

FAST BUNDLE-LEVEL METHODS AND APPLICATIONS TO IMAGE RECONSTRUCTION*

YUNMEI CHEN[†], GUANGHUI LAN[‡], YUYUAN OUYANG[§], AND WEI ZHANG[¶]

Abstract. In this paper, we study a new class of first-order methods, namely fast bundle-level type methods, which can utilize historical first-order information through cutting plane models to significantly accelerate the solutions of imaging problems. Recently, it has been shown in [20] that an accelerated prox-level (APL) method and its variant, the uniform smoothing level (USL) method, have optimal iteration complexity for solving black-box and structured convex programming problems without requiring input of any smoothness information. However, these algorithms require the assumption on the boundedness of the feasible set and their efficiency relies on the solutions of two involved subproblems. These hindered the applicability of these algorithms in solving large-scale and unconstrained optimization problems. In this paper, we first present two new variants of level methods, i.e., the fast APL (FAPL) method and the fast USL (FUSL) method, for solving large scale black-box and structured convex programming problems respectively. Both FAPL and FUSL enjoy the same optimal iteration complexity as APL and USL, while the number of subproblems in each iteration is reduced from two to one, and an exact method is presented to solve the only subproblem in these algorithms. Then, we introduced a generic algorithmic framework to extend these uniformly optimal level methods for solving unconstrained problems. As a result, the proposed FAPL and FUSL methods have improved the performance of APL and USL in practice significantly in terms of both computational time and solution quality. Our numerical results on solving some large-scale least squares problems and total variation based image reconstructions have shown great advantages of these new bundle-level type methods over APL, USL, and some other state-of-the-art first-order methods.

Keywords: convex programming, first-order, optimal method, bundle-level method, total variation, image reconstruction

AMS 2000 subject classification: 90C25, 90C06, 90C22, 49M37

1. Introduction. Many data analysis problems are often modeled as the following broad class of convex programming (CP) problems:

$$f^* := \min_{x \in X} f(x), \tag{1.1}$$

where $X \subseteq \mathbb{R}^n$ is a closed convex set, and $f : X \rightarrow \mathbb{R}$ is a convex function. For example, the classic ridge regression model (namely, Tikhonov regularization) in statistical learning estimates parameters β by

$$\min_{\beta \in \mathbb{R}^n} \|y - X\beta\|^2 \text{ subject to } \|\beta\| \leq \lambda, \tag{1.2}$$

where $\|\cdot\|$ denotes the Euclidean norm, y describes the observed outcome, X are the predictors in the observed data, and λ is a regularization parameter. The above model can be viewed as a special case of (1.1) with $x = \beta$, $f(x) = \|y - X\beta\|^2$, and $X = \{x \in \mathbb{R}^n : \|\beta\| \leq \lambda\}$. Another important example is the two-dimensional total variation based image reconstruction problem given by:

$$\min_{u \in \mathbb{R}^n} \frac{1}{2} \|Au - b\|^2 + \lambda \|u\|_{TV}, \tag{1.3}$$

where A is the measurement matrix, u is the n -vector form of a two-dimensional image to be constructed, b represents the observed data, $\|\cdot\|_{TV}$ is the discrete TV semi-norm, and λ is the regularization parameter. Problem (1.3) can also be casted as (1.1) by setting $x = u$, $f(x) = \|Au - b\|^2/2 + \lambda \|u\|_{TV}$, and $X = \mathbb{R}^n$. It is worth noting that while problem (1.2) has an Euclidean ball constraint, problem (1.3) is an unconstrained CP problem. Moreover, the objective function in (1.2) is smooth, while the objective function in (1.3) is defined as the summation of a smooth term $\|Au - b\|^2/2$ and a nonsmooth term $\lambda \|u\|_{TV}$.

*December, 2014. This research was partially supported by NSF grants CMMI-1254446, DMS-1319050, and ONR grant N00014-13-1-0036.

[†]Department of Mathematics, University of Florida (yun@math.ufl.edu).

[‡]Department of Industrial and System Engineering, University of Florida (glan@ise.ufl.edu).

[§]Department of Industrial and System Engineering, University of Florida (ouyang@ufl.edu).

[¶]Department of Mathematics, University of Florida (weizhang657@ufl.edu).

Due to the high dimensionality of x for many applications in data analysis and imaging, much recent research effort has been directed to the development of efficient first-order methods for solving (1.1). First-order methods use gradients (or subgradients) of f exclusively and hence possess significantly reduced iteration cost than second-order methods. The efficiency of these algorithms are often measured by their iteration complexity in terms of the number of (sub)gradient evaluations required to find an approximate solution of (1.1). In view of the classic complexity theory [24], for any first-order methods the number of (sub)gradient evaluations required to find an ϵ -solution of (1.1) (i.e., a point $x_\epsilon \in \mathbb{R}^n$ satisfying $f(x_\epsilon) - f^* \leq \epsilon$) cannot be smaller than $\mathcal{O}(1/\epsilon^2)$ if f is nonsmooth. This can be achieved, for example, by traditional subgradient methods. For a smooth f , the optimal iteration complexity is $\mathcal{O}(1/\sqrt{\epsilon})$, which can be achieved, for example, by Nesterov's accelerated gradient (AG) algorithms [27, 25, 26]. Recently, by adapting Nesterov's AG schemes and smoothing technique [27], several popular classes of first-order methods have been developed to significantly improve their iteration complexity bounds. For instance, the accelerated primal dual (APD) algorithm [6] and accelerated hybrid proximal extragradient algorithm [12], which exhibit the optimal iterative complexity for solving a broad class of saddle point problems, and several variants of accelerated ADMM (AADMM)[13, 9, 10, 22, 11, 32], have improved the iteration complexity regarding to the smooth component.

In this paper, we focus on a different class of first-order methods, i.e., bundle-level (BL) type methods, which can utilize historical first-order information through cutting plane models to significantly accelerate the numerical performance of the gradient descent type methods as mentioned above. In particular, we will present two new fast BL type methods that are able to uniformly solve smooth, weakly smooth, nonsmooth CPs and a broad class of saddle-point problems with optimal complexity bounds without requiring the input of any problem parameters. In addition, all these methods have cheap computational cost per iteration and can achieve moderate and even high solution accuracy for large-scale CPs in practice. We first discuss the ball-constrained problem, i.e., X is a given Euclidean ball in \mathbb{R}^n , then propose a generic framework to solve unconstrained problem, i.e., $X = \mathbb{R}^n$, through solutions to a series of ball-constrained CPs, and also demonstrate the advantages of these methods for the solutions of imaging/data analysis problems. Let us start by reviewing a few existing BL type methods.

1.1. Cutting plane, bundle and bundle-level methods. The bundle-level method originated from the well-known Kelley's cutting-plane method in 1960 [14]. Consider the convex programming problem

$$f_X^* := \min_{x \in X} f(x), \quad (1.4)$$

where X is a compact convex set and f is a closed convex function. The fundamental idea of the cutting plane method is to generate a sequence of piecewise linear functions to approximate f on X . In particular, given $x_1, x_2, \dots, x_k \in X$, this method approximates f by

$$m_k(x) := \max\{h(x_i, x), 1 \leq i \leq k\}, \quad (1.5)$$

and computes the iterate x_{k+1} by

$$x_{k+1} \in \text{Argmin}_{x \in X} m_k(x), \quad (1.6)$$

where

$$h(z, x) := f(z) + \langle f'(z), x - z \rangle, \quad (1.7)$$

and $f'(x) \in \partial f(x)$, where $\partial f(x)$ denotes the subdifferential of f at x . Clearly, the functions m_i , $i = 1, 2, \dots$, satisfy $m_i(x) \leq m_{i+1}(x) \leq f(x)$ for any $x \in X$, and are identical to f at those search points x_i , $i = 1, \dots, k$. However, the inaccuracy and instability of the piecewise linear approximation m_k over the whole feasible set X may affect the selection of new iterates, and the above scheme converges slowly both theoretically and practically [24, 26]. Some important improvements of Kelley's method have been made in 1990s under the name of bundle methods (see, e.g., [15, 16, 21]). In particular, by incorporating the level sets into Kelley's method, Lemaréchal, Nemirovskii and Nesterov [21] proposed in 1995 the classic bundle-level (BL) method by performing a series of projections over the approximate level sets.

Given x_1, x_2, \dots, x_k , the classic BL iteration consists of the following three steps:

- a) Set $\bar{f}_k := \min\{f(x_i), 1 \leq i \leq k\}$ and compute a lower bound on f_X^* by $\underline{f}_k = \min_{x \in X} m_k(x)$.
- b) Set the level $l_k = \beta \underline{f}_k + (1 - \beta) \bar{f}_k$ for some $\beta \in (0, 1)$.
- c) Set $X_k := \{x \in X : m_k(x) \leq l_k\}$ and determine the new iterate by

$$x_{k+1} = \operatorname{argmin}_{x \in X_k} \|x - x_k\|^2. \quad (1.8)$$

In the BL method, the localizer X_k is used to approximate the level set $L_k := \{x : f(x) \leq l_k\}$, because the projection over L_k is often too difficult to compute. Intuitively, as k increases, the value of l_k will converge to f_X^* , and consequently both L_k and X_k will converge to the set of optimal solutions to problem (1.4). It is shown in [21] that the number of BL iterations required to find an ϵ -solution to problem (1.4), i.e., a point $\hat{x} \in X$ s.t. $f(\hat{x}) - f_X^* \leq \epsilon$, can be bounded by $\mathcal{O}(1/\epsilon^2)$, which is optimal for general non-smooth convex optimization in the black-box model.

Observe that for the above BL methods, the localizer X_k accumulates constraints, and hence the subproblem in Step c) becomes more and more expensive to solve. In order to overcome this difficulty, some restricted memory BL algorithms have been developed in [16, 4]. In particular, Ben-Tal and Nemirovski [4] introduced the non-Euclidean restricted memory level (NERML) method, in which the number of extra linear constraints in X_k can be as small as 1 or 2, without affecting the optimal iteration complexity. Moreover, the objective function $\|\cdot\|^2$ in (1.8) is replaced by a general Bregman distance $d(\cdot)$ for exploiting the geometry of the feasible set X . NERML is often regarded as state-of-the-art for large-scale non-smooth convex optimization as it substantially outperforms subgradient type methods in practice. Some more recent development of inexact proximal bundle methods and BL methods could be found in [34, 29, 28, 33, 18, 17, 8, 19].

1.2. Accelerated bundle-level type methods. While the classic BL method was optimal for solving non-smooth CP problems only, Lan [20] recently significantly generalized this method so that it can optimally solve any black-box CP problems, including non-smooth, smooth and weakly smooth CP problems. In particular, for problem (1.4) over compact feasible set X , the two new BL methods proposed in [20], i.e., the accelerated bundle-level (ABL) and accelerated prox-level (APL) methods, can solve these problems optimally without requiring any information on problem parameters. The ABL method can be viewed as an accelerated version of the classic BL method. Same as the classic BL method, the lower bound on f_X^* is estimated from the cutting plane model m_k in (1.5), the upper bound on f_X^* is given by the best objective value found so far. The novelty of the ABL method exists in that three different sequences, i.e., $\{x_k^l\}, \{x_k\}$ and $\{x_k^u\}$, are used for updating the lower bound, prox-center, and upper bound respectively, which leads to its accelerated iteration complexity for smooth and weakly smooth problems. The APL method is a more practical, restricted memory version of the ABL method, which also employs non-Euclidean prox-functions to explore the geometry of the feasible set X .

We now provide a brief description of the APL method. This method consists of outer-inner loops. Each outer iteration calls an inner loop, a gap reduction procedure, to reduce the gap between the lower and upper bounds on f_X^* by a constant factor. More specifically, at outer iteration s , given the input lower bound lb_s and current best point p_s which provides the input upper bound by $\text{ub}_s = f(p_s)$, the APL method first checks whether $\text{ub}_s - \text{lb}_s > \epsilon$. If so, it calls the gap reduction procedure in order to reduce the gap between the upper and lower bounds. The APL gap reduction procedure (or the inner loop) works as follows. It first sets $\underline{f}_0 = \text{lb}_s, \bar{f}_0 = \text{ub}_s, l = \beta \cdot \text{lb}_s + (1 - \beta) \text{ub}_s, X_0 = X, x_0^u = p_s, x_0 \in X$ be arbitrarily chosen, then iteratively performs the following steps for $k \geq 1$:

- a) Set $x_k^l = (1 - \alpha_k)x_{k-1}^u + \alpha_k x_{k-1}$ and $\underline{f}_k := \max\{\underline{f}_{k-1}, \min\{l, h_k\}\}$, where

$$h_k := \min_{x \in X_{k-1}} h(x_k^l, x). \quad (1.9)$$

- b) Update the prox-center x_k by

$$x_k = \operatorname{argmin}_{x \in X_k} d(x), \quad (1.10)$$

where $X_k := \{x \in X_{k-1} : h(x_k^l, x) \leq l\}$.

- c) Set $\bar{f}_k = \min\{\bar{f}_{k-1}, f(\tilde{x}_k^u)\}$, where $\tilde{x}_k^u = \alpha_k x_k + (1 - \alpha_k)x_{k-1}^u$, and x_k^u is chosen as either \tilde{x}_k^u or x_{k-1}^u such that $f(x_k^u) = \bar{f}_k$.
- d) Choose X_k such that $\underline{X}_k \subseteq X_k \subseteq \bar{X}_k$, where $\bar{X}_k := \{x \in X : \langle \nabla d(x_k), x - x_k \rangle \geq 0\}$.

Here $d(x)$ is the chosen prox-function and $\nabla d(\cdot)$ denotes its gradient. Observe that the parameters α_k and β are fixed a priori, and do not depend on any problem parameters. Moreover, the localizer X_k is chosen between \underline{X}_k and \bar{X}_k , so that the numbers of linear constraints in the two subproblems (1.9) and (1.10) can be fully controlled. It is shown in [20] that the APL method achieves the optimal iteration complexity uniformly for smooth, weakly smooth and non-smooth convex functions for solving problem (1.4). Moreover, by incorporating Nesterov's smoothing technique [27] into the APL method, Lan also presented in [20] that the uniform smoothing level (USL) method which can achieve the optimal complexity for solving an important class of non-smooth structured saddle point (SP) problems without requiring input of any problem parameters (see Subsection 2.2 for more details).

1.3. Contribution of this paper. One crucial problem associated with most existing BL type methods, including APL and USL, is that each iteration of these algorithms involves solving two optimization problems: first a linear programming problem to compute the lower bound, and then a constrained quadratic programming problem to update the prox-center or new iterate. In fact, the efficiency of these algorithms relies on the solutions to these two involved subproblems which are similar to (1.9) and (1.10), and the latter one is often more complicated than the projection subproblem in the gradient projection type methods. Moreover, most existing BL type methods require the assumption that the feasible set is bounded due to the following two reasons. Firstly, the feasible set has to be bounded to compute a meaningful lower bound by solving the aforementioned linear programming problem. Secondly, the iteration complexity analysis of those methods, such as the classical BL method [21], NERML [4], ABL, APL and USL [20], relies on the assumption that the feasible set is compact. These issues have significantly hindered the applicability of existing BL type methods. Our contribution in this paper mainly consists of the following three aspects.

Firstly, we propose the FAPL and FUSL methods that greatly reduced the computational cost per iteration of the APL and USL methods for solving ball-constrained CPs. The improvement is achieved mainly by eliminating the linear optimization subproblem for computing the lower bound and removing the ball constraint in the quadratic subproblem by properly choosing the prox-functions. More importantly we are able to show that with these simplifications the FAPL and FUSL methods can also achieve the optimal iteration complexity uniformly for smooth, nonsmooth and weakly smooth functions. It should be noted that there exist some variants of BL methods [5, 8, 3, 34] for solving non-smooth CP problems in which the computation of the subproblem for updating the lower bound is skipped, so that the feasible set X is allowed to be unbounded. For instance, the level bundle method in [3] updates the level parameter directly when the distance from the stability center to the newly generated iterate becomes larger than a chosen parameter. However, there is no iteration complexity analysis provided in any of these work.

Secondly, we introduce a simple exact approach to solve the dual problem of the only subproblem in our algorithms. By this, the accuracy of the FAPL and FUSL methods is significantly increased and the total number of iterations required to compute an ϵ -solution is greatly decreased comparing to the original APL and USL methods and other first-order methods. Also when the number of linear constraints is fixed and small, the computational cost for solving the subproblem depends on the dimension of the problem only linearly, since the cost for computing the vector inner product will dominate that of solving a few auxiliary linear systems. This feature is very important for solving large-scale CP problems.

Thirdly, we propose a novel algorithmic framework for solving unconstrained CP problems. The proposed framework solves unconstrained CP problems through solutions to a series of ball-constrained CPs and achieves the same order of the iteration complexity as the corresponding ball-constrained CPs. In particular, if there exists a uniformly optimal method (e.g., the APL and USL methods) that solves ball-constrained black-box or structured CPs, then the proposed algorithm solves unconstrained black-box or structured CPs with optimal complexity without requiring the input of any additional problem parameters. To the best of our knowledge, this is the first time in the literature that the complexity analysis has been performed for BL type methods to solve unconstrained CP problems (see Sections 3.2 and 3.3 in [30] for more details).

