

Diagonal Quadratic Approximation for Parallelization of Analytical Target Cascading

Yanjing Li

Ph.D. Candidate
Department of Electrical Engineering and
Computer Science,
Stanford University,
Stanford, CA 94305
e-mail: yanjingl@stanford.edu

Zhaosong Lu

Assistant Professor
Department of Mathematics,
Simon Fraser University,
Burnaby, British Columbia V5A 1S6 Canada
e-mail: zhaosong@sfu.ca

Jeremy J. Michalek

Assistant Professor
Department of Mechanical Engineering,
and Department of Engineering
and Public Policy,
Carnegie Mellon University,
Pittsburgh, PA 15213
e-mail: jmichalek@cmu.edu

Analytical target cascading (ATC) is an effective decomposition approach used for engineering design optimization problems that have hierarchical structures. With ATC, the overall system is split into subsystems, which are solved separately and coordinated via target/response consistency constraints. As parallel computing becomes more common, it is desirable to have separable subproblems in ATC so that each subproblem can be solved concurrently to increase computational throughput. In this paper, we first examine existing ATC methods, providing an alternative to existing nested coordination schemes by using the block coordinate descent method (BCD). Then we apply diagonal quadratic approximation (DQA) by linearizing the cross term of the augmented Lagrangian function to create separable subproblems. Local and global convergence proofs are described for this method. To further reduce overall computational cost, we introduce the truncated DQA (TDQA) method, which limits the number of inner loop iterations of DQA. These two new methods are empirically compared to existing methods using test problems from the literature. Results show that computational cost of nested loop methods is reduced by using BCD, and generally the computational cost of the truncated methods is superior to the nested loop methods with lower overall computational cost than the best previously reported results.

[DOI: 10.1115/1.2838334]

Introduction

When designing complex systems, generally it is not possible or desirable to have a single decision maker in charge of all decisions because of the need to manage problem complexity. Instead, such systems are routinely decomposed hierarchically into subsystems and components, and various design groups interact to coordinate their decisions and achieve a feasible and consistent system solution. For each system in such a hierarchy, target specifications are chosen for the subsystems below such that the system can meet targets set by the supersystem above. If targets cannot be met, then negotiation and rebalancing is necessary to ensure that the final system solution is consistent and achieves system goals. Ford Motor Company refers to this process as target cascading, and the analogous model-based, computational process for such hierarchical systems has been termed analytical target cascading (ATC) [1]. In ATC, top-level design targets are propagated to lower levels, which are optimized to meet the targets. The resulting responses are rebalanced at higher levels to achieve consistency. The optimal system solution is obtained through an iterative process until target/response consistency is achieved globally.

ATC approaches this target-setting and target-matching process through formal mathematical decomposition methods, and thus, it has similarities to many of the multidisciplinary design optimization (MDO) methods that have been developed to coordinate complex analysis models from various disciplines during optimization, such as collaborative optimization (CO), concurrent subspace optimization (CSSO), and bilevel integrated system synthesis (BLISS). In particular, Allison et al. [2] compare and contrast ATC and CO. Apart from the difference in initial motivation, the formulation of ATC also differs in that it is defined for an arbitrarily large hierarchy of subsystems, and formal convergence proofs ensure that the method will reach an optimal system solution under

typical assumptions. More recently, methods for solving nonhierarchical quasi-separable (or block-angular) problems with proven convergence properties have also emerged [3–5]. In this paper, we will focus on hierarchical problems in the spirit of ATC; however, as we will discuss, ATC works by translating a description of a hierarchy of systems and subsystems linked by target-response interactions into a general quasi-separable problem,¹ through relaxation of target-response relationships between systems and their subsystems. The methods we pose for solving hierarchical ATC systems could also be used to approach general quasi-separable systems. ATC has been applied to complex systems such as automotive design [6], architectural design [7], and multidisciplinary product development [8,9], and these studies demonstrate scalability of ATC on large and computationally onerous case studies.

In this paper, we summarize existing methods for coordination of subsystems in ATC and present two new methods that enable parallelization and overcome some of the drawbacks of existing methods. We start by providing an overview of ATC in Sec. 2. In Sec. 3, we describe existing coordination methods, and in Sec. 4, we describe two new methods, including a discussion of convergence properties. In Sec. 5, we compare approaches empirically using test problems from the literature, and we discuss results in Sec. 6 and conclude in Sec. 7.

Overview of ATC

System Structure. ATC is applicable for problems that have a hierarchical structure so that the top-level design is a supersystem that consists of a number of systems, each of which may consist of its own subsystems. For example, an automobile may be composed of powertrain, body, and chassis, and the power train may be composed of engine and transmission, etc. This model is general enough to account for any number of levels in the hierarchy [10]. The objective function for the overall system can be de-

Contributed by the Design Automation Committee of ASME for publication in the JOURNAL OF MECHANICAL DESIGN. Manuscript received October 13, 2006; final manuscript received June 26, 2007; published online March 25, 2008. Review conducted by Shapour Azarm.

¹A quasi-separable problem is nearly separable except for a few linking variables that appear in multiple subsystems, as we will define rigorously later.

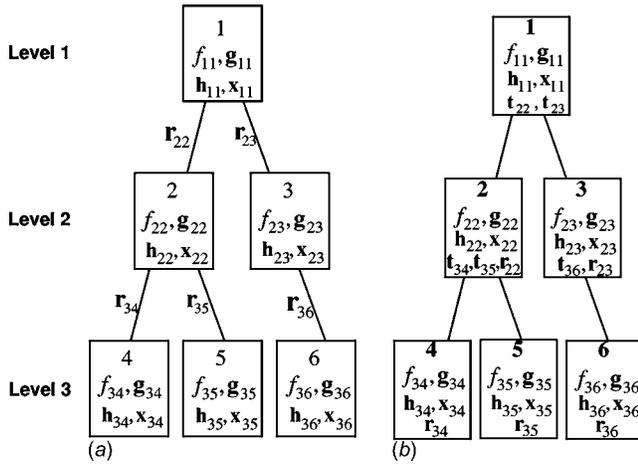


Fig. 1 (a) Hierarchical problem structure and variable allocation for ATC; (b) variable allocation for ATC after introducing target copies

scribed as a sum of the objective functions of its components. Typically, the objective function is entirely at the system level. Moreover, the subproblems are nearly separable except for a few linking variables. Figure 1(a) shows an example. Specifically, a parent and child are connected by a response variable, which represents a child's response to the design specification that its parent imposes. This "response variable" may or may not be a variable in the original predecomposition formulation: It is typically the output of the subsystem simulator and the input to the system simulator, and it may be treated as an intermediate calculation of the original formulation. The effects of subsystem response on system behavior is what prevents subsystems from being designed independently.

Notation. Different notations are used in describing and defining ATC, depending on the application [1,10–12]. In this paper, we adopt the notational system of Tosserams et al. [12] for simplicity. Consider a system that can be decomposed into N levels and M elements. The subscript ij is used to denote the j th element of the system in the i th level. f_{ij} is the scalar objective function, and $\mathbf{g}_{ij} \leq \mathbf{0}$ and $\mathbf{h}_{ij} = \mathbf{0}$ are the inequality and equality constraints, respectively. Local variables of element j are denoted by \mathbf{x}_{ij} . \mathbf{r}_{ij} is the response of element j calculated by analysis model \mathbf{a}_{ij} . Generally speaking, this model is an engineering simulation or a set of equations predicting the behavior of the subsystem. \mathcal{E}_i is the set of elements at level i , and \mathcal{C}_{ij} is the set of children of element j . The system in Fig. 1 is shown corresponding to this notation.

Mathematical Formulation. By the assumption of the problem structure, and using the notation described above, the hierarchical problem before decomposition, also known as the all-in-one (AIO) formulation, can be described as

$$\begin{aligned} \min_{\mathbf{x}_{ij}, \forall j \in \mathcal{E}_i, i=1, \dots, N} f(\mathbf{x}_{ij}, \forall j \in \mathcal{E}_i, i=1, \dots, N) &= \sum_{i=1}^N \sum_{j \in \mathcal{E}_i} f_{ij}(\bar{\mathbf{x}}_{ij}) \\ \text{s.t. } \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) &\leq \mathbf{0}, \\ \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) &= \mathbf{0} \\ \text{where } \bar{\mathbf{x}}_{ij} &= [\mathbf{x}_{ij}, \mathbf{r}_{(i+1)k}, \forall k \in \mathcal{C}_{ij}] \\ \mathbf{r}_{ij} &= \mathbf{a}_{ij}(\bar{\mathbf{x}}_{ij}) \\ \forall j \in \mathcal{E}_i, i &= 1, \dots, N \end{aligned} \quad (1)$$

Note that the response \mathbf{r}_{ij} of each element j depends on the response of its children, which prevents the objective function and the constraint sets from being separable. In order to separate the

set of variables governed by each subsystem, target variables \mathbf{t}_{ij} are created for each shared variable. And the consistency constraint

$$\mathbf{t}_{ij} - \mathbf{r}_{ij} = \mathbf{0} \quad (2)$$

is added to ensure target/response consistency.

We rewrite the problem as

$$\begin{aligned} \min_{\bar{\mathbf{x}}_{11}, \dots, \bar{\mathbf{x}}_{NM}} \sum_{i=1}^N \sum_{j \in \mathcal{E}_i} f_{ij}(\bar{\mathbf{x}}_{ij}) \\ \text{s.t. } \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) &\leq \mathbf{0} \\ \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) &= \mathbf{0} \\ \mathbf{t}_{ij} - \mathbf{r}_{ij} &= \mathbf{0} \\ \text{where } \bar{\mathbf{x}}_{ij} &= [\mathbf{x}_{ij}, \mathbf{t}_{(i+1)k}, \forall k \in \mathcal{C}_{ij}] \\ \mathbf{r}_{ij} &= \mathbf{a}_{ij}(\bar{\mathbf{x}}_{ij}) \quad \forall j \in \mathcal{E}_i, i=1, \dots, N \end{aligned} \quad (3)$$

Figure 1(b) shows the system structure and variable allocations after introducing the target variables. Note that the problem is almost separable except for the consistency constraint $\mathbf{t}_{ij} - \mathbf{r}_{ij} = \mathbf{0}$. In order to make the constraint sets separable, the consistency constraint can be relaxed using penalty functions or Lagrangian relaxation. In general, the problem can be relaxed via a consistency constraint relaxation function π . Alternate methods for consistency constraint relaxation are discussed in Sec. 3. For a general π , the resulting formulation is the relaxed AIO problem

$$\begin{aligned} \min_{\bar{\mathbf{x}}_{11}, \dots, \bar{\mathbf{x}}_{NM}} \sum_{i=1}^N \sum_{j \in \mathcal{E}_i} f_{ij}(\bar{\mathbf{x}}_{ij}) + \sum_{i=2}^N \sum_{j \in \mathcal{E}_i} \pi(\mathbf{t}_{ij} - \mathbf{r}_{ij}) \\ \text{s.t. } \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) &\leq \mathbf{0} \\ \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) &= \mathbf{0} \\ \text{where } \bar{\mathbf{x}}_{ij} &= [\mathbf{x}_{ij}, \mathbf{t}_{(i+1)k}, \forall k \in \mathcal{C}_{ij}], \\ \mathbf{r}_{ij} &= \mathbf{a}_{ij}(\bar{\mathbf{x}}_{ij}) \\ \forall j \in \mathcal{E}_i, i &= 1, \dots, N \end{aligned} \quad (4)$$

