

# A Joint Row and Column Action Method for Cone-Beam Computed Tomography

Yushan Gao  and Thomas Blumensath 

**Abstract**—The inversion of linear systems is fundamental in computed tomography (CT) reconstruction. Computational challenges arise when trying to invert large linear systems, as limited computing resources mean that only a part of the system can be kept in computer memory at any one time. In linear tomographic inversion problems, such as X-ray tomography, even a standard scan can produce millions of individual measurements and the reconstruction of X-ray attenuation profiles typically requires the estimation of a million attenuation coefficients. To deal with the large data sets encountered in real applications and to efficiently utilize modern graphics processing unit based computing architectures, combinations of iterative reconstruction algorithms and parallel computing schemes are increasingly applied. Whilst different parallel methods have been proposed, individual computations currently need to access either the entire set of observations or estimated X-ray absorptions, which can be prohibitive in many realistic applications. We present a fully parallelizable CT image reconstruction algorithm where each computation node works on arbitrary partial subsets of the data and the reconstructed volume. We further develop a non-homogeneously randomized selection criterion which guarantees that submatrices of the system matrix are selected more frequently if they are dense, thus maximizing information flow through the algorithm. We compare our algorithm with block alternating direction method of multipliers and show that our method is significantly faster for CT reconstruction.

**Index Terms**—CT image reconstruction, parallel computing, gradient descent, coordinate descent, linear inverse problems.

## I. INTRODUCTION

IN TRANSMISSION computed tomography (CT), standard scan trajectories, such as rotation based or helical trajectories, allow the use of efficient analytical reconstruction techniques such as the filtered backprojection algorithm (FBP) [1], [2] and the Feldkamp Davis Kress (FDK) [3], [4] method. However, in low signal to noise settings, if scan angles are under-sampled or if nonstandard trajectories are used, then less efficient, algebraic reconstruction methods can provide significantly better reconstructions [5]–[8]. These methods model the

Manuscript received April 23, 2018; revised July 4, 2018; accepted July 12, 2018. Date of publication July 18, 2018; date of current version November 6, 2018. This work was supported in part by EPSRC under Grants EP/K029150/1 and EP/R002495/1, in part by the University of Southampton PGR scholarship, in part by the Faculty of Engineering and the Environment Lancaster Studentship, and in part by the China Scholarship Council. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Ali Bilgin. (Corresponding author: Yushan Gao.)

The authors are with the Faculty of Engineering and Environment, University of Southampton, Southampton, SO17 1BJ, U.K. (e-mail: yg3n15@soton.ac.uk; Thomas.Blumensath@soton.ac.uk).

Digital Object Identifier 10.1109/TCI.2018.2857446

x-ray system as a linear system of equations [9]–[11]:

$$\begin{bmatrix} y_1 \\ \vdots \\ y_r \end{bmatrix} = \begin{bmatrix} a_{11} & \dots & a_{1c} \\ \vdots & \vdots & \vdots \\ a_{r1} & \dots & a_{rc} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_c \end{bmatrix} + \begin{bmatrix} e_1 \\ \vdots \\ e_r \end{bmatrix}, \quad (1)$$

where  $\mathbf{y} = [y_1, \dots, y_r]^T$ ,  $\mathbf{x} = [x_1, \dots, x_c]^T$  and  $\mathbf{e} = [e_1, \dots, e_r]^T$  are projection data, reconstructed image vector and measurement noise respectively. The system matrix  $A \in \mathbb{R}^{r \times c}$ , with non-negative elements  $a_{cd}$  ( $1 \leq c \leq r$ ,  $1 \leq d \leq c$ ), can be computed by Siddon's method [12]. In reality, industrial CT is often used in cases where  $\mathbf{y}$  and  $\mathbf{x}$  can have millions of entries each [13] and where  $A$ , even though it is a sparse matrix, can have billions of entries. In these situations, solving Eq. 1 directly by calculating  $A^{-1}$  is infeasible. More feasible approaches use  $A$  and  $A^T$  directly (instead of their inverse) to iteratively find an approximate solution of Eq. 1.

There are many mature and efficient algorithms to solve Eq. 1, including the conjugate gradient (CG) method and LSQR. These methods can be applied in small scale CT reconstructions when the system matrix  $A$  can be stored in computer memory [14], [15]. When  $A$  is too large to be kept in memory, it is often more efficient to re-compute it on the fly during each iteration, which can be done relatively efficiently using modern graphical processor unit (GPU). However, GPUs have limited internal memory.  $A$  thus has to be broken into smaller subsets so that the GPU only operates on a subset of the data at a time. Whilst it is possible to sequentially process all data in this fashion, this requires constant data transfer to the GPU's internal memory. An alternative is the use of optimisation algorithms that only work on subsets of the matrix  $A$  at any one time. The CG and LSQR methods, whilst having their own block forms [16], [17], do not work with a single subset of  $A$  per iteration. As a result, other algorithms for CT reconstructions have been proposed. Currently, most of these methods can be divided into two categories: *row action methods*, which operate on subsets of the observations  $\mathbf{y}$  at a time and *column action methods*, which operate on subsets of the voxels  $\mathbf{x}$  at a time [18]–[22].

**Row action methods** divides the matrix  $A$  into several row blocks and the system equation thus becomes

$$\begin{bmatrix} \mathbf{y}_{I_1} \\ \vdots \\ \mathbf{y}_{I_M} \end{bmatrix} \approx \begin{bmatrix} A_{I_1} \\ \vdots \\ A_{I_M} \end{bmatrix} \mathbf{x}, \quad (2)$$

where  $A_{I_i} \in \mathbb{R}^{m_i \times c}$  is the row block of system matrix  $A$  and  $\mathbf{y}_{I_i} \in \mathbb{R}^{m_i \times 1}$  is the block of projection  $\mathbf{y}$ . The total block

number is  $M$  and  $\sum_{i=1}^M m_i = r$ . The general iteration scheme in row action method can be summarised as

$$\mathbf{x}^k = \mathbf{x}^{k-1} + R(\mathbf{y}_I - A_I \mathbf{x}), \quad (3)$$

where  $\mathbf{x}^k$  is the  $k^{\text{th}}$  iteration result,  $R$  is a relaxation matrix to control the iteration step length and  $I \in \{I_i\}_{i=1}^M$ . In CT reconstruction, most mature methods are row action methods. These include the classical Kaczmarz family of algorithms (also known as the Algebraic Reconstruction Technique (ART))[7], [23]–[25], the Simultaneous Algebraic Reconstruction Technique (SART) [19], [26], [27], the Simultaneous Iterative Reconstruction Technique (SIRT), [28]–[31] and component averaging (CAV) and its block form (BICAV) [32]–[34].

**Column action methods**, as a counterpart of row action method, divide the system matrix  $A$  into several column blocks. Eq. 1 thus is divided as

$$\mathbf{y} \approx [A^{J_1} \quad \cdots \quad A^{J_N}] \begin{bmatrix} \mathbf{x}_{J_1} \\ \vdots \\ \mathbf{x}_{J_N} \end{bmatrix}, \quad (4)$$

where  $A^{J_j} \in \mathbb{R}^{r \times n_j}$  is the column block of system matrix  $A$  and  $\mathbf{x}_{J_j} \in \mathbb{R}^{n_j \times 1}$  is the block of reconstructed vector  $\mathbf{x}$ . The total block number is  $N$  and  $\sum_{j=1}^N n_j = c$ . The general iteration scheme in column action methods is of the form

$$\begin{aligned} \mathbf{x}_J^k &= \mathbf{x}_J^{k-1} + R\mathbf{r} \\ \mathbf{r} &= \mathbf{r} - A(\mathbf{x}_J^k - \mathbf{x}_J^{k-1}), \end{aligned} \quad (5)$$

where  $J \in \{J_j\}_{j=1}^N$  and the initial  $\mathbf{r} = \mathbf{y}$ . The application of column action methods in CT reconstruction can be traced back to 1990s [35]–[37]. The preliminary column action method updates one voxel each time and it is easy to extend the preliminary method into group form [38]–[43]. Similar to row action methods, the selection criteria on which column blocks to update can be randomised either uniformly [44], [45] or based on specific probability [46], [47].

Row and column action methods also allow for parallel computation. For parallel row action methods, each processor (or node) needs to store the whole reconstructed image vector  $\mathbf{x}$  since each update is on the entire image, whilst for column action methods, each node needs to store all of  $\mathbf{y}$ . As a result, the largest volume they can reconstruct is limited by the storage capacity of the computation nodes. An important exception to this is discussed in [48], where a row action method SIRT is discussed in which each node only requires parts of the reconstructed image vector  $\mathbf{x}$ . However, this method only works for circular scan trajectories and its scalability is limited due to the requirement that overlap between projections of adjacent image blocks should be small.

**General sub-block methods** are proposed in machine learning (ML). There has been interest in the development of methods that use more general updates, updating subsets of  $\mathbf{x}$  with only subsets of  $\mathbf{y}$  at a time [49], [50]. In this paper, these algorithms will be called ‘‘stochastic block coordinate descent’’ (SBCD). In particular, [51] and [52] independently proposed similar SBCD algorithms and both explored the application of variance reduction technique [53] to further accelerate the

convergence rate. [54] proposed a semi-stochastic coordinate descent method to combine the stochastic gradient method and coordinate descent method to minimise a strongly convex problem. [55] mathematically proved the convergence rate of SBCD when the step length is decreasing and when the update strategy adopts a Gauss-Seidel type approach (updating the current column block depends on the previously updated column block), showing that the SBCD method has the same convergence rate as stochastic gradient methods when the objective function is convex. [50], [56] separately proposed the optimal sampling method in the SBCD method by randomly selecting column blocks and selecting row blocks based on a calculated probability.

