{\left(\delta_{w}^{k,r-m+m'} - \left(H_{k,r}^{m'}\right)^{-1} \tilde{\delta}_{h}^{k,r-m+m'}\right) \left(\hat{\delta}_{w}^{k,r-m+m'}\right)^{\top} \left(H_{k,r}^{m'}\right)^{-1}}{\left(\hat{\delta}_{w}^{k,r-m+m'}\right)^{\top} \left(H_{k,r}^{m'}\right)^{-1} \tilde{\delta}_{h}^{k,r-m+m'}}$$

$$(39)$$

for  $m'=0,\ldots,m-1$  with  $H^m_{k,r}=H_{k,r}$ . To avoid the rank deficiency  $\Delta_w$ , a restart checking strategy is used, where the memory is cleared when the Gram-Schmidt orthogonalization becomes ill conditioned  $(\left\|\hat{\delta}^{k,r}_w\right\|<\eta_w\left\|\delta^{k,r}_w\right\|$  for some  $\eta_w\in(0,1)$ ) or the memory exceeds a maximum M; otherwise the memory is allowed to grow. Hence the Anderson acceleration is well-conditioned.

**Lemma 5** (Well-conditioning of Anderson acceleration, Zhang et al. (2018)) *Using the regularization and restart checking techniques, it is guaranteed that* 

$$\left\| H_{k,r}^{-1} \right\|_{2} \le \theta^{-M} \left[ 3 \left( 1 + \theta + \eta_{w} \right)^{M} \eta_{w}^{-N} - 2 \right]^{N-1} < +\infty$$
 (40)

where M is the maximum number of steps in the memory and N is the dimension of w.



A well-conditioned Anderson acceleration is not yet sufficient to guarantee the convergence. Hence we employ a safeguarding technique modified from Zhang et al. (2018) which aims at suppressing a too large increase in the barrier augmented Lagrangian by rejecting such acceleration steps. When Anderson acceleration suggests an update from  $w = (\bar{x}, z)$  to  $\tilde{w} = (\tilde{x}, \tilde{z})$  under the current value of Ax, the resulting Lagrangian increase black, if positive, must not exceed an upper bound:

$$\Delta L_{\text{max}} = \tilde{L}_0 \frac{\eta_L}{\left(R_+ + 1\right)^{1+\sigma}} = \left(\beta^k \left\| B \bar{x}^{k,1} - B \bar{x}^{k,0} \right\|^2 + \frac{\beta^k}{2} \left\| z^{k,1} - z^{k,0} \right\|^2\right) \frac{\eta_L}{\left(R_+ + 1\right)^{1+\sigma}}$$
(41)

where  $\tilde{L}_0$  is the expected Lagrangian decrease after the first non-accelerated iteration after initialization according to Lemma 1, used as a scale for the change in the barrier augmented Lagrangian,  $\eta_L$ ,  $\sigma > 0$  are hyperparameters, and  $R_+$  is the number of already accepted acceleration steps. With safeguarding, it can be guaranteed that the barrier augmented Lagrangian always stays bounded, since  $\sum_{R_+=0}^{\infty} \left(R_+ + 1\right)^{-(1+\sigma)} < +\infty$ . We also require that the acceleration should not lead to a drastic change in w:

$$\|\Delta w\|_{\max}^2 = \frac{\tilde{L}_0}{\beta^k} \frac{\eta_{\bar{w}}}{\sqrt{1 + R_+}},$$
 (42)

where  $\eta_{\bar{w}} > 0$  is a hyperparameter.  $1/\sqrt{1+R_+}$  reflects an expected change according to the plain ADMM iteration, which is used to suppress disproportionate large deviations due to Anderson acceleration.

Finally, the accelerated algorithm using the Anderson acceleration technique for fixed-point iteration of  $(\bar{x}, z)$  is summarized as Algorithm 3. This is our final ELLADA algorithm, whose convergence can now be guaranteed by the following lemma, the proof of which is given in "Appendix 4".

**Lemma 6** (Convergence under Anderson acceleration) Suppose that Assumptions 1–7 hold. Under regulated and safe-guarded Anderson acceleration, Algorithm 3 finds within a finite number of inner iterations r a point satisfying (33). The convergence of outer iterations to an approximate stationary point satisfying (34) hence follows.

Summarizing the conclusions of all the previous lemmas in this section, we have arrived at the following theorem.

#### **Theorem 1** *Suppose that the following assumptions hold:*

- 1. Function f is lower bounded on X;
- 2. Function g is convex and lower bounded on  $\bar{\mathcal{X}}$ ;



- 3. Initialization of outer iterations allows a uniform upper bound of the augmented Lagrangian, e.g., a feasible solution is known a priori;
- 4. Minimization of  $g(\bar{x}) + \frac{\rho}{2} ||B\bar{x} + v||^2$  with respect to  $\bar{x}$  allows an oracle  $G(B, v, \rho)$  returning a unique solution for any v of appropriate dimension and  $\rho > 0$ ;
- 5. Functions f,  $\phi$ , and  $\psi$  are continuously differentiable, and the constraints  $(\phi, \psi)$  are strictly feasible;
- 6. There exists a solver for equality-constrained NLP to any specified tolerances of KKT conditions.

Then given any tolerances  $\epsilon_1, \dots, \epsilon_6 > 0$ , the ELLADA algorithm (Algorithm 3) gives an  $(\epsilon_1, \dots, \epsilon_6)$ -approximate KKT point satisfying the conditions (34).



```
1 Set: Dual bounds \left[\underline{\lambda}, \overline{\lambda}\right], outer iteration parameters \omega \in [0, 1), \gamma > 1, \left\{\epsilon_{1,2,3,4,6}^k\right\}_{k=1}^{\infty} \downarrow 0,
         \{b^k\}_{k=1}^{\infty} \downarrow 0, final tolerances \epsilon_{1,2,3,4,6} > 0, function \pi, acceleration parameters \theta \in (0,1),
         \sigma > 0, \eta_{\epsilon} > 0, \eta_{w} \in (0, 1), \eta_{L} > 0, \eta_{\tilde{w}} > 0, M \in \mathbb{N}. Let \epsilon_{5} = \pi(\epsilon_{4});
2 Initialization: Starting points x^0, \bar{x}^0, z^0, \lambda^1 \in [\underline{\lambda}, \overline{\lambda}], penalty parameter \beta^1 > 0, \epsilon_5^0 = \pi \left(\epsilon_4^0\right);
     Outer iteration count k \leftarrow 0;
      while \epsilon_4^k \ge \epsilon_4 or \epsilon_5^k \ge \epsilon_5 or b^k \ge \epsilon_6 or stationarity criterion (25) is not met do
               Set: Initial tolerances \epsilon_4^{k,0}, \epsilon_5^{k,0} = \pi(\epsilon_4^{k,0}), penalty \rho^k = 2\beta^k, Jacobian estimate H_{k,0}^{-1} = I;
               Inner iteration count r \leftarrow 0, count of accelerated steps R_+^k = 0, memory length m \leftarrow 0;
 6
               Initialization: x^{k,0}, \bar{x}^{k,0}, z^{k,0}, y^{k,0} satisfying \lambda^k + \beta^k z^{k,0} + y^{k,0} = 0;
 7
               while \epsilon_4^{k,r} \ge \epsilon_4^k or \epsilon_5^{k,r} \ge \epsilon_5^k or stopping criterion (22) is not met do
                        x^{k,r+1} = F\left(x^{k,r}; \chi^{k,r}, \psi, \epsilon_4^{k,r}, \epsilon_5^{k,r}\right), where \chi^{k,r} is given by (32);
 q
                        \bar{x}^{k,r+1} = G(B, Ax^{r+1} + z^{k,r} + y^{k,r}/\rho^k, \rho^k), where G is given by (26);
10
                        z^{k,r+1} = -\frac{1}{\rho^k} \left(Ax^{k,r+1} + B\bar{x}^{k,r+1} + \frac{y^{k,r'}}{\rho^k}\right) - \frac{1}{\rho^k + \beta^k} \lambda^k;
11
                        y^{k,r+1} = y^{k,r} + \rho^k \left( A x^{k,r+1} + B \bar{x}^{k,r+1} + z^{k,r+1} \right);
12
                        \tilde{\mathbf{v}}^{k,r} = -\lambda^k - \beta^k \tilde{\mathbf{z}}^{k,r}
13
                        \tilde{x}^{k,r+1} = F\left(x^{k,r}; \tilde{\chi}^{k,r}, \psi, \epsilon_4^{k,r}, \epsilon_5^{k,r}\right), \text{ with } \tilde{\chi} \text{ in (32) with } \bar{x}, z, y \text{ replaced by } \tilde{x}, \tilde{z}, \tilde{y};
14
                        \tilde{x}^{k,r+1} = G\left(B, A\tilde{x}^{r+1} + \tilde{z}^{k,r} + \tilde{y}^{k,r}/\rho^k, \rho^k\right);
15
                        \tilde{z}^{k,r+1} = -\frac{\rho^k}{\rho^k + \beta^k} \left( A \tilde{x}^{k,r+1} + B \tilde{\bar{x}}^{k,r+1} + \frac{\tilde{y}^{k,r}}{\rho^k} \right) - \frac{1}{\rho^k + \beta^k} \lambda^k;
16
17
                                 \tilde{w}^{k,1} \leftarrow (\bar{x}^{k,1}, z^{k,1}), and calculate \tilde{L}_0 by (41);
18
                                  \delta_w^{k,r-1} = \tilde{w}^{k,r} - w^{k,r-1}, \, \delta_h^{k,r} = \tilde{w}_-^{k,r} - \tilde{w}^{k,r+1} - w^{k,r-1} + w^{k,r}, \, m \leftarrow m+1;
20
                                  \hat{\delta}_{w}^{k,r-1} = \delta_{w}^{k,r-1} - \sum_{m'=2}^{m} \frac{\left(\hat{\delta}_{w}^{k,r-m'}\right)^{\top} \delta_{w}^{k,r-1}}{\left\|\hat{\delta}_{w}^{k,r-m'}\right\|^{2}} \hat{\delta}_{w}^{k,r-m'};
                                 if m = M + 1 or \left\| \frac{\delta_w^{k,r-1}}{\delta_w^{k,r-1}} \right\| < \eta_w then m \leftarrow 0, \delta_w^{k,r-1} \leftarrow \delta_w^{k,r-1}, and H_{k,r-1}^{-1} \leftarrow I;
22
                                 Obtain \hat{\delta}_{h}^{k,r-1} by (36) with m' = m-1 and \theta_{k,r-1} = \varphi\left(\frac{\left(\hat{\delta}_{w}^{k,r-1}\right)^{l} H_{k,r-1}^{-1} \delta_{h}^{k,r-1}}{\|\hat{\delta}_{k,r}^{k,r-1}\|^{2}}; \theta\right);
23
                                 \text{Update } H_{k,r}^{-1} = H_{k,r-1}^{-1} + \frac{\left(\delta_{w}^{k,r-1} - H_{k,r-1}^{-1} \hat{\delta}_{h}^{k,r-1}\right) \left(\hat{\delta}_{w}^{k,r-1}\right)^{\mathsf{T}} H_{k,r-1}^{-1}}{\left(\hat{\delta}_{w}^{k,r-1}\right)^{\mathsf{T}} H_{k,r-1}^{-1} \hat{\delta}_{w}^{k,r-1}}, \text{ and suggest}
24
                                    \tilde{w}^{k,r+1} = w^{k,r} - H_{k,r}^{-1} \left( w^{k,r} - w^{k,r+1} \right);
                                 if criteria (41) and (42) hold then w^{k,r+1} \leftarrow \tilde{w}^{k,r+1}, y^{k,r+1} \leftarrow -\lambda^k - \beta^k \tilde{z}^{k,r+1};
25
                        end
26
                        \epsilon_4^{k,r+1} = \left\| \rho^k A^\top \left( B \bar{x}^{k,r+1} - B \bar{x}^{k,r} + z^{k,r+1} - z^{k,r} \right) \right\|, \ \epsilon_5^{k,r+1} = \pi \left( \epsilon_4^{k,r+1} \right);
27
28
29
               (x^{k+1}, \bar{x}^{k+1}, z^{k+1}, y^{k+1}) \leftarrow (x^{k,r}, \bar{x}^{k,r}, z^{k,r}, y^{k,r});
30
               Update \lambda^{k+1} and \beta^{k+1} by (19);
31
               k \leftarrow k + 1;
32
33 end
```