For a general π , consider only the subset of the decision variables that are nonconstant in subsystem j to obtain the general subproblem corresponding to each element:

$$\begin{aligned} \min_{\bar{\mathbf{x}}_{ij}} f_{ij}(\bar{\mathbf{x}}_{ij}) + \pi(\mathbf{t}_{ij} - \mathbf{r}_{ij}) + \sum_{k \in \mathcal{C}_{ij}} \pi(\mathbf{t}_{(i+1)k} - \mathbf{r}_{(i+1)k}) \\ \text{s.t. } \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) &\leq \mathbf{0} \\ \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) &= \mathbf{0} \\ \text{where } \bar{\mathbf{x}}_{ij} &= [\mathbf{x}_{ij}, \mathbf{t}_{(i+1)k}, \forall k \in \mathcal{C}_{ij}] \\ \mathbf{r}_{ij} &= \mathbf{a}_{ij}(\bar{\mathbf{x}}_{ij}) \\ \forall j \in \mathcal{E}_i, i &= 1, \dots, N \end{aligned} \quad (5)$$

Note that in the above formulation, the variables \mathbf{t}_{ij} and $\mathbf{r}_{(i+1)k}$ for $k \in \mathcal{C}_{ij}$ are constants with respect to element j . The constraint sets are now separable, and depending on the consistency relaxation function π , the subproblems may or may not be separable. If subproblems are separable, they can be solved in parallel. Otherwise, sequential computation of each subproblem is required.

It can be shown that by sequentially and iteratively solving each subproblem as specified in (5) in any cyclic order, convergence is guaranteed. This algorithm is called *block coordinate descent* (BCD), and the convergence result is applicable for any general relaxation consistency function π because the constraint sets are independent. The following theorem summarizes the convergence result of block coordinate descent, and the proof is given in Proposition 2.7.1 of [13], p. 268.

THEOREM 1. Suppose that for all i and j , f_{ij} is continuously differentiable and the constraint sets are independent, closed, and convex. Let $\bar{\mathbf{x}}_{ij}^k$ denote the solution of (5) in the k th iteration with $\mathbf{t}_{ij} = \mathbf{t}_{ij}^k$ and $\mathbf{r}_{(i+1)k} = \mathbf{r}_{(i+1)k}^{k-1}$ for $k \in \mathcal{C}_{ij}$. Furthermore, suppose that $\bar{\mathbf{x}}_{ij}^k$ is uniquely attained for all iterations. Then $\{\bar{\mathbf{x}}_{ij}^k\}$ is the sequence generated by the block coordinate descent method, and every limit point of $\{\bar{\mathbf{x}}_{ij}^k\}$ is a stationary point.

In the ATC literature, nonseparable subproblems are handled in nested loop coordination schemes that are carried out in a bilevel fashion, and the original convergence proof for ATC [10] is based on the nested schemes. However, by the theorem above, the nested schemes are not necessary, and the block coordinate descent method alone can produce convergence for any sequence of subproblem solutions.

The inconsistency constraint relaxation function π has been implemented in three ways: a quadratic penalty function [1,10,14], an ordinary Lagrangian function [15], or an augmented Lagrangian function [12,16]. Both the quadratic penalty and augmented Lagrangian approaches do not allow separability of subproblems, and block coordinate descent is required to achieve convergence, which limits efficiency. The ordinary Lagrangian approach does produce separable subproblems. However, the method is not robust when duality gaps exist [17]. These methods will be discussed in more detail in the following sections.

Prior Consistency Constraint Relaxation Methods for ATC

Quadratic Penalty Method. In early ATC literature, a quadratic penalty term was used as the constraint relaxation function π . The basic idea in penalty methods is to eliminate the consistency constraint and add to the objective function a penalty term that prescribes a high cost to infeasible points. Associated with this term are the penalty parameters $\mathbf{w} = [\mathbf{w}_{ij}, \forall i, j]$ that determine the severity of the penalty and, as a consequence, the extent to which the resulting unconstrained problem approximates the original constrained problem.

The general quadratic penalty function is defined as

$$\pi_Q(\mathbf{t}_{ij} - \mathbf{r}_{ij}) = \|\mathbf{w}_{ij} \circ (\mathbf{t}_{ij} - \mathbf{r}_{ij})\|_2^2 \quad (6)$$

where \circ is used to denote the Hadamard product, so that $(\mathbf{A} \circ \mathbf{B})_i = \mathbf{A}_i \mathbf{B}_i$. With π_Q , (4) can be rewritten as

$$\begin{aligned} \min_{\bar{\mathbf{x}}_{11}, \dots, \bar{\mathbf{x}}_{NM}} \quad & \sum_{i=1}^N \sum_{j \in \mathcal{E}_i} f_{ij}(\bar{\mathbf{x}}_{ij}) + \sum_{i=2}^N \sum_{j \in \mathcal{E}_i} \|\mathbf{w}_{ij} \circ (\mathbf{t}_{ij} - \mathbf{r}_{ij})\|_2^2 \\ \text{s.t.} \quad & \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) \leq \mathbf{0} \\ & \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) = \mathbf{0} \\ & \text{where } \bar{\mathbf{x}}_{ij} = [\mathbf{x}_{ij}, \mathbf{t}_{(i+1)k}] \quad \forall k \in \mathcal{C}_{ij} \\ & \mathbf{r}_{ij} = \mathbf{a}_{ij}(\bar{\mathbf{x}}_{ij}) \\ & \forall j \in \mathcal{E}_i, i = 1, \dots, N \end{aligned} \quad (7)$$

Proper selection of weighting coefficients heavily impacts the performance of this method. As \mathbf{w}_{ij} becomes larger, the resulting inconsistency between \mathbf{t}_{ij} and \mathbf{r}_{ij} pairs decreases. Convergence properties of the quadratic penalty method have been studied thoroughly and can be summarized in the following theorem.

THEOREM 2. Assume that the functions f_{ij} , \mathbf{g}_{ij} , and \mathbf{h}_{ij} for all i, j are continuous. Suppose that $\{\bar{\mathbf{x}}_{ij}^k\}$ for all i and j is an exact global minimizer of (7) for each fixed \mathbf{w}^k and that $\mathbf{w}^k \uparrow \infty$. Then every limit point of $\{\bar{\mathbf{x}}^k\}$ is a solution of Eq. (3).

Its proof is similar to that of Theorem 17.1 of [18], p. 496. Indeed, by letting $X = \{\bar{\mathbf{x}} \in \mathfrak{R}^n : \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) \leq \mathbf{0}, \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) = \mathbf{0}, \forall j \in \mathcal{E}_i, i = 1, \dots, N\}$, Eq. (3) can be generalized to

$$\begin{aligned} \min_x \quad & f(x) \\ \text{s.t.} \quad & c_i(x) = 0, i \in \mathcal{E} \\ & x \in X \end{aligned}$$

Under the assumption that \mathbf{g}_{ij} and \mathbf{h}_{ij} for all i, j are continuous on \mathfrak{R}^n , the set X is closed. Thus, replacing \mathfrak{R}^n by X in the proof of Theorem 17.1 of [18] yields a valid proof for Theorem 2. Theorem 2 is also a special case of Proposition 4.2.1 of [13], p. 391 with $\boldsymbol{\lambda}^k = \mathbf{0}$ for all k .

For each subproblem, we have

$$\begin{aligned} \min_{\bar{\mathbf{x}}_{ij}} \quad & f_{ij}(\bar{\mathbf{x}}_{ij}) + \|\mathbf{w}_{ij} \circ (\mathbf{t}_{ij} - \mathbf{r}_{ij})\|_2^2 + \sum_{k \in \mathcal{C}_{ij}} \|\mathbf{w}_{(i+1)k} \circ (\mathbf{t}_{(i+1)k} - \mathbf{r}_{(i+1)k})\|_2^2 \\ \text{s.t.} \quad & \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) \leq \mathbf{0} \\ & \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) = \mathbf{0} \\ & \text{where } \bar{\mathbf{x}}_{ij} = [\mathbf{x}_{ij}, \mathbf{t}_{(i+1)k}] \quad \forall k \in \mathcal{C}_{ij} \\ & \mathbf{r}_{ij} = \mathbf{a}_{ij}(\bar{\mathbf{x}}_{ij}) \\ & \forall j \in \mathcal{E}_i, i = 1, \dots, N \end{aligned} \quad (8)$$

With this formulation, constraint sets are separable, but the objective function is not. Subproblems are solved sequentially throughout the hierarchical chain for fixed \mathbf{w} as an inner loop. As noted before, either the nested coordination scheme or the block coordinate descent method (BCD) can be used as the inner loop. Upon convergence of all inner loops, an outer loop is used to update \mathbf{w} . Figure 7 shows the procedure using nested and BCD schemes. Prior approaches to update \mathbf{w} are either setting the penalty weights through trial and error or initializing them to a small value and then increasing their value by a linear update method, i.e., multiplying the current weights by a constant. However, the trial and error approach can be difficult for large-scale problems. Michalek and Papalambros [14] observed that although any positive \mathbf{w} can ensure convergence to a consistent solution for problems that have attainable targets, no finite \mathbf{w} will lead to perfect consistency for problems with unattainable targets. They proposed an efficient weight update method (WUM) for finding values of \mathbf{w} that achieve solutions within user-specified inconsistency tolerances. It is demonstrated that the WUM achieves better performance than a constant \mathbf{w} .

In practice, large penalty weights can cause ill-conditioning [13]. Moreover, the quadratic term prevents each subproblem from being separable, preventing the use of parallel computing.