The partition of  $A$  in SBCD is the same as our method. To define this, we partition  $A$  into  $M \times N$  blocks. Let  $I_i (1 \leq i \leq M)$  be an index set that indexes  $m_i$  ( $\sum_{i=1}^M m_i = r$ ) rows in  $A$  and  $J_j (1 \leq j \leq N)$  be an index set that indexes  $n_j$  ( $\sum_{j=1}^N n_j = c$ ) columns in  $A$ . The matrix  $A_{I_i}^{J_j}$  thus is a sub-matrix of  $A$  with row indexes  $I_i$  and column indexes  $J_j$ . Thus we can divide the linear system into  $M \times N$  blocks:

$$\begin{bmatrix} \mathbf{y}_{I_1} \\ \vdots \\ \mathbf{y}_{I_M} \end{bmatrix} \approx \begin{bmatrix} A_{I_1}^{J_1} & \cdots & A_{I_1}^{J_N} \\ \vdots & \vdots & \vdots \\ A_{I_M}^{J_1} & \cdots & A_{I_M}^{J_N} \end{bmatrix} \begin{bmatrix} \mathbf{x}_{J_1} \\ \vdots \\ \mathbf{x}_{J_N} \end{bmatrix} \equiv \begin{bmatrix} A_{I_1} \\ \vdots \\ A_{I_M} \end{bmatrix} \mathbf{x}. \quad (6)$$

Note that the index sets can be arbitrary partitions of the columns and rows and do not necessarily have to be consecutive. To facilitate the latter discussion, we also define residual  $\mathbf{r} = \mathbf{y} - A\mathbf{x}$  and let block residual  $\mathbf{r}_I$  be the subset of  $\mathbf{r}$  defined as  $\mathbf{y}_I - A_I \mathbf{x}$ , where  $I \in \{I_i\}_{i=1}^M$ .

With this notation, the general update scheme is

$$\mathbf{x}_J = \mathbf{x}_J - \mu \cdot \nabla_J f_I(\mathbf{x}), \quad (7)$$

where  $J \in \{J_j\}_{j=1}^N$ ,  $\mu$  is the step length and  $f_I(\mathbf{x}) = \|\mathbf{y}_I - A_I \mathbf{x}\|_2^2$ , where the  $\|\cdot\|$  is the  $l_2$  norm of a vector. Thus the gradient  $\nabla_J f_I(\mathbf{x}) = -2(A_I^J)^T (\mathbf{y}_I - A_I \mathbf{x})$ .

There are two difficulties when applying SBCD in large scale CT reconstructions. The first one is that the step length  $\mu$  is determined by calculating the Lipschitz constant of the gradient  $\nabla_J f_I(\mathbf{x})$ , which is computationally challenging. The second issue is that the computation of the gradient in SBCD requires the calculation of the block residue  $\mathbf{r}_I = \mathbf{y}_I - A_I \mathbf{x}$ , which needs the whole of  $\mathbf{x}$ . In other words, the system matrix  $A$  is not actually separable in column direction due to the need to calculate  $\mathbf{r}_I$ . This drawback makes the SBCD algorithm difficult to apply in a totally distributed network where each computation node only has partial access to both  $\mathbf{x}$  and  $\mathbf{y}$ .

Another method that can be used for linear systems is the multisplitting (MS) method [57]. If MS method only partitions the  $A$  into column blocks [58], it always converges as long as  $A$  is of full column rank. However, If MS method partition the  $A$  into both row and column blocks, it becomes only applicable when  $A$  meets certain conditions. For example, when  $A$  is an ‘‘H-matrix’’[59] or is a positive definite matrix [60]. The system matrix  $A$  in CT reconstruction does not meet these requirements and MS dividing  $A$  in both dimensions did not work in initial experiments we conducted for CT reconstruction.

Our CT reconstruction setting leads to a convex optimization problem. The alternating direction method of multipliers (ADMM) is a popular tool in the parallelization of convex optimization problems as it allows the decomposition into several smaller sub-problems [61]. Several versions of ADMM were designed to work on small subsets of either  $\mathbf{x}$  or  $\mathbf{y}$  [62]. In [63], a block ADMM method was proposed that works by breaking the problems into small subsets of  $\mathbf{x}$  and  $\mathbf{y}$ . This version of block ADMM introduces three sets of additional nuisance variables and each sub problem requires the solution to a smaller linear inverse problem, that can often be solved approximately using a small number of CG iterations. The new nuisance variables increase the overall storage requirements and the requirement for each node to solve a linear inverse problem in each iteration leads to relatively slow convergence and a high computation burden in large scale CT reconstruction, which will be illustrated in our paper.

## II. INTRODUCTION OF CSGD

In this paper, we consider large scale CT reconstruction problems in a distributed networks, where each computation node only has limited access to both  $\mathbf{y}$  and  $\mathbf{x}$ . Inspired by SBCD and block ADMM, our goal here is to develop a parallelisable algorithm that works with generic x-ray tomographic scanning trajectories and that is faster and more memory efficient than block ADMM. Our novel algorithm, called coordinate-wise stochastic gradient descent (CSGD), also introduces nuisance variables, but we use fewer variables than block ADMM. Furthermore, CSGD simplifies the step length calculation as well as the residue update scheme of SBCD type algorithms and thus enables a full column decomposition for  $A$ . The algorithm is scalable so that it can be run on a range of computing platforms, including low memory GPU clusters and high performance CPU based clusters.

### A. Derivation of the Algorithm

The proposed CSGD is similar to SBCD. The main goals of CSGD are to find a simpler strategy to compute the step length  $\mu$  and to efficiently approximate the residue  $\mathbf{r}_I$  without having to compute the product  $A_I \mathbf{x}$  in each iteration.

Similar to SBCD, after selecting a row block  $I \in \{I_i\}_{i=1}^M$ , CSGD operates on the object function

$$f_I(\mathbf{x}) = \|\mathbf{y}_I - A_I \mathbf{x}\|^2, \quad (8)$$

with gradient

$$\mathbf{g} = \nabla f_I(\mathbf{x}) = -2(A_I)^T(\mathbf{y}_I - A_I \mathbf{x}) = -2(A_I)^T \mathbf{r}_I. \quad (9)$$

A coordinate descend update scheme is adopted and only those elements whose indices are in the set  $J \in \{J_j\}_{j=1}^N$  are updated, so that the descent operator is

$$\tilde{\mathbf{g}} = \begin{bmatrix} \mathbf{g}_J \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} (A_I^J)^T \mathbf{r}_I \\ \mathbf{0} \end{bmatrix}. \quad (10)$$

Along with this new modified direction, the update on voxels becomes

$$\mathbf{x}^k = \mathbf{x}^{k-1} + \mu \tilde{\mathbf{g}} \rightarrow \begin{cases} \mathbf{x}_J^k = \mathbf{x}_J^{k-1} + \mu \mathbf{g}_J \\ \mathbf{x}_{\hat{J}}^k = \mathbf{x}_{\hat{J}}^{k-1}, \end{cases} \quad (11)$$

where  $\mu$  is the gradient step length,  $\hat{J}$  is the complement to the set  $J$  and  $\mathbf{g}_J$  is a sub set of  $-\frac{1}{2}\mathbf{g}$ . The steepest descent idea is to make the direction of  $\nabla f_I(\mathbf{x}^k)$  perpendicular to the direction of  $\tilde{\mathbf{g}}$ , i.e.

$$((A_I)^T(\mathbf{y}_I - A_I \mathbf{x}^k))^T \tilde{\mathbf{g}} = 0. \quad (12)$$

Use the fact that  $x^k = x^{k-1} + \mu \tilde{\mathbf{g}}$  and  $A_I \tilde{\mathbf{g}} = A_I^J \mathbf{g}_J$ , the  $\mu$  leading to the maximum descent is

$$\mu = \frac{\mathbf{g}_J^T (A_I^J)^T \mathbf{r}_I}{\mathbf{g}_J^T (A_I^J)^T A_I^J \mathbf{g}_J} = \frac{\mathbf{g}_J^T \mathbf{g}_J}{\mathbf{g}_J^T (A_I^J)^T A_I^J \mathbf{g}_J}. \quad (13)$$

This calculation does not require the corresponding computation node to have access to the whole row or column block of matrix  $A$  and does not have to calculate the Lipschitz constant to determine the step length. When the matrix  $A$  cannot be stored and need to be generated on the fly, CSGD iteration only requires to use a sub matrix  $A_I^J$  and its transpose. Furthermore, GPUs can calculate the forward projection  $A_I^J \mathbf{g}_J$  and backward projection  $(A_I^J)^T \mathbf{r}_I$  very efficiently to achieve GPU-accelerated projection operations [64], [65].

The step size derived from Eq. 13 is chosen to reduce  $\|\mathbf{r}_I\|$  instead of  $\|\mathbf{r}\|$ . We observed experimentally that our choice of  $\mu$  can lead to instability in the algorithm. One method to avoid this is to introduce an additional relaxation parameter  $\beta < 1$ , which led CSGD to converge to a precision level similar to that of SIRT and CAV. The alternative of using a constant step length  $\mu$  also led to convergence, but required a careful choice of step length (thus making parameter tuning difficult) and led to relatively slow convergence.