#### 3.4 Parameter tuning

We have so far proved the theoretical convergence of the proposed ELLADA algorithm (Algorithm 3). It is not surprising to expect that such an algorithm that combines multiple techniques aiming at improving the computational performance will require some parameter tuning effort, although these parameters do not alter the ultimate convergence. Three main groups of parameters are involved here.

The first group of parameters is related to the outer iterations, namely the threshold of slack variable decay  $\omega$  and the triggered penalty amplifying ratio  $\gamma$ . Although the bounds of the dual variables  $\left[\frac{\lambda}{\lambda},\overline{\lambda}\right]$  also appear as parameters, they are expected to play a theoretical role in establishing convergence rather than practically affecting the numerical performance. Generally, smaller  $\omega$  and larger  $\gamma$  values result in more radical increase of the penalty parameter, which on one hand reduces the number of outer iterations needed for the slack variables to converge to zero, but on the other hand, is prone the making the problem ill-conditioned and increasing the number of inner iterations. Hence,  $\omega$  and  $\gamma$  should be at neither too large nor too small values, which can be tuned by fixing the tolerance parameters at some conservative levels (to be improved by the next step), and examining when the ELL algorithm is the most efficient.

The second groups of parameters are the tolerances associated with the approximate NLP solution and inner layer ADMM iterations. Since  $B\bar{x}^{k,r+1} - B\bar{x}^{k,r}$  and  $z^{k,r+1} - z^{k,r}$  are related to the first two terms in the linear constraint  $Ax + B\bar{x} + z = 0$  and hence of similar magnitudes,  $\epsilon_1$  and  $\epsilon_2$  as well as  $\epsilon_1^k$  and  $\epsilon_2^k$  during the outer iterations k should be similar. According to Assumption 7,  $\epsilon_4^{k,r}$  should be chosen accordingly with  $\epsilon_1^{k,r}$ , and  $\epsilon_4^k$  accordingly with  $\epsilon_1$ . Across the outer iterations, these tolerances are required to be exponentially decaying so that only a few outer iterations should be needed. To obtain a suitable relation  $\pi$  between stationarity tolerances  $\epsilon_4^{k,r}$ ,  $\epsilon_4^k$  and equality constraint tolerances  $\epsilon_5^{k,r}$ ,  $\epsilon_5^k$  of the NLP solver, we simply let  $\pi$  to be a linear function. That is,

$$\epsilon_{1} = \epsilon_{2} = \epsilon_{4}, 
\epsilon_{1}^{k} = \epsilon_{2}^{k} = \epsilon_{4}^{k} = \max\left(c_{1}/a_{1}^{k-1}, \epsilon_{1}\right), \epsilon_{5}^{k} = \frac{\epsilon_{5}}{\epsilon_{4}} \epsilon_{1}^{k}, 
\epsilon_{4}^{k,r} = \max\left(c_{4}\left(\epsilon_{1}^{k,r}\right)^{a_{4}}, \epsilon_{4}^{k}\right), \epsilon_{5}^{k,r} = \frac{\epsilon_{5}}{\epsilon_{4}} \epsilon_{4}^{k,r}.$$
(43)

The primal and dual residuals  $\epsilon_1$  and  $\epsilon_3$  should also be set in proportion so that the obtained solution is an approximate stationary point rather than a stagnant point. For the barrier constant, we need to guarantee that its diminishing should be in pace with the outer iterations, and hence can be set based on  $\epsilon_3^k$ . We also require that the barrier constants should be clamped within a specified range. That is,

$$\epsilon_3^k = \frac{\epsilon_3}{\epsilon_1} \epsilon_1^k, b^k = \min \left( b^{\max}, \max \left( \epsilon_6, \epsilon_6 \left( \epsilon_1^k / \epsilon_1 \right)^{a_6} \right) \right). \tag{44}$$



All the tolerance parameters thus reduce to  $c_1, a_1, c_4, a_4, a_6$  and the limits,  $\epsilon_1, \epsilon_3, \epsilon_4, \epsilon_5, \epsilon_6, b^{\text{max}}$ . The limits should be determined at the values at which the solution is satisfactory while not causing excessive computational cost. The exponents  $a_4$  and  $a_6$  can be chosen between 1 or 2. Then there exists only three parameters  $c_1, a_1, c_4$  to be determined empirically so that the numerical performance of the resulting ELLA algorithm best improves that of ELL.

The remaining 6 parameters are related to Anderson acceleration and its auxiliary regulation and safeguarding techniques  $-(\eta_{\theta},\eta_{w},M,\sigma,\eta_{\tilde{w}},\eta_{L})$ , which can be classified into 3 groups, namely  $\eta_{\theta}$  involved in regularization,  $(\eta_{w},M)$  in memory clean-up, and  $(\sigma,\eta_{\tilde{w}},\eta_{L})$  in rejecting inappropriate acceleration candidates. Since  $\eta_{\theta}$  aims at guaranteeing the invertibility of the Jacobian estimate with a perturbation, its value should be reasonably small (e.g.,  $\eta_{\theta}=0.05$  or 0.01), and the same holds for  $\eta_{w}$ . The maximum allowed memory M should be large enough to obtain a meaningful "average" curvature, but should not be too large to include too many obsolete previous points. Empirically, nevertheless, it was found that the variation of  $\eta_{\theta}$ ,  $\eta_{w}$ , and M does not have a strong impact. For the last three safeguarding-related parameters, the decrease of  $\sigma$  or the increase of  $\eta_{\tilde{w}}$  or  $\eta_{L}$  allows more accepted Anderson acceleration steps. The effects on the overall computational performance, however, are non-intuitive, since more accelerations may imply not only faster convergence but also possibly more overshoot and oscillations. Therefore, careful tuning is needed to obtain an appropriate extent of acceleration.

# 4 Nonlinear optimization benchmark

To test the proposed ELLADA algorithm as a generic method of nonconvex constrained distributed optimization, we consider a problem from the NLP benchmark library in this section, while the implementation and a case study on distributed nonlinear MPC will be discussed in the next sections. The benchmark is the camshape problem from the COPS (constrained optimization problem set) library (Dolan and Moré 2002), described as follows:

$$\min_{r_{1},...,r_{n}} \sum_{i=1}^{n} r_{i}$$
s.t.  $r_{\min} \leq r_{i} \leq r_{\max}, i = 1, ..., n$ 

$$2r_{i-1}r_{i+1}\cos\theta \leq r_{i}(r_{i-1} + r_{i+1}), i = 0, 1, ..., n, n + 1$$

$$-\alpha \leq \frac{r_{i+1} - r_{i}}{\theta} \leq \alpha, i = 0, 1, ..., n$$
(45)

where  $r_{\min} = 1$ ,  $r_{\max} = 2$ ,  $\alpha = 1.5$ ,  $\theta = 2\pi/5(n+1)$ ,  $r_{-1} = r_0 = r_{\min}$ ,  $r_{n+1} = r_{\max}$ ,  $r_{n+2} = r_n$ . We consider solving this problem with  $n = 4n_0 + 2$ ,  $n_0 = 100$ , in a distributed manner with 4 subproblems. The variables  $r_{in_0+1}$  and  $r_{in_0+2}$ , i = 1, 2, 3 are the shared variables between the neighboring subproblems, and hence  $\bar{x}$  are created for these variables, and equality constraints are imposed (for example, the last component of  $x_1$  and the second component of  $x_2$  should both be equal to the second



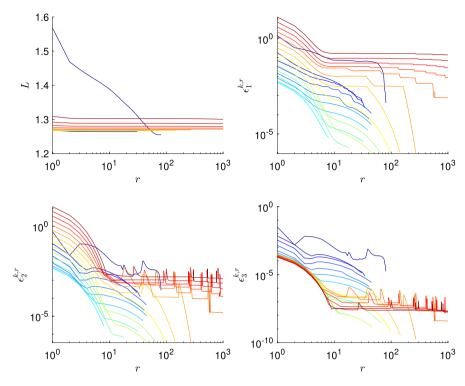


Fig. 2 Results of the ELL algorithm for the camshape problem. Colors from blue to red correspond to increasing outer iterations. (Color figure online)

component of  $\bar{x}$ ). Thus, formulations in the form of (20) and subsequently of (21) are obtained, with the dimensions of  $\bar{x}$  and z being 6 and 12, respectively. Each subproblem contains 102 variables, 100 or 101 nonlinear inequality constraints, and 202 linear inequality constraints. The subproblems are handled with the IPOPT solver used in OPTI Toolbox of Matlab. The gradient of the objective function, the Jacobian matrix of the nonlinear constraints and its structure are provided to the solver.