Ordinary Lagrangian Method. An alternative choice for the constraint relaxation function π is an ordinary Lagrangian function [15]. This method (OL) is based on Lagrangian duality theory [13,18]. Let $\boldsymbol{\lambda} = [\boldsymbol{\lambda}_{ij}, \forall i, j]$ be the vector of Lagrangian multipliers. The consistency constraint relaxation function is specified as

$$\pi_L(\mathbf{t}_{ij} - \mathbf{r}_{ij}) = \boldsymbol{\lambda}_{ij}^T (\mathbf{t}_{ij} - \mathbf{r}_{ij}) \quad (9)$$

With π_L , (4) can be rewritten as

$$\begin{aligned} \Lambda(\boldsymbol{\lambda}) = \min_{\bar{\mathbf{x}}_{11}, \dots, \bar{\mathbf{x}}_{NM}} \quad & \sum_{i=1}^N \sum_{j \in \mathcal{E}_i} f_{ij}(\bar{\mathbf{x}}_{ij}) + \sum_{i=2}^N \sum_{j \in \mathcal{E}_i} \boldsymbol{\lambda}_{ij}^T (\mathbf{t}_{ij} - \mathbf{r}_{ij}) \\ \text{s.t.} \quad & \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) \leq \mathbf{0} \\ & \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) = \mathbf{0} \\ & \text{where } \bar{\mathbf{x}}_{ij} = [\mathbf{x}_{ij}, \mathbf{t}_{(i+1)k}] \quad \forall k \in \mathcal{C}_{ij} \\ & \mathbf{r}_{ij} = \mathbf{a}_{ij}(\bar{\mathbf{x}}_{ij}) \\ & \forall j \in \mathcal{E}_i, i = 1, \dots, N \end{aligned} \quad (10)$$

And the dual problem for (10) is

$$\max_{\boldsymbol{\lambda}} \Lambda(\boldsymbol{\lambda}) \quad (11)$$

Under convexity and constraint qualification assumptions, the strong duality theorem holds and there is no duality gap. As a result, an iterative process can be carried out by first solving (10) with a fixed $\boldsymbol{\lambda}$, then updating $\boldsymbol{\lambda}$ via the subgradient method until convergence. Note that this process is based on solving the dual problem (11), with (10) as its nested problem, to obtain the overall solution of (3). This procedure is theoretically justified by the following theorems, which are presented in Proposition 5.3.2 of [13], p. 514. Again, readers are referred to the text for proofs.

THEOREM 3 (STRONG DUALITY THEOREM). *Assume that f_{ij} and \mathbf{g}_{ij} are convex, and \mathbf{h}_{ij} is linear for $\forall j \in \mathcal{E}_i, i=1, \dots, N$. Suppose that $\{\bar{\mathbf{x}}: \mathbf{t}-\mathbf{r}=\mathbf{0}, \mathbf{g}(\bar{\mathbf{x}})<\mathbf{0}, \mathbf{h}(\bar{\mathbf{x}})=\mathbf{0}\}$ is nonempty. Assume that the optimal value of problem (3) is finite. Then there is no duality gap and there exists at least one Lagrange multiplier.*

THEOREM 4. *Assume that the strong duality theorem is satisfied for Eq. (3). If $\boldsymbol{\lambda}^*$ is the optimal solution to Eq. (11) and if $\bar{\mathbf{x}}^* = [\bar{\mathbf{x}}_{ij}^* \forall i, j]$ solves Eq. (10) with $\boldsymbol{\lambda} = \boldsymbol{\lambda}^*$, then $\bar{\mathbf{x}}^*$ also solves Eq. (3).*

For each subsystem, consider a fixed $\boldsymbol{\lambda}_{ij}$. Also note that since \mathbf{t}_{ij} and $\mathbf{r}_{(i+1)k}$ are constants with respect to element j , the terms $\boldsymbol{\lambda}_{ij}\mathbf{t}_{ij}$ and $\boldsymbol{\lambda}_{(i+1)k}^T\mathbf{r}_{(i+1)k}$ can be eliminated from the objective function, resulting in the following subproblem:

$$\begin{aligned} \min_{\bar{\mathbf{x}}_{ij}} & f_{ij}(\bar{\mathbf{x}}_{ij}) + \boldsymbol{\lambda}_{ij}^T(-\mathbf{r}_{ij}) + \sum_{k \in \mathcal{C}_{ij}} \boldsymbol{\lambda}_{(i+1)k}^T \mathbf{t}_{(i+1)k} \\ \text{s.t.} & \quad \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) \leq \mathbf{0} \\ & \quad \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) = \mathbf{0} \\ \text{where} & \quad \bar{\mathbf{x}}_{ij} = [\mathbf{x}_{ij}, \mathbf{t}_{(i+1)k} \forall k \in \mathcal{C}_{ij}] \\ & \quad \mathbf{r}_{ij} = \mathbf{a}_{ij}(\bar{\mathbf{x}}_{ij}) \\ & \quad \forall j \in \mathcal{E}_i, i=1, \dots, N \end{aligned} \quad (12)$$

In this method, each subproblem is separable and can be solved in parallel, with fixed Lagrange multipliers. Starting with arbitrary Lagrange multiplier estimates, we update $\boldsymbol{\lambda}$ using the subgradient method, defined as follows:

$$\boldsymbol{\lambda}^{\kappa+1} = \boldsymbol{\lambda}^{\kappa} + \tau^{\kappa}(\mathbf{t}^{\kappa} - \mathbf{r}^{\kappa}) \quad (13)$$

The superscript κ denotes the number of iterations, and τ is a scalar representing the step size. The subgradient method is proven to converge, under assumptions shown in the following theorem [15].

THEOREM 5. *If the feasible set of the solution of (10) is bounded and a sequence of step size τ^{κ} satisfies the conditions*

1. $\tau^{\kappa} \rightarrow 0$ as $\kappa \rightarrow \infty$
2. $\sum_{\kappa=1}^{\infty} \tau^{\kappa} = \infty$

then either the iterative process of the subgradient method terminates finitely at $\kappa = \kappa^$ with $\Lambda(\boldsymbol{\lambda}^{\kappa^*}) = \Lambda(\boldsymbol{\lambda}^*)$, or else an infinite sequence of iterates is generated such that $\Lambda(\boldsymbol{\lambda}^{\kappa^*}) \rightarrow \Lambda(\boldsymbol{\lambda}^*)$ as $\kappa \rightarrow \infty$.*

Figure 7 shows the OL procedure. By the theorems above, it is necessary that the strong duality condition holds in order to obtain global convergence. This is the main drawback of this method because duality gaps may exist [17], which can cause instability and limit application of the method for practical problems.

Augmented Lagrangian Method. An improved method, the augmented Lagrangian method (AL) applied to ATC formulations by Tossers et al. [12] combines both the quadratic penalty term and the Lagrangian term, which overcomes some drawbacks of both QP and OL. The Lagrangian term is used to avoid ill-

conditioning. The quadratic term also reduces duality gaps. The augmented Lagrangian function is defined as follows:

$$\pi_{\text{AL}}(\mathbf{t}_{ij} - \mathbf{r}_{ij}) = \boldsymbol{\lambda}_{ij}^T(\mathbf{t}_{ij} - \mathbf{r}_{ij}) + \|\mathbf{w}_{ij} \circ (\mathbf{t}_{ij} - \mathbf{r}_{ij})\|_2^2 \quad (14)$$

With π_{AL} , the general problem is

$$\begin{aligned} \hat{\Lambda}(\boldsymbol{\lambda}) = \min_{\bar{\mathbf{x}}_{11}, \dots, \bar{\mathbf{x}}_{NM}} & \sum_{i=1}^N \sum_{j \in \mathcal{E}_i} f_{ij}(\bar{\mathbf{x}}_{ij}) + \sum_{i=2}^N \sum_{j \in \mathcal{E}_i} (\boldsymbol{\lambda}_{ij}^T(\mathbf{t}_{ij} - \mathbf{r}_{ij}) \\ & + \|\mathbf{w}_{ij} \circ (\mathbf{t}_{ij} - \mathbf{r}_{ij})\|_2^2) \\ \text{s.t.} & \quad \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) \leq \mathbf{0} \\ & \quad \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) = \mathbf{0} \\ \text{where} & \quad \bar{\mathbf{x}}_{ij} = [\mathbf{x}_{ij}, \mathbf{t}_{(i+1)k} \forall k \in \mathcal{C}_{ij}] \\ & \quad \mathbf{r}_{ij} = \mathbf{a}_{ij}(\bar{\mathbf{x}}_{ij}) \\ & \quad \forall j \in \mathcal{E}_i, i=1, \dots, N \end{aligned} \quad (15)$$

The dual problem for (15) is

$$\max_{\boldsymbol{\lambda}} \hat{\Lambda}(\boldsymbol{\lambda}) \quad (16)$$

The following theorem is the basic convergence result of the augmented Lagrangian function. This theorem is a specialization of Proposition 4.2.1 of [13], p. 391. let $X = \{\bar{\mathbf{x}} \in \mathfrak{R}^n: \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) \leq \mathbf{0}, \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) = \mathbf{0}, \forall j \in \mathcal{E}_i, i=1, \dots, N\}$. Under the assumption that \mathbf{g}_{ij} and \mathbf{h}_{ij} for all i, j are continuous on \mathfrak{R}^n , the set X is closed. Thus, Eq. (3) satisfies all assumptions of Proposition 4.2.1 of [13]. Readers are referred to the text for the proof.

THEOREM 6. *Assume that the objective and constraint functions of Eq. (3) are continuous. For $\kappa=0, 1, \dots$, let $\bar{\mathbf{x}}^{\kappa}$ be the global minimum of the problem $\hat{\Lambda}(\boldsymbol{\lambda}^{\kappa})$ with \mathbf{w}^{κ} , where $\boldsymbol{\lambda}^{\kappa}$ is bounded, $0 < \mathbf{w}^{\kappa} < \mathbf{w}^{\kappa+1}$ for all κ , and $\mathbf{w}^{\kappa} \rightarrow \infty$. Then every limit point of the sequence $\{\bar{\mathbf{x}}^{\kappa}\}$ is a global minimum of Eq. (3).*

For each subproblem, we have

$$\begin{aligned} \min_{\bar{\mathbf{x}}_{ij}} & f_{ij}(\bar{\mathbf{x}}_{ij}) + \boldsymbol{\lambda}_{ij}^T(-\mathbf{r}_{ij}) + \|\mathbf{w}_{ij} \circ (\mathbf{t}_{ij} - \mathbf{r}_{ij})\|_2^2 \\ & + \sum_{k \in \mathcal{C}_{ij}} (\boldsymbol{\lambda}_{(i+1)k}^T \mathbf{t}_{(i+1)k} + \|\mathbf{w}_{(i+1)k} \circ (\mathbf{t}_{(i+1)k} - \mathbf{r}_{(i+1)k})\|_2^2) \\ \text{s.t.} & \quad \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) \leq \mathbf{0} \\ & \quad \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) = \mathbf{0} \\ \text{where} & \quad \bar{\mathbf{x}}_{ij} = [\mathbf{x}_{ij}, \mathbf{t}_{(i+1)k} \forall k \in \mathcal{C}_{ij}] \\ & \quad \mathbf{r}_{ij} = \mathbf{a}_{ij}(\bar{\mathbf{x}}_{ij}) \\ & \quad \forall j \in \mathcal{E}_i, i=1, \dots, N \end{aligned} \quad (17)$$

Duality theory also applies to the augmented Lagrangian formulation so we can use the same procedure as OL. However, the subproblems as shown in (17) are not separable due to the quadratic penalty term. As a result, instead of solving all subproblems in parallel as in OL, either the nested coordination scheme or the block coordinate descent method is carried out as an inner loop for the AL relaxation formulation, similar to the QP approach. In order to achieve convergence, we must update the Lagrange multipliers so they approach to the optimal values. The augmented Lagrangian function allows the use of *method of multipliers*. The scheme for selecting new terms $\boldsymbol{\lambda}$ from loop iterate κ to $(\kappa+1)$ is given by the following formula²:

$$\boldsymbol{\lambda}^{(\kappa+1)} = \boldsymbol{\lambda}^{\kappa} + 2\mathbf{w}^{\kappa} \circ \mathbf{w}^{\kappa} \circ (\mathbf{t}_{ij}^{\kappa} - \mathbf{r}_{ij}^{\kappa}) \quad (18)$$

Before addressing convergence properties of the method of multipliers, we introduce a new concept called generalized

²Note that the penalty weights are squared in the definition of multipliers in this case because they are squared in the definition of the quadratic penalty term.