Finally, we apply the FAPL and FUSL methods to solve large-scale least squares problems and total variation based image reconstruction problems. Our experimental results show that these algorithms significantly outperform APL, NERML, Nesterov’s optimal method [27], the accelerated primal dual (APD) method [6], Nesterov’s smoothing (NEST-S) gradient method [27, 2] and the accelerated linearized alternating direction of multipliers (AL-ADMM) with line-search method [32], and even the powerful MATLAB solver for linear systems, especially when the dimension and/or the Lipschitz constant of the problem is large. Moreover, by using the FAPL and FUSL methods, we could achieve much more accurate solutions to the corresponding CPs, which results in better reconstructed image qualities and lower acquisition rates.

1.4. Organization of the paper. The paper is organized as follows. In Section 2, the new FAPL and FUSL methods are proposed followed by their convergence analysis, then an exact approach is introduced to solve the subproblem in these algorithms. In Section 3, we present a general scheme to extend the optimal methods for ball-constrained CPs to unconstrained CPs. The applications to quadratic programming and image processing problems are presented in Section 4.

2. Fast prox-level type methods for ball-constrained problems. In this section, we discuss the following ball-constrained CP:

$$f_{\bar{x},R}^* := \min_{x \in B(\bar{x},R)} f(x), \quad (2.1)$$

where $B(\bar{x}, R) := \{x \in \mathbb{R}^n : \|x - \bar{x}\| \leq R\}$ denotes the Euclidean ball centered at \bar{x} with radius R , and $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex function satisfying

$$f(y) - f(x) - \langle f'(x), y - x \rangle \leq \frac{M}{1 + \rho} \|y - x\|^{1+\rho}, \quad \forall x, y \in B(\bar{x}, R), \quad (2.2)$$

for some $M > 0$ and $\rho \in [0, 1]$. This $f(\cdot)$ can be non-smooth ($\rho = 0$), smooth ($\rho = 1$) and weakly smooth ($0 < \rho < 1$).

This section contains four subsections. We first present a much simplified APL method, referred to the fast APL (FAPL) method, for solving ball-constrained black-box CPs in Subsection 2.1, and then present the fast USL (FUSL) method for solving a special class of ball-constrained structured CPs in Subsection 2.2. We show how to solve the subproblems in these two algorithms in Subsection 2.3.

2.1. FAPL for ball-constrained black-box problems. Our goal in this subsection is to present the FAPL method, which can significantly reduce the iteration cost for the APL method applied to problem (2.1). In particular, we show that only one subproblem, rather than two subproblems (see (1.9) and (1.10)) as in the APL method, is required in the FAPL method for defining a new iterate (or prox-center) and updating lower bound. We also demonstrate that the ball constraint in (2.1) can be eliminated from the subproblem by properly specifying the prox-function.

Similarly to the APL method, the FAPL method consists of outer-inner loops, and in each outer iteration, an inner loop, the FAPL gap reduction procedure denoted by \mathcal{G}_{FAPL} , is called to reduce the gap between the upper and lower bounds on $f_{\bar{x},R}^*$ in (2.1) by a constant factor.

We start by describing the FAPL gap reduction procedure in Procedure 1. This procedure differs from the gap reduction procedure used in the APL method in the following several aspects. Firstly, the localizers Q_k and \bar{Q}_k in procedure \mathcal{G}_{FAPL} (see Steps 1 and 4) only contain linear constraints and hence are possibly unbounded, while the localizers in the APL method must be compact. Secondly, we eliminate the subproblem that updates the lower bound on $f_{\bar{x},R}^*$ in the APL method. Instead, in the FAPL method, the lower bound is updated to the level l directly whenever $Q_k = \emptyset$ or $\|x_k - \bar{x}\| > R$. Thirdly, we choose a specific prox-function $d(x) = \frac{1}{2}\|x - \bar{x}\|^2$, and as a result, all the three sequences $\{x_k\}$, $\{x_k^l\}$ and $\{x_k^u\}$ will reside in the ball $B(\bar{x}, R)$. At last, as we will show in next subsection, since the subproblem (2.6) only contains a limited number of linear constraints (depth of memory), we can solve it very efficiently, or even exactly if the depth of memory is small, say, ≤ 10 .

Procedure 1 The FAPL gap reduction procedure: $(x^+, \text{lb}^+) = \mathcal{G}_{FAPL}(\hat{x}, \text{lb}, R, \bar{x}, \beta, \theta)$

- 0: Set $k = 1$, $\bar{f}_0 = f(\hat{x})$, $l = \beta \cdot \text{lb} + (1 - \beta)\bar{f}_0$, $Q_0 = \mathbb{R}^n$, and let $x_0^u = \hat{x}$, $x_0 = \bar{x}$.
 1: Update the cutting plane model:

$$x_k^l = (1 - \alpha_k)x_{k-1}^u + \alpha_k x_{k-1}, \quad (2.3)$$

$$h(x_k^l, x) = f(x_k^l) + \langle f'(x_k^l), x - x_k^l \rangle, \quad (2.4)$$

$$Q_k = \{x \in Q_{k-1} : h(x_k^l, x) \leq l\}. \quad (2.5)$$

- 2: Update the prox-center and lower bound:

$$x_k = \operatorname{argmin}_{x \in Q_k} \left\{ d(x) := \frac{1}{2} \|x - \bar{x}\|^2 \right\}. \quad (2.6)$$

If $Q_k = \emptyset$ or $\|x_k - \bar{x}\| > R$, then **terminate** with outputs $x^+ = x_{k-1}^u$, $\text{lb}^+ = l$.

- 3: Update the upper bound: set

$$\tilde{x}_k^u = (1 - \alpha_k)x_{k-1}^u + \alpha_k x_k, \quad (2.7)$$

$$x_k^u = \begin{cases} \tilde{x}_k^u, & \text{if } f(\tilde{x}_k^u) < \bar{f}_{k-1}, \\ x_{k-1}^u, & \text{otherwise,} \end{cases} \quad (2.8)$$

and $\bar{f}_k = f(x_k^u)$. If $\bar{f}_k \leq l + \theta(\bar{f}_0 - l)$, then **terminate** with $x^+ = x_k^u$, $\text{lb}^+ = \text{lb}$.

- 4: Choose any polyhedral set Q_k satisfying $Q_k \subseteq Q_{k-1} \subseteq \bar{Q}_k$, where

$$\bar{Q}_k := \{x \in \mathbb{R}^n : \langle x_k - \bar{x}, x - x_k \rangle \geq 0\}. \quad (2.9)$$

Set $k = k + 1$ and go to Step 1.

We now add a few more remarks about the technical details of Procedure 1. Firstly, Procedure 1 is terminated at Step 2 if $Q_k = \emptyset$ or $\|x_k - \bar{x}\| > R$, which can be checked automatically when solving the subproblem (2.6) (see Subsection 2.3 for more details). Secondly, in Step 4, while Q_k can be any polyhedral set between Q_k and \bar{Q}_k , in practice we can simply choose Q_k to be the intersection of the half-space $\{x \in \mathbb{R}^n : \langle x_k - \bar{x}, x - x_k \rangle \geq 0\}$ and a few most recently generated half-spaces, each of which is defined by $\{x \in \mathbb{R}^n : h(x_\tau^l, x) \leq l\}$ for some $1 \leq \tau \leq k$. Finally, in order to guarantee the termination of procedure \mathcal{G}_{FAPL} and the optimal iteration complexity, the parameters $\{\alpha_k\}$ used in this procedure need to be properly chosen. One set of conditions that $\{\alpha_k\}$ should satisfy to guarantee the convergence of procedure \mathcal{G}_{FAPL} is given as follows:

$$\alpha_1 = 1, \quad 0 < \alpha_k \leq 1, \quad \alpha_k \leq \frac{c}{k} \quad \text{and} \quad \frac{1 - \alpha_{k+1}}{\alpha_{k+1}^{1+\rho}} \leq \frac{1}{\alpha_k^{1+\rho}}, \quad \forall k \geq 1, \quad (2.10)$$

for some constants $c > 0$ and $\forall \rho \in [0, 1]$.

The following lemma provides two examples for the selection of $\{\alpha_k\}$.

LEMMA 2.1.

- a) If $\alpha_k = 2/(k+1)$, $k = 1, 2, \dots$, then the condition (2.10) is satisfied with $c = 2$.
 b) If $\{\alpha_k\}$ is recursively defined by

$$\alpha_1 = 1, \quad \alpha_{k+1}^2 = (1 - \alpha_{k+1})\alpha_k^2, \quad \forall k \geq 1, \quad (2.11)$$

then the condition (2.10) holds with $c = 2$.

Proof. part **a)** follows immediately by plugging $\alpha_k = 2/(k+1)$, $k = 1, 2, \dots$, into (2.10). For part **b)**, it is easy to prove by induction that $\alpha_k \in (0, 1], \alpha_{k+1} < \alpha_k, \forall k \geq 1$. Then since

$$\frac{1}{\alpha_{k+1}} - \frac{1}{\alpha_k} = \frac{\alpha_k - \alpha_{k+1}}{\alpha_k \alpha_{k+1}} = \frac{\alpha_k}{\alpha_k + \alpha_{k+1}} \geq \frac{1}{2}, \forall k \geq 1,$$

and notice that $\alpha_1 = 1$, we have $\frac{1}{\alpha_k} - 1 \geq \frac{k-1}{2}$. Therefore $\alpha_k \leq \frac{2}{k+1} < \frac{2}{k}, \forall k \geq 1$. Moreover, for any $k \geq 1$, $1 - \alpha_{k+1} = \left(\frac{\alpha_{k+1}}{\alpha_k}\right)^2$ by (2.11), and $\left(\frac{\alpha_{k+1}}{\alpha_k}\right)^2 \leq \left(\frac{\alpha_{k+1}}{\alpha_k}\right)^{1+\rho}$ due to $\alpha_{k+1} < \alpha_k$ and $\rho \in [0, 1]$. Therefore $\frac{1-\alpha_{k+1}}{\alpha_{k+1}^{1+\rho}} \leq \frac{1}{\alpha_k^{1+\rho}}, \forall k \geq 1$, hence (2.10) holds. \square

The following lemma describes some important observations regarding the execution of procedure \mathcal{G}_{FAPL} .

LEMMA 2.2. Let $\mathcal{E}_f(l) := \{x \in B(\bar{x}, R) : f(x) \leq l\}$, the following statements hold for procedure \mathcal{G}_{FAPL} .

- a) If $\mathcal{E}_f(l) \neq \emptyset$, then $\mathcal{E}_f(l) \subseteq \underline{Q}_k \subseteq Q_k \subseteq \bar{Q}_k$ for any $k \geq 1$.
- b) If $\underline{Q}_k \neq \emptyset$, then problem (2.6) in Step 2 has a unique solution. Moreover, if procedure \mathcal{G}_{FAPL} terminates at Step 2, then $l \leq f_{\bar{x}, R}^*$.

Proof. To prove part **a)**, it suffices to prove that $\mathcal{E}_f(l) \subseteq Q_k$ and $Q_k \subseteq \bar{Q}_k$ for all $k \geq 1$, since from the definition of Q_k we already have $\underline{Q}_k \subseteq Q_k \subseteq \bar{Q}_k$. we first use induction to prove the former relation. As Q_0 is set to \mathbb{R}^n , we have $\mathcal{E}_f(l) \subseteq Q_0$. Moreover, if $\mathcal{E}_f(l) \subseteq Q_{k-1}$ for some $k \geq 1$, then from the definition of \underline{Q}_k in (2.5) and the observation that $h(x_k^l, z) \leq f(z) \leq l$ for all $z \in \mathcal{E}_f(l)$, we have $\mathcal{E}_f(l) \subseteq \underline{Q}_k$. To prove the relation that $Q_k \subseteq \bar{Q}_k$, observe that the definition of \bar{Q}_k in (2.9) implies that $\bar{Q}_k = \{x \in \mathbb{R}^n : d(x) \geq d(x_k)\}$. Combining such observation with the definition of x_k in (2.6), we have $Q_k \subseteq \bar{Q}_k$ and conclude part a).

We now provide the proof of part **b)**. From the definition of Q_k in Step 4 and the definition of \underline{Q}_k in (2.5), we can see that \underline{Q}_k is the intersection of half-spaces, hence it is convex and closed. Therefore, the subproblem (2.6) always has a unique solution as long as \underline{Q}_k is non-empty.

To finish the proof it suffices to show that $\mathcal{E}_f(l) = \emptyset$ when either $\underline{Q}_k = \emptyset$ or $\|x_k - \bar{x}\| > R$, which can be proved by contradiction. Firstly, if $\underline{Q}_k = \emptyset$ but $\mathcal{E}_f(l) \neq \emptyset$, then by part **a)** proved above, we have $\mathcal{E}_f(l) \subseteq \underline{Q}_k$, which contradicts the assumption that \underline{Q}_k is empty. On the other hand, suppose that $\|x_k - \bar{x}\| > R$ and $\mathcal{E}_f(l) \neq \emptyset$, let $x_R^* \in \text{Argmin}_{x \in B(\bar{x}, R)} f(x)$, it is clear that $x_R^* \in \mathcal{E}_f(l) \subseteq \underline{Q}_k$ by **a)**, however $\|x_R^* - \bar{x}\| \leq R < \|x_k - \bar{x}\|$ which contradicts the definition of x_k in (2.6). \square

Let $\text{ub} := f(\hat{x}), \text{ub}^+ := f(x^+)$ be the input and output upper bounds on $f_{\bar{x}, R}^*$ in procedure \mathcal{G}_{FAPL} , respectively. The following lemma shows that whenever procedure \mathcal{G}_{FAPL} terminates, the gap between the upper and lower bounds on $f_{\bar{x}, R}^*$ is reduced by a constant factor.

LEMMA 2.3. Whenever procedure \mathcal{G}_{FAPL} terminates, we have $\text{ub}^+ - \text{lb}^+ \leq q(\text{ub} - \text{lb})$, where

$$q := \max\{\beta, 1 - (1 - \theta)\beta\}. \quad (2.12)$$

Proof. By the definition of x_k^u in (2.7) and the definition of \bar{f}_k in Step 3 of procedure \mathcal{G}_{FAPL} , we have $\bar{f}_k \leq \bar{f}_{k-1}, \forall k \geq 1$, which implies $\text{ub}^+ \leq \text{ub}$. Procedure \mathcal{G}_{FAPL} could terminate at either Step 2 or Step 3. We first assume that it terminates at Step 2 at K^{th} iteration. The termination condition gives $\text{lb}^+ = l = \beta \cdot \text{lb} + (1 - \beta)\text{ub}$, then we have

$$\text{ub}^+ - \text{lb}^+ \leq \text{ub} - \beta \text{lb} - (1 - \beta)\text{ub} = \beta(\text{ub} - \text{lb}).$$

Next suppose that Procedure 1 terminates at Step 3 at K^{th} iteration. We have $\text{ub}^+ = f(x^+) = \bar{f}_K \leq l + \theta(\text{ub} - l)$, since $\text{lb}^+ \geq \text{lb}$ and $l = \beta \cdot \text{lb} + (1 - \beta)\text{ub}$, we conclude that

$$\text{ub}^+ - \text{lb}^+ \leq l + \theta(\text{ub} - l) - \text{lb} = [1 - (1 - \theta)\beta](\text{ub} - \text{lb}).$$

The lemma is proved by combining the above two cases. \square

We now provide a bound on the number of iterations performed by procedure \mathcal{G}_{FAPL} . Note that the proof of this result is similar to Theorem 3 in [20].

PROPOSITION 2.4. *If the stepsizes $\{\alpha_k\}_{k \geq 1}$ are chosen such that (2.10) holds, then the number of iterations performed by procedure \mathcal{G}_{FAPL} does not exceed*

$$N(\Delta) := \left(\frac{c^{1+\rho} M R^{1+\rho}}{(1+\rho)\theta\beta\Delta} \right)^{\frac{2}{1+3\rho}} + 1, \quad (2.13)$$

where $\Delta := \text{ub} - \text{lb}$.