For ELL, the final tolerances  $\epsilon_1 = \epsilon_2 = 10^{-6}$ ,  $\epsilon_3 = 10^{-4}$  are set as equal to the default values in the IPOPT solver, and the parameters are tuned as  $\omega = 0.75$ ,  $\gamma = 2$ ,  $\epsilon_1^k = \epsilon_2^k = 10^{-3}/2^{k-1}$ ,  $\epsilon_3^k = 10^{-1}/2^{k-1}$ . The initial penalty parameter  $\beta = 4$ . The maximum number of outer iterations is set as 20, and the maximum number of inner iterations for each outer iteration is set as 1000. The intermediate calculated results of augmented Lagrangian L and  $\epsilon_1^{k,r}$ ,  $\epsilon_2^{k,r}$ ,  $\epsilon_3^{k,r}$  are shown in Fig. 2. The expected monotonic decreasing property of L is observed, leading to the convergence of the asymptotic inner iterations. However, for all the outer iterations after the 15th, the maximum allowed number of inner iterations is reached, and the maximum allowed number of outer iterations is also reached, resulting in a total number of 6002 inner iterations and a computational time of 761.40 s (0.1269 s for each inner iteration on average). At the end of the 20 outer iterations,  $||z||_{\infty} = 5.69 \times 10^{-4}$  remains high. This is due to the large value of the penalty parameter values (above  $2^{10}$ ), which fundamentally limits the speed of convergence.



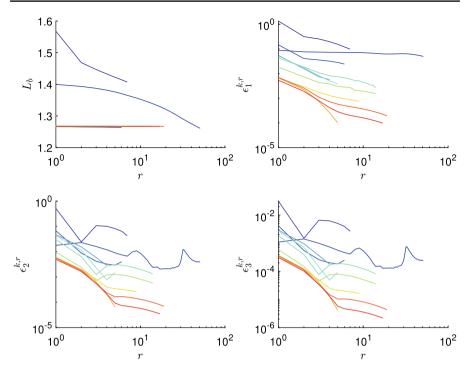


Fig. 3 Results of the ELLA algorithm for the camshape problem. Colors from blue to red correspond to increasing outer iterations. (Color figure online)

By observing Fig. 2, for ELLA we decide to reset the final tolerances to coarser values  $\epsilon_1 = \epsilon_2 = \epsilon_4 = 10^{-4}$ ,  $\epsilon_3 = 10^{-3}$  since it appears difficult for the solver to practically reach lower precisions. We let  $\epsilon_1^{k,r} = \epsilon_2^{k,r} = 10^{-1}/2^{k-1}$ ,  $\epsilon_5 = 10^{-4}$ ,  $\epsilon_6 = 10^{-4}$ . The barrier parameter is set as  $b^k = \epsilon_6 (\epsilon_3^k/\epsilon_3)^2$ . The total number of inner iterations needed to converge is 151 in 11 outer iterations, with a computational time of 14.38 s. In contrast, ELL achieves these tolerances in 7 outer iterations, which takes 317 inner iterations and 33.67 s. Hence, ELLA needs 47.6% less iterations and 57.3% less computational time than ELL to reach solutions with the same level of tolerances. On average for each inner iteration, 10.3% of computational time is saved compared to ELL. In ELLA, the behaviors of the barrier augmented Lagrangian  $L_b$  and  $\epsilon_{1.2.3}^{k,r}$ , as shown in Fig. 3, are similar to that of ELL.

Finally we seek to decrease the number of inner iterations needed by using ELLADA. The Anderson acceleration-related parameters are tuned as  $\eta_w = 0.001$ ,  $\theta = 0.001$ ,  $\sigma = 0.1$ ,  $\eta_L = 5$ ,  $\eta_{\bar{w}} = 1$ , and M = 3. After acceleration, the total number of inner iterations needed becomes 84, 44.4% lower than ELLA. Within these inner iterations, 72 iterations are under Anderson acceleration. The resulting computational time is 8.1010 s, 43.7% lower than that of ELLA. It is noted from Fig. 4 that under ELLADA, the curves of  $e_{1,2,3}^{k,r}$  becomes more kinky and larger deviations from monotonicity. Such phenomena highlight the necessity of safeguarding techniques to curb the behavior of Anderson acceleration as well as careful tuning of the relevant parameters.



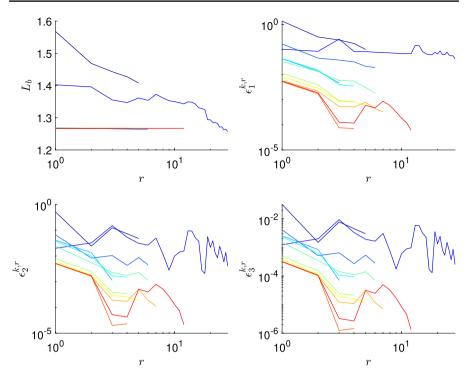


Fig. 4 Results of the ELLADA algorithm for the camshape problem. Colors from blue to red correspond to increasing outer iterations. (Color figure online)

## 5 Implementation on distributed nonlinear MPC

Consider a nonlinear discrete-time dynamical system

$$x(t+1) = f(x(t), u(t))$$
(46)

where  $x(t) \in \mathbb{R}^n$  and  $u(t) \in \mathbb{R}^m$  are the vectors of states and inputs, respectively, for t = 0, 1, 2, ..., and  $f : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}^n$ . Suppose that at time t we have the current states x = x(t), then in MPC, the control inputs are determined by the following optimal control problem:

min 
$$J = \sum_{\tau=t}^{t+T-1} \ell(\hat{x}(\tau), \hat{u}(\tau)) + \ell^{f}(\hat{x}(t+T))$$
s.t. 
$$\hat{x}(\tau+1) = f(\hat{x}(\tau), \hat{u}(\tau)), \tau = t, \dots, t+T-1$$

$$p(\hat{x}(\tau), \hat{u}(\tau), \tau) \le 0, \tau = t, \dots, t+T-1$$

$$q(\hat{x}(\tau), \hat{u}(\tau), \tau) = 0, \tau = t, \dots, t+T-1$$

$$\hat{x}(t) = x.$$
(47)



In the above formulation, the optimization variables  $\hat{x}(\tau)$  and  $\hat{u}(\tau)$  represent the predicted states and inputs in a future horizon  $\{t, t+1, \ldots, t+T\}$  with length  $T \in \mathbb{N}$ . The predicted trajectory is constrained by the dynamics (46) as well as some additional path constraints p, q such as the bounds on the inputs and states or Lyapunov descent to enforce stability. Functions  $\ell$  and  $\ell^f$  are called the stage cost and terminal cost, respectively. By solving (47), one executes  $u(t) = \hat{u}(t)$ . For simplicity it is assumed here that the states are observable; otherwise, the states can be estimated using an optimization formulation such as moving horizon estimation (MHE). For continuous-time systems, collocation techniques can be used to discretize the resulting optimal control problem into a finite-dimensional one.

Now suppose that the system (46) is large-scale with its states and outputs decomposed into n subsystems:  $x = \begin{bmatrix} x_1^\mathsf{T}, x_2^\mathsf{T}, \dots, x_n^\mathsf{T} \end{bmatrix}^\mathsf{T}$ ,  $u = \begin{bmatrix} u_1^\mathsf{T}, u_2^\mathsf{T}, \dots, u_n^\mathsf{T} \end{bmatrix}^\mathsf{T}$ , and that the optimal control problem should be solved by the corresponding n agents, each containing the model of its own subsystem:

$$x_{i}(\tau+1) = f_{i}\left(x_{i}(\tau), u_{i}(\tau), \left\{x_{ji}(\tau), u_{ji}(\tau)\right\}_{j \in \mathcal{P}(i)}\right). \tag{48}$$

where  $\{x_{ji}, u_{ji}\}$  stands for the states and inputs in subsystem j (i.e., components of  $x_j$  and  $u_j$ ) that appear in the arguments of  $f_i$ , which comprise of the components of f corresponding to the ith subsystem.  $\mathcal{P}_i$  is the collection of subsystems j that has some inputs and outputs influencing subsystem i. We assume that the cost functions and the path constraints are separable:

$$\ell(\hat{x}, \hat{u}) = \sum_{i=1}^{n} \ell_{i}(\hat{x}_{i}, \hat{u}_{i}), \ell^{f}(\hat{x}) = \sum_{i=1}^{n} \ell_{i}^{f}(\hat{x}_{i}),$$

$$p(\hat{x}, \hat{u}, \tau) = \left[ p_{1}(\hat{x}_{1}, \hat{u}_{1}, \tau)^{\mathsf{T}}, \dots, p_{n}(\hat{x}_{n}, \hat{u}_{n}, \tau)^{\mathsf{T}} \right]^{\mathsf{T}},$$

$$q(\hat{x}, \hat{u}, \tau) = \left[ q_{1}(\hat{x}_{1}, \hat{u}_{1}, \tau)^{\mathsf{T}}, \dots, q_{n}(\hat{x}_{n}, \hat{u}_{n}, \tau)^{\mathsf{T}} \right]^{\mathsf{T}}.$$
(49)

## 5.1 Formulation on directed and bipartite graphs

Graphs or networks are systematic tools to visualize distributed control/optimization problems and allow the practitioners of distributed MPC to construct and configure their specific problem of interest in a structured manner, by using network decomposition methods to automatically generate subsystems (Daoutidis et al. 2019) and associating the subsystems and their interactions with a graph topology (Jalving et al. 2019). To better illustrate the application of ELLADA on the distributed solution of the optimal control problem (47) reformulated into the form of (21), we introduce some graph-theoretic representations of optimization problems. For problem (47), we first define a directed graph (digraph), which is a straightforward characterization of the mutual impacts among the subsystem models.



**Definition 1** (*Digraph*) The digraph of system (46) under the decomposition  $x = \begin{bmatrix} x_1^\mathsf{T}, x_2^\mathsf{T}, \dots, x_n^\mathsf{T} \end{bmatrix}^\mathsf{T}$  and  $u = \begin{bmatrix} u_1^\mathsf{T}, u_2^\mathsf{T}, \dots, u_n^\mathsf{T} \end{bmatrix}^\mathsf{T}$  is  $\mathcal{G}_1 = \{\mathcal{V}_1, \mathcal{E}_1\}$  with nodes  $\mathcal{V}_1 = \{1, 2, \dots, n\}$  and edges  $\mathcal{E}_1 = \{(j, i) | j \in \mathcal{P}(i)\}$ . If  $(i, j) \in \mathcal{E}_1$ , i.e.,  $j \in \mathcal{P}(i)$ , we say that j is a parent of i and i is a child of j (denoted as  $i \in \mathcal{C}(j)$ ).

