

Consensus Equilibrium Framework for Super-Resolution and Extreme-Scale CT Reconstruction

Xiao Wang
Harvard Medical School
Boston Children's Hospital

Venkatesh Sridhar
Purdue University

Zahra Ronaghi
NVIDIA Corporation

Rollin Thomas
Jack Deslippe
Dilworth Parkinson
Lawrence Berkeley Laboratory

Gregory T. Buzzard
Samuel P. Midkiff
Charles A. Bouman
Purdue University

Simon K. Warfield
Harvard Medical School
Boston Children's Hospital

ABSTRACT

Computed tomography (CT) image reconstruction is a crucial technique for many imaging applications. Among various reconstruction methods, Model-Based Iterative Reconstruction (MBIR) enables super-resolution with superior image quality. MBIR, however, has a high memory requirement that limits the achievable image resolution, and the parallelization for MBIR suffers from limited scalability. In this paper, we propose Asynchronous Consensus MBIR (AC-MBIR) that uses Consensus Equilibrium (CE) to provide a super-resolution algorithm with a small memory footprint, low communication overhead and a high scalability. Super-resolution experiments show that AC-MBIR has a 6.8 times smaller memory footprint and 16 times more scalability, compared with the state-of-the-art MBIR implementation, and maintains a 100% strong scaling efficiency at 146880 cores. In addition, AC-MBIR achieves an average bandwidth of 3.5 petabytes per second at 587520 cores.

CCS CONCEPTS

• **Computing methodologies** → **Massively parallel algorithms; Image processing;**

ACM Reference format:

Xiao Wang, Venkatesh Sridhar, Zahra Ronaghi, Rollin Thomas, Jack Deslippe, Dilworth Parkinson, Gregory T. Buzzard, Samuel P. Midkiff, Charles A. Bouman, and Simon K. Warfield. 2019. Consensus Equilibrium Framework for Super-Resolution and Extreme-Scale CT Reconstruction. In *Proceedings of The International Conference for High Performance Computing, Networking, Storage, and Analysis, Denver, CO, USA, November 17–22, 2019 (SC '19)*, 13 pages.

<https://doi.org/10.1145/3295500.3356142>

1 INTRODUCTION

A computed tomography (CT) system is an imaging modality to produce a 3D reconstruction of a scanned object. As one of the most

popular imaging systems, CT is heavily used for medical imaging, with 72 million scans per year in the US alone, and is also widely used in security imaging [12], scientific imaging [18, 24, 25, 27, 28], and industrial imaging for inspecting manufacturing parts [19].

An important goal for CT imaging is super-resolution, namely to improve image resolution beyond the limit of detector resolution. CT super-resolution is not only used to detect early-stage cancers [13], but is also used to study microscopic structure, such as brain neurons and synthetic material micro-structures [8]. To accomplish super-resolution, interpolation methods produce a Filtered-back Projection (FBP) image through an inverse radon transform, and then interpolate the FBP image to super-resolution. Interpolation-based methods, however, fail to provide more image detail than the uninterpolated FBP images [17]. Example-based methods, such as deep learning and dictionary learning [11], are another popular choice for super-resolution. They extract image details from similar example images and apply these details to new images. In many cases, unfortunately, example-based methods require an unattainably large training database. More importantly in the presence of unusual image outliers such as metal, image details learned from example images are not always consistent with the object's taken X-ray measurements and example-based methods fail to truthfully reconstruct.

In contrast, Model-Based Iterative Reconstruction (MBIR) [23, 29] is a Bayesian based method that not only achieves good super-resolution quality consistent with X-ray measurements [9, 21], but also reduces the radiation doses by up to 80% for certain applications [14]. At a high level, MBIR has better image quality because it incorporates complex mathematical models to preserve X-ray acquisition information, while estimating image details and suppressing noise from image spatial property.

Despite of the benefits described above, the realistic but complex models for MBIR have a very large memory requirement to store a physics-related geometry system matrix, whose size grows linearly with image size. Any naive way to shrink or distribute the models among distributed nodes corrupts the CT physics-geometry system matrix. Therefore, the high memory requirement restricts the image resolution that MBIR can achieve, as such geometry system matrix information has to be broadcast to and stored on all nodes.