$\mathbf{r}_I$  plays an important role in CSGD since it determines the update direction of  $\mathbf{x}$ . If we had access to all  $\mathbf{x}_j$ , then we could simply compute  $\mathbf{r}_I = \mathbf{y}_I - \sum_j \mathbf{x}_j$ . In our setting, the node only has access to one  $\mathbf{x}_j$ , thus the above computation is not possible. Instead, each node computes a quantity  $\mathbf{z}_I^j = \mathbf{A}_I^J \mathbf{x}_{j_j}$ . However, as discussed later, we might not compute all  $\mathbf{z}_I^j$  in each iteration. In this case, we use an older versions of  $\mathbf{z}_I^j$  to approximate  $\mathbf{r}_I = \mathbf{y}_I - \sum_{j=1}^N \mathbf{z}_I^j$ . For smoothly varying cost functions, using old residual information when calculating gradients is similar to the use of old gradients in parallel methods that can often be shown to converge [66], [67].

The pseudo-code of the basic computation blocks is shown in the Fig. 1.

To parallelize the algorithm, we do not have to update all sub-blocks before updating  $\mathbf{r}$ . We define percentages  $\alpha$  and  $\gamma$  to specify the fraction of row and column blocks to be updated. The algorithm is shown in Algo.2. In the following discussion, when operations from line 4 to line 16 are performed for one time, the algorithm is said to perform for ‘‘one epoch’’.

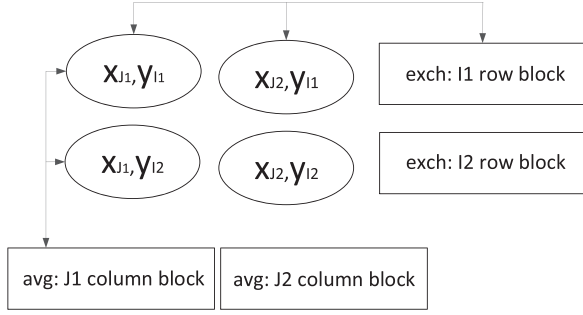


Fig. 1. One parallel implementation of block ADMM. Computation nodes (ovals) compute  $\Pi$  projection and master nodes (square boxes) compute “avg” and “exch”.

---

### Algorithm 1: The Algorithm for the Computation Performed at Each Node.

---

Initialization: select system matrix's row index  $I \in \{I_i\}_{i=1}^M$  and column index  $J \in \{J_j\}_{j=1}^N$   
 $\mathbf{g}_J = (A_I^J)^T \mathbf{r}_I$   
 $\mu = \beta \frac{(\mathbf{g}_J)^T \mathbf{g}_J}{(\mathbf{g}_J)^T (A_I^J)^T A_I^J \mathbf{g}_J}$   
 $\mathbf{x}_J^k = \mathbf{x}_J^{k-1} + \mu \mathbf{g}_J$   
 $\mathbf{z}_I^k = A_I^J \mathbf{x}_J^k$

---



---

### Algorithm 2: CSGD Algorithm which Parallelize Both Row and Column Blocks.

---

1: Input:  $\mathbf{y}$ ,  $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\{I_i\}_{i=1}^M$  and  $\{J_j\}_{j=1}^N$ .  
2: Initialisation:  $\mathbf{x} = \mathbf{0}$ , (i.e.,  $\{\mathbf{x}_{J_j}\}_{j=1}^N = \mathbf{0}$ ),  $\{\hat{\mathbf{x}}^i\}_{i=1}^M = \mathbf{x}$ ,  $\{\mathbf{z}^j\}_{j=1}^N = \mathbf{0}$  and  $\mathbf{r} = \mathbf{y}$ .  
3: **while** stopping criterion is not met **do**  
4:   **for**  $\tilde{j} = 1, 2, \dots, \gamma N$  **in parallel** (J loop) **do**  
5:     randomly draw  $J$  from  $\{J_j\}_{j=1}^N$  with replacement  
6:      $\hat{\mathbf{x}}^i (1 \leq i \leq M) = \mathbf{0}$   
7:     **for**  $\tilde{i} = 1, 2, \dots, \alpha M$  **in parallel** (I loop) **do**  
8:       randomly draw  $I$  from  $\{I_i\}_{i=1}^M$  with replacement  
9:        $\mathbf{g}_J = (A_I^J)^T \mathbf{r}_I$   
10:        $\mu_j^i = \beta \frac{(\mathbf{g}_J)^T \mathbf{g}_J}{(\mathbf{g}_J)^T (A_I^J)^T A_I^J \mathbf{g}_J}$   
11:        $\hat{\mathbf{x}}_J^i = \mathbf{x}_J + \mu_j^i \mathbf{g}_J$   
12:        $\mathbf{z}_I^j = A_I^J \hat{\mathbf{x}}_J^i$   
13:     **end for**  
14:   **end for**  
15:    $\mathbf{r} = \mathbf{y} - \sum_j \mathbf{z}^j$   
16:   for all updated  $J$ ,  $\mathbf{x}_J = \sum_i (\hat{\mathbf{x}}_J^i) / (\alpha M)$   
17: **end while**  
18:  $\mathbf{x}_{solution} = [\mathbf{x}_{J_1}, \dots, \mathbf{x}_{J_N}]^T$

---

### B. Comparison CSGD With Block ADMM

Block ADMM has the same parallel computing architecture as CSGD. To facilitate the comparison, we briefly present the basic operation of block ADMM algorithm [63], as shown in Algo.3. Here avg is the element-wise averaging operator. The

---

### Algorithm 3: Block ADMM Iteration.

---

1: Input:  $\mathbf{y}$ ,  $\rho$ ,  $\{I_i\}_{i=1}^M$  and  $\{J_j\}_{j=1}^N$ .  
2: Initialisation:  $\mathbf{x} = \tilde{\mathbf{x}} = \mathbf{0}$ ,  $\mathbf{z} = \tilde{\mathbf{z}} = \mathbf{0}$ ,  $\{\mathbf{x}^i\}_{i=1}^M = \mathbf{0}$ ,  $\{\mathbf{z}^j\}_{j=1}^N = \mathbf{0}$   
3: **for**  $k = 1, 2, \dots, k_{max}$  **do**  
4:    $\mathbf{x}^{k+\frac{1}{2}} = \mathbf{x}^k - \tilde{\mathbf{x}}^k$   
5:    $\mathbf{z}^{k+\frac{1}{2}} = \frac{1}{1+\frac{\rho}{2}} \mathbf{y} + \frac{1}{\frac{\rho}{2}+1} (\mathbf{z}^k - \tilde{\mathbf{z}}^k)$   
6:   **for**  $\tilde{j} = 1, 2, \dots, N$  **in parallel** (J loop) **do**  
7:     randomly draw  $J$  from  $\{J_j\}_{j=1}^N$  with replacement  
8:     **for**  $\tilde{i} = 1, 2, \dots, M$  **in parallel** (I loop) **do**  
9:       randomly draw  $I$  from  $\{I_i\}_{i=1}^M$  with replacement  
10:        $(\mathbf{x}_J^{i, k+\frac{1}{2}}, \mathbf{z}_I^{j, k+\frac{1}{2}}) = \Pi_{A_I^J} (\mathbf{x}_J^k - (\tilde{\mathbf{x}}_J^i)^k, \mathbf{z}_I^k + \tilde{\mathbf{z}}_I^k)$   
11:     **end for**  
12:   **end for**  
13:    $\mathbf{x}^{k+1} = \text{avg}(\mathbf{x}^{k+\frac{1}{2}}, \{\mathbf{x}^{i, k+\frac{1}{2}}\}_{i=1}^M)$   
14:    $(\mathbf{z}^{k+1}, \{\mathbf{z}^j\}_{j=1}^N) = \text{exch}(\mathbf{z}^{k+\frac{1}{2}}, \{\mathbf{z}^j\}_{j=1}^N)$   
15:    $\tilde{\mathbf{x}}^{k+1} = \tilde{\mathbf{x}}^k + \mathbf{x}^{k+\frac{1}{2}} - \mathbf{x}^{k+1}$   
16:    $\tilde{\mathbf{x}}^i)^{k+1} = \tilde{\mathbf{x}}^i)^k + \mathbf{x}^{j, k+\frac{1}{2}} - \mathbf{x}^{k+1}$   
17:    $\tilde{\mathbf{z}}^{k+1} = \tilde{\mathbf{z}}^k + \mathbf{z}^{k+\frac{1}{2}} - \mathbf{z}^{k+1}$   
18: **end for**  
19:  $\mathbf{x}_{solution} = \mathbf{x}^{k_{max} + \frac{1}{2}}$

---

$\Pi$  projection is a linear operator

$$\Pi_{A_I^J}(c, d) = \begin{bmatrix} I_d & (A_I^J)^T \\ A_I^J & -I_d \end{bmatrix}^{-1} \begin{bmatrix} I_d & (A_I^J)^T \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} c \\ d \end{bmatrix} \quad (14)$$

where  $I_d$  is an identity matrix whose size is determined by the size of  $A_I^J$ . To solve this equation, the CG method can be used. Standard techniques to speed up block ADMM include the early termination of the CG iteration and a variable  $\rho$ -update scheme [61]. The exchange operator  $\text{exch}(c, \{c_j\}_{j=1}^N)$  is given by

$$\mathbf{z}_{I_i}^j = c_j + \frac{(c - \sum_{j=1}^N c_j)}{N+1}$$

$$\mathbf{z}_{I_i} = c - \frac{(c - \sum_{j=1}^N c_j)}{N+1} \quad (15)$$

When applying block ADMM and CSGD in a distributed network, they share the same architecture. For  $M = N = 2$ , parallel block ADMM is shown diagrammatically in Fig. 1. Here, each computation node (ovals) stores one image block  $\mathbf{x}_J$  ( $J \in \{J_j\}_{j=1}^N$ ) and one data block  $\mathbf{y}_I$  ( $I \in \{I_i\}_{i=1}^M$ ). There are two types of master node (square boxes). One type performs the “avg” procedure for column blocks and the other type performs the “exch” procedure for row blocks. One set of master nodes stores  $\mathbf{x}_J^{k+\frac{1}{2}}$ ,  $\mathbf{x}_J^k$  and  $\tilde{\mathbf{x}}_J^k$  and the other set of master nodes stores  $\mathbf{y}_I^{k+\frac{1}{2}}$ ,  $\mathbf{y}_I^k$  and  $\tilde{\mathbf{y}}_I^k$ . The same distributed architecture of Fig. 1 is also suitable for CSGD with  $\alpha = \gamma = 1$ . For each  $J$ , CSGD requires the exchange and summation of  $\mathbf{z}_I^j$  (i.e. over columns in Fig. 1) and for each  $I$ , CSGD requires the exchange and summation of  $\hat{\mathbf{x}}_J^i$  (i.e. over rows in Fig. 1). These operations

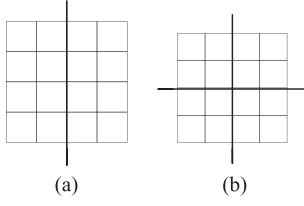


Fig. 2. (a) Shows the partition that along with one dimension for a 16-pixel image. (b) Shows a partition that along both image dimensions.