Lagrange multiplier.

DEFINITION 1 (GENERALIZED LAGRANGE MULTIPLIER). Consider the following problem:

$$\begin{aligned} \min f(x) \\ \text{s.t. } c(x) = 0, \quad x \in X \end{aligned} \quad (19)$$

where $f: \mathfrak{R}^n \rightarrow \mathfrak{R}$, $c: \mathfrak{R}^n \rightarrow \mathfrak{R}^m$ are continuously differentiable, and X is a given subset in \mathfrak{R}^n . For the set X , the set of feasible directions at $x^* \in X$ is defined as $D_X(x^*) = \{0 \neq p \in \mathfrak{R}^n : x^* + tp \in X, \forall t \in [0, \delta]\}$ for some $\delta > 0$. $\lambda^* \in \mathfrak{R}^m$ is called a generalized Lagrange multiplier for the constraint $c(x) = 0$ at $x = x^*$ if

$$p^T [\nabla f(x^*) + \lambda^T \nabla c(x^*)] \geq 0, \quad \forall p \in D_X(x^*)$$

holds.

Clearly, when $X = \mathfrak{R}^n$, we have $D_X(x^*) = \mathfrak{R}^n \setminus \{0\}$, and the generalized Lagrange multiplier λ^* becomes a standard Lagrange multiplier that satisfies

$$\nabla f(x^*) + \lambda^T \nabla c(x^*) = 0$$

Let

$$L_w(x, \lambda) = f(x) + \lambda^T c(x) + \frac{w}{2} \|c(x)\|^2$$

We next address convergence properties of the method of multipliers for Eq. (19).

PROPOSITION 1. Assume that $f(x)$ and $c(x)$ are continuously differentiable, and that X is a closed convex set. For $\kappa = 0, 1, \dots$, let x^κ be a global minimum of the problem

$$\begin{aligned} \min L_{w^\kappa}(x, \lambda^\kappa) \\ \text{s.t. } x \in X \end{aligned}$$

where $\{\lambda^\kappa\}$ is bounded, $0 < w^\kappa < w^{\kappa+1}$ for all κ , and $w^\kappa \rightarrow \infty$. Assume that a subsequence $\{x^{\kappa_j}\}_K$ converges to a vector x^* and $\{\lambda^{\kappa_j} + w^{\kappa_j} c(x^{\kappa_j})\}_K \rightarrow \lambda^*$. Then x^* is a minimum of Eq. (19) and λ^* is a generalized Lagrange multiplier for the constraint $c(x) = 0$ at $x = x^*$.

Proof. It immediately follows from Proposition 4.2.1 of [13], p. 391 that x^* is a minimum of Eq. (19). We next show that λ^* is a generalized Lagrange multiplier for the constraint $c(x) = 0$ at $x = x^*$. Without loss of generality, we assume that the entire sequence $\{x^\kappa\}$ and $\{\lambda^\kappa + w^\kappa c(x^\kappa)\}$ converges to x^* and λ^* , respectively. Let

$$\tilde{\lambda}^\kappa = \lambda^\kappa + w^\kappa c(x^\kappa), \quad \forall \kappa \quad (20)$$

By definition of x^κ and the first-order optimality condition, we have

$$p^T [\nabla f(x^\kappa) + (\tilde{\lambda}^\kappa)^T \nabla c(x^\kappa)] \geq 0, \quad \forall p \in D_X(x^\kappa) \quad (21)$$

For any $p \in D_X(x^*)$, there exists a $t > 0$ such that $x^* + tp \in X$. This together with the fact that X is convex and $x^\kappa \in X$ implies that $\alpha(x^* + tp) + (1 - \alpha)x^\kappa \in X$ for $\alpha \in [0, 1]$. Thus, $x^* + \alpha(x^* + tp - x^\kappa) \in X, \forall \alpha \in [0, 1]$. It follows that $x^* + tp - x^\kappa \in D_X(x^\kappa)$. Using this relation, (20) and (21), we have

$$(x^* + tp - x^\kappa)^T [\nabla f(x^\kappa) + (\tilde{\lambda}^\kappa)^T \nabla c(x^\kappa)] \geq 0, \quad \forall p \in D_X(x^*)$$

Upon letting $\kappa \rightarrow \infty$, and using the assumption that $x^\kappa \rightarrow x^*$ and $\tilde{\lambda}^\kappa \rightarrow \lambda^*$, we further obtain that

$$tp^T [\nabla f(x^*) + (\lambda^*)^T \nabla c(x^*)] \geq 0, \quad \forall p \in D_X(x^*)$$

Noting that $t > 0$, we have

$$p^T [\nabla f(x^*) + (\lambda^*)^T \nabla c(x^*)] \geq 0, \quad \forall p \in D_X(x^*)$$

and hence, λ^* is a generalized Lagrange multiplier for the con-

straint $c(x) = 0$ at $x = x^*$. ■

We observe that Eq. (3) lies in the class of problems described in (19) with the set X defined as

$$X = \{\bar{x} \in \mathfrak{R}^n : \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) \leq \mathbf{0}, \quad \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) = \mathbf{0}, \quad \forall j \in \mathcal{E}_i, \quad i = 1, \dots, N\} \quad (22)$$

Convergence properties of the method of multipliers applied to Eq. (3) are presented as the following theorem. Its proof is similar to that of Proposition 1.

THEOREM 7. Assume that the objective function of Eq. (3) is continuously differentiable, and the set X defined in (22) is closed convex. Let

$$\begin{aligned} f_{\text{AL}}(\bar{\mathbf{x}}, \boldsymbol{\lambda}, \mathbf{w}) = & \sum_{i=1}^N \sum_{j \in \mathcal{E}_i} f_{ij}(\bar{\mathbf{x}}_{ij}) + \sum_{i=2}^N \sum_{j \in \mathcal{E}_i} (\boldsymbol{\lambda}_{ij}^T (\mathbf{t}_{ij} - \mathbf{r}_{ij}) + \|\mathbf{w}_{ij} \circ (\mathbf{t}_{ij} \\ & - \mathbf{r}_{ij})\|_2^2) \end{aligned}$$

where $\bar{\mathbf{x}} = [\bar{\mathbf{x}}_{11}, \dots, \bar{\mathbf{x}}_{NM}]$ and $\boldsymbol{\lambda} = [\boldsymbol{\lambda}_{11}, \dots, \boldsymbol{\lambda}_{NM}]$. Let $\bar{\mathbf{x}}^*$ be a minimum of Eq. (15) and suppose that $\{\lambda^\kappa\}$ is bounded, and $\{w^\kappa\}$ satisfies

$$0 < w^\kappa < w^{\kappa+1} \quad \forall \kappa, \quad w^\kappa \rightarrow \infty$$

Assume that a subsequence $\{\bar{\mathbf{x}}^{\kappa_j}\}_K$ converges to a vector $\bar{\mathbf{x}}^*$ and $\{\lambda^{\kappa_j} + 2w^{\kappa_j} \circ (\mathbf{t}^{\kappa_j} - \mathbf{r}^{\kappa_j})\}_K \rightarrow \lambda^*$. Then $\bar{\mathbf{x}}^*$ is a minimum of Eq. (3) and λ^* is a generalized Lagrange multiplier corresponding to the constraints $\mathbf{t}_{ij} - \mathbf{r}_{ij} = \mathbf{0}$ at $\bar{\mathbf{x}} = \bar{\mathbf{x}}^*$ for all $j \in \mathcal{E}_i, i = 1, \dots, N$.

There are other ways of updating $\boldsymbol{\lambda}$ when using the augmented Lagrangian method. For example, the subgradient update scheme is used in [16]. In AL, the penalty weight \mathbf{w} is also updated. A linear scheme with constant scale β for updating \mathbf{w} is used in [16]

$$\mathbf{w}^{(\kappa+1)} = \beta \mathbf{w}^\kappa, \quad \beta \geq 1 \quad (23)$$

Because the quadratic term prevents subproblems from being separable, the nested coordination scheme or BCD is used in the inner loop for finding the optimal solution with respect to fixed $\boldsymbol{\lambda}$ and \mathbf{w} . The outer loop applies (18) and (23) for updating the Lagrange multipliers and penalty weights. In implementing the AL method, the update of \mathbf{w} can also be separate from the update of $\boldsymbol{\lambda}$. \mathbf{w} can be updated only if the improvement of the current iteration is not large enough [5]. Figure 7 shows the procedure using nested and BCD schemes. Although the augmented Lagrangian approach shows stable convergence properties and a superior performance compared to QP, the inner loop can still induce large computational cost. Moreover, throughput is low due to nonseparability of subproblems.

Augmented Lagrangian With Alternating Direction Method of Multipliers. To reduce computational cost of the augmented Lagrangian approach, Tossierams et al. [12] applied the augmented Lagrangian with alternating direction method of multipliers (ALAD). The key observation is that all elements in the odd levels of the hierarchy only depend on the elements in the even levels and vice versa; thus, it is possible to first solve all odd-level elements in parallel, then all even-level elements, for a fixed number of iterations. The Lagrange multipliers are then updated using the method of multipliers. The penalty weight can also be updated. Contrary to quadratic penalty method (QP) or AL, a large penalty weight may have negative effect on convergence, but a small penalty weight may result in unbounded subproblems [19]. In [12], Tossierams et al. keep a constant \mathbf{w} . Furthermore, the inner loop is solved for only one iteration to minimize computational effort. This method is demonstrated to be effective on all example problems with superior computational properties compared to QP and AL. It has been proven to converge under the assumption that the problem is convex, the feasibility sets are nonempty, and the constraint sets are bounded [19]. In practice, however, ALAD shows good numerical convergence behavior for nonconvex problems as well [5]. Figure 7 shows the ALAD procedure. This method is

promising because it demonstrates both good convergence properties and low computational cost through the truncated inner loops and partial parallelization [12].

Summary. Except for the ordinary Lagrangian approach, all other prior methods for consistency constraint relaxation prevent fully parallel computation of subproblems. The ALAD method requires sequential computation of odd levels and even levels, and the QP and AL methods require sequential computation of all levels. As parallel and distributed computing becomes more popular, it is desirable to have fully separable subproblems so that each subproblem can be solved concurrently and computational throughput can be improved. Because the ordinary Lagrangian approach has convergence difficulties and other approaches prevent separability, we are motivated to explore alternate approaches.

Proposed Consistency Constraint Relaxation Methods for ATC

Diagonal Quadratic Approximation Method. In order for the subproblems of the augmented Lagrangian approach to be separable, we apply the diagonal quadratic approximation (DQA) originally proposed in [20], where it is used to solve block-angular structure problems. Proofs on convergence and convergence rate are available for this method [20]. There are other variations of solving similar problems using approximation techniques, for example, the method discussed in [17].