In addition to an excessive memory requirement, the parallel scalability of MBIR is another performance bottleneck. With its complex models, MBIR often has much longer reconstruction times and is computationally much more expensive than other methods

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SC '19, November 17–22, 2019, Denver, CO, USA

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-6229-0/19/11...\$15.00

<https://doi.org/10.1145/3295500.3356142>

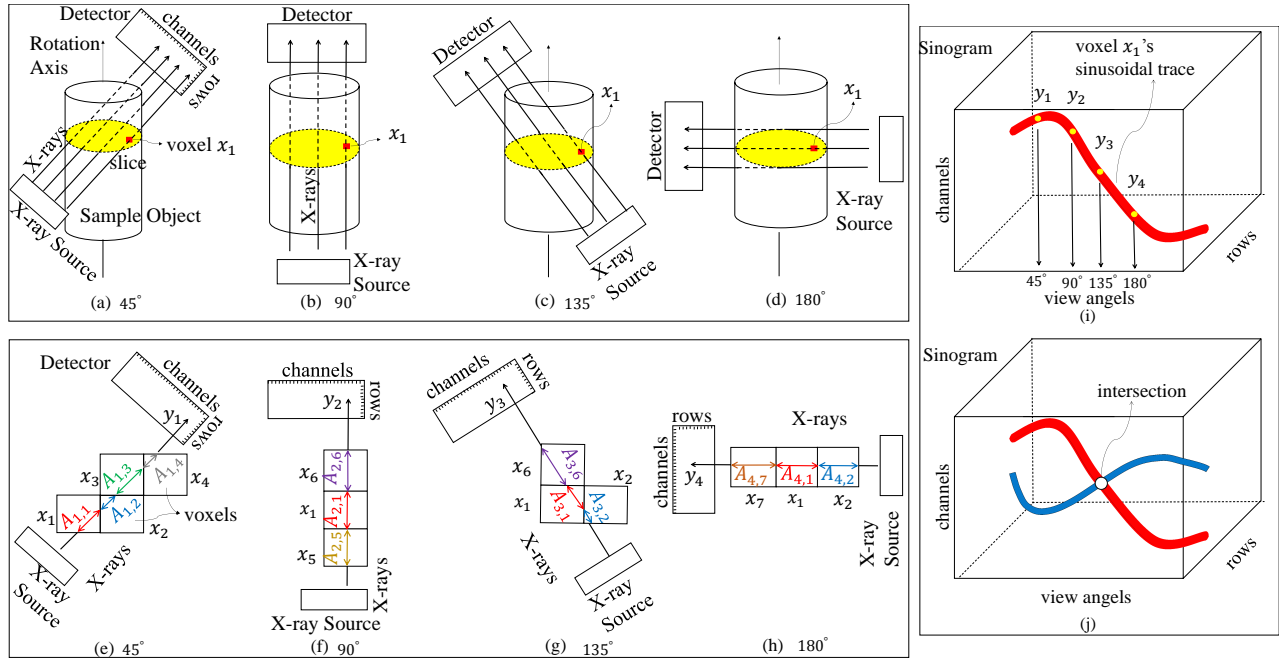


Figure 1: (a) A CT system at 45° view angle with a slice in yellow and a voxel x_1 in red. (b), (c) and (d) CT system rotates to 90°, 135° and 180° respectively. (e) An X-ray at 45° view angle passes through x_1 and its neighbor voxels and hits the detector with a measurement y_1 . The intersection between the X-ray and each voxel is shown as a line segment with a different color for each. (f), (g) and (h) X-rays pass through x_1 at 90°, 135° and 180° respectively. (i) Shows the sinusoidal memory layout for x_1 in Sinogram. (j) Demonstrate an intersection between two voxel traces.

by several orders of magnitude. Therefore, four parallelization levels have been proposed to reduce computation times: (1) parallelizing SIMD operations for a voxel update [27], where a voxel is a three-dimensional pixel in the reconstruction volume; (2) parallelizing across voxels [6, 31]; (3) parallelizing across super-voxels (SVs), where an SV is a group of spatially-close voxels in the shape of a rectangular cuboid [26, 27]; and (4) parallelizing across slices [3, 7, 18], where a slice is defined as a cross-section of the scanned object with a fixed thickness. Unfortunately, these four levels of parallelization are insufficient to fully utilize a large supercomputer. With all four parallelization levels, the scalability for the state-of-the-art MBIR implementation is still limited to a small portion of a supercomputer [27].