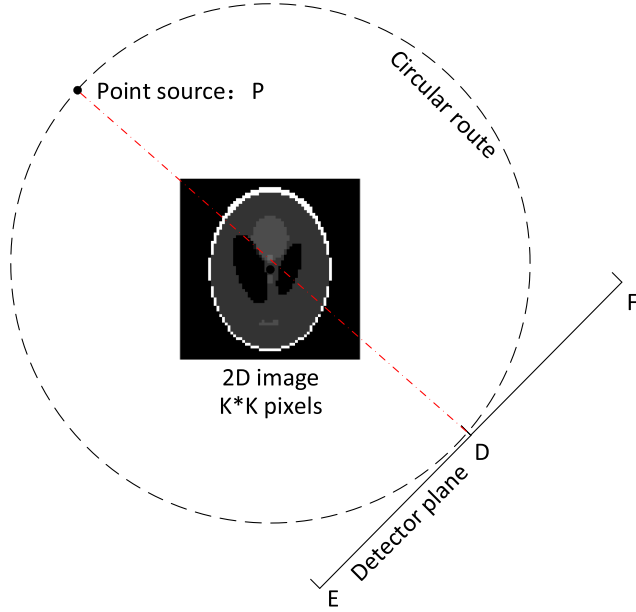


Fig. 3. A standard 2D scanning geometry with a Shepp-Logan phantom, where the P is the x-ray source, O is the centre of the object and the rotation centre. D is the centre of the detector. Source and detector rotate around the centre and take measurements at different angles. The linear detector is evenly divided into sub-areas DE and DF, which will be used in importance sampling discussed later.

are similar to the “avg” and “exch” procedure of ADMM respectively. Note that these operations can be implemented efficiently using a message passing approach.

The storage demand is different between block ADMM and CSGD. Overall, CSGD requires storage of vectors  $\mathbf{y} \in \mathbb{R}^{r \times 1}$ ,  $\mathbf{r} \in \mathbb{R}^{r \times 1}$  and  $\mathbf{x} \in \mathbb{R}^{c \times 1}$  and of the set of N vectors  $\{\mathbf{z}^j\}_{j=1}^N \in \mathbb{R}^{r \times 1}$ , so the total memory requirement for CSGD is  $(N + 2)r + c$ . Block ADMM requires more storage, as it requires storage of  $\mathbf{x}$ ,  $\tilde{\mathbf{x}}$ ,  $\mathbf{z}$ ,  $\tilde{\mathbf{z}}$ ,  $\mathbf{y}$  and the M vectors  $\{\mathbf{x}^i\}_{i=1}^M$  and  $\{\tilde{\mathbf{x}}^i\}_{i=1}^M$  and the N vectors  $\{\mathbf{z}^j\}_{j=1}^N$ , leading to a total storage demand of  $(N + 3)r + (2M + 2)c$ .

### C. Partition Methods in CT Case

In our CT reconstruction problem, the partition of  $A$  along columns is equivalent to the partition of the image. Two examples of column partition are 1) to cut the image along one dimension and 2) to cut the image along all dimensions. Two 2D examples are shown in Fig. 2.

To understand how partitioning row blocks of  $A$  relates to the tomographic imaging setting, we take the 2D scanning model shown in Fig. 3 as an example. When detector and source are at a given location, we obtain projections along a range of

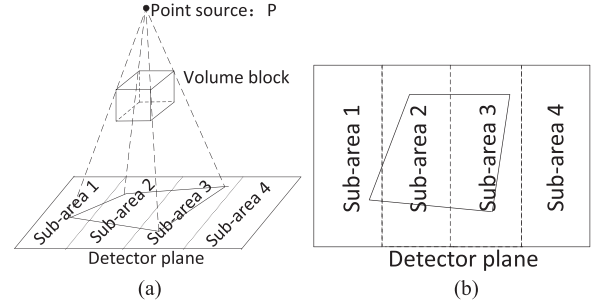


Fig. 4. (a) Shows a 3D model of the cone beam setup with one block of the volume being projected on a detector plane. (b) Shows the projection area of the volume block.

paths from the source to the different detector elements. For one source/detector location, these measurements will be called one “projection”. We could form the row blocks of  $A$  with random projections, but then keep them fixed during different epochs. A non-random version of this will be called “deterministic partitioning” in which one row block contains sequential projections from successive projection angles. We can further divide the detector into several sub-areas and treat projections from each sub-area as a “sub-projection”. An illustration of this in 3D is shown in Fig. 4, where one projection is divided into 4 sub-projections. When forming the row blocks, different sub-projections from different projection angles can be grouped together. This will be seen to be advantageous in the importance sampling strategy discussed next.

### D. Importance Sampling

When slicing the detector into several sub-areas, we can also form row blocks dynamically. It means that the sub-areas forming row blocks are changing for different epochs. Based on the proposed algorithm, it is straightforward to develop a random sampling strategy that goes through all projection views and all detector sub-areas arbitrarily. Looking at Algo.2, for one image block  $\mathbf{x}_{J_j}$  ( $1 \leq j \leq N$ ), we could randomly select sub-areas to form one row block  $I_i$  ( $1 \leq i \leq M$ ), with each sub-areas being chosen with equal probability as long as the sub-area receives x-rays passing through  $\mathbf{x}_{J_j}$ . This sample method will be called “uniform sampling”. Considering the sparsity of  $A$ , we further develop an importance sampling strategy. Fig. 4 illustrates that the projection of a sub-image block only intersects a small part of the detector. When sampling sub-projections for each sub-image block  $\mathbf{x}_{J_j}$  ( $1 \leq j \leq N$ ), a vector  $P_j$  representing the probability of choosing each detector sub-area is calculated. The calculation is based on projection areas of  $\mathbf{x}_{J_j}$  on each sub-detector area at different projection angles. Then  $\alpha * 100\%$  of the sub-areas are sampled with probability  $P_j$  and are grouped into  $\alpha M$  row blocks. These row blocks are assigned to different nodes to perform line 9–12 in Algo.2. In this case, the detector sub-areas receiving more projections from the current sub-image block  $\mathbf{x}_{J_j}$  are more likely to be chosen.

### E. Computational Complexity

There are several important aspects when comparing computational efficiency of the methods. The methods are designed to

allow parallel computation. We envisage this to be performed in a distributed network of computation nodes<sup>1</sup>. Computation nodes produce two outputs, as displayed in line 11 and line 12 in Algo.2. These are then either sent to larger, but slow storage or directly to other nodes, where they are eventually used to compute line 15 and line 16 in Algo.2, which can be performed efficiently using message passing interface reduction methods.

Three aspects affect performance:

- 1) Computational complexity in terms of multiply-adds.
  - 2) Data transfer requirements between data storage and a processing unit as well as between different processing units.
  - 3) Data storage requirements, both in terms of fast access RAM and in slower access (e.g. disk based) data storage.
- Each of these costs are dominated by different properties:

- 1) Computational complexity is dominated by the computation of matrix vector products involving  $A_I^j$  and its transpose, especially as  $A$  is not generally stored but might have to be re-computed every time it is needed. The computational complexity is thus  $O(|I| * |J|)$ , though computations performed on highly parallel architectures, such as modern GPUs, are able to perform millions of these computations in parallel.
- 2) Data transfer requirements are dominated by the need for each of the parallel computation nodes to receive  $\mathbf{r}_I$  and  $\mathbf{x}_J$  and transmit  $\hat{\mathbf{x}}_J^i$  and  $\mathbf{z}_J^j$ . Note that the size of the required input and output vectors are the same, the data transfer requirement is thus  $O(|I| + |J|)$ .
- 3) Central data storage requirements are dominated by the need to store the original data and the current estimate of  $\mathbf{x}$ . We also need to compute and store averages over  $\hat{\mathbf{x}}_J^i$  and  $\mathbf{z}_J^j$ . These computations can be performed efficiently using parallel data reduction techniques. Our approach would mean that each node would thus require  $O(|I| + |J|)$  local memory.

### III. SIMULATIONS

The simulations are divided into two parts, the first part explores the performance of CSGD in CT reconstruction. The second part compares CSGD with block ADMM.