Then under the decomposition, (47) can be written as

$$\min \sum_{i \in \mathcal{V}_{1}} J_{i} = \sum_{i \in \mathcal{V}} \sum_{\tau=t}^{t+T-1} \ell_{i}(\hat{x}_{i}(\tau), \hat{u}_{i}(\tau)) + \ell_{i}^{f}(\hat{x}_{i}(t+T))$$
s.t. 
$$\hat{x}_{i}(\tau+1) = f_{i}(\hat{x}_{i}(\tau), \hat{u}_{i}(\tau), \left\{\hat{x}_{ji}(\tau), \hat{u}_{ji}(\tau)\right\}_{j \in \mathcal{P}(i)}),$$

$$p_{i}(\hat{x}_{i}(\tau), \hat{u}_{i}(\tau), \tau) \leq 0, \tau = t, \dots, t+T-1, i \in \mathcal{V}_{1}$$

$$q_{i}(\hat{x}_{i}(\tau), \hat{u}_{i}(\tau), \tau) = 0, \tau = t, \dots, t+T-1, i \in \mathcal{V}_{1}$$

$$\hat{x}_{i}(t) = x_{i}, i \in \mathcal{V}_{1},$$

$$(50)$$

We denote the variables of the ith agent as

$$\xi_{i} = \left[\hat{x}_{i}(t)^{\mathsf{T}}, \hat{u}_{i}(t)^{\mathsf{T}}, \dots, \hat{x}_{i}(t+T-1)^{\mathsf{T}}, \hat{u}_{i}(t+T-1)^{\mathsf{T}}, \hat{x}_{i}(t+T)^{\mathsf{T}}, \left\{\hat{x}_{ji}(t)^{\mathsf{T}}, \hat{u}_{ji}(t)^{\mathsf{T}}\right\}_{j \in \mathcal{P}(i)}, \dots, \left\{\hat{x}_{ji}(t+T-1)^{\mathsf{T}}, \hat{u}_{ji}(t+T-1)^{\mathsf{T}}\right\}_{j \in \mathcal{P}(i)}\right]^{\mathsf{T}},$$
(51)

in which the variables related to the *j*th subsystem are denoted as  $\xi_{ji}$ . Since  $\xi_{ji}$  is a part of the predicted states and inputs from subsystem *j*, i.e., some components of  $\xi_{j}$ , the interactions between the parent *j* and the child *i* be captured by a matrix  $\overrightarrow{D}_{ji}$  with exactly one unit entry ("1") on every row:  $\xi_{ji} = \overrightarrow{D}_{ji}\xi_{j}$ , where the right arrow represents the impact of the parent subsystem *j* on the child subsystem *i*. By denoting the model and path constraints in agent *i* as  $\xi_{i} \in \Xi_{i}$ , the optimal control problem (46) is expressed in a compact way as follows:

$$\min \sum_{i \in \mathcal{V}_1} J_i(\xi_i) \quad \text{s.t.} \xi_i \in \Xi_i, i \in \mathcal{V}_1, \xi_{ji} = \overrightarrow{D}_{ji}\xi_j, (j,i) \in \mathcal{E}_1.$$
 (52)

This is an optimization problem defined on a *directed graph*. An illustration for a simple case when  $\mathcal{E}_1 = \{(1, 2), (2, 3), (3, 1)\}$  is shown in Fig. 5a.

Although it is natural to represent the interactions among the subsystems in a digraph, performing distributed optimization on digraphs where the agents communicate among themselves without a coordinator can be challenging. For example, it is known that the ADMM algorithm, which behaves well for distributed optimization with 2 blocks of variables, can become divergent when directly extended to multi-block problems (Chen et al. 2016). Hence we construct such a 2-block architecture by using a *bipartite graph*.

**Definition 2** (Bipartite graph) The bipartite graph of system (46)  $\mathcal{G}_2$  is constructed from the digraph  $\mathcal{G}_1$  by taking both the nodes and edges as the new nodes, and adding an edge between  $i \in \mathcal{V}_1$  and  $e \in \mathcal{E}_1$  if i is the



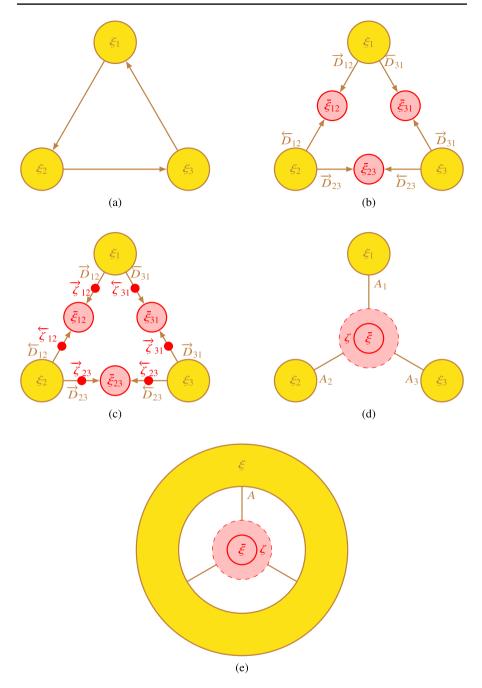


Fig. 5 Graphical illustrations of the problem structure of distributed MPC

head or tail of e in the digraph, i.e.,  $\mathcal{G}_2 = \left(\mathcal{V}_2, \mathcal{E}_2\right)$  with  $\mathcal{V}_2 = \mathcal{V}_1 \cup \mathcal{E}_1$ ,  $\mathcal{E}_2 = \left\{(i,e)|i\in\mathcal{V}_1,e\in\mathcal{E}_1,e=(i,j),j\in\mathcal{C}(i)\text{ or }e=(j,i),j\in\mathcal{P}(i)\right\}$ .



Such a graph is bipartite since any edge is between a node of  $\mathcal{V}_1$  and a node of  $\mathcal{E}_1$ . We note that the last line of (52) corresponds to the digraph edges  $\mathcal{E}_1$ . In the bipartite graph, these edges should become nodes and hence new groups of variables should be associated with them. For this purpose, we simply need to pull out  $\xi_{ji}$  as overlapping variables  $\bar{\xi}_{ji}$ , and add the constraint that  $\bar{\xi}_{ji}$  are some selected components of  $\xi_i$ :  $\xi_{ji} = \overline{D}_{ji}\xi_i$ :

$$\min \sum_{i \in \mathcal{V}_1} J_i(\xi_i) \quad \text{s.t. } \xi_i \in \Xi_i, i \in \mathcal{V}_1, \bar{\xi}_{ji} = \overrightarrow{D}_{ji}\xi_j = \overleftarrow{D}_{ji}\xi_j, (j, i) \in \mathcal{E}_1$$
 (53)

In (53), variables  $\xi_i$  ( $i \in \mathcal{V}_1$ ) and  $\bar{\xi}_{ji}$  ( $(j,i) \in \mathcal{E}_1$ ) are defined on the nodes of the bipartite graph, and the constraints captured by the matrices  $\vec{D}_{ji}$  and  $\vec{D}_{ji}$  correspond to the bipartite edges (j,(j,i)) and (i,(j,i)), respectively. We may also write the last line of (53) as

$$\bar{\xi}_e = D_{ie}\xi_i, (i, e) \in \mathcal{E}_2. \tag{54}$$

Therefore (53) is an optimization problem on the bipartite graph. An illustration is given in Fig. 5b. Under this reformulation, the problem structure becomes a 2-block one—distributed agents i = 1, ..., N manage the decision variables  $\xi_i$ ,  $\mathcal{V}_1$  in parallel without interference, and the coordinator regulates the agents by using overlapping variables  $\bar{\xi}_e$ ,  $e \in \mathcal{E}_1$ .

#### 5.2 Reformulation with slack variables

It is known that a key condition for distributed optimization in the context of the ADMM algorithm to converge is that one block of variables can always be made feasible given the other block (Wang et al. 2019). Unfortunately this condition is not always met by the problem (53). For example, given  $\xi_1$  and  $\xi_2$ , there may not be a  $\bar{\xi}_{12}$  satisfying both  $\bar{\xi}_{12} = \vec{D}_{12}\xi_1$  and  $\bar{\xi}_{12} = \vec{D}_{12}\xi_2$ . To deal with this issue, it was proposed to associate with each linear constraint in (53), namely each edge in the bipartite graph, a slack variable  $\zeta_{ie}$  (e.g., Sun and Sun 2019):

$$\begin{aligned} & \min \quad \sum_{i \in \mathcal{V}_1} J_i(\xi_i) \\ & \text{s.t.} \quad \xi_i \in \Xi_i, i \in \mathcal{V}_1 \\ & \quad D_{ie} \xi_i - \bar{\xi}_e + \zeta_{ie} = 0, (i, e) \in \mathcal{E}_2 \\ & \quad \zeta_{ie} = 0, (i, e) \in \mathcal{E}_2. \end{aligned} \tag{55}$$

Similar to the notation for D, we write  $\zeta_{ie}$  as  $\vec{\zeta}_{ij}$  if e = (i,j) and  $\vec{\zeta}_{ij}$  if e = (j,i). Such a problem structure is graphically illustrated in Fig. 5c.

Finally, we stack all the subscripted variables into  $\xi$ ,  $\bar{\xi}$ ,  $\zeta$  in a proper ordering of  $i \in \mathcal{V}_1$ ,  $e \in \mathcal{E}_1$ , and  $(i, e) \in \mathcal{E}_2$ . The matrices  $D_{ie}$  are stacked in a block diagonal pattern in the same ordering of  $(i, e) \in \mathcal{E}_2$  into A. The appearance of  $\bar{\xi}_e$  in the equality constraints is represented by a matrix B (satisfying  $B^TB = 2I$ ). We write the



objective function as  $J(\xi)$ , and the set constraints  $\Xi_i$  are lumped into a Cartesian product  $\Xi = \times_{i \in \mathcal{V}_i} \Xi_i$ . Finally, we reach a compact formulation for (55):

$$\min J(\xi) \quad \text{s.t.} \xi \in \Xi, A\xi + B\bar{\xi} + \zeta = 0, \zeta = 0$$
 (56)

Such an architecture is represented by the graphs in Fig. 5d, e. The variables  $\bar{\xi}$  and  $\zeta$  belong to the coordinator (marked in red), and  $\xi$  is in the distributed agents. Clearly, the optimal control problem formulated as (55) with the afore-mentioned graph representations is now a special form of (21), with  $\xi$ ,  $\bar{\xi}$  and  $\zeta$  rewritten as x,  $\bar{x}$  and z, respectively, and  $g(\bar{x}) = 0$ ,  $\bar{\mathcal{X}}$  equal to the entire Euclidean space.

#### 5.3 Implementation of ELLADA

As long as the cost function J is lower bounded (e.g., a quadratic cost), it follows from Theorem 1 that Algorithm 3 is applicable to (55), where the operations on  $\bar{x}$ , z, y are performed by the coordinator, and the operations on x is handled by the distributed agents. Specifically,

- The update steps of  $\bar{x}$ , z, y (Lines 10–13, 15, 16) and the entire Anderson acceleration (Lines 17–26) belong to the coordinator. The updates of penalty parameters and outer-layer dual variables  $\lambda$  (Lines 31) should also be performed by the coordinator. The conditions for  $\epsilon_1^k$ ,  $\epsilon_2^k$ ,  $\epsilon_3^k$  and  $\epsilon_1$ ,  $\epsilon_2$ ,  $\epsilon_3$  are checked by the coordinator.
- The distributed agents are responsible for carrying out a trial x-update step for the Anderson acceleration (Line 9) as well as the plain x-update (Line 14). The conditions and updates for  $\epsilon_4^{k,r}$ ,  $\epsilon_5^{k,r}$ ,  $\epsilon_4^k$ ,  $\epsilon_5^k$ , and  $\epsilon_4$ ,  $\epsilon_5$ ,  $\epsilon_6$  are checked by the agents.