This paper proposes the *Asynchronous Consensus MBIR Algorithm* (AC-MBIR) that significantly reduces MBIR’s memory footprint and achieves better scalability. Without AC-MBIR, MBIR has a limited super-resolution and cannot scale to a large supercomputer. AC-MBIR allows each node to perform an individual reconstruction using a subset of measurements, and then uses a lossless *Consensus Equilibrium Method* (CE) to merge the individual reconstructions asynchronously into a consensus solution. But unlike the previous work on CE that has no memory footprint reduction or distributed node implementation [4, 22], AC-MBIR uses a new technique, *grouped-view partition*, to optimize memory footprint reduction. Furthermore, parallelism for AC-MBIR is not only achieved by dividing a reconstruction volume into slices and SVs, but also by dividing models and measurements into data subsets without

corrupting their information, and performing parallel individual reconstructions with a low communication overhead.

In summary, this paper makes three contributions: (1) proposing a novel grouped-view partition technique to distribute the CT geometry model and measurements among nodes, and optimizing memory footprint reduction; (2) introducing the AC-MBIR algorithm that enables parallel asynchronous consensus updates for a high scalability on a supercomputer with a low communication overhead; and (3) applying the CE Method to AC-MBIR to ensure a convergence to the global optimum. Super-resolution experiments show that AC-MBIR has a 6.8 times smaller memory footprint and 16 times more scalability, compared with the state-of-the-art MBIR implementation [27], and maintains a 100% strong scaling efficiency at 146880 cores. In addition, AC-MBIR achieves an average bandwidth of 3.5 petabytes per second at 587520 cores.

2 BACKGROUND

2.1 CT System Setup

A CT system places the sample object in the center and the CT system has a gantry that rotates around the sample object with an x-ray source at one end of the gantry and an X-ray detector at the other end. At each view angle of a rotation, the X-ray source emits X-rays that pass through the sample object and project onto an X-ray detector, whose horizontal sensor units are channels and vertical sensor units are rows. Each received CT measurement on the X-ray detector corresponds to the integral of radio density of the sample

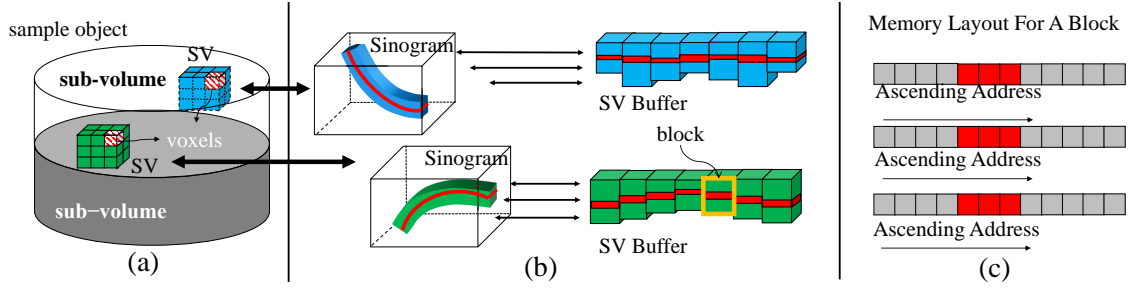


Figure 2: (a) Two nodes update two sub-volumes simultaneously and each core of a node updates a different SV in its assigned sub-volume. (b) Two private SV buffers for efficient data operations and each SV buffer consists of blocks of measurements. (c) Shows the memory layout for a block with completely regular access pattern.