#### A. CSGD in CT Reconstruction

We used two criteria to evaluate the performance. 1) signal to noise ratio (SNR):  $20\log_{10}(\|\mathbf{x}_{true}\|/\|\mathbf{x}_{true} - \mathbf{x}_{est}\|)$ , 2) relative distance (RD):  $20\log_{10}(\|\mathbf{x}_{lsq}\|/\|\mathbf{x}_{lsq} - \mathbf{x}_{est}\|)$ . The  $\mathbf{x}_{true}$ ,  $\mathbf{x}_{est}$  and  $\mathbf{x}_{lsq}$  are the true phantom image vector, the reconstructed image vector and the least square solution respectively.

In the following simulations, we used the Shepp-Logan phantom and, unless stated otherwise,  $K$  in Fig. 3 is 16 and the side length of each pixel is 1mm, the length of OP and OD are always 100 mm. The rotation interval for source and detector is  $10^\circ$ .

<sup>1</sup>A serial version running on a single computation node where each computation is done independently, but one after the other, is also possible and this is how many of the simulations reported here were computed.

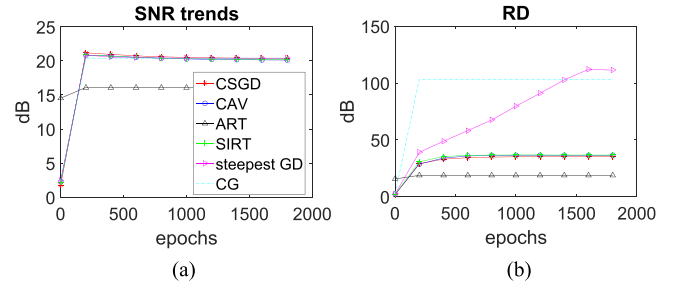


Fig. 5. Comparison of CSGD with other methods. The performance of CSGD, including the final SNR and RD level, is similar to the SIRT and CAV methods.

The detector contains 30 pixels of size 1mm. The  $\mathbf{e}$  in Eq. 1 is Gaussian white noise with variance  $\sigma$  of 0.1 to the simulated observations. The SNR of projection data  $\mathbf{y}$  ( $20\log_{10}(\|\mathbf{y}\|/\|\mathbf{e}\|)$ ) is 25.8 dB. The default partition method of projection data  $\mathbf{y}$  uses the static type “deterministic partitioning”. The default partition in the image domain uses the method shown in Fig. 2a.

#### B. Comparison to Other Methods

We start by comparing our method to a range of other algorithms popular in CT reconstruction. For CSGD we used  $M = 8$ ,  $N = 4$ ,  $\beta = 0.23$  and  $\alpha = \gamma = 1$ . We compared our method to CG, steepest gradient descent (GD), SIRT, ART and CAV. The results are shown in Fig. 5. It can be seen that for noisy observations, the CSGD and the other method (except for ART method) all obtains nearly the same SNR. However, whilst CG and steepest GD converge to the least square solution, CSGD, SIRT, ART and CAV do not. This indicates that similar to SIRT, ART and CAV, CSGD also iterates the  $\mathbf{x}$  towards to a weighted least square solution rather than the least square solution itself (see our fixed point analysis in the Appendix).

#### C. Influence of Algorithm Parameters

Our algorithm has three parameters that have to be set,  $\alpha$ ,  $\beta$  and  $\gamma$ . Due to space constraints we only discuss  $\alpha$  and  $\beta$  and set the  $\gamma \equiv 1$ . Smaller  $\gamma$  lead to similar results to those observed with smaller  $\alpha$ . We start by an empirical evaluation of  $\beta$ . To determine the range of suitable parameters  $\beta$ , we explored the performance by varying  $M$ ,  $N$  and  $\beta$ . To show the largest value of  $\beta$  that can be used for different values of  $N$  and  $M$ , we run the algorithm for 800 epochs for different values of  $\beta$ . The results are shown in Fig. 6 where we see that the largest  $\beta$  value leading to highly accurate solutions is approximately  $\frac{1}{N}$ . Increasing  $N$  reduces the  $\beta$  range, and there is also a relationship on  $M$  which is less clear. We next turn to the influence of the  $\alpha$ . Reducing  $\alpha$  reduces the computation of CSGD within one epoch as only some of the updates are computed. We set  $M = 12$  and  $N = 2$  and generated the data as before. Simulation results show that reducing  $\alpha$  increase the acceptable  $\beta$  range, as shown in Fig. 7.

Convergence speed is influenced by  $\alpha$ ,  $\beta$ ,  $M$  and  $N$ . In general, convergence is slower for larger values of  $N$  and  $M$  and for smaller  $\beta$  and  $\alpha$ . The comparisons of  $\beta$  and  $M$  are shown in Fig. 8. It can be seen that in general, increasing  $\beta$  was found to lead to faster convergence, up to a point where the algorithms

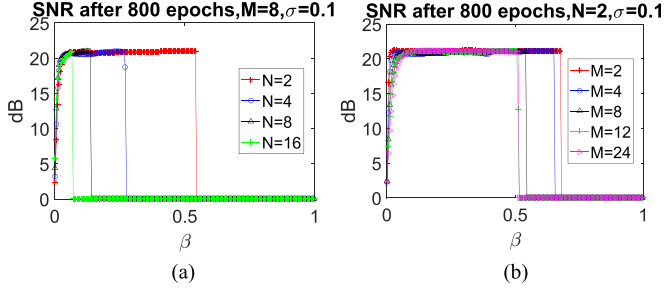


Fig. 6. SNR value after 800 epochs for different values of  $\beta$  and for different  $N$  and  $M$ . Increasing the  $M$  or  $N$  reduces the acceptable  $\beta$  range and the range is influenced more by  $N$  value instead of  $M$ . In this simulation, the  $\alpha \equiv 1$ .

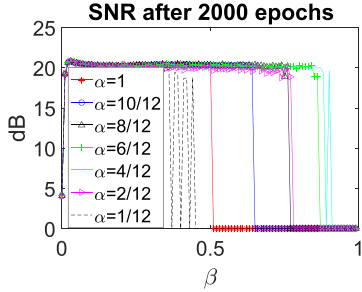


Fig. 7. When  $M = 12$ ,  $N = 2$ , the  $\beta$  value range for different  $\alpha$ . Unless that  $\alpha$  is too small ( $\frac{1}{12}$ ), decreasing  $\alpha$  generally enlarges the range of acceptable  $\beta$ .

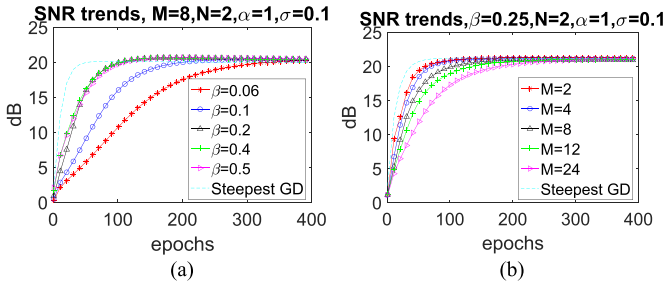


Fig. 8. (a) Shows that within the acceptable  $\beta$  range, increasing the  $\beta$  is found to lead to faster convergence, up to a point where the algorithms is about to diverge. (b) Shows that when  $N$  is fixed, increasing  $M$  slows down the total convergence rate.

started to diverge. When the  $\beta$  is constant, increasing the  $M$  (or  $N$ , which is not provided here) slows down the convergence speed. We here also provide the steepest GD results as a reference to show that CSGD can obtain the same precision level as steepest GD does.

Before discussing the influence on convergence speed caused by  $\alpha$ , we introduce a new concept “effective epoch”. The number of sub-matrices  $A_{I_i}^{J_j}$  (in Eq. 6) that CSGD used per epoch is proportional to the parameter  $\alpha$ . As the computational speed is dominated by matrix vector products, when we compare the difference in convergence speed for methods using different  $\alpha$ s (for example,  $\alpha = 1$  and  $\alpha = 0.5$ ), we took account of the reduction in computational effort when using smaller  $\alpha$ . The epoch count is multiplied by  $\alpha$  and is called as the *effective epoch*. As a result, for different  $\alpha$ , the computation amount after one “effective epoch” is the same with each other. We need to point out that the “effective epoch” does not consider the data transfer influence. When  $\alpha$  is small, one “effective epoch” includes more

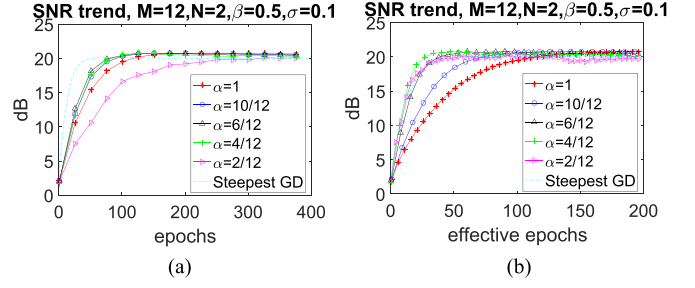


Fig. 9. (a) Shows that when  $\beta$  is fixed, reducing the  $\alpha$  to around 0.5 increases the convergence speed. However, when the  $\alpha$  is reduced too small as  $\frac{2}{12}$ , the convergence speed is slowed down. (b) Shows that from “effective epoch” point of view, reducing the  $\alpha$  is always helpful to increase the convergence speed and the increased convergence speed can be similar to the steepest GD when  $\alpha \leq 0.5$ .