In the original paper, DQA is used to solve problems in the following form:

$$\min_{\mathbf{x}} f(\mathbf{x}) = \sum_{i=1}^L f_i(\mathbf{x}_i) \quad (24)$$

$$\text{s.t. } \sum_{i=1}^L \mathbf{A}_i \mathbf{x}_i = \mathbf{b} \quad (25)$$

$$\mathbf{x}_i \in X_i, \quad i = 1, 2, \dots, L \quad (26)$$

where f is convex, and X_1, X_2, \dots, X_L are nonempty closed convex sets. Our ATC formulation, as in (3), falls into this form for convex problems. First, the objective function is separable with respect to individual subproblems. Second, the consistency constraint function is in the form of (25). Third, the inequality and equality constraints \mathbf{g} and \mathbf{h} , now fully separable, form sets X_i . In DQA, a linear approximation is applied on the cross term $\mathbf{t}_{ij} \circ \mathbf{r}_{ij}$ included in the quadratic penalty term $\|\mathbf{t}_{ij} - \mathbf{r}_{ij}\|_2^2$ in (15).

$$\|\mathbf{t}_{ij} - \mathbf{r}_{ij}\|_2^2 = \|\mathbf{t}_{ij} \circ \mathbf{t}_{ij} + \mathbf{r}_{ij} \circ \mathbf{r}_{ij} - 2(\mathbf{t}_{ij} \circ \mathbf{r}_{ij})\|_1 \quad (27)$$

Using the Taylor expansion for multiple variable scalar functions up to the first order, a linearization at the point $\mathbf{t}_{ij}^k, \mathbf{r}_{ij}^k$ gives

$$\begin{aligned} \mathbf{t}_{ij} \circ \mathbf{r}_{ij} &\cong \mathbf{t}_{ij}^k \circ \mathbf{r}_{ij}^k + \mathbf{r}_{ij}^k \circ (\mathbf{t}_{ij} - \mathbf{t}_{ij}^k) + \mathbf{t}_{ij}^k \circ (\mathbf{r}_{ij} - \mathbf{r}_{ij}^k) \\ &= \mathbf{r}_{ij}^k \circ \mathbf{t}_{ij} + \mathbf{t}_{ij}^k \circ \mathbf{r}_{ij} - \mathbf{t}_{ij}^k \circ \mathbf{r}_{ij}^k \end{aligned} \quad (28)$$

Combining (27) and (28), we have

$$\begin{aligned} \|\mathbf{t}_{ij} - \mathbf{r}_{ij}\|_2^2 &\cong \|\mathbf{t}_{ij} \circ \mathbf{t}_{ij} + \mathbf{r}_{ij} \circ \mathbf{r}_{ij} - 2(\mathbf{r}_{ij}^k \circ \mathbf{t}_{ij} + \mathbf{t}_{ij}^k \circ \mathbf{r}_{ij} - \mathbf{t}_{ij}^k \circ \mathbf{r}_{ij}^k)\|_1 \\ &= \|\mathbf{t}_{ij}^k \circ \mathbf{t}_{ij} + \mathbf{r}_{ij} \circ \mathbf{r}_{ij} - 2(\mathbf{t}_{ij}^k \circ \mathbf{r}_{ij})\|_1 + \|\mathbf{t}_{ij} \circ \mathbf{t}_{ij} + \mathbf{r}_{ij}^k \circ \mathbf{r}_{ij}^k \\ &\quad - 2(\mathbf{t}_{ij} \circ \mathbf{r}_{ij}^k)\|_1 + \|2(\mathbf{t}_{ij}^k \circ \mathbf{r}_{ij}^k) - \mathbf{t}_{ij}^k \circ \mathbf{t}_{ij} - \mathbf{r}_{ij}^k \circ \mathbf{r}_{ij}^k\|_1 \\ &= \|\mathbf{t}_{ij}^k - \mathbf{r}_{ij}^k\|_2^2 + \|\mathbf{t}_{ij} - \mathbf{r}_{ij}^k\|_2^2 + \text{const} \end{aligned} \quad (29)$$

By substituting (29) and (15) where \mathbf{t}_{ij}^k and \mathbf{r}_{ij}^k are solutions obtained from the previous iterations and are constant with respect to the problem of the current iteration, we have derived the overall problem after applying DQA. Note that we have omitted the constant term because the solution of the objective function will not be affected by a constant.

$$\begin{aligned} \min_{\bar{\mathbf{x}}_{11}, \dots, \bar{\mathbf{x}}_{NM}} & \sum_{i=1}^N \sum_{j \in \mathcal{E}_i} f_{ij}(\bar{\mathbf{x}}_{ij}) + \sum_{i=2}^N \sum_{j \in \mathcal{E}_i} (\boldsymbol{\lambda}_{ij}^T (\mathbf{t}_{ij} - \mathbf{r}_{ij}) + \|\mathbf{w}_{ij} \circ (\mathbf{t}_{ij}^k - \mathbf{r}_{ij})\|_2^2 \\ & + \|\mathbf{w}_{ij} \circ (\mathbf{t}_{ij} - \mathbf{r}_{ij}^k)\|_2^2) \\ \text{s.t. } & \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) \leq \mathbf{0} \\ & \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) = \mathbf{0} \\ \text{where } & \bar{\mathbf{x}}_{ij} = [\mathbf{x}_{ij}, \mathbf{t}_{(i+1)k}] \quad \forall k \in \mathcal{C}_{ij}] \\ & \mathbf{r}_{ij} = \mathbf{a}_{ij}(\bar{\mathbf{x}}_{ij}) \\ & \forall j \in \mathcal{E}_i, i = 1, \dots, N \end{aligned} \quad (30)$$

And for each subproblem,

$$\begin{aligned} \min_{\bar{\mathbf{x}}_{ij}} & f_{ij}(\bar{\mathbf{x}}_{ij}) + \boldsymbol{\lambda}_{ij}^T (-\mathbf{r}_{ij}) + \|\mathbf{w}_{ij} \circ (\mathbf{t}_{ij}^k - \mathbf{r}_{ij})\|_2^2 \\ & + \sum_{k \in \mathcal{C}_{ij}} (\boldsymbol{\lambda}_{(i+1)k}^T \mathbf{t}_{(i+1)k} + \|\mathbf{w}_{(i+1)k} \circ (\mathbf{t}_{(i+1)k} - \mathbf{r}_{(i+1)k}^k)\|_2^2) \\ \text{s.t. } & \mathbf{g}_{ij}(\bar{\mathbf{x}}_{ij}) \leq \mathbf{0} \\ & \mathbf{h}_{ij}(\bar{\mathbf{x}}_{ij}) = \mathbf{0} \\ \text{where } & \bar{\mathbf{x}}_{ij} = [\mathbf{x}_{ij}, \mathbf{t}_{(i+1)k}] \quad \forall k \in \mathcal{C}_{ij}] \\ & \mathbf{r}_{ij} = \mathbf{a}_{ij}(\bar{\mathbf{x}}_{ij}) \\ & \forall j \in \mathcal{E}_i, i = 1, \dots, N \end{aligned} \quad (31)$$

The DQA approach consists of an inner loop and an outer loop. The inner loop is used to improve linearization while the outer loop is used to implement the method of multipliers. We concurrently update the penalty weight using the linear update method, similar to the AL approach. The general DQA algorithm consists of the following steps:

1. Initialize $\bar{\mathbf{x}}$, $\boldsymbol{\lambda}$, and \mathbf{w} , and set $\kappa=0$, where κ denotes the number of outer loop iteration.
2. Given $\bar{\mathbf{x}}^\kappa$, the final solution of the previous $(\kappa-1)$ st outer loop iteration upon inner loop convergence, set $s=0$, where s is the inner loop iteration, and $\bar{\mathbf{x}}^{\kappa+1,0} = \bar{\mathbf{x}}^\kappa$, where $\bar{\mathbf{x}}^{\kappa+1,s}$ is the solution of the s th inner loop iteration and the current outer loop iteration.
3. For each element, solve for $\bar{\mathbf{x}}_{ij}$ in (31) in parallel, and obtain $\bar{\mathbf{x}}^{\kappa+1,s+1}$.
4. If $\max(\|\mathbf{t}^{\kappa+1,s+1} - \mathbf{t}^{\kappa+1,s}\|, \|\mathbf{r}^{\kappa+1,s+1} - \mathbf{r}^{\kappa+1,s}\|) \leq \sigma_{\text{inner}}$, where σ_{inner} is the inner loop consistency deviation tolerance, set $\bar{\mathbf{x}}^{\kappa+1} = \bar{\mathbf{x}}^{\kappa+1,s+1}$, and go to step 5. Otherwise, set $\bar{\mathbf{x}}^{\kappa+1,s+1} = \bar{\mathbf{x}}^{\kappa+1,s} + \tau(\bar{\mathbf{x}}^{\kappa+1,s+1} - \bar{\mathbf{x}}^{\kappa+1,s})$, where τ is the step size, set $s = s+1$, and go to step 3.
5. If $\max(\|\mathbf{t}^{\kappa+1} - \mathbf{t}^\kappa\|, \|\mathbf{r}^{\kappa+1} - \mathbf{r}^\kappa\|) \leq \sigma_{\text{outer}}$, where σ_{outer} is the outer loop consistency deviation tolerance, then stop and set the optimal solution $\bar{\mathbf{x}}^*$ to be $\bar{\mathbf{x}}^{\kappa+1}$; otherwise, set $\kappa = \kappa+1$ and update the Lagrange multipliers using the method of multipliers by setting $\boldsymbol{\lambda}^{\kappa+1} = \boldsymbol{\lambda}^\kappa + \mathbf{w}^\kappa \circ (\mathbf{t}^\kappa - \mathbf{r}^\kappa)$, $\mathbf{w}^{(\kappa+1)} = \beta \mathbf{w}^\kappa$, and go to step 2.

This procedure is shown in Fig. 7. Note that the consistency deviation tolerance for both the inner loop σ_{inner} and outer loop σ_{outer} should be significantly smaller than the step size τ to prevent premature convergence.

Convergence is proven for the above algorithm with $\sigma_{\text{inner}}=0$, as stated in [20]. For practical purposes, however, we allow σ_{inner} to be small but nonzero.

THEOREM 8. *Suppose that the constraint sets are bounded, $\sigma_{\text{inner}}=0$, and the step size $\tau \in (0, 1)$ is significantly small. The following statements hold*

1. $\lim_{s \rightarrow \infty} (\mathbf{t}^{\kappa,s+1} - \mathbf{t}^{\kappa,s}) = \mathbf{0}$, $\lim_{s \rightarrow \infty} (\mathbf{r}^{\kappa,s+1} - \mathbf{r}^{\kappa,s}) = \mathbf{0}$

2. Each limit point of the sequence $\{\bar{\mathbf{x}}^{\kappa,s}\}_{s=0}^{\infty}$ is a solution of $\hat{\Lambda}(\boldsymbol{\lambda}^{\kappa})$.