When executing the updates, the agents need the values of  $B\bar{x} + z + y/\rho$  to add to Ax, and the coordinator needs the value of Ax from the agents. When the variables x are distributed into agents  $x_1, \dots, x_n$ , and the equality constraints between the agents and the coordinator is expressed on a bipartite graph:

$$D_{ie}x_i - \bar{x}_e + z_{ie} = 0, (i, e) \in \mathcal{E}_2,$$
 (57)

the communication of Ax and  $B\bar{x}+z+y/\rho$  takes place in a distributed and parallel way, i.e., the ith agent obtains the information of  $-\bar{x}_e+z_{ie}+y_{ie}/\rho$  for all e such that  $(i,e)\in\mathcal{E}_2$  from the coordinator. The coordinator, based on inter-subsystem edges e in the digraph, obtains the information of  $D_{ie}x_i$  for all related agents i. When the objective function and  $\mathcal{X}$  are separable  $f(x)=\sum_{i=1}^n f_i(x_i),\ \mathcal{X}=\mathcal{X}_1\times\cdots\times\mathcal{X}_n$ , based on such distributed and parallel communication, the optimization problem

min 
$$f(x) - b \sum_{c=1}^{C_{\phi}} \ln\left(-\phi_c(x)\right) + \frac{\rho}{2} \left\| Ax + B\bar{x} + z + \frac{y}{\rho} \right\|^2$$
  
s.t.  $\psi(x) = 0$  (58)

in an x-update step can be solved in a distributed and parallel manner:



$$\min_{x_{i}} f_{i}(x_{i}) - b \sum_{c=1}^{C_{\phi,i}} \ln \left( -\phi_{c,i}(x_{i}) \right) + \frac{\rho}{2} \sum_{\{e \mid (i,e) \in \mathcal{E}_{2}\}} \left\| D_{ie}x_{i} - \bar{x}_{e} + z_{ie} + \frac{y_{ie}}{\rho} \right\|^{2} \right\} / / \text{ for } i.$$
s.t. 
$$\psi_{i}(x_{i}) = 0 \tag{59}$$

Similarly, the  $\bar{x}$ -update with the G-mapping is in parallel for its components e, if  $\bar{\mathcal{X}}$  is separable, i.e., if  $\bar{\mathcal{X}}$  is a closed hypercube (whether bounded or unbounded), and if g is also separable. That is,  $\bar{x}$ -update can be expressed as

$$\begin{aligned} & \min_{\bar{x}_i} & g_i(\bar{x}_i) + \frac{\rho}{2} \sum_{\left\{i \mid (i,e) \in \mathcal{E}_2\right\}} \left\| D_{ie} x_i - \bar{x}_e + z_{ie} + \frac{y_{ie}}{\rho} \right\|^2 \\ & \text{s.t.} & \bar{x}_i \in \bar{\mathcal{X}}_i \end{aligned}$$

$$(60)$$

The z and y updates are in parallel for the edges (i, e) on the bipartite graph.

In Algorithm 3, the procedures are written such that in each iteration, the update steps are carried out in sequence. This requires a synchronization of all the agents i and the coordinating elements e and (i, e). For example, for the x-update, every distributed agent needs to create a "finish" signal after solving  $x_i$  in (59) and send it to the coordinator. Only after the coordinator receives the "finish" signals from all the distributed agents can the  $\bar{x}$ -update be carried out. Due to the possible computational imbalance among the agents and the coordinator, such synchronization implies that faster updates must idle for some time to wait for slower ones. In fact, the convergence properties of the ELLADA algorithm do not rely on the synchronization. Even when the inner iterations are asynchronous, the update steps still contribute to the convergence of the barrier augmented Lagrangian and hence result in convergence to KKT conditions. The only exception is that under Anderson acceleration, the steps for generating the candidate of accelerated updates are allocated to another coordinator and another set of distributed agents, and they should communicate to make the decision on executing the accelerations.

# 6 Application to a quadruple tank process

The quadruple tank process is a simple benchmark process for distributed model predictive control (Johansson 2000) with 4 states (water heights in the 4 tanks) and 2 inputs (flow rates from the reservoir). The dynamic model is written as follows:

$$\dot{h}_{1} = -\frac{a_{1}}{A_{1}}\sqrt{h_{1}} + \frac{a_{3}}{A_{1}}\sqrt{h_{3}} + \frac{\gamma_{1}k_{1}}{A_{1}}v_{1}, \dot{h}_{2} = -\frac{a_{2}}{A_{2}}\sqrt{h_{2}} + \frac{a_{4}}{A_{2}}\sqrt{h_{4}} + \frac{\gamma_{2}k_{2}}{A_{2}}v_{2}$$

$$\dot{h}_{3} = -\frac{a_{3}}{A_{3}}\sqrt{h_{3}} + \frac{(1 - \gamma_{2})k_{2}}{A_{3}}v_{2}, \qquad \dot{h}_{4} = -\frac{a_{4}}{A_{4}}\sqrt{h_{4}} + \frac{(1 - \gamma_{1})k_{1}}{A_{4}}v_{1}.$$
(61)

Other parameter values and the nominal steady state are given in Table 1.

The process is considered to have 2 subsystems, one containing tanks 1 and 4 and the other containing tanks 2 and 3. Each subsystem has 2 states, 1 input and 1



Parameter	Value	Parameter	Value
$A_1, A_3$	28	$a_{1}, a_{3}$	3.145
$A_2, A_4$	32	$a_2, a_4$	2.525
$\gamma_1$	0.43	$k_1$	3.14
$\gamma_2$	0.34	$k_2$	3.29
Input	Value	Input	Value
$\overline{v_1}$	3.15	$v_2$	3.15
State	Value	State	Value
$h_1$	12.44	$h_2$	13.17
$h_3$	4.73	$h_4$	4.99

Table 1 Parameters and nominal steady state

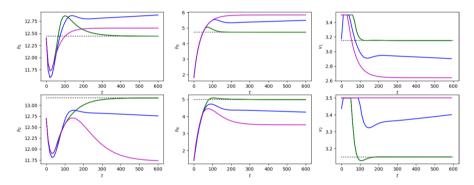


Fig. 6 Closed-loop trajectories under traditional MPC controllers for the quadruple tank process

upstream state. We first design a centralized MPC with quadratic objective function for each tank, and bounds on the inputs  $2.5 \le v_1, v_2 \le 3.5$ . We first decide through the simulation of centralized MPC that a receding horizon of T=400 with sampling time  $\delta t=10$  is appropriate. The computations are performed using the Python module pyomo.dae with an IPOPT solver (Nicholson et al. 2018).

The closed-loop trajectories under the traditional MPC controllers, including a centralized MPC (black), a semi-centralized MPC where the inputs are iteratively updated based on predictions over the entire process (green), a decentralized MPC (blue), and a distributed MPC with only state feedforwarding among the agents (purple), are shown in Fig. 6. It was observed that a semi-centralized MPC based on system-wide prediction maintains the control performance, yielding trajectories overlapping with those of the centralized MPC. However, the state-feedforward distributed MPC without sufficient coordination accounting for the state interactions results in unsatisfactory control performance, whose ultimate deviation from the steady state is even larger than the decentralized MPC without any communication between the controllers.



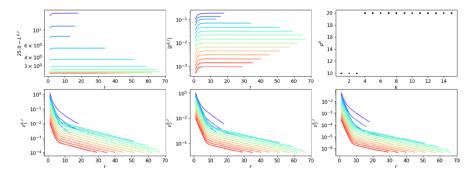


Fig. 7 Solution results of the ELL algorithm for the quadruple tank process

Next we use the proposed ELLADA algorithm for distributed nonlinear MPC of the process. We first examine the basic ELL algorithm (Algorithm 1) by solving the corresponding distributed MPC problem at a state with  $h_1 = 12.6$ ,  $h_2 = 12.4$ ,  $h_3 = 5.0$ ,  $h_4 = 4.5$ , where we set  $\omega = 0.75$ ,  $\gamma = 2$ ,  $\lambda = -\lambda = 10$  (in an element-wise sense) and  $\varepsilon_1^k = \varepsilon_2^k = 10^{-2}/2^{k-1}$ ,  $\varepsilon_3^k = 10^{-1}/2^{k-1}$ ,  $\varepsilon_1 = \varepsilon_2 = 10^{-4}$ ,  $\varepsilon_3 = 10^{-3}$  as a baseline that seeks to terminate only when a highly precise solution is reached. The solution results in terms of the variation of the augmented Lagrangian  $L^{k,r}$ , the violations to the KKT conditions  $\varepsilon_{1,2,3}^{k,r}$ , and penalty parameters  $\rho^k$  throughout the inner and outer iterations are presented in Fig. 7, where the rainbow colormap from blue to red colors stand for increasing outer iteration number. In accordance to the conclusion of Lemma 1, the augmented Lagrangian is monotonically decreasing in each outer iterations and remains upper bounded, which guarantees the convergence of the algorithm. Using the ELL algorithm for the afore-mentioned closed-loop MPC simulation, the resulting trajectories are found identical to those of the centralized control, which corroborates the theoretical property of the algorithm of converging to the set of stationary solutions.

With the preserved control performance of the ELL algorithm, we seek to improve its computational efficiency with the ELLA and ELLADA algorithms (Algorithms 2 and 3). In ELLA, the tolerances for approximate NLP solution are tuned to  $\epsilon_1 = \epsilon_2 = \epsilon_4 = 10^3 \epsilon_3 = 1, \qquad \epsilon_1^k = \epsilon_2^k = 10^3 \epsilon_3^k = \epsilon_4^k = 100/2^{k-1}, \\ \epsilon_4^{k,r} = 10^3 \epsilon_5^{k,r} = \max\left(\epsilon_4^k, 40\left(\epsilon_1^{k,r}\right)^2\right).$  The barrier constants are updated throughout outer iterations according to  $\|z\|$  according to  $b^{k+1} = \min\left(10^{-1}, \max\left(10^{-4}, 25\left(\epsilon_3^k\right)^2\right)\right).$  Compared to ELL, the accumulated number of iterations and computational time of ELLA are reduced by over an order of magnitude. To seek for better computational performance, we apply the ELLADA algorithm, where we set M=10,  $\sigma=1$ ,  $\eta_L=\eta_{\tilde{w}}=0.01$ ,  $\eta_\theta=0.5$ ,  $\eta_w=0.05$ . This further reduces the number of iterations and computational time. These results are shown in Fig. 8.