Proof. Suppose that procedure \mathcal{G}_{FAPL} does not terminate at the K^{th} iteration for some $K > 0$. Then observing that $x_k \in \underline{Q}_k \subseteq Q_{k-1} \subseteq \overline{Q}_{k-1}$ for any $2 \leq k \leq K$ due to (2.5) and (2.6), and $x_1 \in Q_0, x_0 = \arg\min_{x \in Q_0} d(x)$, we have $\langle \nabla d(x_{k-1}), x_k - x_{k-1} \rangle \geq 0, \forall 1 \leq k \leq K$. Since $d(x)$ is strongly convex with modulus 1, we also have

$$d(x_k) \geq d(x_{k-1}) + \langle \nabla d(x_{k-1}), x_k - x_{k-1} \rangle + \frac{1}{2} \|x_k - x_{k-1}\|^2.$$

Combining the two relations above, it implies $\frac{1}{2} \|x_k - x_{k-1}\|^2 \leq d(x_k) - d(x_{k-1})$. Summing up these inequalities for $1 \leq k \leq K$, we conclude that

$$\frac{1}{2} \sum_{k=1}^K \|x_k - x_{k-1}\|^2 \leq d(x_K) - d(x_0) = \frac{1}{2} \|x_K - \bar{x}\|^2 \leq \frac{1}{2} R^2. \quad (2.14)$$

By (2.2), (2.7), (2.8) and the convexity of f , we have for any $1 \leq k \leq K$,

$$f(x_k^u) \leq f(\tilde{x}_k^u) \leq h(x_k^l, \tilde{x}_k^u) + \frac{M}{1+\rho} \|\tilde{x}_k^u - x_k^l\|^{1+\rho} \quad (2.15)$$

$$= (1 - \alpha_k) h(x_k^l, x_{k-1}^u) + \alpha_k h(x_k^l, x_k) + \frac{M}{1+\rho} \|\tilde{x}_k^u - x_k^l\|^{1+\rho}. \quad (2.16)$$

Since $h(x_k^l, x_{k-1}^u) \leq f(x_{k-1}^{u-1})$, $h(x_k^l, x_k) \leq l$ due to (2.5) and (2.6), and $\tilde{x}_k^u - x_k^l = \alpha_k(x_k - x_{k-1})$ due to (2.3) and (2.7), we have

$$f(x_k^u) - l \leq (1 - \alpha_k)(f(x_{k-1}^u) - l) + \frac{\alpha_k^{1+\rho} M}{1+\rho} \|x_k - x_{k-1}\|^{1+\rho} \quad (2.17)$$

Dividing both sides of (2.17) by $\alpha_k^{1+\rho}$, and then summing up these inequalities for $1 \leq k \leq K$, by (2.10) and the fact $f(x_k^u) - l \geq 0, \forall 1 \leq k \leq K$, we have

$$f(x_K^u) - l \leq \frac{\alpha_K^{1+\rho} M}{1+\rho} \sum_{k=1}^K \|x_k - x_{k-1}\|^{1+\rho} \quad (2.18)$$

Apply Hölder's inequality, and use (2.14), we have

$$\sum_{k=1}^K \|x_k - x_{k-1}\|^{1+\rho} \leq K^{\frac{1-\rho}{2}} \left(\sum_{k=1}^K \|x_k - x_{k-1}\|^2 \right)^{\frac{1+\rho}{2}} \leq K^{\frac{1-\rho}{2}} R^{1+\rho},$$

and $\alpha_K^{1+\rho} \leq c^{1+\rho} K^{-(1+\rho)}$ due to (2.10). Therefore (2.18) gives

$$f(x_K^u) - l \leq \frac{M R^{1+\rho}}{(1+\rho)} \cdot \frac{c^{1+\rho}}{K^{\frac{1+3\rho}{2}}} \quad (2.19)$$

In view of Step 3 in procedure 1, and using the fact that $l = \beta \cdot \text{lb} + (1 - \beta) \text{ub}$ in Step 0, we also have

$$f(x_K^u) - l > \theta(\text{ub} - l) = \theta\beta\Delta.$$

Combining the above two relations, and using (2.10) and (2.14), we obtain

$$\theta\beta\Delta < \frac{MR^{1+\rho}}{(1+\rho)} \cdot \frac{c^{1+\rho}}{K^{\frac{1+3\rho}{2}}}, \quad (2.20)$$

which implies that

$$K < \left(\frac{c^{1+\rho}MR^{1+\rho}}{(1+\rho)\theta\beta\Delta} \right)^{\frac{2}{1+3\rho}}. \quad (2.21)$$

□

In view of Lemma 2.3 and Proposition 2.4, we are now ready to describe the FAPL method, which performs a sequence of calls to procedure \mathcal{G}_{FAPL} until an approximate solution with sufficient accuracy is found.

Algorithm 1 The fast accelerated prox-level (FAPL) method

- 0: Given ball $B(\bar{x}, R)$, choose initial point $p_0 \in B(\bar{x}, R)$, tolerance $\epsilon > 0$ and parameters $\beta, \theta \in (0, 1)$.
 - 1: Set $p_1 \in \text{Argmin}_{x \in B(\bar{x}, R)} h(p_0, x)$, $\text{lb}_1 = h(p_0, p_1)$, $\text{ub}_1 = \min\{f(p_0), f(p_1)\}$, let \hat{x}_1 be either p_0 or p_1 such that $f(\hat{x}_1) = \text{ub}_1$, and $s = 1$.
 - 2: If $\text{ub}_s - \text{lb}_s \leq \epsilon$, **terminate** and **output** approximate solution \hat{x}_s .
 - 3: Set $(\hat{x}_{s+1}, \text{lb}_{s+1}) = \mathcal{G}_{FAPL}(\hat{x}_s, \text{lb}_s, R, \bar{x}, \beta, \theta)$ and $\text{ub}_{s+1} = f(\hat{x}_{s+1})$.
 - 4: Set $s = s + 1$ and go to Step 2.
-

For the sake of simplicity, each iteration of Procedure \mathcal{G}_{FAPL} is also referred to as an iteration of the FAPL method. The following theorem establishes the complexity bounds on the numbers of gap reduction procedures \mathcal{G}_{FAPL} and total iterations performed by the FAPL method, its proof is similar to that of Theorem 4 in [20].

THEOREM 2.5. *For any given $\epsilon > 0$, if the stepsizes $\{\alpha_k\}$ in procedure \mathcal{G}_{FAPL} are chosen such that (2.10) holds, then the following statements hold for the FAPL method to compute an ϵ -solution to problem (2.1).*

a) *The number of gap reduction procedures \mathcal{G}_{FAPL} performed by the FAPL method does not exceed*

$$S := \left\lceil \max \left\{ 0, \log_{\frac{1}{q}} \left(\frac{(2R)^{1+\rho}M}{(1+\rho)\epsilon} \right) \right\} \right\rceil. \quad (2.22)$$

b) *The total number of iterations performed by the FAPL method can be bounded by*

$$N(\epsilon) := S + \frac{1}{1 - q^{\frac{2}{1+3\rho}}} \left(\frac{c^{1+\rho}MR^{1+\rho}}{(1+\rho)\theta\beta\epsilon} \right)^{\frac{2}{1+3\rho}}, \quad (2.23)$$

where q is defined in (2.12).

Proof. We first prove part a). Let $\Delta_s := \text{ub}_s - \text{lb}_s$, without loss of generality, we assume that $\Delta_1 > \epsilon$. In view of Step 0 in Algorithm 1 and (2.2), we have

$$\Delta_1 \leq f(p_1) - h(p_0, p_1) = f(p_1) - f(p_0) - \langle f'(p_0), p_1 - p_0 \rangle \leq \frac{(2R)^{1+\rho}M}{1+\rho}. \quad (2.24)$$

Also, by Lemma 2.3 we can see that $\Delta_{s+1} \leq q\Delta_s$ for any $s \geq 1$, which implies that

$$\Delta_{s+1} \leq q^s \Delta_1, \quad \forall s \geq 0.$$

Moreover, if an ϵ -solution is found after performing S gap reduction procedures \mathcal{G}_{FAPL} , then we have

$$\Delta_S > \epsilon \geq \Delta_{S+1}. \quad (2.25)$$

Combining the above three inequalities, we conclude that

$$\epsilon < q^{S-1} \Delta_1 \leq q^{S-1} \frac{(2R)^{1+\rho} M}{1+\rho}, \quad (2.26)$$

and part a) follows immediately from the above inequality. We are now ready to prove part b). In view of Lemma 2.3 and (2.25), we have $\Delta_s \geq \epsilon q^{s-S}$ for any $1 \leq s \leq S$. Using this estimate, part a), and Proposition 2.4, we conclude that the total number of iterations performed by the FAPL method is bounded by

$$N(\epsilon) = \sum_{s=1}^S N_s = S + \left(\frac{c^{1+\rho} M R^{1+\rho}}{(1+\rho)\theta\beta\epsilon} \right)^{\frac{2}{1+3\rho}} \sum_{s=1}^S q^{\frac{2(S-s)}{1+3\rho}} < S + \frac{1}{1 - q^{\frac{2}{1+3\rho}}} \left(\frac{c^{1+\rho} M R^{1+\rho}}{(1+\rho)\theta\beta\epsilon} \right)^{\frac{2}{1+3\rho}}, \quad (2.27)$$

where N_s denotes the number of the iterations performed by the s^{th} gap reduction procedure \mathcal{G}_{FAPL} for any $1 \leq s \leq S$. \square

In view of Theorem 2.5, the FAPL method achieves the optimal iteration complexity for solving non-smooth, weakly smooth, and smooth CP problems, which is the same as the convergence properties of the ABL and APL methods (see [24, 26] for the discussions on the complexity theories for solving CPs, and [20] for the convergence properties of the ABL and APL methods).

2.2. FUSL for ball-constrained structured problems. In this subsection, we still consider the ball-constrained problem (2.1), but assume that its objective function is given by

$$f(x) := \hat{f}(x) + F(x), \quad (2.28)$$

where \hat{f} is a smooth convex function, i.e., $\exists L_{\hat{f}} > 0$ s.t.

$$\hat{f}(y) - \hat{f}(x) - \langle \nabla \hat{f}(x), y - x \rangle \leq \frac{L_{\hat{f}}}{2} \|y - x\|^2, \quad (2.29)$$

and

$$F(x) := \max_{y \in Y} \{ \langle Ax, y \rangle - \hat{g}(y) \}. \quad (2.30)$$

Here, $Y \subseteq \mathbb{R}^m$ is a compact convex set, $\hat{g} : Y \rightarrow \mathbb{R}$ is a relatively simple convex function, and $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a linear operator. We will integrate Nesterov's smoothing technique for minimizing non-smooth functions [27] into the FAPL method to solve problem (2.1)-(2.28) (i.e., problem (2.1) with f defined in (2.28)), which can significantly reduce the iteration cost of the USL method [20].

Let $v : Y \rightarrow \mathbb{R}$ be a prox-function with modulus σ_v and denoting $c_v := \operatorname{argmin}_{v \in Y} v(y)$, by employing the idea in the important work of Nesterov [27], we can approximate F in (2.30) by the smooth function

$$F_\eta(x) := \max_{y \in Y} \{ \langle Ax, y \rangle - \hat{g}(y) - \eta V(y) \}, \quad (2.31)$$

where $\eta > 0$ is called the smoothing parameter, and $V(\cdot)$ is the Bregman divergence defined by

$$V(y) := v(y) - v(c_v) - \langle \nabla v(c_v), y - c_v \rangle. \quad (2.32)$$

It was shown in [27] that the gradient of $F_\eta(\cdot)$ given by $\nabla F_\eta(x) = A^* y^*(x)$ is Lipschitz continuous with constant

$$L_\eta := \|A\|^2 / (\eta \sigma_v), \quad (2.33)$$

where $\|A\|$ is the operator norm of A , A^* is the adjoint operator, and $y^*(x) \in Y$ is the solution to the optimization problem (2.31). Moreover, the "closeness" of $F_\eta(\cdot)$ to $F(\cdot)$ depends linearly on the smoothing parameter η , i.e.,

$$F_\eta(x) \leq F(x) \leq F_\eta(x) + \eta D_{v,Y}, \quad \forall x \in X, \quad (2.34)$$

where

$$D_{v,Y} := \max_{y,z \in Y} \{v(y) - v(z) - \langle \nabla v(z), y - z \rangle\}. \quad (2.35)$$

Therefore, if we denote

$$f_\eta(x) := \hat{f}(x) + F_\eta(x), \quad (2.36)$$

then

$$f_\eta(x) \leq f(x) \leq f_\eta(x) + \eta D_{v,Y}. \quad (2.37)$$

Applying an optimal gradient method to minimize the smooth function f_η in (2.36), Nesterov proves in [27] that the iteration complexity for computing an ϵ -solution to problem (2.1)-(2.28) is bounded by $\mathcal{O}(\frac{1}{\sqrt{\epsilon}} + \frac{1}{\epsilon})$. However, the values of quite a few problem parameters, such as $\|A\|, \sigma_v$ and $D_{v,Y}$, are required for the implementation of Nesterov's smoothing scheme.

By incorporating Nesterov's smoothing technique [27] into the APL method, Lan developed in [20] a new bundle-level type method, namely the uniform smoothing level (USL) method, to solve structured problems given in the form of (2.1)-(2.28). While the USL method achieves the same optimal iteration complexity as Nesterov's smoothing scheme in [27], one advantage of the USL method over Nesterov's smoothing scheme is that the smoothing parameter η is adjusted dynamically during the execution, and an estimate of $D_{v,Y}$ is obtained automatically, which makes the USL method problem parameter free. However, similar to the APL method, each iteration of the USL method involves the solutions to two subproblems. Based on the USL method in [20] and our analysis of the FAPL method in Section 2.1, we propose a fast USL (FUSL) method that solves problem (2.1)-(2.28) with the same optimal iteration complexity as the USL method, but requiring only to solve one simpler subproblem in each iteration.

Similar to the FAPL method, the FUSL method has an outer-inner loops structure, and each outer iteration calls an inner loop, a gap reduction procedure denoted by \mathcal{G}_{FUSL} , to reduce the gap between the upper and lower bounds on $f_{\bar{x},R}^*$ in (2.1)-(2.28) by a constant factor unless an ϵ -solution is found. We start by describing procedure \mathcal{G}_{FUSL} .

Procedure 2 The FUSL gap reduction procedure: $(x^+, D^+, \text{lb}^+) = \mathcal{G}_{FUSL}(\hat{x}, D, \text{lb}, R, \bar{x}, \beta, \theta)$

0: Let $k = 1, \bar{f}_0 = f(\hat{x}), l = \beta \cdot \text{lb} + (1 - \beta)\bar{f}_0, Q_0 = \mathbb{R}^n, x_0^u = \hat{x}, x_0 = \bar{x}$, and

$$\eta := \theta(\bar{f}_0 - l)/(2D). \quad (2.38)$$

1: Update the cutting plane model: Set x_k^l to (2.3), Q_k to (2.5), and

$$h(x_k^l, x) = h_\eta(x_k^l, x) = f_\eta(x_k^l) + \langle f'_\eta(x_k^l), x - x_k^l \rangle. \quad (2.39)$$

2: Update the prox-center: Set x_k to (2.6). If $Q_k = \emptyset$ or $\|x_k - \bar{x}\| > R$, then **terminate** with outputs $x^+ = x_{k-1}^u, D^+ = D, \text{lb}^+ = l$.

3: Update the upper bound and the estimate of $D_{v,Y}$: Set \tilde{x}_k^u to (2.7), x_k^u to (2.8), and $\bar{f}_k = f(x_k^u)$. Check the following conditions:

3a) if $f(x_k^u) \leq l + \theta(\bar{f}_0 - l)$, then **terminate** with outputs $x^+ = x_k^u, D^+ = D, \text{lb}^+ = \text{lb}$;

3b) if $f(x_k^u) > l + \theta(\bar{f}_0 - l)$ and $f_\eta(x_k^u) \leq l + \frac{\theta}{2}(\bar{f}_0 - l)$, then **terminate** with outputs $x^+ = x_k^u, D^+ = 2D, \text{lb}^+ = \text{lb}$.

4: Choose Q_k as same as Step 4 in \mathcal{G}_{FAPL} . Set $k = k + 1$, and go to step 1.

A few remarks about procedure \mathcal{G}_{FUSL} are in place. Firstly, since the non-smooth objective function f is replaced by its smoothed approximation f_η , we replace the cutting plane model in (1.7) with the one for f_η (see (2.39)). Also note that for the USL method in [20], \hat{f} is assumed to be a simple Lipschitz continuous convex

function, and only F_η is approximated by the linear estimation. However, in the FUSL method, we assume \hat{f} is general smooth and convex, and linearize both \hat{f} and F_η in (2.39). Secondly, the smoothing parameter η is specified as a function of the parameter D , \bar{f}_0 and l , where D is an estimator of $D_{v,Y}$ defined by (2.35) and given as an input parameter for procedure \mathcal{G}_{FUSL} . Thirdly, same as the FAPL method, the parameters $\{\alpha_k\}$ are chosen according to (2.10). Such conditions are required to guarantee the optimal iteration complexity of the FUSL method for solving problem (2.1)-(2.28). Fourthly, similar to the FAPL method, the localizers $\underline{Q}_k, Q_k, \bar{Q}_k$ only contain a limited number of linear constraints, and there is only one subproblem (i.e., (2.6)) involved in procedure \mathcal{G}_{FUSL} , which can be solved exactly when the depth of memory is small.

The following lemma provides some important observations about procedure \mathcal{G}_{FUSL} , which are similar to those for the USL gap reduction procedure in [20].