The proof of the theorem is based on the estimation of the difference between the approximation and the augmented Lagrangian function. This estimation shows that a progress in the former also introduces progress in the later. The step size τ plays an important role in the DQA method. Because linearization is only accurate in a neighborhood of the point at which the linearization takes place, we use a small step size to obtain a good approximation. As noted in [20], the step size τ is related to the number of linking variables. Proof of the above theorem requires smaller τ for problems with a larger number of linking variables. However, empirical results indicate that a step size close to 1 can bring convergence to all example problems presented in this paper.

Local convergence is also discussed in [20], and this property of DQA is based on the quadratic growth condition of the augmented Lagrangian function, defined as follows:

DEFINITION 2 (QUADRATIC GROWTH CONDITION). Let

$$f_{\text{AL}} = \sum_{i=1}^N \sum_{j \in \mathcal{E}_i} f_{ij}(\bar{\mathbf{x}}_{ij}) + \sum_{i=2}^N \sum_{j \in \mathcal{E}_i} (\boldsymbol{\lambda}_{ij}^T (\mathbf{t}_{ij} - \mathbf{r}_{ij}) + \|\mathbf{w}_{ij} \circ (\mathbf{t}_{ij} - \mathbf{r}_{ij})\|_2^2)$$

and $\bar{\mathbf{x}}^*$ is the solution of (15), where $\bar{\mathbf{x}} = [\bar{\mathbf{x}}_{11}, \dots, \bar{\mathbf{x}}_{NM}]$. The augmented Lagrangian function is said to satisfy the quadratic growth condition if there exists $\gamma > 0$ such that for every $\bar{\mathbf{x}}$,

$$f_{\text{AL}}(\bar{\mathbf{x}}) - f_{\text{AL}}(\bar{\mathbf{x}}^*) \geq \gamma \|\bar{\mathbf{x}} - \bar{\mathbf{x}}^*\|_2^2$$

Under the quadratic growth rate of the augmented Lagrangian function, the rate of convergence for DQA is linear. However, other factors, such as a large number of linking variables can slow down the progress [20]. As a variation, it is possible to linearize at the midpoint of the solution from the previous iteration $[(\mathbf{t}^{\kappa} + \mathbf{r}^{\kappa})/2, (\mathbf{t}^{\kappa} + \mathbf{r}^{\kappa})/2]$. We obtain similar results as linearization at end points in all our test examples.

Truncated Diagonal Quadratic Approximation Method. As will be discussed in the next section, DQA performs well on test problems in terms of throughput, but the overall computational cost is still high. From the experiments, we have found out that much of the computational effort is spent on the inner linearization loop. We have also observed that usually the inner loop progresses slowly and introduces a high cost to reach the desirable inner loop convergence tolerance. However, when the Lagrange multipliers are not optimal, high accuracy of the subproblem solutions is not necessary, and the computational effort is wasted. It is more desirable to quickly update the Lagrange multiplier to move toward its optimal value. This can be achieved by limiting the total number of inner loop iterations in DQA by treating it as a user-specified parameter in a way that is similar to the ALAD approach, which reduces the computational cost for solving the inner loop. Because some inner loop calculation is omitted in DQA, this method is called the truncated diagonal quadratic approximation method (TDQA). The outer loop of TDQA implements the method of multipliers. Similar to ALAD, the penalty weight is held fixed.

In order to minimize overall cost, only one iteration of computation is carried out in the inner loop. The TDQA algorithm is given as follows:

1. Initialize $\bar{\mathbf{x}}^0, \boldsymbol{\lambda}$, and \mathbf{w} , and set $\kappa=0$, where κ denotes the number of loop iterations.
2. For each element, solve for $\bar{\mathbf{x}}_{ij}$ in (31) in parallel, and obtain $\bar{\mathbf{x}}_{ij}^{\kappa+1} = \bar{\mathbf{x}}_{ij}^{\kappa} + \tau(\bar{\mathbf{x}}_{ij}^{\kappa+1} - \bar{\mathbf{x}}_{ij}^{\kappa})$ where τ is the step size.
3. If $\max(\|\mathbf{t}^{\kappa+1} - \mathbf{t}^{\kappa}\|, \|\mathbf{r}^{\kappa+1} - \mathbf{r}^{\kappa}\|) \leq \sigma$, where σ is the outer loop consistency deviation tolerance, then stop, and set the optimal solution $\bar{\mathbf{x}}^*$ to be $\bar{\mathbf{x}}^{\kappa+1}$; otherwise, set $\kappa = \kappa + 1$, update the Lagrange multipliers using the method of multipliers by

setting $\boldsymbol{\lambda}^{\kappa+1} = \boldsymbol{\lambda}^{\kappa} + \mathbf{w}^{\kappa} \circ (\mathbf{t}^{\kappa} - \mathbf{r}^{\kappa})$, $\mathbf{w}^{(\kappa+1)} = \beta \mathbf{w}^{\kappa}$, and go to step 2.

This procedure is shown in Fig. 7. The outer loop tolerance σ should be significantly smaller than the step size τ to prevent premature convergence. Furthermore, we can also linearize at the midpoint of the solution obtained from the previous iteration. Again, empirical results show similar performance as linearization at end points.

Empirical results show promising results of the TDQA method. The intuition is that TDQA can be considered as an approximation of the ALAD method. As long as the approximation is accurate, solution of TDQA at each iteration is close to that of the ALAD method, and convergence follows by the convergence property of the ALAD method. We can also impose a strategy for increasing the total number of inner loop iterations when the improvement of the actual function is not large enough. In the extreme case when we allow the limit of the number of inner loop iterations to achieve infinity, the method becomes the DQA method and convergence can be obtained. A proof of convergence for finite inner loop truncation is left for future work.

DQA and TDQA With the Trust Region Technique. From the discussion above, it is possible to have a linearization at the current point that is a poor approximation to the augmented Lagrangian function. A small step size τ can ensure accuracy of the approximation. However, if the step size is too small, convergence will be slow. As a result, some trial and error is required for finding a good value of the step size for practical applications.

An alternative approach for finding a good step size is to use a trust region. To determine whether the current linearization is accurate, we estimate the ratio ρ of the actual reduction (the reduction of the original augmented Lagrangian function) to the predicted reduction (the reduction of the diagonal quadratic approximation of the augmented Lagrangian function) between current solution $\bar{\mathbf{x}}_{ij}^{\kappa+1}$ and previous solution $\bar{\mathbf{x}}_{ij}^{\kappa}$. If ρ is close to 1, the linearization is accurate. On the other hand, whenever the ratio ρ is far from 1, it indicates that the step size for the current solution is too big, the linearization is not accurate, and it is desirable to use a smaller step size. The algorithm of DQA and TDQA can be modified according to this method in the following way. Instead of using a predetermined step size, we calculate the reduction ratio ρ after each iteration. If $\rho \geq \gamma$, where γ is the ratio threshold, then accept the solution. Otherwise, shrink the step size by half. The recommended value of γ is 0.25 in the standard trust region method. This procedure can also be carried out iteratively.

Results and Comparisons

For the purpose of a clear comparison among all ATC methods, we use the four examples discussed in [12], in the same order. Readers are referred to [12] for a more in-depth discussion on the structures and properties of each problem. To make sure that the comparison is fair, we use the same condition on terminating the update of the Lagrange multipliers for all methods, as shown in the following:

$$\max(\|\mathbf{t}^{\kappa} - \mathbf{t}^{\kappa-1}\|, \|\mathbf{r}^{\kappa} - \mathbf{r}^{\kappa-1}\|) \leq \sigma_{\text{outer}} \quad (32)$$

This is different from [12], where the termination condition is

$$\max(\|(\mathbf{t}^{\kappa} - \mathbf{r}^{\kappa}) - (\mathbf{t}^{\kappa-1} - \mathbf{r}^{\kappa-1})\|) \leq \sigma_{\text{outer}} \quad (33)$$

We use a new criterion because (33) might result in premature convergence.³

³Generally speaking, three convergence criteria are widely used in the implementation of optimization methods: The gradient of the Lagrangian function is close to 0, objective function value stops changing, and solution point stops changing. We are using the third method here, and for practical purposes, a small nonzero convergence tolerance is given to determine when the solution stops changing.

In this paper, three quantities are used to evaluate performance: The solution error, the total number of function evaluations, and the overall computational throughput. Solution error ξ is defined as follows, the same as [12]:

$$\xi = \|\mathbf{x}^* - \mathbf{x}^{\text{ATC}}\|_{\infty},$$

where \mathbf{x}^* is the actual optimal solution and \mathbf{x}^{ATC} is the solution found by ATC. The magnitude of the solution error is controlled by the convergence tolerance. The tighter the tolerance, the more accurate the solution. The total number of function evaluations, also defined in [12], is used to measure the overall computational cost, and is reported by the subproblem solver, in this case, TOMLAB NPSOL. Computational throughput is defined as the number of tasks that can be performed for a given time period. Because each problem is considered a task, for simplicity we use the inverse of throughput, or latency in CPU seconds, to compare the speed for solving each problem of each method. In other words, latency measures the total CPU time required to solve the overall problem, taking into account parallelization.

We compare the performance of each method in the following two ways: function evaluations versus solution error and computational latency versus solution error.⁴ The methods being compared are the quadratic penalty method with the nested coordination scheme and WUM (QP); the quadratic penalty method with BCD and WUM (QP-BCD); the augmented Lagrangian method with nested coordination scheme (AL); the augmented Lagrangian method with BCD (AL-BCD); ALAD; DQA; and TDQA. The value of β for updating penalty weights is set to 2.0 for AL and DQA, and 1.0 for ALAD and TDQA. w step size for DQA is 0.9 for all problems. For TDQA, we invoke the trust region once to ensure the accuracy of the linearization and the initial step size is set to 0.7. Note that OL is not included because it encounters convergence difficulties for the last three problems. All tests are carried out on the same machine with the same applications; thus, comparison is fair.⁵

Example 1. The first example is a two-level decomposition of the geometric programming problem. It has a unique solution $\mathbf{z}^* = [2.15, 2.06, 1.32, 0.76, 1.07, 1.00, 1.47]$ with all constraints active.

$$\begin{aligned} \min_{z_1, \dots, z_7} f &= f_1 + f_2 = z_1^2 + z_2^2 \\ \text{s.t.} \quad g_1 &= (z_3^{-2} + z_4^{-2})z_5^{-2} - 1 \leq 0 \\ g_2 &= (z_5^2 + z_6^{-2})z_7^{-2} - 1 \leq 0 \\ h_1 &= (z_3^2 + z_4^{-2} + z_5^2)z_1^{-2} - 1 = 0 \\ h_2 &= (z_5^2 + z_6^2 + z_7^2)z_2^{-2} - 1 = 0 \\ z_1, z_2, \dots, z_7 &\geq 0 \end{aligned} \quad (34)$$

Figure 2 shows the structure of the problem and Fig. 3 shows results for each method. Note that since this is a two-level problem, QP and QP-BCD are identical, as are AL and AL-BCD. It is clear that the truncated methods ALAD and TDQA perform better than all the nested loop methods in terms of both function evaluations and latency. TDQA outperforms ALAD. DQA requires more function evaluations than AL but has a lower latency due to parallelization. QP does not perform as well as the other methods.