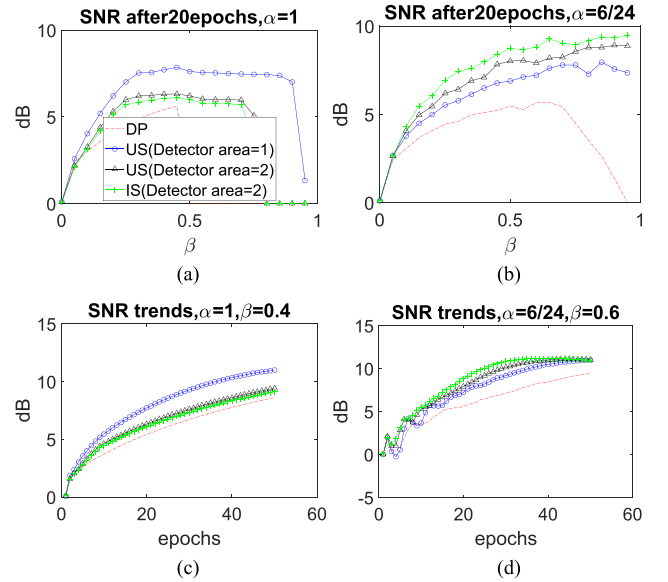


Fig. 10. Comparisons of different sample methods. The  $M = 24$ ,  $N = 2$ ,  $\sigma = 1$ . “DP” is deterministic partitioning (defined in II-C), “US” is uniform sampling method (defined in II-D) and “IS” is importance sampling method (defined in II-D). The uniform sample or importance sampling both show faster convergence speed and wider acceptable  $\beta$  range than previous deterministic partitioning. Importance sampling methods outperform the others when  $\alpha$  is  $\frac{6}{24}$ .

epochs, which means more data transfer between nodes. Here we mainly focus on cases when the data transfer is much more efficient than calculating matrix-vector multiplications and thus is negligible. The comparison of convergence speed with different  $\alpha$  is shown in Fig. 9. Interestingly, by looking at effective epochs, we see that CSGD can converge faster than gradient descent for smaller  $\alpha$ .

#### D. Partitioning the Image and the Detector

In this section, larger simulations are conducted. The CT scanning geometry still uses the form shown in Fig. 3, but with  $K = 64$ . The length of OP and OD is set to 115 mm. The rotation interval for the point source and the detector is  $3^\circ$ . The detector contains 130 pixels. The  $e$  in Eq. 1 is Gaussian white noise with variance  $\sigma$  of 1 to the simulated observations. The SNR of projection data  $\mathbf{y}$  ( $20 \log_{10}(\|\mathbf{y}\|/\|\mathbf{e}\|)$ ) is 33.8 dB. In

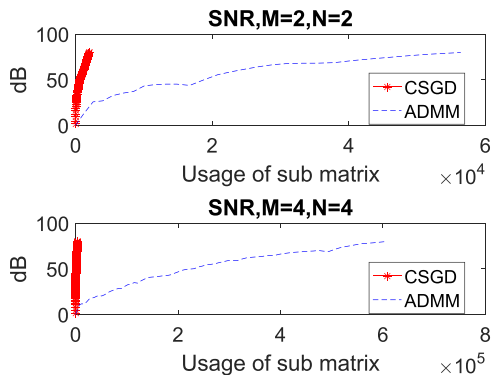


Fig. 11. Under noise free assumption, CSGD uses much fewer matrix-vector multiplications to achieve predefined SNR. The  $\beta$  used in CSGD is  $\frac{1}{2N}$ .

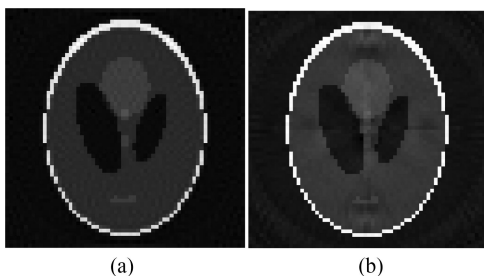


Fig. 12. (a) Comparison of reconstruction results using CSGD and (b) block ADMM. Although the SNR of both images is the same, (b) shows much clearer inner artefacts than (a).

the following section, the partition of the image follows Fig. 2b. We here compare the partition of the linear detector into one or two areas, as shown in Fig. 3. We also compare uniform sampling and importance sampling. The algorithm is performed for 20 epochs and the SNR value for different  $\beta$  values at this point is shown in Fig. 10 (a)-(b) and two typical SNR trend in Fig. 10 (c)-(d). Under different  $\alpha$ , the uniform sampling strategy always obtains higher SNR after the same epochs than deterministic partitioning. Besides, the uniform sampling method has an even wider  $\beta$  range than deterministic partitioning, which is easier for parameter tuning. The uniform sampling method can be further improved when the detector is sliced into two sub-areas. This method avoids choosing the sub-detector areas which do not receive any projections from the currently selected sub-image block. However, when  $\alpha = 1$ , the two sub-detector situation does not increase the convergence rate. This is because that those sub-detectors which do not receive projections from the current sub-block still have to be selected since the  $\alpha$  is too large. When the  $\alpha$  decrease to a value (e.g.  $\frac{6}{24}$ ) that guarantees all selected sub-detectors are those who receive projections, slicing the detector area into two sub-areas obtains higher SNR than not slicing situation. Furthermore, when  $\alpha$  is small, the importance sampling strategy can obtain higher SNR than uniform sampling method at the initial iterations.

#### E. Comparison of CSGD and Block ADMM

As discussed previously, block ADMM allows for the same partitioning of  $A$ . To compare CSGD and block ADMM, a random system matrix  $A_{256 \times 128}$  and a random vector  $\mathbf{x}_{128 \times 1}$  are

used to reduce the computation requirements in our simulations. The linear system is a noise free model here.

As each CSGD iteration and each block ADMM iteration require different amounts of computational effort, we here do not compare SNR after each epoch. Instead, we plot SNR against the number of times the algorithm has computed a matrix-vector product involving the sub-matrices of  $A$ , as this is the dominating computational cost here. Both CSGD and block ADMM use all subsets of the matrix  $A$  and are stopped when their solutions reach 80 dB SNR. For block ADMM, we determined the best parameter values for variable  $\rho$  and also stopped the CG method after as few iterations as possible to allow algorithm convergence to the required level.

The convergence comparison for CSGD and block ADMM are shown in Fig. 11. The results demonstrate the significant speed advantage offered by CSGD, which requires significantly fewer matrix-vector multiplications compared to block ADMM. Furthermore, the fully distributed form of CSGD used here, i.e.  $\alpha = 1$ , is not the optimal use of CSGD. According to the previous simulations, setting  $\alpha < 1$  can further increase convergence speed.

We have also compared CSGD and block ADMM on the CT simulation of the previous simulation data with  $K = 64$  and a detector with 130 pixels. The image was partitioned again as in Fig. 2b. Both methods were stopped once the SNR of the reconstructed image had reached 20 dB. The reconstructed images under noise free situation are shown in Fig. 12. Although the SNR of both images is the same, the ADMM reconstruction in Fig. 12 b shows much clearer artefacts along with the boundaries of adjacent sub-image blocks whilst the CSGD results in Fig. 12 a do not show these effects.

#### IV. CONCLUSION

CSGD was designed for a distributed reconstruction of cone beam CT data under arbitrary scan trajectories. In the distributed network, each node is assumed to have limited storage capacity and thus all nodes operate with limited access to the projection data and reconstructed volume. Whilst the method does not converge to the least squares solution, the solution is found empirically to be comparable in quality to those found with other common tomographic reconstruction algorithms, such as SIRT and CAV, but at a significant computational advantage. The parallel architecture is the same as that of block ADMM, which is a general algorithm for separable convex optimization. However, for large scale CT reconstruction, block ADMM is less attractive compared to CSGD. One advantage of CSGD is that it requires less storage compared to block ADMM. Another significant advantage of CSGD is that it converges with significantly fewer matrix vector products as it avoids the calculation of matrix inverses. This means that CSGD converges much faster compared to block ADMM.

We have furthermore developed an importance sampling strategy, that has been shown to further increase initial convergence. A theoretical analysis of the algorithm's convergence properties is ongoing and so is the inclusion of regularisation terms in the method.



## APPENDIX

Whilst we do not have a formal convergence proof of CSGD yet, it is instructive to analyse the fixed points of the algorithm. We here look at the deterministic version of the algorithm with  $\alpha = 1$ .

The algorithm updates two quantities,  $\mathbf{x}$  and  $\mathbf{z}$ . Let  $(\mathbf{x}^{k+1}, \mathbf{z}^{k+1}) = T(\mathbf{x}^k, \mathbf{z}^k)$  define one iteration of the algorithm. Let  $\mathbf{x}^*$  and  $\mathbf{z}^*$  be fixed points of the operator  $T(\mathbf{x}, \mathbf{z})$  defined by  $(\mathbf{x}^*, \mathbf{z}^*) = T(\mathbf{x}^*, \mathbf{z}^*)$ . Similar results to the once derived here for the deterministic algorithm can also be obtained for the randomised versions and for  $\alpha < 1$  if we look at points for which  $(\mathbf{x}^*, \mathbf{z}^*) = \mathbb{E}\{T_r(\mathbf{x}^*, \mathbf{z}^*)\}$ , where  $\mathbb{E}\{\cdot\}$  is the expectation with respect to the random iteration operator  $T_r(\mathbf{x}, \mathbf{z})$ , given the current state. In the following demonstration,  $I$  and  $J$  are arbitrarily selected from  $\{I_i\}_{i=1}^M$  and  $\{J_j\}_{j=1}^N$ .