LEMMA 2.6. *The following statements hold for procedure \mathcal{G}_{FUSL} .*

- a) *If procedure \mathcal{G}_{FUSL} terminates at Steps 2 or 3a), then we have $\text{ub}^+ - \text{lb}^+ \leq q(\text{ub} - \text{lb})$, where q is defined by (2.12) and $\text{ub} := f(\hat{x}), \text{ub}^+ := f(x^+)$.*
- b) *If procedure \mathcal{G}_{FUSL} terminates at Step 3b), then $D < D_{v,Y}$ and $D^+ < 2D_{v,Y}$.*

Proof. The proof of part a) is as the same as that of Lemma 2.3, and we only prove part b) here. By the termination condition at Step 3b), we have

$$f(x_k^u) - f_\eta(x_k^u) > \frac{\theta}{2}(\bar{f}_0 - l).$$

We conclude from the above relation, (2.37), and (2.38) that

$$D_{v,Y} \geq \frac{f(x_k^u) - f_\eta(x_k^u)}{\eta} > \frac{\theta(\bar{f}_0 - l)}{2\eta} = D.$$

Finally, $D^+ < 2D_{v,Y}$ follows immediately from the above relation and the definition of D^+ in Step 3b). \square

The following result provides a bound on the number of iterations performed by procedure \mathcal{G}_{FUSL} .

PROPOSITION 2.7. *Suppose that $\{\alpha_k\}_{k \geq 1}$ in procedure \mathcal{G}_{FUSL} are chosen such that (2.10) holds. Then, the number of iterations performed by this procedure does not exceed*

$$\bar{N}(\Delta, D) := cR \left(\sqrt{\frac{L_{\hat{f}}}{\theta\beta\Delta}} + \frac{\sqrt{2}\|A\|}{\theta\beta\Delta} \sqrt{\frac{D}{\sigma_v}} \right) + 1, \quad (2.40)$$

where $\Delta := f(\hat{x}) - \text{lb}$.

Proof. It is easy to see that the gradient of f_η in (2.28) has Lipschitz continuous gradient with constant $L = L_{\hat{f}} + L_\eta$, where L_η and $L_{\hat{f}}$ are defined in (2.33) and (2.29), respectively. Suppose that procedure \mathcal{G}_{FUSL} does not terminate at iteration K . As the \mathcal{G}_{FUSL} procedure could be viewed as applying the \mathcal{G}_{FAPL} procedure to f_η , similar to the discussion on (2.19), and notice that $\rho = 1$ as f_η is smooth, we have

$$f_\eta(x_K^u) - l \leq \frac{c^2 L d(x_K)}{K^2} \leq \frac{c^2 L R^2}{2K^2}, \quad (2.41)$$

where c is defined in (2.10). Also, since procedure \mathcal{G}_{FUSL} does not terminate at iteration K , in view of the termination condition in Step 3b) and the definition of l in Step 0, we have

$$f_\eta(x_K^u) - l > \frac{\theta\beta\Delta}{2}. \quad (2.42)$$

By (2.41), (2.42), (2.33) and (2.38), we conclude that

$$K < \sqrt{\frac{c^2 L R^2}{\theta\beta\Delta}} \leq cR \left(\sqrt{\frac{L_{\hat{f}}}{\theta\beta\Delta}} + \frac{\sqrt{2}\|A\|}{\theta\beta\Delta} \sqrt{\frac{D}{\sigma_v}} \right). \quad (2.43)$$

Algorithm 2 The fast uniform smoothing level (FUSL) method

- 0: Given ball $B(\bar{x}, R)$, choose initial point $p_0 \in B(\bar{x}, R)$, prox-function $v(\cdot)$ for the smoothing function F_η in (2.31) and (2.32), initial guess D_1 on the size $D_{v,Y}$ in (2.35), tolerance $\varepsilon > 0$, and parameters $\beta, \theta \in (0, 1)$.
- 1: Set $p_1 \in \text{Argmin}_{x \in B(\bar{x}, R)} h(p_0, x)$, $\text{lb}_1 = h(p_0, p_1)$, $\text{ub}_1 = \min\{f(p_0), f(p_1)\}$, let \hat{x}_1 be either p_0 or p_1 such that $f(\hat{x}_1) = \text{ub}_1$, and $s = 1$.
- 2: If $\text{ub}_s - \text{lb}_s \leq \varepsilon$, **terminate** and **output** approximate solution \hat{x} .
- 3: Set $(\hat{x}_{s+1}, D_{s+1}, \text{lb}_{s+1}) = \mathcal{G}_{FUSL}(\hat{x}_s, D_s, \text{lb}_s, R, \bar{x}, \beta, \theta)$ and $\text{ub}_{s+1} = f(\hat{x})$.
- 4: Set $s = s + 1$ and go to Step 2.
-

□

We are now ready to describe the FUSL method which iteratively calls procedure \mathcal{G}_{FUSL} to solve the structured saddle point problem (2.1)-(2.28).

For simplicity, we say that a phase (i.e., an outer iteration) of the FUSL method occurs when s increases by 1. More specifically, similar to the USL method, we classify two types of phases in the FUSL method. A phase is called *significant* if the corresponding \mathcal{G}_{FUSL} procedure terminates at Steps 2 or 3a), otherwise it is called *non-significant*. Clearly, if the value of $D_{v,y}$ is provided, which is the assumption made in Nesterov's smoothing scheme [27], then we can set $D_1 = D_{v,Y}$ in the schemes of both the USL and FUSL methods, and consequently, all the phases of both methods become significant.

As the same as the FAPL method, an iteration of procedure \mathcal{G}_{FUSL} is also referred to an iteration of the FUSL method. The following result establishes a bound on the total number of iterations performed by the FUSL method to find an ε -solution to problem (2.1)-(2.28). Note that the proof of this result is similar to that of Theorem 7 in [20].

THEOREM 2.8. *Suppose that $\{\alpha_k\}$ in procedure \mathcal{G}_{FUSL} are chosen such that (2.10) holds. Then, the total number of iterations performed by the FUSL method for computing an ε -solution to problem (2.1)-(2.28) is bounded by*

$$\bar{N}(\varepsilon) := S_1 + S_2 + \left(\frac{2}{\sqrt{2}-1} + \frac{\sqrt{2}}{1-q} \right) \frac{cR\|A\|}{\theta\beta\varepsilon} \sqrt{\frac{\tilde{D}}{\sigma_v}} + \left(S_1 + \frac{1}{1-\sqrt{q}} \right) cR \sqrt{\frac{L_{\hat{f}}}{\theta\beta\varepsilon}}, \quad (2.44)$$

where q and $D_{v,Y}$ are defined by (2.12) and (2.35) respectively, and

$$\tilde{D} := \max\{D_1, 2D_{v,Y}\}, S_1 := \max\left\{ \left\lceil \log_2 \frac{D_{v,Y}}{D_1} \right\rceil, 0 \right\} \quad \text{and} \quad S_2 := \left\lceil \log_{\frac{1}{q}} \frac{4\sqrt{2}R\|A\| \sqrt{\frac{D_{v,Y}}{\sigma_v}} + 2R^2 L_{\hat{f}}}{\varepsilon} \right\rceil. \quad (2.45)$$

Proof. We prove this result by estimating the numbers of iterations performed within non-significant and significant phases separately. Suppose that the sets of indices of the non-significant and significant phases are $\{m_1, m_2, \dots, m_{s_1}\}$ and $\{n_1, n_2, \dots, n_{s_2}\}$ respectively. For any non-significant phase m_k , $1 \leq k \leq s_1$, we can easily see from Step 3b) that $D_{m_{k+1}} = 2D_{m_k}$, by part b) in Lemma 2.6, the number of non-significant phases performed by the FUSL method is bounded by S_1 defined above, i.e., $s_1 \leq S_1$.

In addition, since $D_{m_{s_1}} \leq \tilde{D}$, we have $D_{m_k} \leq (1/2)^{s_1-k} \tilde{D}$, where \tilde{D} is defined above. Combining the above estimates on s_1 and D_{m_k} , and in view of the fact $\Delta_{m_k} > \varepsilon$ for all $1 \leq k \leq s_1$, we can bound the number

of iterations performed within the non-significant phases by

$$\begin{aligned}
\bar{N}_1 &= \sum_{k=1}^{s_1} \bar{N}(\Delta_{m_k}, D_{m_k}) \leq \sum_{k=1}^{s_1} \bar{N}\left(\epsilon, \tilde{D}/2^{s_1-k}\right) \\
&\leq S_1 \left(cR \sqrt{\frac{L_{\hat{f}}}{\theta\beta\epsilon}} + 1 \right) + \frac{\sqrt{2}cR\|A\|}{\theta\beta\epsilon} \sqrt{\frac{\tilde{D}}{\sigma_v}} \sum_{k=1}^{s_1} 2^{-\frac{(s_1-k)}{2}} \\
&\leq S_1 \left(cR \sqrt{\frac{L_{\hat{f}}}{\theta\beta\epsilon}} + 1 \right) + \frac{2cR\|A\|}{(\sqrt{2}-1)\theta\beta\epsilon} \sqrt{\frac{\tilde{D}}{\sigma_v}}.
\end{aligned} \tag{2.46}$$

Applying Lemma 8 in [20] and relation (2.29), and in view of the fact that $p_0, p_1 \in B(\bar{x}, R)$ in Algorithm 2, the initial gap is bounded by

$$\Delta_1 := \text{ub}_1 - \text{lb}_1 \leq [F(p_0) - F(p_1) - \langle F'(p_1), p_0 - p_1 \rangle] + \left[\hat{f}(p_0) - \hat{f}(p_1) - \langle \hat{f}'(p_1), p_0 - p_1 \rangle \right] \tag{2.47}$$

$$\leq 4\sqrt{2}R\|A\| \sqrt{\frac{D_{v,Y}}{\sigma_v}} + 2R^2L_{\hat{f}}, \tag{2.48}$$

where $F'(p_1) \in \partial F(p_1)$. Then for the significant phases, similar to the proof of Theorem 2.5, we have $s_2 \leq S_2$. Moreover, for any n_k , $1 \leq k \leq s_2$, using Lemmas 2.3, 2.6, we have $D_{n_k} \leq \tilde{D}$, $\Delta_{n_{k+1}} \leq q\Delta_{n_k}$, and $\Delta_{n_{s_2}} > \epsilon$, which implies $\Delta_{n_k} > \epsilon/q^{s_2-k}$. Combine these estimates on D_{n_k}, Δ_{n_k} and bound on s_2 , we can see that the total number of iterations performed within the significant phases is bounded by

$$\begin{aligned}
\bar{N}_2 &= \sum_{k=1}^{s_2} \bar{N}(\Delta_{n_k}, D_{n_k}) \leq \sum_{k=1}^{s_2} \bar{N}(\epsilon/q^{s_2-k}, \tilde{D}) \\
&\leq S_2 + cR \sqrt{\frac{L_{\hat{f}}}{\theta\beta\epsilon}} \sum_{k=1}^{s_2} q^{\frac{s_2-k}{2}} + \frac{\sqrt{2}cR\|A\|}{\theta\beta\epsilon} \sqrt{\frac{\tilde{D}}{\sigma_v}} \sum_{k=1}^{s_2} q^{s_2-k} \\
&\leq S_2 + \frac{cR}{1-\sqrt{q}} \sqrt{\frac{L_{\hat{f}}}{\theta\beta\epsilon}} + \frac{\sqrt{2}cR\|A\|}{\theta\beta\epsilon(1-q)} \sqrt{\frac{\tilde{D}}{\sigma_v}}.
\end{aligned} \tag{2.49}$$

Finally, the total number of iterations performed by the FUSL method is bounded by $\bar{N}_1 + \bar{N}_2$, and thus (2.44) holds. \square

From (2.44) in the above theorem, we can see that the iteration complexity of the FUSL method for solving problem (2.1)-(2.28) is bounded by

$$\mathcal{O} \left(\sqrt{\frac{L_{\hat{f}}}{\epsilon}} + \frac{\|A\|}{\epsilon} \right). \tag{2.50}$$

The above iteration complexity is the same as that of the Nesterov smoothing scheme in [27] and the USL method in [20]. However, both the USL and FUSL methods improve Nesterov's smoothing scheme in that both of them are problem parameter free. In addition, as detailed in Subsection 2.3 below, the FUSL method further improves the USL method by significantly reducing its iteration cost and improving the accuracy for solving its subproblems.

2.3. Solving the subproblems of FAPL and FUSL. In this section, we introduce an efficient method to solve the subproblems (2.6) in the FAPL and FUSL methods, which are given in the form of

$$x_c^* := \operatorname{argmin}_{x \in Q} \frac{1}{2} \|x - p\|^2. \tag{2.51}$$

Here, Q is a closed polyhedral set described by m linear inequalities, i.e.,

$$Q := \{x \in \mathbb{R}^n : \langle A_i, x \rangle \leq b_i, i = 1, 2, \dots, m\},$$

where $A_i \in \mathbb{R}^n$ and $b_i \in \mathbb{R}$ for $1 \leq i \leq m$.

Now let us examine the Lagrange dual of (2.51) given by

$$\max_{\lambda \geq 0} \min_{x \in \mathbb{R}^n} \frac{1}{2} \|x - p\|^2 + \sum_{i=1}^m \lambda_i [\langle A_i, x \rangle - b_i]. \quad (2.52)$$

It can be checked from the theorem of alternatives that problem (2.52) is solvable if and only if $Q \neq \emptyset$. Indeed, if $Q \neq \emptyset$, it is obvious that the optimal value of (2.52) is finite. On the other hand, if $Q = \emptyset$, then there exists $\bar{\lambda} \geq 0$ such that $\bar{\lambda}^T A = 0$ and $\bar{\lambda}^T b < 0$, which implies that the optimal value of (2.52) goes to infinity. Moreover, if (2.52) is solvable and λ^* is one of its optimal dual solutions, then

$$x_c^* = p - \sum_{i=1}^m \lambda_i^* A_i. \quad (2.53)$$

It can also be easily seen that (2.52) is equivalent to

$$\max_{\lambda \geq 0} -\frac{1}{2} \lambda^T M \lambda + C^T \lambda, \quad (2.54)$$

where $M_{ij} := \langle A_i, A_j \rangle$, $C_i := \langle A_i, p \rangle - b_i$, $\forall i, j = 1, 2, \dots, m$. Hence, we can determine the feasibility of (2.51) or compute its optimal solution by solving the relatively simple problem (2.54).

Many algorithms are capable of solving the above nonnegative quadratic programming in (2.54) efficiently. Due to its low dimension (usually less than 10 in our practice), we propose a brute-force method to compute the exact solution of this problem. Consider the Lagrange dual associated with (2.54):

$$\min_{\lambda \geq 0} \max_{\mu \geq 0} \mathcal{L}(\lambda, \mu) := \frac{1}{2} \lambda^T M \lambda - (C^T + \mu) \lambda,$$

where the dual variable is $\mu := (\mu_1, \mu_2, \dots, \mu_m)$. Applying the KKT condition, we can see that $\lambda^* \geq 0$ is a solution to problem (2.54) if and only if there exists $\mu^* \geq 0$ such that

$$\nabla_{\lambda} \mathcal{L}(\lambda^*, \mu^*) = 0 \quad \text{and} \quad \langle \lambda^*, \mu^* \rangle = 0. \quad (2.55)$$

Note that the first identity in (2.55) is equivalent to a linear system:

$$(M \quad -I) \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_m \\ \mu_1 \\ \vdots \\ \mu_m \end{pmatrix} = \begin{pmatrix} C_1 \\ C_2 \\ \vdots \\ C_m \end{pmatrix}, \quad (2.56)$$

where I is the $m \times m$ identity matrix. The above linear system has $2m$ variables and m equations. But for any $i = 1, \dots, m$, we have either $\lambda_i = 0$ or $\mu_i = 0$, and hence we only need to consider 2^m possible cases on the non-negativity of these variables. Since m is rather small in practice, it is possible to exhaust all these 2^m cases to find the exact solution to (2.55). For each case, we first remove the m columns in the matrix $(M - I)$ which correspond to the m variables assumed to be 0, and then solve the remaining determined linear system. If all variables of the computed solution are non-negative, then solution (λ^*, μ^*) to (2.55) is found, and the exact solution x_c^* to (2.51) is computed by (2.53), otherwise, we continue to examine the next case. It is interesting to observe that these different cases can also be considered in parallel to take the advantages of high performance computing techniques.

3. Solving unconstrained CP problems through ball-constrained CPs. In this section we discuss the following unconstrained CP:

$$f^* := \min_{x \in \mathbb{R}^n} f(x), \quad (3.1)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, and for any closed sets $\Omega \in \mathbb{R}^n$, there exists $M(\Omega) > 0$ and $\rho(\Omega) \in [0, 1]$, such that

$$f(y) - f(x) - \langle f'(x), y - x \rangle \leq \frac{M(\Omega)}{1 + \rho(\Omega)} \|y - x\|^{1+\rho(\Omega)}, \quad \forall x, y \in \Omega. \quad (3.2)$$

The above assumption on $f(\cdot)$ covers unconstrained non-smooth ($\rho(\Omega) = 0$), smooth ($\rho(\Omega) = 1$) and weakly smooth ($0 < \rho(\Omega) < 1$) CPs.