⁴We did not use the number of redesigns (i.e., the number of times each subproblem must be solved) as a metric because it can be misleading. In some iterations the subproblem may take many function evaluations to solve, whereas in many of the later iterations it will be very fast because the starting point is close to the solution. Thus, the redesigns metric does not seem to be an accurate or easy-to-interpret metric of computational cost.

⁵Application: MATLAB Version 7.0 with TOMLAB NPSOL SOLVER Version 5.3; OS: SUSE Linux; Processor: Intel(R) Xeon(TM); CPU: 2.80 GHz. Also, for the test results we used the time required to complete the longest running subproblem in each iteration to measure the effect of imperfect load balancing. We did not explicitly capture the communication overhead of multiprocessors, but this aspect will be very small in the examples compared to the computation time required at each iteration.

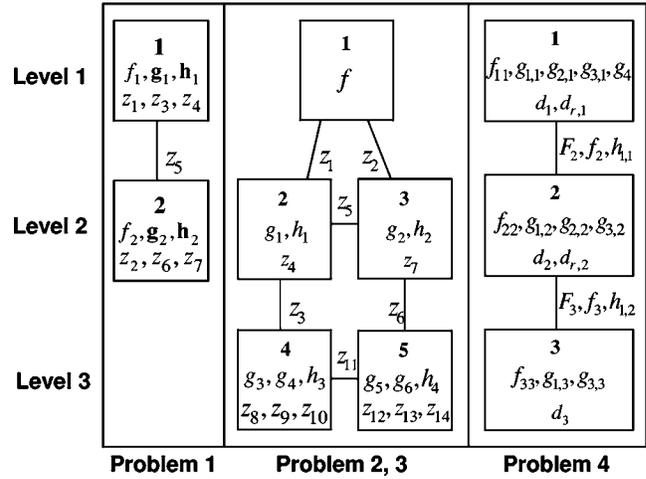


Fig. 2 Example problem structures

Example 2. The second example problem is a three-level decomposition of a posynomial geometric programming problem. Its unique solution is $\mathbf{z}^* = [2.84, 3.09, 2.36, 0.76, 0.87, 2.81, 0.94, 0.97, 0.87, 0.80, 1.30, 0.84, 1.76, 1.55]$ with all constraints active

$$\begin{aligned} \min_{z_1, \dots, z_{14}} f &= f_1 + f_2 = z_1^2 + z_2^2 \\ \text{s.t.} \quad g_1 &= (z_3^{-2} + z_4^{-2})z_5^{-2} - 1 \leq 0 \\ g_2 &= (z_5^2 + z_6^{-2})z_7^{-2} - 1 \leq 0 \\ g_3 &= (z_8^2 + z_9^{-2})z_{11}^{-2} - 1 \leq 0 \\ g_4 &= (z_8^{-2} + z_{10}^2)z_{11}^{-2} - 1 \leq 0 \\ g_5 &= (z_{11}^2 + z_{12}^{-2})z_{13}^{-2} - 1 \leq 0 \\ g_6 &= (z_{11}^2 + z_{12}^2)z_{14}^{-2} - 1 \leq 0 \\ h_1 &= (z_3^2 + z_4^{-2} + z_5^2)z_1^{-2} - 1 = 0 \\ h_2 &= (z_5^2 + z_6^2 + z_7^2)z_2^{-2} - 1 = 0 \\ h_3 &= (z_8^2 + z_9^{-2} + z_{10}^{-2} + z_{11}^2)z_3^{-2} - 1 = 0 \\ h_4 &= (z_{11}^2 + z_{12}^2 + z_{13}^2 + z_{14}^2)z_6^{-2} - 1 = 0 \\ z_1, z_2, \dots, z_{14} &\geq 0 \end{aligned} \quad (35)$$

Figure 2 shows the structure of the problem, and Fig. 4 shows the results.

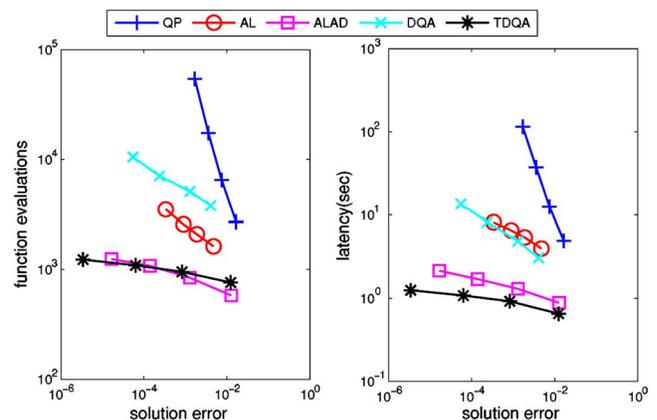


Fig. 3 Example 1: Computation cost and latency versus solution accuracy

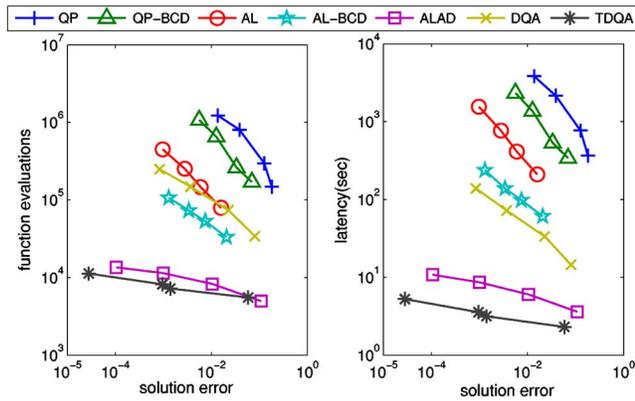


Fig. 4 Example 2: Computational cost and latency versus solution accuracy

In this example, the advantage of BCD is clear because QP-BCD and AL-BCD outperform QP and AL, respectively. Moreover, AL outperforms QP and QP-BCD. DQA outperforms AL and AL-BCD in computational latency due to parallelization. Similar to example 1, the truncated methods are superior to all nested loop methods, and TDQA outperforms DQA in terms of both function evaluations and latency.

Example 3. The third example is a slight variation of the second, changing only the objective function so the targets are attainable. The objective function is now $f=f_1+f_2=(z_1-2.9)^2+(z_2-3.1)^2$. In this example, only $z_1=2.9$ and $z_2=3.1$ are unique. The other variables can have different values as long as all constraints are satisfied.

This example is a special case. As discussed in [12], the Lagrange multipliers were originally set to zero, which happens to be the optimal value for this problem. As a result, the outer loops of all nested loop methods are theoretically unnecessary. However, due to numerical errors, the outer loops are updated several times in practice.

Similar to example 2, QP-BCD and AL-BCD outperform QP and AL, respectively. However, because λ is at its optimal value originally, the advantage of the Lagrangian term cannot be observed, and the QP methods outperform the AL methods in this example. For the truncated methods ALAD and TDQA, which outperform all nested loop methods, ALAD outperforms TDQA in terms of both function evaluations and latency. These results are shown in Fig. 5. The total computational cost, and the variation among methods, is small for this problem because updating λ and w is unnecessary.

Example 4. The fourth example is a structured optimization

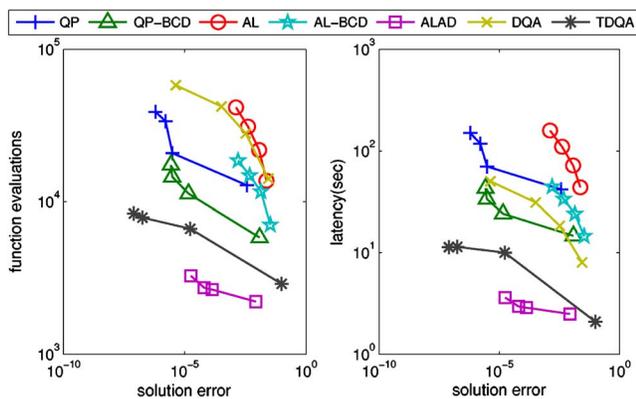


Fig. 5 Example 3: Computational cost and execution latency versus solution accuracy

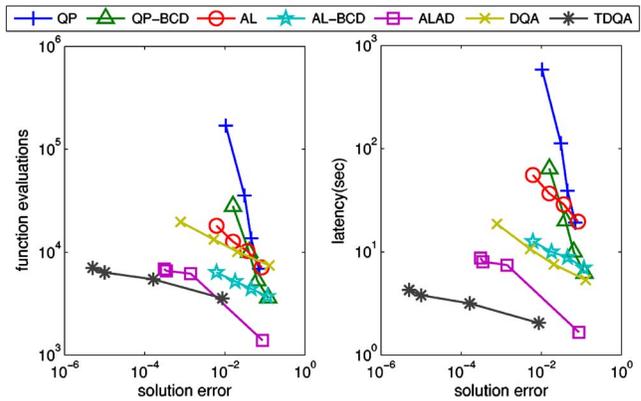


Fig. 6 Example 4: Computation cost and latency versus solution accuracy

problem based on the analytical mass allocation problem. It has a unique solution $\mathbf{z}^*=[0.0346, 0.0349, 0.0294, 0.0046, 0.0028]$