The deterministic version of the algorithm computes updates of the form

$$\begin{aligned} \mathbf{x}_J^{k+1} &= \frac{1}{M} \sum_{i=1}^M \mathbf{x}_J^{i, k+1} \\ &= \frac{1}{M} \sum_{i=1}^M (\mathbf{x}_J^k + \mu_J^i \mathbf{g}_J^i) \\ &= \mathbf{x}_J^k + \frac{1}{M} \sum_{i=1}^M \mu_J^i (A_{I_i}^J)^T (\mathbf{y}_{I_i} - \mathbf{z}_{I_i}^k), \end{aligned} \quad (16)$$

and

$$\begin{aligned} \mathbf{z}_I^{k+1} &= \sum_{j=1}^N \mathbf{z}_I^{j, k} = \sum_{j=1}^N A_I^{J_j} (\mathbf{x}_{J_j}^k + \mu_{J_j}^i \mathbf{g}_{J_j}^i) \\ &= \sum_{j=1}^N A_I^{J_j} (\mathbf{x}_{J_j}^k + \mu_{J_j}^i (A_{I_i}^{J_j})^T \mathbf{r}_I^k) \\ &= A_I \mathbf{x}^k + \sum_{j=1}^N \mu_{J_j}^i A_I^{J_j} (A_{I_i}^{J_j})^T (\mathbf{y}_{I_i}^k - \mathbf{z}_{I_i}^k) \\ &= A_I \mathbf{x}^k + S_I (\mathbf{y}_I^k - \mathbf{z}_I^k), \end{aligned} \quad (17)$$

where  $S_I = \sum_{j=1}^N \mu_{J_j}^i A_I^{J_j} (A_{I_i}^{J_j})^T$ .

This implies that, at the fixed point  $\mathbf{x}^*$  and  $\mathbf{z}^*$ , we have

$$\mathbf{z}_I^* = (I + S_I)^{-1} A_I \mathbf{x}^* + (I + S_I)^{-1} S_I \mathbf{y}_I \quad (18)$$

and

$$\sum_{i=1}^M \mu_J^i (A_{I_i}^J)^T (\mathbf{y}_{I_i} - \mathbf{z}_I^*) = \mathbf{0}. \quad (19)$$

Eq. 18 can be expanded into the whole  $\mathbf{z}$

$$\mathbf{z}^* = (I + S_T)^{-1} A \mathbf{x}^* + (I + S_T)^{-1} S_T \mathbf{y}, \quad (20)$$

where  $S_T$  is a block diagonal matrix:

$$S_T = \begin{bmatrix} S_{I_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & S_{I_2} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & S_{I_M} \end{bmatrix} \quad (21)$$

Note that  $(I + S_T)$  is a positive semi-definite matrix if the  $\mu_J^i$  are positive. Define a diagonal matrix  $D_j$  with diagonal entries  $\mu_J^i$  where  $D_j(h, h) = \mu_J^i$  if  $h \in I_i$ . Eq. 19 can thus be written as

$$(D_j A^J)^T (\mathbf{y} - \mathbf{z}^*) = \mathbf{0} \quad (22)$$

Combine Eq. 20 and Eq. 22 gives

$$\begin{aligned} \mathbf{0} &= (D_j A^J)^T (\mathbf{y} - (I + S_T)^{-1} A \mathbf{x}^* - (I + S_T)^{-1} S_T \mathbf{y}) \\ &= (D_j A^J)^T (I + S_T)^{-1} (\mathbf{y} - A \mathbf{x}^*). \end{aligned} \quad (23)$$

This equation has to hold for all  $J$ . As  $D_j$  is a diagonal matrix, this implies that

$$\mathbf{0} = A^T (I + S_T)^{-1} (\mathbf{y} - A \mathbf{x}^*), \quad (24)$$

which shows

$$\mathbf{x}^* = (A^T (I + S_T)^{-1} A)^{-1} A^T (I + S_T)^{-1} \mathbf{y} \quad (25)$$

is a weighted least squares solution.

## REFERENCES

- [1] Y. Sagara, A. K. Hara, and W. Pavlicek, "Abdominal CT: Comparison of low-dose CT with adaptive statistical iterative reconstruction and routine-dose CT with filtered back projection in 53 patients," *Amer. J. Roentgenol.*, vol. 195, no. 3, pp. 713–719, Sep. 2010.
- [2] E. J. Hoffman, S. C. Huang, and M. E. Phelps, "Quantitation in positron emission computed tomography: 1. Effect of object size," *J. Comput. Assisted Tomography*, vol. 3, no. 3, pp. 299–308, Jun. 1979.
- [3] T. Rodet, F. Noo, and M. Defrise, "The cone-beam algorithm of feldkamp, davis, and kress preserves oblique line integrals," *Med. Phys.*, vol. 31, no. 7, pp. 1972–1975, Jun. 2004.
- [4] L. Feldkamp, L. Davis, and J. Kress, "Practical cone-beam algorithm," *J. Opt. Soc. Amer. A*, vol. 1, no. 6, pp. 612–619, Jun. 1984.
- [5] A. Gervaise, B. Osemont, and S. Lecocq, "CT image quality improvement using adaptive iterative dose reduction with wide-volume acquisition on 320-detector CT," *Eur. Radiol.*, vol. 22, no. 2, pp. 295–301, Feb. 2012.
- [6] G. Wang, H. Yu, and B. De Man, "An outlook on X-ray CT research and development," *Med. Phys.*, vol. 35, no. 3, pp. 1051–1064, Feb. 2008.
- [7] J. Deng, H. Yu, and J. Ni, "Parallelism of iterative ct reconstruction based on local reconstruction algorithm," *J. Supercomput.*, vol. 48, no. 1, pp. 1–14, Apr. 2009.
- [8] M. J. Willemink, P. A. Jong, and T. Leiner, "Iterative reconstruction techniques for computed tomography Part 1: Technical principles," *Eur. Radiol.*, vol. 23, no. 6, pp. 1623–1631, Jun. 2013.
- [9] M. Soleimani and T. Pengpen, "Introduction: A brief overview of iterative algorithms in X-ray computed tomography," *Philos. Trans. A Math. Phys. Eng. Sci.*, vol. 373, no. 2043, pp. 20140399, Mar 2015.
- [10] X. Guo, "Convergence studies on block iterative algorithms for image reconstruction," *Appl. Math. Comput.*, vol. 273, pp. 525–534, Jan. 2016.
- [11] M. Beister, D. Kolditz, and W. A. Kalender, "Iterative reconstruction methods in X-ray CT," *Phys. Medica*, vol. 28, no. 2, pp. 94–108, Apr. 2012.
- [12] F. Jacobs, E. Sundermann, B. De Sutter, and M. Christiaens, "A fast algorithm to calculate the exact radiological path through a pixel or voxel space," *J. Comput. Inf. Technol.*, vol. 6, no. 1, pp. 89–94, Mar. 1998.
- [13] J. Ni, X. Li, T. He, and G. Wang, "Review of parallel computing techniques for computed tomography image reconstruction," *Current Med. Imag. Rev.*, vol. 2, no. 4, pp. 405–414, Nov. 2006.
- [14] M. Zibetti, C. Lin, and G. Herman, "Total variation superiorized conjugate gradient method for image reconstruction," *Inverse Probl.*, vol. 34, no. 3, pp. 1–26, Jan. 2018.
- [15] M. Chillaron, V. Vidal, D. Segrelles, I. Blanquer, and G. Verdu, "Combining grid computing and docker containers for the study and parametrization of CT image reconstruction methods," *Procedia Comput. Sci.*, vol. 108, pp. 1195–1204, Jun. 2017.
- [16] S. Karimi and F. Toutounian, "The block least squares method for solving nonsymmetric linear systems with multiple right-hand sides," *Appl. Math. Comput.*, vol. 177, no. 2, pp. 852–862, Jun. 2006.