We first present a generic algorithmic framework to solve unconstrained CP problems through solutions to a series of ball-constrained CPs in Subsection 3.1, then extend the FAPL and FUSL methods to unconstrained CPs with iteration complexity analysis in Subsection 3.2.

3.1. Expansion algorithm for unconstrained CP problems. In order to analyze the computational complexity of BL methods, it is commonly assumed that the feasible set of interest is compact (e.g., [21, 4, 20]). While the compactness assumption is essential for performing complexity analysis, the applicability of existing BL methods on unconstrained CP may be impaired. In practice, one possible workaround for BL methods to solve (3.1) may involve an assumption on the distance from a point \bar{x} to an optimal solution x^* , namely, $x^* \in B(\bar{x}, R) := \{x \in \mathbb{R}^n : \|x - \bar{x}\| \leq R\}$ for some \bar{x} and R . With such an assumption, we can solve (3.1) by considering a ball-constrained CP

$$f_{\bar{x}, R}^* := \min_{x \in B(\bar{x}, R)} f(x). \quad (3.3)$$

It should be noted that, while the above equivalent reformulation seems straightforward, its computational complexity relies almost exclusively on the radius R . In particular, if R is close to the distance from \bar{x} to X^* , the optimal solution set of (3.1), i.e., $R - D^*$ is small enough, where

$$x^* := \operatorname{argmin}_x \{\|\bar{x} - x\| : x \in X^*\} \quad \text{and} \quad D^* := \|\bar{x} - x^*\|, \quad (3.4)$$

then the computational complexity for computing approximate solutions to the ball-constrained CP (3.3) and the original unconstrained CP (3.1) are close, and it is definitely reasonable to solve (3.3) instead. However, if R is severely overestimated, the computational complexity for solving the ball-constrained CP (3.3) may become much higher than the optimal complexity bound that depends on D^* .

Based on the above discussion, we can conclude that a good BL method should tackle the unconstrained CP (3.1) from two perspectives. Firstly, without any satisfiable knowledge regarding D^* , such BL method should still be able to solve (3.1) with optimal complexity bound that depends on D^* . Secondly, if there exists an approximate distance R that is close to D^* , a good BL method should solve the ball-constrained problem (3.3) efficiently. We will consider a generic framework that follows the former perspective in this subsection, and then extend the BL method to solve (3.1) in the next subsection.

In this subsection, our goal is to design a generic algorithmic framework that follows the former perspective in the above discussion, and solve the unconstrained CP (3.1) through solutions to a series of ball-constrained CPs. It should be noted that such concept indeed follows naturally from the two perspectives discussed above: if a BL method is computationally efficient for ball-constrained CPs of form (3.3), starting from certain R for (3.3) and enlarging it by two each time, we will eventually reach a close enough estimate of D^* after logarithmic amount of times.

Given $\bar{x} \in \mathbb{R}^n$, $R > 0$, $\epsilon > 0$, let us assume that there exists a first-order algorithm, denoted by $\mathcal{A}(\bar{x}, R, \epsilon)$, which can find an ϵ -solution to (3.3). In other words, we assume that each call to $\mathcal{A}(\bar{x}, R, \epsilon)$ will compute a point $z \in B(\bar{x}, R)$ such that $f(z) - f_{\bar{x}, R}^* \leq \epsilon$. Moreover, throughout this section, we assume that the number

of (sub)gradient evaluations required by $\mathcal{A}(\bar{x}, R, \epsilon)$ for finding an ϵ -solution to (3.3) is bounded by

$$N_{\bar{x}, R, \epsilon} := \frac{C_1(\bar{x}, R, f)R^{\alpha_1}}{\epsilon^{\beta_1}} + \frac{C_2(\bar{x}, R, f)R^{\alpha_2}}{\epsilon^{\beta_2}}, \quad (3.5)$$

where $\alpha_1 \geq \beta_1 > 0$ and $\alpha_2 \geq \beta_2 > 0$. $C_1(\bar{x}, R, f)$ and $C_2(\bar{x}, R, f)$ are constants that depend on f in (3.3) and nondecreasing with respect to R . For example, if f is a smooth convex function, ∇f is Lipschitz continuous in \mathbb{R}^n with constant L , i.e., (??) holds with $\rho(\mathbb{R}^n) = 1$ and $M(\mathbb{R}^n) = L$, and we apply the APL method to (3.3), then we have only one term with $\alpha_1 = 1$, $\beta_1 = 1/2$, and $C_1(\bar{x}, R, f) = c\sqrt{L}$ in (3.5), where c is a universal constant. Observe that the two complexity terms in (3.5) will be useful for analyzing some structured CP problems in Section 2.2. It should also be noted that a more accurate estimate of $C_1(\bar{x}, R, f)$ is $cM(B(\bar{x}, R))$, since the Lipschitz constant $L = M(\mathbb{R}^n)$ throughout \mathbb{R}^n is larger than or equal to the local Lipschitz constant on $B(\bar{x}, R)$.

Let \mathcal{A} be the algorithm that can find an ϵ -solution to ball-constrained problem (3.3). By utilizing a novel guess and check procedure, we present an expansion algorithm for unconstrained convex optimizations as follows:

Algorithm 3 Expansion algorithm for unconstrained CP problems

Choose an arbitrary $r_1 > 1$ and compute the initial gap $\Delta_1 := f(\bar{x}) - \min_{x \in B(\bar{x}, r_1)} h(\bar{x}, x)$.

For $k = 1, 2, \dots$,

Step 1. Solve $\bar{x}'_k = \mathcal{A}(\bar{x}, r_k, \Delta_k)$, then use \bar{x}'_k as the initial point to solve $\bar{x}''_k = \mathcal{A}(\bar{x}, 2r_k, \Delta_k)$.

Step 2. If $f(\bar{x}'_k) - f(\bar{x}''_k) > \Delta_k$, update $r_k \leftarrow 2r_k$ and go to Step 1.

Step 3. Otherwise, output $\bar{x}^*_k = \bar{x}''_k$, and let $\Delta_{k+1} = \Delta_k/2$ and $r_{k+1} = r_k$.

Note that in order to reflect the flexibility we choose r_1 as any positive constant in Algorithm 3. However, from the theoretical complexity point of view, we prefer to starting with a smaller r_1 (i.e., a lower bound on D^*). In fact, given the target accuracy ϵ one can derive a theoretically viable selection of r_1 given as follows. For simplicity, assume that f has Lipschitz continuous gradient with constant L , we have $f(x_0) - f^* \leq \frac{L}{2} \|x_0 - x^*\|^2$. If $\|x_0 - x^*\| \leq \sqrt{2\epsilon/L}$, then x_0 already satisfies $f(x_0) - f^* \leq \epsilon$. Therefore, we can set $r_1 = \sqrt{2\epsilon/L}$. Such an estimate requires some prior information of L . It should be noted that an overestimate on L does not hurt the complexity bound, since r_1 is updated exponentially fast and will approach D^* quickly.

Steps 1 and 2 in Algorithm 3 constitute a loop for finding a pair of solutions $(\bar{x}'_k, \bar{x}''_k)$ satisfying $0 \leq f(\bar{x}'_k) - f(\bar{x}''_k) \leq \Delta_k$ for any $k \geq 1$. Since \bar{x}'_k and \bar{x}''_k are Δ_k -optimal solutions to $\min_{x \in B(\bar{x}, r_k)} f(x)$ and $\min_{x \in B(\bar{x}, 2r_k)} f(x)$ respectively, this loop must terminate in finite time, because it will terminate whenever $r_k \geq D^*$.

For simplicity, we call it an expansion if we double the radius in Step 2. Each iteration may contain several expansions before outputting solution \bar{x}^*_k in Step 3. Note that, for the implementation of Algorithm 3, most computational work exists in Step 1, which involves a sequence of calls to algorithm \mathcal{A} . However, notice that the gaps (Δ_k) and radiuses (r_k) for these calls are monotonically decreasing and increasing, respectively, numerous computational cost could be saved by using results from previous expansions and iterations. For instance, the output solutions of previous calls to \mathcal{A} associated with larger gaps or smaller radiuses could always be used as the starting point for the current call to \mathcal{A} . For successive expansions, if the previous execution of Step 1 called $\mathcal{A}(\bar{x}, r_k, \Delta_k)$ and $\mathcal{A}(\bar{x}, 2r_k, \Delta_k)$ and the current execution of Step 1 calls $\mathcal{A}(\bar{x}, 2r_k, \Delta_k)$ and $\mathcal{A}(\bar{x}, 4r_k, \Delta_k)$, the computation of call $\mathcal{A}(\bar{x}, 2r_k, \Delta_k)$ could be saved. Moreover, if \mathcal{A} is referred to some specific algorithms, such as BL type methods, more previous results, like the lower bounds and prox-centers, could also be utilized to save computational cost.

Note that there are proximal bundle methods incorporating with trust region techniques for updating the new iterate [1, 31]. In [1] an quadratic cutting plane model based method and in [31] the idea of Chebychev center were used to generate the trust regions. These trust region methods restrict the iterates in the trust region for better convergence to the optimal solution, while the approximate solutions in the searching balls generated by our expansion algorithm are used only for checking whether or not the current searching ball needs to be expanded in order to get a better estimate of D^* .

Before analyzing the iteration complexity of Algorithm 3, we discuss some important observations related to the aforementioned expansions.

LEMMA 3.1. *Let $\bar{x} \in \mathbb{R}^n$ and $r > 0$ be a constant, \bar{x}_1 and \bar{x}_2 be Δ -solutions to CP problems*

$$f_{\bar{x},r}^* := \min_{x \in B(\bar{x},r)} f(x) \quad \text{and} \quad f_{\bar{x},2r}^* := \min_{x \in B(\bar{x},2r)} f(x), \quad (3.6)$$

respectively. If $0 \leq f(\bar{x}_1) - f(\bar{x}_2) \leq \Delta$, then we have

$$f(\bar{x}_2) - f^* \leq \left(3 + \frac{2D^*}{r}\right) \Delta, \quad (3.7)$$

where f^* and D^* are defined in (3.1) and (3.4) respectively.

Proof. Clearly, by definition, we have $\|\bar{x}_1 - \bar{x}\| \leq r$, $\|\bar{x}_2 - \bar{x}\| \leq 2r$, $0 \leq f(\bar{x}_1) - f_{\bar{x},r}^* \leq \Delta$ and $0 \leq f(\bar{x}_2) - f_{\bar{x},2r}^* \leq \Delta$. It suffices to consider the case when $f_{\bar{x},2r}^* > f^*$ and $\|x^* - \bar{x}\| > 2r$, since otherwise (3.7) holds trivially. Suppose x_1^* and x_2^* are the solutions to the first and second problems in (3.6) respectively, let \hat{x} be the intersection of the line segment (x^*, x_1^*) with the ball $B(\bar{x}, 2r)$, and denote $R_1 := \|\hat{x} - x_1^*\|$ and $R_2 := \|x^* - x_1^*\|$. Clearly, $\hat{x} = (1 - \frac{R_1}{R_2})x_1^* + \frac{R_1}{R_2}x^*$. By the convexity of $f(\cdot)$, we have

$$f(\hat{x}) \leq \left(1 - \frac{R_1}{R_2}\right)f(x_1^*) + \frac{R_1}{R_2}f(x^*), \quad (3.8)$$

which implies that

$$\frac{R_1}{R_2}[f(x_1^*) - f(x^*)] \leq f(x_1^*) - f(\hat{x}), \quad (3.9)$$

and that $f(\hat{x}) \leq f(x_1^*)$ due to the fact that $f(x^*) \leq f(x_1^*)$. Also, we have $f(\hat{x}) \geq f(x_2^*)$ since $\hat{x} \in B(\bar{x}, 2r)$. In addition,

$$f(x_1^*) - f(x_2^*) = [f(x_1^*) - f(\bar{x}_1)] + [f(\bar{x}_1) - f(\bar{x}_2)] + [f(\bar{x}_2) - f(x_2^*)] \quad (3.10)$$

$$\leq 0 + \Delta + \Delta = 2\Delta. \quad (3.11)$$

Combining the previous inequalities, we obtain

$$f(x_1^*) - 2\Delta \leq f(x_2^*) \leq f(\hat{x}) \leq f(x_1^*), \quad (3.12)$$

which implies that $f(x_1^*) - f(\hat{x}) \leq 2\Delta$. Using (3.9), and the fact that $R_1 \geq r$ and $R_2 \leq D^* + r$, we have

$$f(x_1^*) - f(x^*) \leq \frac{2R_2\Delta}{R_1} \leq \left(2 + \frac{2D^*}{r}\right) \Delta.$$

Therefore,

$$f(\bar{x}_2) - f(x^*) \leq f(\bar{x}_1) - f(x^*) \leq [f(\bar{x}_1) - f(x_1^*)] + [f(x_1^*) - f(x^*)] \leq \left(3 + \frac{2D^*}{r}\right) \Delta.$$

□

We are now ready to prove the iteration complexity of Algorithm 3 for solving the unconstrained CP problem (3.1).

THEOREM 3.2. *Suppose that the number of (sub)gradient evaluations required by $\mathcal{A}(\bar{x}, R, \epsilon)$ for finding an ϵ -solution to (3.3) is bounded by (3.5). For any $k \geq 1$, denote $\epsilon_k := f(\bar{x}_k^*) - f^*$ for the output \bar{x}_k^* in Algorithm 3. Then we have*

- a) $r_k \leq \bar{r} := \max\{r_1, 2D^*\}$ for all $k \geq 1$, where D^* is defined in (3.4);
- b) $\lim_{k \rightarrow \infty} \epsilon_k = 0$;

c) The total number of (sub)gradient evaluations performed by Algorithm 3 for finding the ϵ_k -solution \bar{x}_k^* to problem (3.1) is bounded by

$$\mathcal{O}\left(\frac{C_1(\bar{x}, 2\bar{r}, f)\bar{r}^{\alpha_1}}{\epsilon_k^{\beta_1}} + \frac{C_2(\bar{x}, 2\bar{r}, f)\bar{r}^{\alpha_2}}{\epsilon_k^{\beta_2}}\right). \quad (3.13)$$

Proof. We start by proving a), if $r_1 \geq D^*$, then $f(\bar{x}'_k) - f(\bar{x}''_k) \leq \Delta_k$ for any $k \geq 1$ and no expansion takes place, hence $r_k = r_1 = \bar{r}$. If $r_1 \leq D^*$, from Algorithm 3, we see that expansion occurs at Step 2 if and only if $f(\bar{x}'_k) - f(\bar{x}''_k) > \Delta_k$. Hence, if $r_k \geq D^*$, this condition is not satisfied and no more expansion is performed. This implies $r_k < 2D^*$.

To prove b), observe that \bar{x}'_k is used as the initial point for computing \bar{x}''_k in Algorithm 3 and hence $f(\bar{x}''_k) \leq f(\bar{x}'_k)$. Combining this observation with the condition in Step 3, we have $0 \leq f(\bar{x}'_k) - f(\bar{x}''_k) \leq \Delta_k$. Applying Lemma 3.1 implies

$$\epsilon_k = f(\bar{x}''_k) - f^* \leq \left(3 + \frac{2D^*}{r_k}\right) \Delta_k. \quad (3.14)$$

Note that the total number of expansion is bounded by

$$\bar{S}_1 := \left\lceil \log_2 \frac{\bar{r}}{r_1} \right\rceil, \quad (3.15)$$

hence Δ_k decreases to 0 as k increases, we have $\lim_{k \rightarrow \infty} \epsilon_k = 0$.