Compared to the basic ELL algorithm, ELLADA achieves acceleration by approximately 18 times in terms of iterations and 19 times in computational time for the entire simulation time span. These improvements are more significant when the



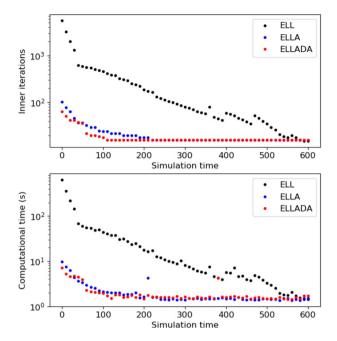


Fig. 8 Iteration and computational time under ELL, ELLA and ELLADA algorithms for the quadruple tank process

states are far from the target steady state (43 and 45 times, respectively, for the first 1/6 of the simulation). We note that the improvement from ELLA to ELLADA by using the Anderson scheme is not an order-of-magnitude one mainly because each outer iteration needs only a few number of inner iterations, leaving little space for further acceleration (e.g., for the first sampling time, 12 outer iterations including only 102 inner iterations are needed in ELLA, and in ELLADA, 61 inner iterations are needed). Under the accelerations, ELLADA returns the identical solution to the centralized optimization, thus preserving the control performance of the centralized MPC.

#### 7 Conclusions and discussions

We have proposed a new algorithm for distributed optimization allowing nonconvex constraints, which simultaneously guarantees convergence under mild assumptions and achieves fast computation. Specifically, convergence is established by adopting a two-layer architecture. In the outer layer, the slack variables are tightened using the method of multipliers, and the inequalities are handled using a barrier technique. In the inner layer, ADMM iterations are performed in a distributed and coordinated manner. Approximate NLP solution and Anderson acceleration techniques are integrated into inner iterations for computational acceleration. A large-scale optimization problem from a NLP benchmark library is used to show the numerical



advantages of using the proposed ELLADA algorithm compared to the basic convergent algorithms ELL (Algorithm 1) and ELLA (Algorithm 2).

Distributed nonlinear MPC, as an challenging problem in process control, is an important application of such an algorithm. The advantages of applying ELLADA to distributed nonlinear MPC include:

- Arbitrary input and state couplings among subsystems are allowed. No specific pattern is required a priori.
- The convergence property of the algorithm towards a stationary point is theoretically guaranteed, and its performance can be monitored throughout iterations.
- Equality-constrained NLP solvers can be used only as a subroutine. No internal
  modification of solvers is needed, and the choice of any appropriate solver is
  flexible.
- Asynchronous updates are allowed without affecting the convergence properties.
- Although motivated with a nominal optimal control problem, the algorithm could be suitable for more intricate MPC formulations such as stochastic/robust MPC or sensitivity-based advance-step MPC.

The application of the ELLADA algorithm on the distributed nonlinear MPC of a quadruple tank process has already shown its improved computational performance compared to ELL and ELLA, and improved control performance compared to the decentralized MPC and distributed MPC without accounting for state interactions. Of course, due to the small size of the specific benchmark process, the control can be realized easily with a centralized MPC. A truly large-scale control problem is more suitable to demonstrate the effectiveness of our algorithm, and this shall be presented in an upcoming separate paper.

**Acknowledgements** This work was supported by National Science Foundation (NSF-CBET). The authors would also like to thank Prof. Qi Zhang for his constructive opinions.

# Appendix 1: Proof of Lemma 1

We first prove that

$$L(x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r+1}, y^{k,r+1}) \le L(x^{k,r}, \bar{x}^{k,r}, z^{k,r}, y^{k,r})$$
$$- \beta^k \left\| B\bar{x}^{k,r+1} - B\bar{x}^{k,r} \right\|^2 - \frac{\beta^k}{2} \left\| z^{k,r+1} - z^{k,r} \right\|^2$$
(62)

for r = 0, 1, 2, ... First, since  $x^{k,r+1}$  is chosen as the minimizer of the augmented Lagrangian with respect to x (Line 9, Algorithm 1), the update of x leads to a decrease in L:

$$L(x^{k,r+1}, \bar{x}^{k,r}, z^{k,r}, y^{k,r}) \le L(x^{k,r}, \bar{x}^{k,r}, z^{k,r}, z^{k,r}, y^{k,r}). \tag{63}$$

Then consider the decrease resulted from  $\bar{x}$ -update:



$$L(x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r}, y^{k,r}) - L(x^{k,r+1}, \bar{x}^{k,r}, z^{k,r}, y^{k,r})$$

$$= g(\bar{x}^{k,r+1}) - g(\bar{x}^{k,r}) + y^{k,r} (B\bar{x}^{k,r+1} - B\bar{x}^{k,r})$$

$$+ \frac{\rho^{k}}{2} \|Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r}\|^{2} - \frac{\rho^{k}}{2} \|Ax^{k,r+1} + B\bar{x}^{k,r} + z^{k,r}\|^{2}$$

$$= g(\bar{x}^{k,r+1}) - g(\bar{x}^{k,r}) - \frac{\rho^{k}}{2} \|B\bar{x}^{k,r+1} - B\bar{x}^{k,r}\|^{2}$$

$$- \rho^{k} (\bar{x}^{k,r} - \bar{x}^{k,r+1})^{\mathsf{T}} B^{\mathsf{T}} \left(Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r} + \frac{y^{k,r}}{\rho^{k}}\right).$$
(64)

The minimization of  $\bar{x}$  (Line 10, Algorithm 1) should satisfy the optimality condition

$$0 \in \rho^k B^{\mathsf{T}} \left( A x^{k,r+1} + B \bar{x}^{k,r+1} + z^{k,r} + \frac{y^{k,r}}{\rho^k} \right) + \partial g \left( \bar{x}^{k,r+1} \right) + \mathcal{N}_{\bar{\mathcal{X}}} \left( \bar{x}^{k,r+1} \right), \quad (65)$$

i.e., there exist vectors  $v_1 \in \partial g(\bar{x}^{k,r+1})$  and  $v_2 \in \mathcal{N}_{\bar{\mathcal{X}}}(\bar{x}^{k,r+1})$  with

$$\rho^k B^{\top} \left( A x^{k,r+1} + B \bar{x}^{k,r+1} + z^{k,r} + \frac{y^{k,r}}{\rho^k} \right) = -\nu_1 - \nu_2. \tag{66}$$

Since  $v_1 \in \partial g(\bar{x}^{k,r+1})$  and g is convex,  $v_1^{\top}(\bar{x}^{k,r} - \bar{x}^{k,r+1}) \leq g(\bar{x}^{k,r}) - g(\bar{x}^{k,r+1})$ . And  $v_2 \in \mathcal{N}_{\bar{x}}(\bar{x}^{k,r+1})$  implies  $v_2^{\top}(\bar{x}^{k,r} - \bar{x}^{k,r+1}) \leq 0$ . Hence

$$\rho^{k} (\bar{x}^{k,r} - \bar{x}^{k,r+1})^{\top} B^{\top} \left( A x^{k,r+1} + B \bar{x}^{k,r+1} + z^{k,r} + \frac{y^{k,r}}{\rho^{k}} \right)$$

$$= -v_{1}^{\top} (\bar{x}^{k,r} - \bar{x}^{k,r+1}) - v_{2}^{\top} (\bar{x}^{k,r} - \bar{x}^{k,r+1})$$

$$\geq -(g(\bar{x}^{k,r}) - g(\bar{x}^{k,r+1})).$$
(67)

Substituting the above inequality in (64), we obtain

$$L(x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r}, y^{k,r}) \le L(x^{k,r+1}, \bar{x}^{k,r}, z^{k,r}, y^{k,r}) - \frac{\rho^k}{2} \|B\bar{x}^{k,r+1} - B\bar{x}^{k,r}\|^2.$$
(68)

Third, we consider the decrease resulted from *z*- and *y*-updates:



$$L(x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r+1}, y^{k,r+1}) - L(x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r}, y^{k,r})$$

$$= \lambda^{k\top} (z^{k,r+1} - z^{k,r}) + \frac{\beta^k}{2} (\|z^{k,r+1}\|^2 - \|z^{k,r}\|^2)$$

$$+ y^{k,r+1\top} (Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r+1})$$

$$- y^{k,r\top} (Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r})$$

$$+ \frac{\rho^k}{2} \|Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r+1}\|^2$$

$$- \frac{\rho^k}{2} \|Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r}\|^2.$$
(69)

Since  $v(z; \lambda, \beta) = \lambda^{T} z + \frac{\beta}{2} ||z||^{2}$  is a convex function, whose gradient is  $\nabla v(z; \lambda, \beta) = \lambda + \beta z$ ,

$$v(z^{k,r+1};\lambda^k,\beta^k) - v(z^{k,r};\lambda^k,\beta^k) \le (\lambda^k + \beta^k z^{k,r+1})^\top (z^{k,r+1} - z^{k,r}), \tag{70}$$

From Line 11 of Algorithm 1 it can be obtained

$$\lambda^k + \beta z^{k,r+1} = -y^{k,r+1}. (71)$$

Substituting into (69), we obtain

$$L(x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r+1}, y^{k,r+1}) - L(x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r}, y^{k,r})$$

$$\leq (y^{k,r+1} - y^{k,r})^{\mathsf{T}} (Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r})$$

$$+ \frac{\rho^{k}}{2} \|Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r+1}\|^{2} - \frac{\rho^{k}}{2} \|Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r}\|^{2}$$

$$= \frac{\rho^{k}}{2} (Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r+1})^{\mathsf{T}} (Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r})$$

$$+ \frac{\rho^{k}}{2} \|Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r+1}\|^{2} - \frac{\rho^{k}}{2} \|Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r}\|^{2}$$

$$= -\frac{\rho^{k}}{2} \|z^{k,r+1} - z^{k,r}\|^{2} + \rho^{k} \|Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r+1}\|^{2}$$

$$= -\frac{\rho^{k}}{2} \|z^{k,r+1} - z^{k,r}\|^{2} + \rho^{k} \|Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r+1}\|^{2}$$

From (71),

$$Ax^{k,r+1} + B\bar{x}^{k,r+1} + z^{k,r+1} = \frac{1}{\rho_k} (y^{k,r+1} - y^{k,r}) = -\frac{\beta^k}{\rho^k} (z^{k,r+1} - z^{k,r}).$$
 (73)

Then (72) becomes

$$L(x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r+1}, y^{k,r+1}) - L(x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r}, y^{k,r})$$

$$\leq -\left(\frac{\rho^k}{2} - \frac{(\beta^k)^2}{\rho^k}\right) \left\|z^{k,r+1} - z^{k,r}\right\|^2 = -\frac{\beta^k}{2} \left\|z^{k,r+1} - z^{k,r}\right\|^2.$$
(74)

Summing up the inequalities (63), (68) and (74), we have proved the inequality (62).