- [17] P. Dianne and O. Leary, "The block conjugate gradient algorithm and related methods," *Linear Algebra Appl.*, vol. 29, pp. 293–322, Feb. 1980.
- [18] J. Hsieh, B. Nett, and Z. Yu, "Recent advances in CT image reconstruction," *Current Radiol. Rep.*, vol. 1, no. 1, pp. 39–51, Mar. 2013.
- [19] J. Bilbao-Castro, J. Carazo, J. Fernández, and I. García, "Performance of parallel 3-D iterative reconstruction algorithms," *WSEAS Trans. Biol. Biomed.*, vol. 1, no. 1, pp. 112–119, Jan. 2004.
- [20] Y. Censor, "Row-action methods for huge and sparse systems and their applications," *SIAM Rev.*, vol. 23, no. 4, pp. 444–466, Jun. 1981.
- [21] D. W. Watt, "Column-relaxed algebraic reconstruction algorithm for tomography with noisy data," *Appl. Opt.*, vol. 33, no. 20, pp. 4420–4427, Sep. 1994.
- [22] T. Elfving, "Block-iterative methods for consistent and inconsistent linear equations," *Numer. Math.*, vol. 35, no. 1, pp. 1–12, Mar. 1980.
- [23] J. Kole and F. Beekman, "Parallel statistical image reconstruction for cone-beam X-ray CT on a shared memory computation platform," *Phys. Med. Biol.*, vol. 50, no. 6, pp. 1265–1272, Mar. 2005.
- [24] T. Li, T. J. Kao, D. Isaacson, J. C. Newell, and G. J. Saulnier, "Adaptive Kaczmarz method for image reconstruction in electrical impedance tomography," *Physiol. Meas.*, vol. 34, no. 6, pp. 595–608, May 2013.
- [25] Y. Censor, P. P. Eggermont, and D. Gordon, "Strong underrelaxation in Kaczmarz's method for inconsistent systems," *Numer. Math.*, vol. 41, no. 1, pp. 83–92, Feb. 1983.
- [26] A. H. Andersen and A. C. Kak, "Simultaneous algebraic reconstruction technique (SART): A superior implementation of the ART algorithm," *Ultrasonic Imag.*, vol. 6, no. 1, pp. 81–94, Jan. 1984.
- [27] Y. Censor, "Parallel application of block-iterative methods in medical imaging and radiation therapy," *Math. Program.*, vol. 42, no. 1, pp. 307–325, Apr. 1988.
- [28] J. Gregor and T. Benson, "Computational analysis and improvement of SIRT," *IEEE Trans. Med. Imag.*, vol. 27, no. 7, pp. 918–924, Jun. 2008.
- [29] N. D. Tang, N. De Ruyter, J. Mohr, and A. P. Butler, "Using algebraic reconstruction in computed tomography," in *Proc. 27th Conf. Image Vis. Comput.*, Dunedin, New Zealand, Nov. 2012, pp. 216–221.
- [30] T. M. Benson and J. Gregor, "Framework for iterative cone-beam micro-CT reconstruction," *IEEE Trans. Nucl. Sci.*, vol. 52, no. 5, pp. 1335–1340, Oct. 2005.
- [31] H. M. Hudson and R. S. Larkin, "Accelerated image reconstruction using ordered subsets of projection data," *IEEE Trans. Med. Imag.*, vol. 13, no. 4, pp. 601–609, Dec. 1994.
- [32] Y. Censor, D. Gordon, and R. Gordon, "Component averaging: An efficient iterative parallel algorithm for large and sparse unstructured problems," *Parallel Comput.*, vol. 27, no. 5, pp. 777–808, May 2001.
- [33] Y. Censor and T. Elfving, "Block-iterative algorithms with diagonally scaled oblique projections for the linear feasibility problem," *SIAM J. Matrix Anal. Appl.*, vol. 24, no. 1, pp. 40–58, Jul. 2002.
- [34] Y. Censor, D. Gordon, and R. Gordon, "BICAV: A block-iterative parallel algorithm for sparse systems with pixel-related weighting," *IEEE Trans. Med. Imag.*, vol. 20, no. 10, pp. 1050–1060, Oct. 2001.
- [35] K. Sauer and C. Bouman, "A local update strategy for iterative reconstruction from projections," *IEEE Trans. Signal Process.*, vol. 41, no. 2, pp. 534–548, Feb. 1993.
- [36] C. Bouman and K. Sauer, "A unified approach to statistical tomography using coordinate descent optimization," *IEEE Trans. Signal Process.*, vol. 5, no. 3, pp. 480–492, Mar. 1996.
- [37] J. Thibault, K. Sauer, C. Bouman, and J. Hsieh, "A three dimensional statistical approach to improved image quality for multislice helical CT," *Med. Phys.*, vol. 34, no. 11, pp. 4526–4544, Oct. 2007.
- [38] T. M. Benson, B. K. De Man, L. Fu, and J.-B. Thibault, "Block-based iterative coordinate descent," in *Proc. IEEE Nucl. Sci. Symp. Med. Imag. Conf.*, 2010, pp. 2856–2859.
- [39] D. Kim and J. A. Fessler, "Parallelizable algorithms for X-ray CT image reconstruction with spatially nonuniform updates," in *Proc. 2nd Int. Conf. Image Formation X-ray Comput. Tomography*, 2012, pp. 33–36.
- [40] K. Sauer, S. Borman, and C. Bouman, "Parallel computation of sequential pixel updates in statistical tomographic reconstruction," in *Proc. Int. Conf. Image Process.*, Oct. 1995, pp. 93–96.
- [41] J. Fessler, "Grouped-coordinate ascent algorithms for penalized-likelihood transmission image reconstruction," *IEEE Trans. Med. Imag.*, vol. 16, no. 2, pp. 166–175, Apr. 1997.
- [42] J. Zhang, S. Saquib, and K. Sauer, "Parallelizable Bayesian tomography algorithms with rapid, guaranteed convergence," *IEEE Trans. Image Process.*, vol. 9, no. 10, pp. 1745–1759, Oct. 2000.
- [43] J. A. Fessler and K. Donghwan, "Axial block coordinate descent (ABCD) algorithm for X-ray CT image reconstruction," in *Proc. 11th Int. Meeting Fully 3D Image Reconstruct. Radiol. Nucl. Med.*, Jul. 2011, pp. 262–265.
- [44] C. Hsieh, K. Chang, and C. Lin, "A dual coordinate descent method for large-scale linear SVM," in *Proc. 25th Int. Conf. Mach. Learn.*, Jul. 2008, pp. 408–415.
- [45] K. Chang, C. Hsieh, and C. Lin, "Coordinate descent method for large-scale  $l_2$ -loss linear support vector machines," *J. Mach. Learn. Res.*, vol. 9, no. 2, pp. 1369–1398, Jul. 2008.
- [46] Z. Yu, J. B. Thibault, C. A. Bouman, K. D. Sauer, and J. Hsieh, "Non-homogeneous updates for the iterative coordinate descent algorithm," in *Proc. SPIE-IS&T Electron. Imag.*, San Jose, CA, USA, Feb. 2007, pp. 64981B1–64981B12.
- [47] Z. Yu, J.-B. Thibault, C. A. Bouman, and K. D. Sauer, "Fast model-based X-ray CT reconstruction using spatially nonhomogeneous ICD optimization," *IEEE Trans. Image Process.*, vol. 20, no. 1, pp. 161–175, Jan. 2011.
- [48] W. J. Palenstijn, J. Bédorf, and K. J. Batenburg, "A distributed SIRT implementation for the ASTRA toolbox," in *Proc. Fully Three-Dimensional Image Reconstruct. Radiol. Nucl. Med.*, Jun. 2015, pp. 166–169.
- [49] J. Chen and Q. Gu, "Accelerated stochastic block coordinate gradient descent for sparsity constrained nonconvex optimization," in *Proc. 32nd Conf. Uncertainty Artif. Intell.*, Jun. 2016, pp. 132–141.
- [50] A. Zhang and Q. Gu, "Accelerated stochastic block coordinate descent with optimal sampling," in *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining*, Aug. 2016, pp. 2035–2044.
- [51] H. Wang and A. Banerjee, "Randomized block coordinate descent for online and stochastic optimization," *arXiv preprint arXiv:1407.0107*, Jul. 2014.
- [52] T. Zhao, M. Yu, Y. Wang, R. Arora, and H. Liu, "Accelerated mini-batch randomized block coordinate descent method," in *Proc. Neural Inf. Process. Syst.*, Dec. 2014, pp. 3329–3337.
- [53] R. Johnson and T. Zhang, "Accelerating stochastic gradient descent using predictive variance reduction," in *Proc. Neural Inf. Process. Syst.*, Dec. 2013, pp. 315–323.
- [54] J. Konecny, Q. Zheng, and P. Richtarik, "Semi-stochastic coordinate descent," *Optim. Methods Softw.*, vol. 32, no. 5, pp. 993–1005, Mar. 2017.
- [55] Y. Xu and W. Yin, "Block stochastic gradient iteration for convex and nonconvex optimization," *SIAM J. Optim.*, vol. 25, no. 3, pp. 1686–1716, Aug. 2015.
- [56] Z. Lu and L. Xiao, "On the complexity analysis of randomized block-coordinate descent methods," *Math. Program.*, vol. 152, no. 1, pp. 615–642, Aug. 2015.
- [57] M. Neumann and R. Plemmons, "Convergence of parallel multisplitting iterative methods for M-matrices," *Numer. Linear Algebra Appl.*, vol. 88, pp. 559–573, Apr. 1987.
- [58] R. Renault, "A parallel multisplitting solution of the least squares problem," *Numer. Linear Algebra Appl.*, vol. 5, no. 1, pp. 11–31, May 1998.
- [59] A. Frommer and G. Mayer, "Convergence of relaxed parallel multisplitting methods," *Linear Algebra Appl.*, vol. 119, pp. 141–152, Jul. 1989.
- [60] R. Wen and H. Duan, "A parallel multisplitting method with self-adaptive weightings for solving H-matrix linear systems," *J. Inequalities Appl.*, vol. 2017, no. 1, pp. 95–105, Jun. 2017.
- [61] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Found. Trends Mach. Learn.*, vol. 3, no. 1, pp. 1–122, Jul. 2011.
- [62] W. Deng, M. Lai, Z. Peng, and W. Yin, "Parallel multiblock ADMM with  $o(1/k)$  convergence," *J. Sci. Comput.*, vol. 71, no. 2, pp. 712–736, May 2017.
- [63] N. Parikh and S. Boyd, "Block splitting for distributed optimization," *Math. Program. Comput.*, vol. 6, no. 1, pp. 77–102, Mar. 2014.
- [64] W. Aarle *et al.*, "The ASTRA toolbox: A platform for advanced algorithm development in electron tomography," *Ultramicroscopy*, vol. 157, pp. 35–47, May 2015.
- [65] A. Biguri, M. Dosanjh, S. Hancock, and M. Soleimani, "TIGRE: A MATLAB-GPU toolbox for CBCT image reconstruction," *Biomed. Phys. Eng. Express*, vol. 2, no. 5, Sep. 2016, Art. no. 055010.
- [66] J. Langford, A. J. Smola, and M. Zinkevich, "Slow learners are fast," in *Proc. Adv. Neural Inf. Process. Syst.*, 2009, pp. 2331–2339.
- [67] A. Alekh and J. C. Duchi, "Distributed delayed stochastic optimization," in *Proc. Adv. Neural Inf. Process. Syst.*, 2011, pp. 873–881.