To prove c), assume that the number of executions of Step 1 in Algorithm 3 for finding \bar{x}_k^* is K . For any $1 \leq j \leq K$, let $\mathcal{A}(\bar{x}, R_j, \delta_j)$ and $\mathcal{A}(\bar{x}, 2R_j, \delta_j)$ be called in the j^{th} execution of Step 1. By using (3.5) and noting that $C_1(\bar{x}, R, f)$ and $C_2(\bar{x}, R, f)$ are nondecreasing with respect to R , we have the number of (sub)gradient evaluations performed by the j^{th} execution of Step 1 is bounded by

$$N_j := \frac{(1 + 2^{\alpha_1})C_1(\bar{x}, 2R_j, f)R_j^{\alpha_1}}{\delta_j^{\beta_1}} + \frac{(1 + 2^{\alpha_2})C_2(\bar{x}, 2R_j, f)R_j^{\alpha_2}}{\delta_j^{\beta_2}}. \quad (3.16)$$

Let N'_j and N''_j be the first and second terms on the right of (3.16), respectively. The $(j+1)^{\text{th}}$ execution of Step 1 either doubles the radius or reduces the gap by half comparing to the j^{th} execution, i.e., $R_{j+1} = 2R_j$ or $\delta_{j+1} = \delta_j/2$ respectively. Therefore we have either $N'_{j+1} \geq 2^{\alpha_1}N'_j$ and $N''_{j+1} \geq 2^{\alpha_2}N''_j$, or $N'_{j+1} \geq 2^{\beta_1}N'_j$ and $N''_{j+1} \geq 2^{\beta_2}N''_j$. Since $2^{\alpha_1} \geq 2^{\beta_1} > 1$ and $2^{\alpha_2} \geq 2^{\beta_2} > 1$, we can combine these two cases and have

$$N'_j \leq 2^{-\beta_1}N'_{j+1} \text{ and } N''_j \leq 2^{-\beta_2}N''_{j+1}, \text{ for } 1 \leq j \leq K-1, \quad (3.17)$$

which further implies

$$N'_j \leq 2^{-\beta_1(K-j)}N'_K \text{ and } N''_j \leq 2^{-\beta_2(K-j)}N''_K, \text{ for } 1 \leq j \leq K. \quad (3.18)$$

Then the total number of (sub)gradient evaluations performed by these K executions of Step 1 in Algorithm 3 is bounded by

$$N := \sum_{j=1}^K (N'_j + N''_j) \leq N'_K \sum_{j=1}^K 2^{-\beta_1(K-j)} + N''_K \sum_{j=1}^K 2^{-\beta_2(K-j)} \quad (3.19)$$

$$< N'_K \sum_{j=0}^{+\infty} 2^{-\beta_1 j} + N''_K \sum_{j=0}^{+\infty} 2^{-\beta_2 j} \leq \frac{1}{1-2^{-\beta_1}} N'_K + \frac{1}{1-2^{-\beta_2}} N''_K \quad (3.20)$$

$$\leq \frac{(1 + 2^{\alpha_1})C_1(\bar{x}, 2r_k, f)}{1-2^{-\beta_1}} \cdot \frac{r_k^{\alpha_1}}{\Delta_k^{\beta_1}} + \frac{(1 + 2^{\alpha_2})C_2(\bar{x}, 2r_k, f)}{1-2^{-\beta_2}} \cdot \frac{r_k^{\alpha_2}}{\Delta_k^{\beta_2}}. \quad (3.21)$$

The last inequality follows from the observation that the last (i.e., K^{th}) execution of Step 1 before \bar{x}_k^* is found has $R_K = r_k$ and $\delta_K = \Delta_k$. Combining the above inequality with (3.14), we have

$$N < \sum_{i=1}^2 \frac{(1 + 2^{\alpha_i})C_i(\bar{x}, 2r_k, f)}{1 - 2^{-\beta_i}} \cdot \frac{r_k^{\alpha_i} (3 + \frac{2D^*}{r_k})^{\beta_i}}{\epsilon_k^{\beta_i}}. \quad (3.22)$$

Since $\alpha_i \geq \beta_i > 0$, $C_i(\bar{x}, 2r_k, f)r_k^{\alpha_i} (3 + \frac{2D^*}{r_k})^{\beta_i} = C_i(\bar{x}, 2r_k, f)r_k^{\alpha_i - \beta_i} (3r_k + 2D^*)^{\beta_i}$ for $i = 1, 2$ are monotonically increasing with respect to r_k , which, in view of the fact $r_k < 2\bar{r}$ proved by part a), therefore clearly implies

$$N < \sum_{i=1}^2 \frac{(2^{3\beta_i} + 2^{\alpha_i + 3\beta_i})C_i(\bar{x}, 2\bar{r}, f)}{2^{\beta_i} - 1} \cdot \frac{\bar{r}^{\alpha_i}}{\epsilon_k^{\beta_i}}. \quad (3.23)$$

Hence the proof is complete. \square

Note that to solve the unconstrained CP problem (3.1), the termination criterions of most first-order algorithms are based on the residual of the (sub)gradient, which would lead to different complexity analysis. To the best of our knowledge, without any prior information on D^* , there is no verifiable termination criterion based on functional optimality gap that could guarantee the termination of algorithms for finding an ϵ -solution to (3.1). Comparing to Nesterov's optimal gradient method for unconstrained problems in [25], Algorithm 3 only provides efficiency estimates about $\epsilon_k := f(\bar{x}_k^*) - f^*$ when the output \bar{x}_k^* is updated, while the optimal gradient method could have estimates about $\bar{\epsilon}_k := f(x_k) - f^*$ for each iterate x_k . For both methods the efficiency estimates involve D^* . Since Algorithm 3 extends methods for ball-constrained CP problems to solve (3.1), and the iterations in the expansions of Algorithm 3 could be regarded as a guess and check procedure to estimate D^* , it is reasonable that the efficiency estimates are only provided for unexpansive steps, i.e., Step 3 of Algorithm 3, which output \bar{x}_k^* .

It has been shown in [20] that the APL method and its variant, the USL method, achieve the optimal iteration complexities in (3.3) for smooth, non-smooth and weakly smooth CPs and a class of structured saddle point problems respectively, on any convex and compact feasible set. So Algorithm 3 could be incorporated to solve (3.1) with the optimal iteration complexities too. Therefore, the remaining part of this paper will focus on how to improve the efficiency of these BL type methods for solving ball-constrained CPs.

3.2. Extending FAPL and FUSL for unconstrained problems. In this subsection, we study how to utilize the FAPL and FUSL methods to solve the unconstrained problems based on our results in Section 3.

Let us first consider the case when f in (3.1) satisfies (3.2). If the method \mathcal{A} in Step 1 of Algorithm 3 is given as the FAPL method, then by Theorem 2.5, and the fact that only one (sub)gradient of f is computed in each iteration of the FAPL method, the number of (sub)gradient evaluations within one call to $\mathcal{A}(\bar{x}, R, \epsilon)$ is bounded by $N(\epsilon)$ given by (2.30).

Therefore, for the FAPL method, we have $\alpha_1 = \frac{2(1+\rho)}{3+2\rho}$, $\beta_1 = \frac{2}{1+3\rho}$, $C_1(\bar{x}, R, f) = C' M^{\frac{2}{1+3\rho}}$ and $C_2(\bar{x}, R, f) = \alpha_2 = \beta_2 = 0$ in (3.5), where C' is a constant depending on the parameters q, θ, β, ρ and c in the FAPL method. Letting $\epsilon_k := f(\bar{x}_k^*) - f^*$ for $k \geq 1$ in Algorithm 3 and applying Theorem 3.2, we then conclude for finding the ϵ_k -solution \bar{x}_k^* to problem (3.1), the total number of (sub)gradient evaluations performed by Algorithm 3 is bounded by

$$\mathcal{O} \left(\left[\frac{M(2\bar{r})^{1+\rho}}{\epsilon_k} \right]^{\frac{2}{1+3\rho}} \right), \quad (3.24)$$

where $M := M(B(\bar{x}, 2\bar{r}))$, $\rho := \rho(B(\bar{x}, 2\bar{r}))$ and \bar{r} is defined in part a) of Theorem 3.2. It should be noted that the constants M and ρ are local constants that depend on the size of the initial ball, i.e., r_1 , and the distance from \bar{x} and x^* , which are not required for the FAPL method and Algorithm 3, and also generally smaller than the constants $M(\mathbb{R}^n)$ and $\rho(\mathbb{R}^n)$, respectively, for the global Hölder continuity condition.

Moreover, if f in (3.1) is given in the form of (2.28) as a structured non-smooth CP problem, then the FUSL method could be applied to solve the corresponding structured ball-constrained problems in Algorithm 3. By Theorem 2.8, the number of (sub)gradient evaluations of f within one call to $\mathcal{A}(\bar{x}, R, \epsilon)$ is bounded by

$$S_1 + S_2 + C'R\sqrt{\frac{L_{\hat{f}}}{\epsilon}} + \frac{C''R\|A\|}{\epsilon}, \quad (3.25)$$

where C', C'' are some constants depending on the parameters $q, \theta, \beta, \sigma_v, c, D_0$ and $D_{v,Y}$ in the FUSL method.

Applying Theorem 3.2 with $\alpha_1 = \alpha_2 = 1, \beta_1 = \frac{1}{2}, \beta_2 = 1, C_1(\bar{x}, R, f) = C'\sqrt{L_{\hat{f}}}$, and $C_2(\bar{x}, R, f) = C''\|A\|$, we conclude that the total number of (sub)gradient evaluations performed by Algorithm 3 to find the ϵ_k -solution \bar{x}_k^* to problem (3.1)-(2.28) is bounded by

$$\mathcal{O}\left(2C'\bar{r}\sqrt{\frac{L_{\hat{f}}}{\epsilon_k}} + \frac{2C''\bar{r}\|A\|}{\epsilon_k}\right). \quad (3.26)$$

Similar to the FAPL method, here $L_f := L_f(B(\bar{x}, 2\bar{r}))$ is a lower bound of $L_f(\mathbb{R}^n)$.

4. Numerical experiments. In this section we apply the FAPL and FUSL methods to solve a few large-scale CP problems, including the quadratic programming problems with large Lipschitz constants, synthetic and real world total variation based image reconstruction problems, then compare them with some other first-order algorithms. All the algorithms were implemented in MATLAB, Version R2011a and all experiments were performed on a desktop with an Inter Dual Core 2 Duo 3.3 GHz CPU and 8G memory.

4.1. Quadratic programming. The main purpose of this section is to investigate the performance of the FAPL method for solving smooth CP problems especially with large Lipschitz constants and demonstrate the improvement of the FAPL method comparing to some other state-of-the-art BL type methods. And since most state-of-the-art BL type methods require compact feasible sets, we consider the following quadratic programming problem:

$$\min_{\|x\| \leq 1} \|Ax - b\|^2, \quad (4.1)$$

where $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$. We compare the FAPL method with Nesterov's optimal method (NEST) for smooth functions [27], NERML [4], and APL [20]. We also compare the FAPL method with the built-in Matlab linear system solver in view of its good practical performance. In the APL method, the subproblems are solved by MOSEK [23], an efficient software package for linear and second-order cone programming. Two cases with different choices of initial lower bound LB in these experiments are conducted: (1). $LB = 0$ and (2). $LB = -\infty$.

In our experiments, given m and n , two types of matrix A are generated. The first type of matrix A is randomly generated with entries uniformly distributed in $[0,1]$, while the entries of the second type are normally distributed according to $N(0,1)$. We then randomly choose an optimal solution x^* within the unit ball in \mathbb{R}^n , and generate the data b by $b = Ax^*$. We apply all these four methods aforementioned to solve (4.1) with this set of data A and b , and the accuracy of the generated solutions are measured by $e_k = \|Ax_k - b\|^2$. The results are shown in Tables 4.1, 4.2 and 4.3.

In order to investigate the efficiency of solving unconstrained CPs using the proposed expansion algorithm (i.e., Algorithm 3), we conduct two sets of experiments to compare the performance of solving the unconstrained QP problem $\min_{x \in \mathbb{R}^n} \|Ax - b\|^2$ using two different strategies: one starts with small initial feasible ball, then applies the unconstrained FAPL method (i.e., incorporating the expansion algorithm with the FAPL method as subroutine \mathcal{A}), while the other one, under the assumption that a bound on D^* defined in (3.4) is known, applies the ball-constrained FAPL method directly by choosing some large ball that contains at least one optimal solution.

TABLE 4.1
Uniformly distributed QP instances

$A : n = 4000, m = 3000, L = 2.0e6, e_0 = 2.89e4$							
Alg	LB	Iter.	Time	Acc.	Iter.	Time	Acc.
FAPL	0	103	3.06	9.47e-7	142	3.76	8.65e-9
	$-\infty$	277	6.55	5.78e-7	800	19.18	2.24e-11
APL	0	128	37.10	9.07e-7	210	60.85	9.82e-9
	$-\infty$	300	85.65	6.63e-6	800	234.69	2.59e-9
NERML	0	218	58.32	9.06e-7	500	134.62	1.63e-8
	$-\infty$	300	84.01	1.02e-2	800	232.14	1.71e-3
NEST	-	10000	220.1	3.88e-5	20000	440.02	3.93e-6

$A : n = 8000, m = 4000, L = 8.0e6, e_0 = 6.93e4$							
Alg	LB	Iter.	Time	Acc.	Iter.	Time	Acc.
FAPL	0	70	4.67	7.74e-7	95	6.46	6.85e-10
	$-\infty$	149	8.99	6.27e-7	276	16.94	6.10e-10
APL	0	79	71.24	7.79e-7	144	129.52	3.62e-9
	$-\infty$	248	205.48	8.16e-7	416	358.96	8.68e-9
NERML	0	153	128.71	7.30e-7	300	251.79	4.03e-9
	$-\infty$	300	257.54	1.18e-3	800	717.13	9.24e-5
NEST	-	10000	681.03	5.34e-5	20000	1360.52	4.61e-6

FAPL method for large dimension matrix							
Matrix A: $m \times n$	LB	Iter.	Time	Acc.	Iter.	Time	Acc.
10000×20000 L=5.0e7	0	97	36.65	6.41e-11	185	69.31	7.29e-21
	$-\infty$	207	73.70	8.28e-8	800	292.06	2.32e-15
10000×40000 L=1.0e8	0	67	49.95	9.21e-11	122	91.49	7.27e-21
	$-\infty$	130	88.40	7.11e-8	421	295.15	1.95e-16
10000×60000 L=1.5e8	0	52	58.06	7.68e-11	95	106.14	8.43e-21
	$-\infty$	156	160.93	9.84e-8	394	422.4	7.48e-16

Since the performance of both methods would be affected by the distance between the initial point x_0 and the optimal solution set, in both experiments, we set $\bar{x} = 0, D^* \approx 1$, and for any initial ball $B(\bar{x}, R)$, we choose the starting point x_0 randomly within the ball and then normalize x_0 and set $\|x_0\| = \frac{R}{2}$. For both methods, the initial lower bound is set to $-\infty$, and the parameters of the FAPL method are the same. In this first experiment, A is generated randomly with entries uniformly distributed in $[0, 1]$. In the second experiment, we use the worst-case QP instance for first-order methods which are generated by A. Nemirovski (see the construction scheme in [26] and [24]). When we apply the unconstrained FAPL method, the radiuses of the initial balls are chosen as $D/1e5, D/1e4, D/1e3, D/1e2$ and $D/10$, respectively. While the ball-constrained FAPL method is employed, the radiuses of the balls are selected as $1e5 \times D, 1e4 \times D, 1e3 \times D, 1e2 \times D$ and $10 \times D$, respectively. The results are shown in Table 4.4.

The advantages of the FAPL method can be observed from these experiments. Firstly, it is evident that BL type methods have much less iterations than NEST especially when the Lipschitz constant of the objective function is large. Among these three BL type methods, NERML requires much more iterations than APL and FAPL, which have optimal iteration complexity for this problem.

Secondly, comparing to the existing BL type methods (APL and NERML), FAPL has much lower computational cost for each iteration. The computational cost of FAPL for each iteration is just slightly larger than that of NEST. However, the costs of each iteration of APL and NERML are 10 times larger than that of NEST.

Thirdly, consider the difference between the performance of setting the initial lower bound equal to 0 and $-\infty$, it is also evident that FAPL is more robust to the choice of the initial lower bound and it updates the

TABLE 4.2
Gaussian distributed QP instances

$A : n = 4000, m = 3000, L = 2.32e4, e_0 = 2.03e3$							
Alg	LB	Iter.	Time	Acc.	Iter.	Time	Acc.
FAPL	0	105	2.78	8.43e-7	153	4.10	7.84e-10
	$-\infty$	338	8.02	6.86e-7	696	16.58	9.74e-10
APL	0	128	35.12	9.01e-7	172	47.49	9.28e-9
	$-\infty$	639	200.67	7.92e-7	800	258.25	1.03e-7
NERML	0	192	48.44	7.05e-7	276	70.31	1.09e-8
	$-\infty$	300	93.32	3.68e-1	800	257.25	6.41e-2
NEST	-	10000	211.30	7.78e-4	20000	422.78	1.95e-4

$A : n = 8000, m = 4000, L = 2.32e4, e_0 = 2.03e3$							
Alg	LB	Iter.	Time	Acc.	Iter.	Time	Acc.
FAPL	0	49	3.25	8.34e-7	68	4.37	7.88e-10
	$-\infty$	165	9.77	5.17e-7	280	16.18	5.06e-10
APL	0	59	48.91	8.59e-7	78	64.95	1.70e-8
	$-\infty$	300	268.47	9.81e-7	670	637.70	9.42e-10
NERML	0	105	181.23	9.14e-7	133	102.68	1.39e-8
	$-\infty$	300	282.56	9.92e-3	800	760.26	8.32e-4
NEST	-	10000	567.59	3.88e-4	20000	1134.38	9.71e-5

FAPL method for large dimension matrix							
Matrix A: $m \times n$	LB	Iter.	Time	Acc.	Iter.	Time	Acc.
10000×20000 L=5.7e4	0	78	27.88	7.22e-11	145	51.81	6.81e-21
	$-\infty$	228	78.57	9.92e-8	800	280.19	1.37e-15
10000×40000 L=9e4	0	48	34.36	5.97e-11	87	62.24	8.26e-21
	$-\infty$	156	106.12	7.18e-8	390	271.15	4.29e-16
10000×60000 L=1.2e5	0	34	36.30	9.88e-11	65	69.56	7.24e-21
	$-\infty$	98	98.11	9.50e-8	350	361.83	8.34e-16

TABLE 4.3
Comparison to Matlab solver

Matrix A: $m \times n$	Matlab $A \setminus b$		FAPL method		
	Time	Acc.	Iter.	Time	Acc.
Uniform 2000×4000	4.41	5.48e-24	204	3.59	6.76e-23
Uniform 2000×6000	7.12	9.04e-24	155	4.10	9.73e-23
Uniform 2000×8000	9.80	9.46e-24	135	4.45	9.36e-23
Uniform 2000×10000	12.43	1.04e-23	108	4.23	7.30e-23
Gaussian 3000×5000	11.17	5.59e-25	207	6.25	7.18e-23
Gaussian 3000×6000	13.96	1.43e-24	152	5.50	9.59e-23
Gaussian 3000×8000	19.57	1.66e-24	105	4.83	8.17e-23
Gaussian 3000×10000	25.18	1.35e-24	95	5.43	5.81e-23

lower bound more efficiently than the other two BL type methods. Though setting the initial lower bound equal to $-\infty$ increases the numbers of iterations for all these three BL type methods, a close examination reveals that the difference between setting the lower bound to zero and $-\infty$ for FAPL is not so significant as that for APL and NERML, especially for large matrix, for example, the second one in Table 4.1 .