Next, we show that the augmented Lagrangian is lower bounded, and hence is convergent towards some  $\underline{L}^k \in \mathbb{R}$ . We note that  $v(z;\lambda,\beta)$  is a convex function of modulus  $\beta$ , it can be easily verified that

$$v(z^{k,r};\lambda^{k},\beta^{k}) + (\lambda^{k} + \beta^{k}z^{k,r})^{\top}(z' - z^{k,r}) + \frac{\rho^{k}}{2} \|z' - z^{k,r}\|^{2} \ge v(z';\lambda^{k},\beta^{k})$$
(75)

for any z', i.e.,

$$\upsilon(z^{k,r};\lambda^k,\beta^k) + y^{k,r\top}(z^{k,r} - z') \ge \upsilon(z';\lambda^k,\beta^k) - \frac{\rho^k}{2} \|z' - z^{k,r}\|^2. \tag{76}$$

Let  $z' = -(Ax^{k,r} + B\bar{x}^{k,r})$  and remove the last term on the right-hand side. Then

$$v(z^{k,r};\lambda^k,\beta^k) + y^{k,r\top} (Ax^{k,r} + B\bar{x}^{k,r} + z^{k,r})$$

$$\geq v(-(Ax^{k,r} + B\bar{x}^{k,r});\lambda^k,\beta^k). \tag{77}$$

Hence

$$L(x^{k,r}, \bar{x}^{k,r+1}, z^{k,r}, y^{k,r})$$

$$= f(x^{k,r}) + g(\bar{x}^{k,r}) + v(z^{k,r}; \lambda^{k}, \beta^{k})$$

$$+ y^{k,r}(Ax^{k,r} + B\bar{x}^{k,r} + z^{k,r}) + \frac{\rho^{k}}{2} \|Ax^{k,r} + B\bar{x}^{k,r} + z^{k,r}\|^{2}$$

$$\geq f(x^{k,r}) + g(\bar{x}^{k,r}) + v(-(Ax^{k,r} + B\bar{x}^{k,r}); \lambda^{k}, \beta^{k}).$$
(78)

Since  $v(z) = \lambda^T z + \frac{\beta}{2} ||z||^2 \ge -||\lambda||^2/(2\beta)$ ,  $\lambda$  is bounded in  $\left[\underline{\lambda}, \overline{\lambda}\right]$ ,  $\beta^k \ge \beta^1$ , and f and g are bounded below, L has a lower bound.

Taking the limit  $r \to \infty$  on the both sides of inequality (62), it becomes obvious that  $B\bar{x}^{k,r+1} - B\bar{x}^{k,r}$  and  $z^{k,r+1} - z^{k,r}$  converge to 0. Due to (73), we have  $Ax^{k,r} + B\bar{x}^{k,r} + z^{k,r} \to 0$ . Hence there must exist a r such that (22) is met. At this time, the optimality conditions for  $x^{k,r+1}$  is written as

$$0 \in \partial f(x^{k,r+1}) + \mathcal{N}_{\mathcal{X}}(x^{k,r+1}) + A^{\mathsf{T}} y^{k,r} + \rho^k A^{\mathsf{T}} (A x^{k,r+1} + B \bar{x}^{k,r} + z^{k,r}). \tag{79}$$

According to the update rule of  $y^{k,r}$ , the above expression is equivalent to

$$0 \in \partial f(x^{k,r+1}) + \mathcal{N}_{\mathcal{X}}(x^{k,r+1}) + A^{\mathsf{T}} y^{k,r+1} - \rho^k A^{\mathsf{T}} (B\bar{x}^{k,r+1} + z^{k,r+1} - B\bar{x}^{k,r} - z^{k,r}),$$
(80)

i.e.,

$$\rho^{k} A^{\top} (B \bar{x}^{k,r+1} + z^{k,r+1} - B \bar{x}^{k,r} - z^{k,r}) \in \partial f(x^{k,r+1}) + \mathcal{N}_{\mathcal{X}} (x^{k,r+1}) + A^{\top} y^{k,r+1}.$$
(81)



According to the first inequality of (22), the norm of the left hand side above is not larger than  $\epsilon_1^k$ , which directly implies the first condition in (23). In a similar manner, the second condition in (23) can be established. The third one follows from (71) and the fourth condition is obvious.

# Appendix 2: Proof of Lemma 2

We first consider the situation when  $\beta^k$  is unbounded. From (78), we have

$$\overline{L} \ge f(x^{k+1}) + g(x^{k+1}) - \lambda^{k\top} (Ax^{k+1} + B\bar{x}^{k+1}) + \frac{\beta^k}{2} \left\| Ax^{k+1} + B\bar{x}^{k+1} \right\|^2. \tag{82}$$

Since f and g are both lower bounded, as  $\beta^k \to \infty$ , we have  $Ax^{k+1} + B\bar{x}^{k+1} \to 0$ . Combined with the first two conditions of (23) in the limit of  $\epsilon_1^k$ ,  $\epsilon_2^k$ ,  $\epsilon_3^k \downarrow 0$ , we have reached (25).

Then we suppose that  $\beta^k$  is bounded, i.e., the amplification step  $\beta^{k+1} = \gamma \beta^k$  is executed for only a finite number of outer iterations. According to Lines 17–21 of Algorithm 1, expect for some finite choices of k,  $||z^{k+1}|| \le \omega ||z^k||$  always hold. Therefore  $z^{k+1} \to 0$ . Apparently, (25) follows from the limit of (23).

## Appendix 3: Proof of Lemma 3

From Lemma 1 one knows that within R inner iterations

$$\frac{\overline{L} - \underline{L}^k}{\beta^k} \ge \sum_{r=1}^R \left( \left\| B \overline{x}^{k,r+1} - B \overline{x}^{k,r} \right\|^2 + \frac{1}{2} \left\| \overline{z}^{k,r+1} - z^{k,r} \right\|^2 \right). \tag{83}$$

Then

$$\left\|B\bar{x}^{k,R+1} - B\bar{x}^{k,R}\right\|, \left\|z^{k,R+1} - z^{k,R}\right\| \sim \mathcal{O}\left(1/\sqrt{\beta^k R}\right). \tag{84}$$

For the *k*th outer iteration, its inner iterations are terminated when (22) is met, which is translated into the following relations:

$$\mathcal{O}\left(\rho^{k}/\sqrt{\beta^{k}R^{k}}\right) \leq \epsilon_{1}^{k}, \epsilon_{2}^{k} \sim \mathcal{O}(\vartheta^{k}), 
\mathcal{O}\left(1/\sqrt{\beta^{k}R^{k}}\right) \leq \epsilon_{3}^{k} \sim \mathcal{O}(\vartheta^{k}/\beta^{k}).$$
(85)

where the last relation uses (73) with  $\rho^k = 2\beta^k$ . Therefore

$$R^k \sim \mathcal{O}(\beta^k/\vartheta^{2k}).$$
 (86)

At the end of the kth iteration, suppose that Lines 19–20 and Lines 17–18 of Algorithm 1 have been executed for  $k_1$  and  $k_2$  times,



respectively  $(k_1 + k_2 = k)$ . Then the obtained  $z^{k+1}$  satisfies  $||z^{k+1}|| \sim \mathcal{O}(\omega^{k_1})$ , and  $||Ax^{k+1} + B\bar{x}^{k+1} + z^{k+1}|| \leq \epsilon_3^k \sim \mathcal{O}(\vartheta^k/\beta^k)$ , which imply

$$\left\| Ax^{k+1} + B\bar{x}^{k+1} \right\| \le \mathcal{O}(\vartheta^k/\beta^k) + \mathcal{O}(\omega^{k_1}). \tag{87}$$

From (82),

$$\beta^{k} \left\| A x^{k+1} + B \bar{x}^{k+1} \right\|^{2} \sim \beta^{k} \left( \mathcal{O} \left( \vartheta^{k} / \beta^{k} \right) + \mathcal{O} \left( \omega^{k_{1}} \right) \right)^{2} \sim \mathcal{O}(1). \tag{88}$$

Substituting (88) into (86), we obtain

$$R^k \sim \mathcal{O}\left(\frac{1}{\vartheta^{2k}} \frac{1}{\left(\mathcal{O}\left(\vartheta^k/\beta^k\right) + \mathcal{O}\left(\omega^{k_1}\right)\right)^2}\right).$$
 (89)

When  $\vartheta \le \omega$ ,  $\vartheta^k \le \omega^k \le \omega^{k_1} \gamma^{k_2}$ , and hence  $\gamma^{k_2} \vartheta^k \le \omega^{k_1}$ , i.e.,  $\omega^{k_1}$  dominates over  $\vartheta^k/\beta^k$ , leading to

$$R^k \sim \mathcal{O}(1/\vartheta^{2k}\omega^{2k_1}) \sim \mathcal{O}(1/\vartheta^{2k}\omega^{2k}).$$
 (90)

For *K* outer iterations, the total number of inner iterations is

$$R = \sum_{k=1}^{K} R^k \sim \mathcal{O}\left(\sum_{k=1}^{K} \frac{1}{\vartheta^{2k} \omega^{2k}}\right) \sim \mathcal{O}\left(\frac{1}{\vartheta^{2K} \omega^{2K}}\right). \tag{91}$$

The number of outer iterations needed to reach an  $\epsilon$ -approximate stationary point is obviously  $K \sim \mathcal{O}(\log_{\Omega} \epsilon)$ . Then

$$R \sim \mathcal{O}(\epsilon^{-2(1+\varsigma)}).$$
 (92)

# Appendix 4: Proof of Lemma 6

Through the inner iterations, only Anderson acceleration might lead to an increase in the barrier augmented Lagrangian. Combining Assumptions 3, 5, and the safeguarding criterion (41), we obtain

$$L_{b^{k}}\left(x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r+1}, y^{k,r+1}\right) \le \overline{L} + \tilde{L}_{0}\eta_{L} \sum_{r=0}^{\infty} \frac{1}{r^{1+\sigma}} < +\infty, \tag{93}$$

Together with Assumptions 1 and 2,  $L_{b^k}$  is also bounded below. Therefore  $L_{b^k}$  is bounded in a closed interval and must have converging subsequences. Therefore we can choose a subsequence converging to the lower limit  $\underline{L}$ . For any  $\varepsilon > 0$  there exists an index R of inner iteration in this subsequence, such that  $\tilde{L}_0 \eta_L \sum_{r=R}^{\infty} r^{-(1+\sigma)} < \varepsilon/2$  and  $L_{b^k} \left( x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r+1}, y^{k,r+1} \right) < \underline{L} + \varepsilon/2$  for any  $r \geq R$  on this subsequence. It then follows that for any  $r \geq R$ , whether on the subsequence or not, it holds that



$$L_{b^k}(x^{k,r+1},\bar{x}^{k,r+1},z^{k,r+1},y^{k,r+1}) < \underline{L} + \varepsilon. \tag{94}$$

Hence the upper limit is not larger than  $\underline{L} + \varepsilon$ . Due to the arbitrariness of  $\varepsilon > 0$ , the lower limit coincides with the upper limit, and hence the sequence of barrier augmented Lagrangian is convergent.