$$\min_{d_1, d_2, d_{r,1}, d_{r,2}} \sum_{i=1}^3 m_i + \sum_{j=1}^2 m_{r,j}$$

$$\text{s.t. } g_{1,i} = \sigma_{b,i} - \bar{\sigma} \leq 0, i = 1, 2, 3$$

$$g_{2,j} = \sigma_{a,j} - \bar{\sigma} \leq 0, j = 1, 2$$

$$g_{3,i} = F_{t,i} - \bar{F}_t \leq 0, i = 1, 2, 3$$

$$g_4 = f_1 - \bar{f}_1 \leq 0$$

$$h_1 = f_i - f_{i+1} - f_{r,i} = 0, i = 1, 2$$

$$\text{where } m_i = \frac{\pi}{4} d_i^2 L \rho, i = 1, 2, 3$$

$$m_{r,j} = \frac{\pi}{4} d_{r,j}^2 L \rho, j = 1, 2$$

$$\sigma_{b,i} = \frac{32L(F_i - F_{i+1})}{\pi d_i^3}, i = 1, 2, 3$$

$$f_i = \frac{64L^3(F_i - F_{i+1})}{3\pi E d_i^4}, i = 1, 2, 3$$

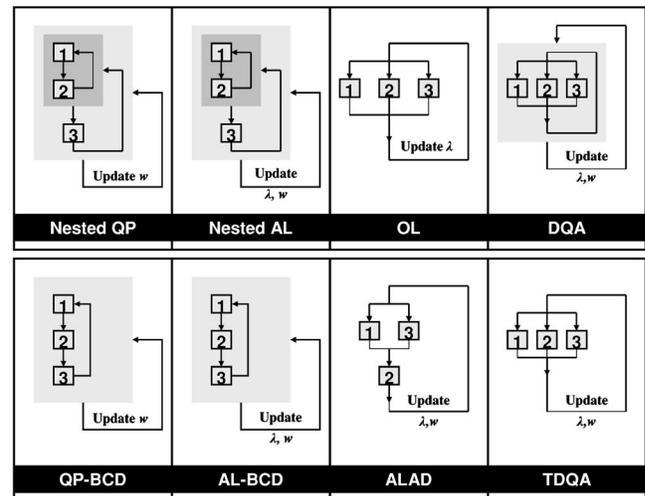


Fig. 7 Flow charts of methods

Table 1 Summary of methods

Consistency Constraint Relaxation Method	Refs	Separability	Nested Loops	Global Convergence Properties	Performance (latency in sec at 10 ⁻⁴ accuracy)			
					Ex 1	Ex 2	Ex 3	Ex 4
QP Quadratic Penalty	[1, 10, 11, 14]	Separable constraints => block coordinate descent	Nested subsystem or BCD loops, w update loop	Optimal solution as $w \rightarrow \infty$, but ill conditioning problems	>114	>2335 (BCD) >3827 (nested)	>70 (BCD) >24 (nested)	>583 (BCD) >63 (nested)
OL Ordinary Lagrangian Relaxation	[15]	Separable subsystems that can be solved in parallel	λ update loop	Optimal solution if no duality gaps, which is not always guaranteed	3	does not converge	does not converge	does not converge
AL Augmented Lagrangian	[5, 12, 16]	Separable constraints => block coordinate descent	Nested subsystem or BCD loops, λ and w update loop	Optimal solution if objective function and constraint functions are continuous, and the sequence of λ is bounded	8	>237 (BCD) 1549 (nested)	>159 (BCD) >44 (nested)	>55 (BCD) >12 (nested)
ALAD Alternating Direction Augmented Lagrangian	[5, 12]	Solve all elements in odd levels in parallel, then all elements in even levels in parallel.	λ and w update loop	Optimal solution if objective function and constraint functions are convex, and w is bounded	2	11	3	8
DQA Diagonal Quadratic Approximation	[17, 20]	Separable subsystems that can be solved in parallel	λ and w update loop	Optimal solution if objective function and constraint functions are convex, feasible region is bounded and closed, and step size is sufficiently small	8	138	31	17
TDQA Truncated Diagonal Quadratic Approximation		Separable subsystems that can be solved in parallel	λ and w update loop	not known	1	4	4	3

$$\sigma_{a,j} = \frac{4F_{j+1}}{\pi d_{rj}^2}, \quad j = 1, 2,$$

$$f_{r,j} = \frac{4F_{j+1}L}{\pi E d_{rj}^2}, \quad j = 1, 2$$

$$\rho = 2700 \text{ kg/m}^3, \quad E = 70 \text{ GPa}, \quad F_1 = 1000 \text{ N} \quad (36)$$

Figure 2 shows the structure of the problem. The results shown in Fig. 6 are very similar to those in example 2. QP-BCD and AL-BCD outperform QP and AL, respectively. The AL methods outperform the QP methods. DQA outperforms AL: It has a larger number of function evaluations than AL-BCD, but it has the smallest latency among all nested loop methods. The truncated methods are superior to all nested loop methods, and TDQA outperforms ALAD.

Discussion

The biggest advantage of both DQA and TDQA is their ability to separate subproblems for parallel computation. This property is highly desirable, especially in large-scale problems. Because each subproblem is separable, nested loops are avoided. Moreover, because each problem can be solved in parallel, throughput is greatly increased. Although we only present DQA and TDQA in the context of hierarchical ATC problems, these methods can be used in a wide variety of decomposition methods. DQA was first proposed in [20] as a method solving general block-angular structured problems, and TDQA can be applied to any problems to which DQA can be applied. For example, in [5], an augmented Lagrangian approach is used in solving nonhierarchical dual-block angular problems, and DQA and TDQA can also be used for solving these problems. This is also an advantage of DQA and TDQA over ALAD. DQA and TDQA do not require problems to be hierarchical, whereas ALAD requires identification of independent levels in hierarchical problems. It is expected that DQA and TDQA will perform similarly in nonhierarchical decomposition schemes as in ATC.

One limitation of DQA and TDQA is that, theoretically, these methods achieve good performance only when the number of linking variables is small, as in the case of quasiseparable prob-

lems. Further investigation is needed to determine their performance in practice on problems that have large number of linking variables after decomposition. Future research is also needed to investigate the global and local convergence properties of the TDQA method.

Conclusion

In this paper, we have summarized the existing methods, compared the block coordinate descent method to nested schemes, and presented the diagonal quadratic approximation method and the truncated diagonal quadratic approximation method used in handling consistency constraint relaxation in ATC. The DQA method is supported by theoretical justification. Fundamental results show that by using BCD instead of the nested coordination scheme, computational cost and throughput are greatly improved. BCD is also theoretically justified. In general, DQA has lower computational cost than most nested loop methods and it has the highest throughput due to parallelization. TDQA has the least total cost and best throughput of all methods in all the test examples, except for Example 3, for which ALAD has achieved the best performance. In this degenerate case, updates of the penalty weights and Lagrange multipliers are unnecessary, and the total computational time is much smaller than the other nondegenerate cases, making distinctions less critical. As a result, we believe that TDQA is more preferable when parallel processing is available. DQA has proofs for local and global convergence, and proofs for TDQA are left for future work. The flowcharts for the algorithms of all methods are shown in Fig. 7, and Table 1 summarises methods and results.

The proposed methods overcome many of the concerns with prior approaches to ATC, such as convergence difficulties, ill-conditioning, and computational cost associated with the coordination strategies. DQA and TDQA enable parallelization of all subsystems in the ATC hierarchy, which can improve computational throughput when parallel processing is available. Given the theoretical benefits and promising empirical results of the DQA and TDQA approaches to ATC, we hope to see these methods utilized in future ATC studies and applications. Additionally, because DQA and TDQA do not require the assumption of hierar-

chical structures, extension to nonhierarchical structures has the potential to extend applicability to a wider range of systems optimization problems.

Acknowledgment

The authors wish to thank Simon Tosserams for sharing his code and both Simon Tosserams and Larry Biegler for their helpful feedback. This work was supported in part by the Pennsylvania Infrastructure Technology Alliance, a partnership of Carnegie Mellon, Lehigh University, and the Commonwealth of Pennsylvania's Department of Community and Economic Development (DCED), and by a SFU President's Research Grant and NSERC Discovery Grant.

References

- [1] Kim, H. M., Michelena, N. F., Papalambros, P. Y., and Jiang, T., 2003, "Target Cascading in Optimal System Design," *ASME J. Mech. Des.*, **125**(3), pp. 474–480.
- [2] Allison, J., Kokkolaras, M., Zawaislak, M., and Papalambros, P., 2005, "On the Use of Analytical Target Cascading and Collaborative Optimization for Complex System Design," *Proceedings of 6th World Congress on Structural and Multidisciplinary Optimization*, Rio de Janeiro, Brazil.
- [3] DeMiguel, A., and Murray, W., 2006, "A Local Convergence Analysis of Bilevel Decomposition Algorithms," *Optim. Eng.*, **7**, pp. 99–133.
- [4] Haftka, R. T., and Watson, L. T., 2005, "Multidisciplinary Design Optimization With Quasiseparable Subsystems," *Optim. Eng.*, **6**(1), pp. 9–20.
- [5] Tosserams, S., Etman, L. F. P., and Rooda, J. E., 2007, "An Augmented Lagrangian Decomposition Method for Quasi-Separable Problems in MDO," *Struct. Multidiscip. Optim.*, **34**, pp. 211–227.
- [6] Kim, H. M., Kokkolaras, M., Louca, L. S., Delagrammatikas, G. J., Michelena, N. F., Filipi, Z. S., Papalambros, P. Y., Stein, J. L., and Assanis, D. N., 2002, "Target Cascading in Vehicle Redesign: A Class VI Truck Study," *Int. J. Veh. Des.*, **29**(3), pp. 199–225.
- [7] Choudhary, R., Malkawi, A., and Papalambros, P. Y., 2005, "Analytic Target Cascading in Simulation-Based Building Design," *Autom. Constr.*, **14**(4), pp. 551–568.
- [8] Michalek, J. J., Ceryan, O., Papalambros, P. Y., and Koren, Y., 2006, "Balancing Marketing and Manufacturing Objectives in Product Line Design," *ASME J. Mech. Des.*, **128**(6), pp. 1196–1204.
- [9] Michalek, J. J., Feinberg, F. M., and Papalambros, P. Y., 2005, "Linking Marketing and Engineering Product Design Decisions via Analytical Target Cascading," *Journal of Product Innovation Management*, **22**, pp. 42–62.
- [10] Michelena, N., Park, H., and Papalambros, P. Y., 2003, "Convergence Properties of Analytical Target Cascading," *AIAA J.*, **41**(5), pp. 897–905.
- [11] Michalek, J. J., and Papalambros, P. Y., 2005, "Weights, Norms, and Notation in Analytical Target Cascading," *ASME J. Mech. Des.*, **127**(3), pp. 499–501.
- [12] Tosserams, S., Etman, L. F. P., and Rooda, J. E., 2006, "An Augmented Lagrangian Relaxation for Analytical Target Cascading Using the Alternating Directions Method of Multipliers," *Struct. Multidiscip. Optim.*, **31**(3), pp. 176–189.
- [13] Bertsekas, D. P., 2003, *Nonlinear Programming*, 2nd ed., Athena Scientific, Belmont, MA.
- [14] Michalek, J. J., and Papalambros, P. Y., 2005, "An Efficient Weighting Update Method to Achieve Acceptable Inconsistency Deviation in Analytical Target Cascading," *J. Qual. Maint. Eng.*, **127**(3), pp. 206–214.
- [15] Lassiter, J. B., Wiecek, M. M., and Andrighetti, K. R., 2005, "Lagrangian Coordination and Analytical Target Cascading: Solving ATC-Decomposed Problems With Lagrangian Duality," *Optim. Eng.*, **6**(3) pp. 361–381.
- [16] Kim, H. M., Chen, W., and Wiecek, M. M., 2006, "Lagrangian Coordination for Enhancing the Convergence of Analytical Target Cascading," *AIAA J.*, **44**(10), pp. 2197–2207.
- [17] Stephanopoulos, G., and Westerberg, A. W., 1975, "The Use of Hestenes' Method of Multipliers to Resolve Dual Gaps in Engineering System Optimization," *J. Optim. Theory Appl.*, **15**(3), pp. 285–309.
- [18] Nocedal, J., and Wright, S. J., 1999, *Numerical Optimization*, Springer Series in Operations Research, Springer, New York.
- [19] Bertsekas, D. P., and Tsitsiklis, J. N., 1989, *Parallel and Distributed Computation*, Prentice-Hall, Englewood Cliffs, NJ.
- [20] Ruzcynski, A., 1995, "On Convergence of an Augmented Lagrangian Decomposition Method for Space Convex Optimization," *Math. Op. Res.*, **20**(3), pp. 634–656.