Fourthly, FAPL needs less number of iterations than APL, especially when the required accuracy is high. A plausible explanation is that exactly solving the subproblems provides better updating for the prox-centers, and

TABLE 4.4
Unconstrained QP instances

Uniform instance: $A = \text{rand}(4000, 8000)$							
FAPL-Unconstrained				FAPL-Ball constrained			
Radius	Iter.	Time	Acc.	Radius	Iter.	Time	Acc.
1e-5	1405	73.38	2.14e-11	1e5	1888	104.16	7.07e-10
1e-4	1187	62.70	2.14e-11	1e4	1414	77.24	5.55e-10
1e-3	1128	59.34	6.72e-11	1e3	1172	64.51	9.69e-11
1e-2	933	55.38	9.38e-11	1e2	812	43.99	3.71e-11
1e-1	835	45.89	5.38e-11	1e1	617	33.01	6.94e-11

Worst-case instance: $A = \text{Bdata}(2062, 4124)$							
FAPL-Unconstrained				FAPL-Ball constrained			
Radius	Iter.	Time	Acc.	Radius	Iter.	Time	Acc.
1e-5	1409	23.63	9.93e-8	1e5	4000	77.17	5.32e-1
1e-4	1262	21.57	9.97e-8	1e4	4000	78.20	3.29e-3
1e-3	1135	22.81	9.98e-8	1e3	3500	68.98	6.38e-5
1e-2	848	15.10	9.99e-8	1e2	2500	46.67	1.01e-6
1e-1	847	15.30	9.89e-8	1e1	1602	28.92	9.89e-8

consequently, more accurate prox-centers improve the efficiency of algorithm significantly. The experiments show that, for APL and NERML, it is hard to improve the accuracy beyond 10^{-10} . However, FAPL can keep almost the same speed for decreasing the objective value from 10^6 to 10^{-21} .

Fifthly, we can clearly see from Table 4.3 that FAPL is comparable to or significantly outperforms the built-in Matlab solver for randomly generated linear systems, even though our code is implemented in MATLAB rather than lower-level languages, such as C or FORTRAN. We can expect that the efficiency of FAPL will be much improved by using C or FORTRAN implementation, which has been used in the MATLAB solver for linear systems.

Finally, from Table 4.4 it is evident that the performance of both the unconstrained FAPL method and the ball-constrained FAPL method are affected by the distance between the starting point x_0 and the optimal solution set. And improper estimations on D^* would increase the computational cost significantly. Comparing the results presented in the same rows of the left and right columns in Table 4.4, one can see that when the estimations are proper (i.e., R is close to D^*), the performance of both methods are comparable, and when the estimations are improper, the unconstrained FAPL method requires much less iterations and CPU time and can achieve higher accuracy than the ball-constrained FAPL method. Especially, for the worst-case instance, the advantages are much more significant.

In summary, due to its low iteration cost and effective usage of the memory of first-order information, the FAPL method is a powerful tool for solving ball-constrained smooth CP problems especially when the number of variables is huge and/or the value of Lipschitz constant is large. And by incorporating the proposed Expansion Algorithm, the unconstrained FAPL method is very efficient for solving unconstrained CPs especially when a proper estimation on the optimal solution set is not available.

4.2. Total variation based image reconstruction. In this subsection, we apply the FUSL method to solve the non-smooth total variation (TV) based image reconstruction problem:

$$\min_{u \in \mathbb{R}^N} \frac{1}{2} \|Au - b\|_2^2 + \lambda \|u\|_{TV}, \quad (4.2)$$

where A is a given matrix, u is the vector form of the image to be reconstructed, b represents the observed data, and $\|\cdot\|_{TV}$ is the discrete TV semi-norm defined by

$$\|u\|_{TV} := \sum_{i=1}^N \|D_i u\|_2, \quad (4.3)$$

where $D_i u \in \mathbb{R}^2$ is the discrete gradient (finite differences along the coordinate directions) of the i^{th} component of u , and N is the number of pixels in the image. Note that $\|u\|_{TV}$ is convex and non-smooth.

One of the approaches to solve this problem is to consider the associated dual or primal-dual formulations of (4.3) based on the dual formulation of the TV norm:

$$\|u\|_{TV} = \max_{p \in Y} \langle p, Du \rangle, \text{ where } Y = \{p = (p_1, \dots, p_N) \in \mathbb{R}^{2N} : p_i \in \mathbb{R}^2, \|p_i\|_2 \leq 1, 1 \leq i \leq N\}. \quad (4.4)$$

Consequently, we can rewrite (4.2) as a saddle-point problem:

$$\min_{u \in \mathbb{R}^n} \max_{p \in Y} \frac{1}{2} \|Au - b\|_2^2 + \lambda \langle p, Du \rangle. \quad (4.5)$$

Note that (4.5) is exactly the form we considered in the USL and FUSL methods if we set $\hat{g}(y) = 0$. Specifically, the prox-function $v(y)$ on Y is simply chosen as $v(y) = \frac{1}{2} \|y\|_2^2$ in these smoothing techniques.

In our experiments, we consider two types of instances depending on how the matrix A is generated. Specifically, for the first case, the entries of A are normally distributed, while for the second one, the entries are uniformly distributed. For both types of instances, first, we generate the matrix $A \in \mathbb{R}^{m \times n}$, then choose some true image x_{true} and convert it to a vector, and finally compute b by $b = Ax_{true} + \epsilon$, where ϵ is the Gaussian noise with distribution $\epsilon = N(0, \sigma)$. We compare the following algorithms: the accelerated primal dual (APD) method [6], Nesterov's smoothing (NEST-S) method [27, 2], and the FUSL method.

For our first experiment, the matrix A is randomly generated of size $4,096 \times 16,384$ with entries normally distributed according to $N(0, \sqrt{4,096})$, the image x_{true} is a 128×128 Shepp-Logan phantom generated by MATLAB. Moreover, we set $\lambda = 10^{-3}$ and the standard deviation $\sigma = 10^{-3}$. The values of Lipschitz constants are provided for APD and NEST-S, and the initial lower bound for FUSL is set to 0. We run 300 iterations for each algorithm, and present the objective values of problem (4.2) and the relative errors defined by $\|x_k - x_{true}\|_2 / \|x_{true}\|_2$ in Figure 4.1. In our second experiment, the matrix A is randomly generated with entries uniformly distributed in $[0, 1]$. We use a 200×200 brain image [7] as the true image x_{true} , and set $m = 20,000, \lambda = 10, \sigma = 10^{-2}$. Other setup is the same as the first experiment, and the results are shown in Figure 4.2.

We make some observations about the results in Figures 4.1 and 4.2. For the first experiment, there is almost no difference between APD and NEST-S, but FUSL outperforms both of them after 5 seconds in terms of both objective value and relative error. The second experiment clearly demonstrates the advantage of FUSL for solving CP problems with large Lipschitz constants. The Lipschitz constant of matrix A in this instance is about 2×10^8 , much larger than the Lipschitz constant (about 5.9) in the first experiment. FUSL still converges quickly and decreases the relative error to 0.05 in less than 100 iterations, while APD and NEST-S converge very slowly and more than 1,000 iterations are required due to the large Lipschitz constants. It seems that FUSL is not so sensitive to the Lipschitz constants as the other two methods. This feature of FUSL makes it more efficient for solving large-scale CP problems which often have big Lipschitz constants.

In summary, for the TV-based image reconstruction problem (4.2), FUSL not only enjoys the completely parameter-free property (and hence no need to estimate the Lipschitz constant), but also demonstrates significant advantages for its speed of convergence and its solution quality in terms of relative error, especially for large-scale problems.

4.3. Partially parallel imaging. In this subsection, we compare the performance of the FUSL method with several related algorithms in reconstruction of magnetic resonance (MR) images from partial parallel imaging (PPI), to further confirm the observations on advantages of the FUSL method. The detailed background and description of PPI reconstruction can be found in [7]. This image reconstruction problem can be modeled as

$$\min_{u \in \mathbb{C}^n} \sum_{j=1}^k \|M\mathcal{F}S_j u - f_j\|^2 + \lambda \sum_{i=1}^N \|D_i u\|_2,$$

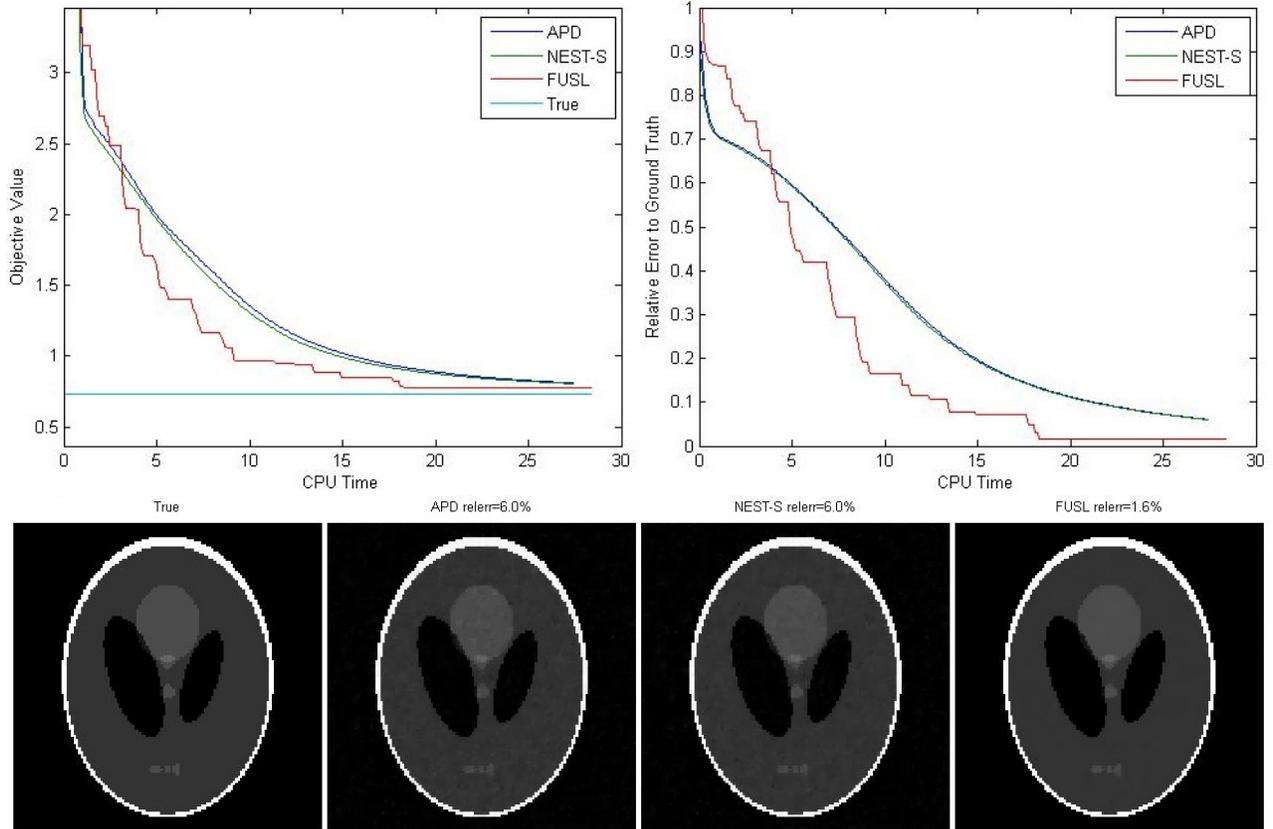


FIG. 4.1. TV-based reconstruction (Shepp-Logan phantom)

where $n = 2$ (we consider two dimensional case), u is the N-vector form of a two-dimensional complex valued image to be reconstructed, k is the number of coils (consider them as sensors) in the magnetic resonance (MR) parallel imaging system. $F \in C^{n \times n}$ is a 2D discrete Fourier transform matrix, $S_j \in C^{n \times n}$ is the sensitivity map of the j -th sensor, and $M \in R^{n \times n}$ is a binary mask describes the scanning pattern. Note that the percentages of nonzero elements in M describes the compression ratio of PPI scan. In our experiments, the sensitivity maps $\{S_j\}_{j=1}^k$ are shown in Figure 4.3, the image x_{true} is of size 512×512 shown in Figures 4.4 and 4.5, and the measurements $\{f_j\}$ are generated by

$$f_j = M(FS_j x_{true} + \epsilon_j^{re}/\sqrt{2} + \epsilon_j^{im}/\sqrt{-2}), \quad j = 1, \dots, k, \quad (4.6)$$

where $\epsilon_j^{re}, \epsilon_j^{im}$ are the noise with entries independently distributed according to $N(0, \sigma)$. We conduct two experiments on this data set with different acquisition rates, and compare the FUSL method to NEST-S method, and the accelerated linearized alternating direction of multipliers (AL-ADMM) with line-search method [32].

For both experiments, set $\sigma = 3 \times 10^{-2}, \lambda = 10^{-5}$, and $\{f_j\}_{j=1}^k$ are generated by (4.6). In the first experiment, we use Cartesian mask with acquisition rate 14%: acquire image in one row for every successive seven rows, while for the second one, we use Cartesian mask with acquisition rate 10%: acquire image in one row for every successive ten rows. The two masks are shown in Figure 4.3. The results of the first and second experiments are shown in Figures 4.4 and 4.5, respectively. These experiments again demonstrate the advantages of the FUSL method over these state-of-the-art techniques for PPI image reconstruction,