The convergence of the barrier augmented Lagrangian implies that as  $r \to \infty$ ,

$$L_{bk}(x^{k,r+1}, \bar{x}^{k,r+1}, z^{k,r+1}, y^{k,r+1}) - L_{bk}(x^{k,r}, \bar{x}^{k,r}, z^{k,r}, y^{k,r}) \to 0.$$
 (95)

Suppose that r is not an accelerated iteration, then since this quantity does not exceed  $-\beta^k \|B\bar{x}^{k,r+1} - B\bar{x}^{k,r}\|^2 - (\beta^k/2)\|z^{k,r+1} - z^{k,r}\|^2$ , we must have  $B\bar{x}^{k,r+1} - B\bar{x}^{k,r} \to 0$  and  $z^{k,r+1} - z^{k,r} \to 0$ . Otherwise if inner iteration r is accelerated, the convergence of  $B\bar{x}^{k,r+1} - B\bar{x}^{k,r}$  and  $z^{k,r+1} - z^{k,r}$  are automatically guaranteed by the second criterion (42) of accepting Anderson acceleration. The convergence properties of these two sequences naturally fall into the paradigm of Lemma 1 for establishing the convergence to approximate KKT conditions of the relaxed problem.

#### References

Anderson DG (1965) Iterative procedures for nonlinear integral equations. J ACM 12(4):547-560

Bertsekas DP (2016) Nonlinear programming, 3rd edn. Athena Scientific, Nashua

Biegler LT, Thierry DM (2018) Large-scale optimization formulations and strategies for nonlinear model predictive control. IFAC-PapersOnLine 51(20):1–15, the 6th IFAC Conference on Nonlinear Model Predictive Control (NMPC)

Boyd S, Parikh N, Chu E, Peleato B, Eckstein J (2011) Distributed optimization and statistical learning via the alternating direction method of multipliers. Found Trend Mach Learn 3(1):1–122

Chen X, Heidarinejad M, Liu J, Christofides PD (2012) Distributed economic MPC: application to a non-linear chemical process network. J Process Control 22(4):689–699

Chen C, He B, Ye Y, Yuan X (2016) The direct extension of ADMM for multi-block convex minimization problems is not necessarily convergent. Math Prog 155(1–2):57–79

Christofides PD, Scattolini R, Muñoz de la Peña D, Liu J (2013) Distributed model predictive control: a tutorial review and future research directions. Comput Chem Eng 51:21–41

Daoutidis P, Tang W, Jogwar SS (2018) Decomposing complex plants for distributed control: perspectives from network theory. Comput Chem Eng 114:43–51

Daoutidis P, Tang W, Allman A (2019) Decomposition of control and optimization problems by network structure: concepts, methods and inspirations from biology. AIChE J 65(10):e16708

Dhingra NK, Khong SZ, Jovanović MR (2019) The proximal augmented Lagrangian method for nonsmooth composite optimization. IEEE Trans Autom Control 64(7):2861–2868

Dolan ED, Moré JJ (2002) Benchmarking optimization software with performance profiles. Math Prog 91(2):201–213

Eckstein J, Bertsekas DP (1992) On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators. Math Prog 55(1–3):293–318

Eckstein J, Yao W (2017) Approximate ADMM algorithms derived from Lagrangian splitting. Comput Optim Appl 68(2):363–405

Eckstein J, Yao W (2018) Relative-error approximate versions of Douglas-Rachford splitting and special cases of the ADMM. Math Prog 170(2):417–444

Fang Hr, Saad Y (2009) Two classes of multisecant methods for nonlinear acceleration. Numer Linear Algebra Appl 16(3):197–221

Farokhi F, Shames I, Johansson KH (2014) Distributed MPC via dual decomposition and alternative direction method of multipliers. In: Distributed model predictive control made easy. Springer, Berlin, pp 115–131



- Fu A, Zhang J, Boyd S (2019) Anderson accelerated Douglas–Rachford splitting. arXiv preprint arXiv:190811482
- Gabay D, Mercier B (1976) A dual algorithm for the solution of nonlinear variational problems via finite element approximation. Comput Math Appl 2(1):17–40
- Giselsson P, Doan MD, Keviczky T, De Schutter B, Rantzer A (2013) Accelerated gradient methods and dual decomposition in distributed model predictive control. Automatica 49(3):829–833
- Glowinski R, Marroco A (1975) Sur l'approximation, par éléments finis d'ordre un, et la résolution, par pénalisation-dualité d'une classe de problèmes de dirichlet non linéaires. Rev Fr Autom Inform Rech Opér, Anal Numér 9(R2):41–76
- Goldstein T, O'Donoghue B, Setzer S, Baraniuk R (2014) Fast alternating direction optimization methods. SIAM J Imaging Sci 7(3):1588–1623
- Hajinezhad D, Hong M (2019) Perturbed proximal primal-dual algorithm for nonconvex nonsmooth optimization. Math Prog 176(1-2):207-245
- He B, Yuan X (2012) On the O(1/n) convergence rate of the Douglas–Rachford alternating direction method. SIAM J Numer Anal 50(2):700–709
- Hong M, Luo ZQ (2017) On the linear convergence of the alternating direction method of multipliers. Math Prog 162(1–2):165–199
- Hong M, Luo ZQ, Razaviyayn M (2016) Convergence analysis of alternating direction method of multipliers for a family of nonconvex problems. SIAM J Optim 26(1):337–364
- Hours JH, Jones CN (2015) A parametric nonconvex decomposition algorithm for real-time and distributed NMPC. IEEE Trans Autom Control 61(2):287–302
- Houska B, Frasch J, Diehl M (2016) An augmented Lagrangian based algorithm for distributed nonconvex optimization. SIAM J Optim 26(2):1101–1127
- Jalving J, Cao Y, Zavala VM (2019) Graph-based modeling and simulation of complex systems. Comput Chem Eng 125:134–154
- Jiang B, Lin T, Ma S, Zhang S (2019) Structured nonconvex and nonsmooth optimization: algorithms and iteration complexity analysis. Comput Optim Appl 72(1):115–157
- Johansson KH (2000) The quadruple-tank process: a multivariable laboratory process with an adjustable zero. IEEE Trans Control Syst Technol 8(3):456–465
- Li G, Pong TK (2015) Global convergence of splitting methods for nonconvex composite optimization. SIAM J Optim 25(4):2434–2460
- Liu J, Chen X, Muñoz de la Peña D, Christofides PD (2010) Sequential and iterative architectures for distributed model predictive control of nonlinear process systems. AIChE J 56(8):2137–2149
- Mota JF, Xavier JM, Aguiar PM, Püschel M (2014) Distributed optimization with local domains: applications in MPC and network flows. IEEE Trans Autom Control 60(7):2004–2009
- Nesterov YuE (1983) A method of solving a convex programming problem with convergence rate  $O(\frac{1}{k^2})$ . Dokl Akad Nauk SSSR 269(3):543–547
- Nicholson B, Siirola JD, Watson JP, Zavala VM, Biegler LT (2018) pyomo.dae: a modeling and automatic discretization framework for optimization with differential and algebraic equations. Math Prog Comput 10(2):187–223
- Nishihara R, Lessard L, Recht B, Packard A, Jordan M (2015) A general analysis of the convergence of ADMM. Proc Mach Learn Res 37:343–352
- Ouyang Y, Chen Y, Lan G, Pasiliao E Jr (2015) An accelerated linearized alternating direction method of multipliers. SIAM J Imaging Sci 8(1):644–681
- Patterson MA, Rao AV (2014) GPOPS-II: a MATLAB software for solving multiple-phase optimal control problems using  $h_p$ -adaptive Gaussian quadrature collocation methods and sparse nonlinear programming. ACM Trans Math Softw (TOMS) 41(1):1–37
- Pulay P (1980) Convergence acceleration of iterative sequences. The case of SCF iteration. Chem Phys Lett 73(2):393–398
- Rawlings JB, Mayne DQ, Diehl MM (2017) Model predictive control: theory, computation, and design, 2nd edn. Nob Hill Publishing, Madison
- Rockafellar RT, Wets RJB (1998) Variational analysis. Springer, Berlin
- Scattolini R (2009) Architectures for distributed and hierarchical model predictive control—a review. J Process Control 19(5):723-731
- Scutari G, Facchinei F, Lampariello L (2016) Parallel and distributed methods for constrained nonconvex optimization—part I: theory. IEEE Trans Signal Process 65(8):1929–1944
- Stewart BT, Venkat AN, Rawlings JB, Wright SJ, Pannocchia G (2010) Cooperative distributed model predictive control. Syst Control Lett 59(8):460–469



- Sun K, Sun XA (2019) A two-level distributed algorithm for general constrained non-convex optimization with global convergence. arXiv preprint arXiv:190207654
- Tang W, Allman A, Pourkargar DB, Daoutidis P (2018) Optimal decomposition for distributed optimization in nonlinear model predictive control through community detection. Comput Chem Eng 111:43–54
- Themelis A, Patrinos P (2020) Douglas-Rachford splitting and ADMM for nonconvex optimization: tight convergence results. SIAM J Optim 30(1):149–181
- Toth A, Kelley C (2015) Convergence analysis for Anderson acceleration. SIAM J Numer Anal 53(2):805-819
- Wächter A, Biegler LT (2005) Line search filter methods for nonlinear programming: motivation and global convergence. SIAM J Optim 16(1):1–31
- Wächter A, Biegler LT (2006) On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming. Math Prog 106(1):25–57
- Wang Y, Boyd S (2009) Fast model predictive control using online optimization. IEEE Trans Control Syst Technol 18(2):267–278
- Wang Z, Ong CJ (2017) Distributed model predictive control of linear discrete-time systems with local and global constraints. Automatica 81:184–195
- Wang Y, Yin W, Zeng J (2019) Global convergence of ADMM in nonconvex nonsmooth optimization. J Sci Comput 78(1):29–63
- Xie J, Liao A, Yang X (2017) An inexact alternating direction method of multipliers with relative error criteria. Optim Lett 11(3):583–596
- Yang Y, Hu G, Spanos CJ (2020) A proximal linearization-based fecentralized method for nonconvex problems with nonlinear constraints. arXiv preprint arXiv:200100767
- Zhang RY, White JK (2018) GMRES-accelerated ADMM for quadratic objectives. SIAM J Optim  $28(4){:}3025{-}3056$
- Zhang J, O'Donoghue B, Boyd S (2018) Globally convergent type-I Anderson acceleration for non-smooth fixed-point iterations. arXiv preprint arXiv:180803971
- Zhang J, Peng Y, Ouyang W, Deng B (2019) Accelerating ADMM for efficient simulation and optimization. ACM Trans Graph 38(6):163

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