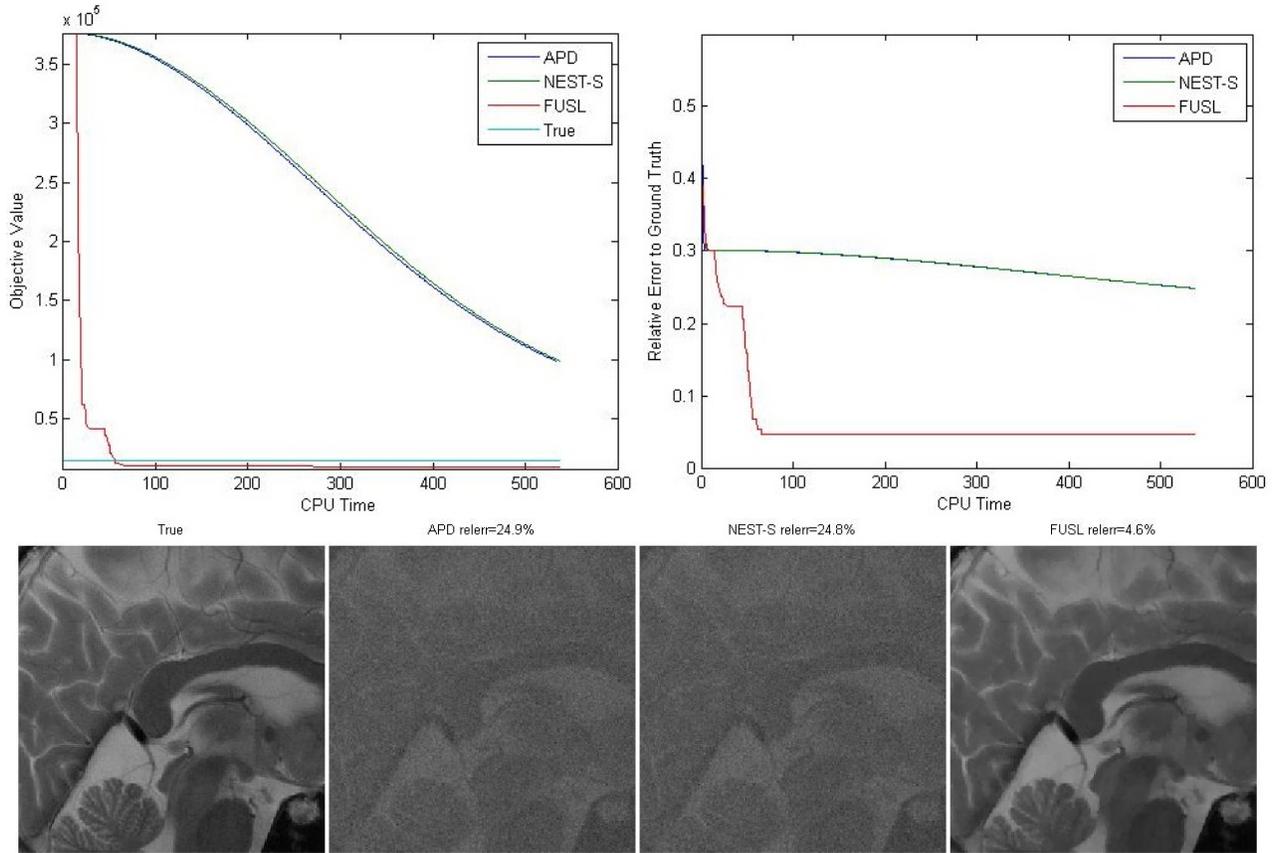


FIG. 4.2. *TV-based reconstruction (brain image)*

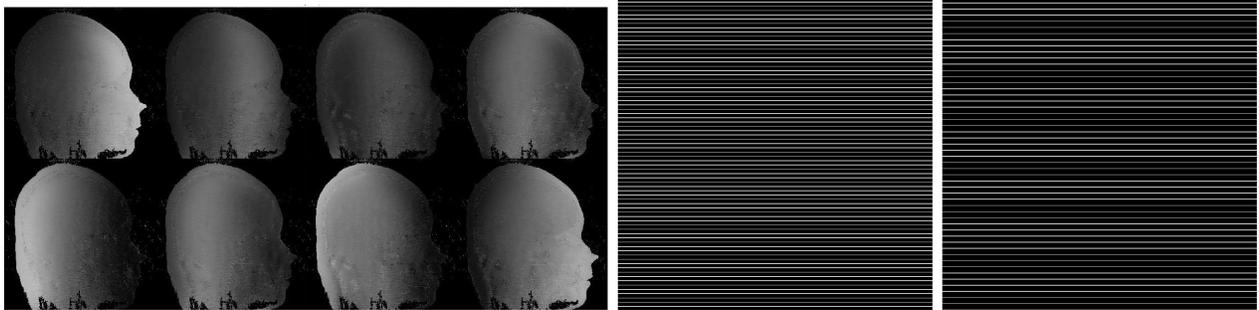


FIG. 4.3. *Sensitivity map and Cartesian masks*

5. Concluding remarks. In this paper, we presented two new BL type methods, the FAPL and FUSL methods, to uniformly solve smooth, non-smooth, and weakly smooth CP problems and a class of structured non-smooth problems with optimal iteration complexities. Because of the use of cutting plane model, technique of restricted memory and an efficient and scalable solver for solving involved subproblem, the proposed methods have lower iteration cost, and can find a solution with higher accuracy within less number of iterations than many gradient descent type methods. Moreover, these BL type methods do not require the input of any problem parameter, or involve any stepsize that could be affected by large Lipschitz constant of the objective function and large dimension of the problem. These built-in features of the proposed methods are very

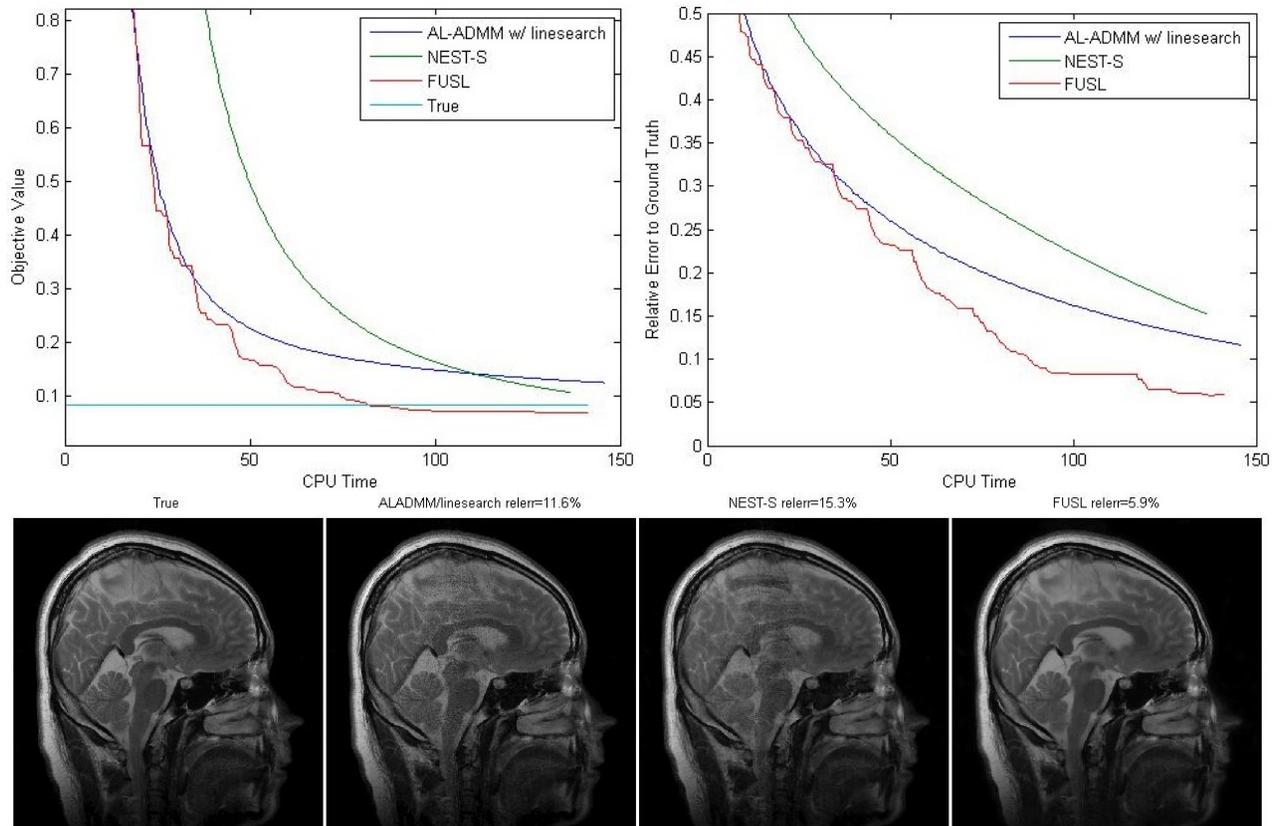


FIG. 4.4. PPI image reconstruction (acquisition rate: 14%)

important for solving large-scale CPs. Our numerical results of least squares problems and total variation based image reconstructions clearly demonstrate the advantages of the FAPL and FUSL methods over the original APL, USL and some other state-of-the-art first-order methods.

REFERENCES

- [1] Annabella Astorino, Antonio Frangioni, Manlio Gaudio, and Enrico Gorgone. Piecewise-quadratic approximations in convex numerical optimization. *SIAM Journal on Optimization*, 21(4):1418–1438, 2011.
- [2] Stephen Becker, Jérôme Bobin, and Emmanuel J Candès. NESTA: a fast and accurate first-order method for sparse recovery. *SIAM Journal on Imaging Sciences*, 4(1):1–39, 2011.
- [3] JY Bello Cruz and W de Oliveira. Level bundle-like algorithms for convex optimization. *Journal of Global Optimization*, pages 1–23, 2013.
- [4] A. Ben-Tal and A. S. Nemirovski. Non-Euclidean restricted memory level method for large-scale convex optimization. *Mathematical Programming*, 102:407–456, 2005.
- [5] Ulf Brännlund, Krzysztof C Kiwiel, and Per Olof Lindberg. A descent proximal level bundle method for convex nondifferentiable optimization. *Operations Research Letters*, 17(3):121–126, 1995.
- [6] Y. Chen, G. Lan, and Y. Ouyang. Optimal primal-dual methods for a class of saddle point problems. *SIAM Journal on Optimization*, 2014. to appear.
- [7] Yunmei Chen, William Hager, Feng Huang, Dzung Phan, Xiaojing Ye, and Wotao Yin. Fast algorithms for image reconstruction with application to partially parallel mr imaging. *SIAM Journal on Imaging Sciences*, 5(1):90–118, 2012.
- [8] W de Oliveira and CLAUDIA Sagastizábal. Level bundle methods for oracles with on-demand accuracy. *Optimization Methods and Software*, (ahead-of-print):1–30, 2014.
- [9] Wei Deng and Wotao Yin. On the global and linear convergence of the generalized alternating direction method of multipliers. *Journal of Scientific Computing*, 66(3):889–916, 2016.
- [10] Donald Goldfarb, Shiqian Ma, and Katya Scheinberg. Fast alternating linearization methods for minimizing the sum of two

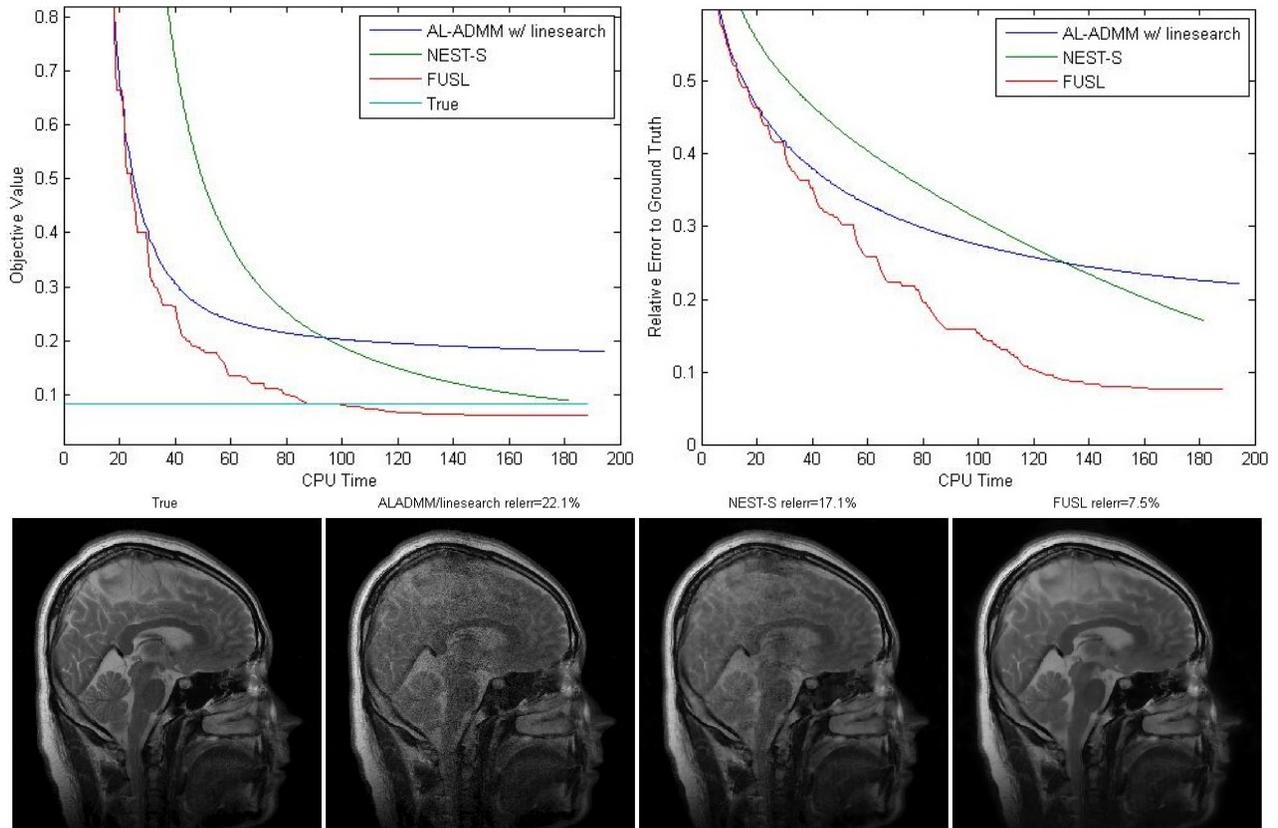


FIG. 4.5. PPI image reconstruction (acquisition rate: 10%)

- convex functions. *Mathematical Programming*, 141(1-2):349–382, 2013.
- [11] Tom Goldstein, Brendan O’Donoghue, Simon Setzer, and Richard Baraniuk. Fast alternating direction optimization methods. *SIAM Journal on Imaging Sciences*, 7(3):1588–1623, 2014.
- [12] Yunlong He and Renato DC Monteiro. An accelerated hpe-type algorithm for a class of composite convex-concave saddle-point problems. *SIAM Journal on Optimization*, 26(1):29–56, 2016.
- [13] Mingyi Hong and Zhi-Quan Luo. On the linear convergence of the alternating direction method of multipliers. *arXiv preprint arXiv:1208.3922*, 2012.
- [14] J.E. Kelley. The cutting plane method for solving convex programs. *Journal of the SIAM*, 8:703–712, 1960.
- [15] K.C. Kiwiel. Proximity control in bundle methods for convex nondifferentiable minimization. *Mathematical Programming*, 46:105–122, 1990.
- [16] K.C. Kiwiel. Proximal level bundle method for convex nondifferentiable optimization, saddle point problems and variational inequalities. *Mathematical Programming, Series B*, 69:89–109, 1995.
- [17] KC Kiwiel. Bundle methods for convex minimization with partially inexact oracles. *Comput. Optim. Appl.*, to appear, 2009.
- [18] Krzysztof C Kiwiel. A proximal bundle method with approximate subgradient linearizations. *SIAM Journal on optimization*, 16(4):1007–1023, 2006.
- [19] Krzysztof C Kiwiel and Claude Lemaréchal. An inexact bundle variant suited to column generation. *Mathematical programming*, 118(1):177–206, 2009.
- [20] Guanghui Lan. Bundle-level type methods uniformly optimal for smooth and nonsmooth convex optimization. *Mathematical Programming*, 149(1-2):1–45, 2015.
- [21] C. Lemaréchal, A. S. Nemirovski, and Y. E. Nesterov. New variants of bundle methods. *Mathematical Programming*, 69:111–148, 1995.
- [22] Renato DC Monteiro and Benar F Svaiter. Iteration-complexity of block-decomposition algorithms and the alternating direction method of multipliers. *SIAM Journal on Optimization*, 23(1):475–507, 2013.
- [23] Mosek. The mosek optimization toolbox for matlab manual. version 6.0 (revision 93). <http://www.mosek.com>.
- [24] A. S. Nemirovski and D. Yudin. *Problem complexity and method efficiency in optimization*. Wiley-Interscience Series in Discrete Mathematics. John Wiley, XV, 1983.

- [25] Y. E. Nesterov. A method for unconstrained convex minimization problem with the rate of convergence $O(1/k^2)$. *Doklady AN SSSR*, 269:543–547, 1983.
- [26] Y. E. Nesterov. *Introductory Lectures on Convex Optimization: a basic course*. Kluwer Academic Publishers, Massachusetts, 2004.
- [27] Y. E. Nesterov. Smooth minimization of nonsmooth functions. *Mathematical Programming*, 103:127–152, 2005.
- [28] W Oliveira, C Sagastizábal, and C Lemaréchal. Bundle methods in depth: a unified analysis for inexact oracles. *Optim Online Rep (Submitted)*, 2013.
- [29] Wellington Oliveira, Claudia Sagastizábal, and Susana Scheimberg. Inexact bundle methods for two-stage stochastic programming. *SIAM Journal on Optimization*, 21(2):517–544, 2011.
- [30] Wellington de Oliveira and Claudia Sagastizábal. Bundle methods in the xxist century: A bird’s-eye view. *Pesquisa Operacional*, 34(3):647–670, 2014.
- [31] Adam Ouorou. A proximal cutting plane method using chebychev center for nonsmooth convex optimization. *Mathematical Programming*, 119(2):239–271, 2009.
- [32] Yuyuan Ouyang, Yunmei Chen, Guanghui Lan, and Eduardo Pasiliao Jr. An accelerated linearized alternating direction method of multipliers. *arXiv preprint arXiv:1401.6607*, 2014.
- [33] Peter Richtárik. Approximate level method for nonsmooth convex minimization. *Journal of Optimization Theory and Applications*, 152(2):334–350, 2012.
- [34] Wim van Ackooij and Claudia Sagastizábal. Constrained bundle methods for upper inexact oracles with application to joint chance constrained energy problems. *SIAM Journal on Optimization*, 24(2):733–765, 2014.