object along the path of X-rays at the current view angle. Figure 1(a) shows the CT system setup at view angle of 45° with a single slice depicted in yellow and a voxel, x_1 , in that slice depicted in red. In addition, detector channels and rows are represented as short dashes in the detector. Figure 1(e) shows an X-ray that passes through voxel x_1 , and x_1 's neighbor voxels, x_2 , x_3 and x_4 , at 45° view angle and hits the right side of the detector with a CT measurement, y_1 . We define the length of intersection between the X-ray in Figure 1(e) and voxel x_2 as $A_{1,2}$, shown as a blue line segment. Similarly, the lengths of intersection for the X-ray in Figure 1(e) and voxels x_1 , x_3 and x_4 are denoted as $A_{1,1}$, $A_{1,3}$ and $A_{1,4}$ respectively. Since y_1 is a line integral of all voxels' radio density on the path of the X-ray, y_1 can be written as $y_1 = A_{1,1}x_1 + A_{1,2}x_2 + A_{1,3}x_3 + A_{1,4}x_4 + e_1$, where e_1 is y_1 's data error and represents the data corruption from scanner noise and detector artifacts.

Figure 1(b), (c) and (d) are the CT system setup rotated to view angles of 90° , 135° and 180° respectively, and Figure 1(f), (g) and (h) show measurements, y_2 , y_3 and y_4 for the X-rays that pass through voxel x_1 at corresponding view angles. Note that the X-rays for y_2 , y_3 and y_4 pass through a different set of voxels with different lengths of intersections at each view angle. For example, $A_{2,1}$ is longer than $A_{3,1}$ and $A_{1,1}$. In addition, X-rays hit different detector channels at different view angles. In this example, a right detector channel receives measurement y_1 at 45° and a center detector channel receives measurement y_4 at 180° . If we use e_2 , e_3 and e_4 to represent the data errors for y_2 , y_3 and y_4 , then voxel x_1 can then be computed as the solution to the following linear equations:

$$y_1 = A_{1,1}x_1 + A_{1,2}x_2 + A_{1,3}x_3 + A_{1,4}x_4 + e_1, \quad (1)$$

$$y_2 = A_{2,5}x_5 + A_{2,1}x_1 + A_{2,6}x_6 + e_2, \quad (2)$$

$$y_3 = A_{3,2}x_2 + A_{3,1}x_1 + A_{3,6}x_6 + e_3, \quad (3)$$

$$y_4 = A_{4,7}x_7 + A_{4,1}x_1 + A_{4,2}x_2 + e_4, \quad (4)$$

If we generalize the above equations for all voxels, then we need to compute the solution to the following linear system:

$$y = Ax + e, \quad (5)$$

where reconstruction x is an $Nx \times Ny \times Nz$ matrix ($Nx \times Ny$ voxels in a slice and Nz slices); measurements y is an $N_r \times N_c \times N_v$ matrix, also called *Sinogram*, with N_r detector rows, N_c detector channels, and N_v view angles; A is a 6 dimensional geometry system matrix and

the size of A increases proportionally with the sizes of x and y . For simplicity, all following equations in this paper vectorize x matrix as an N dimensional vector, where $N = Nx \times Ny \times Nz$, y matrix as an M dimensional vector, where $M = N_r \times N_c \times N_v$, and A as an $M \times N$ system matrix. Each element of A matrix, $A_{i,j}$, represents the length of intersection between the X-ray for measurement y_i and voxel x_j . $A_{i,j}$ is zero when x_j is not on the path of the X-ray for y_i . To show how large A matrix is and why it leads to a large memory footprint, consider an x and an y of size $1000 \times 1000 \times 1000$ for both. In addition, we assume that the A matrix is stored in memory in floating point sparse matrix format with 0.1% non-zero elements. In this example, the entire A matrix then requires a memory size of 3638 terabytes!

In addition, if we trace out the memory locations of each voxel's measurements in Sinogram, the memory locations follow a sinusoidal trace in Sinogram as X-rays that pass through the same voxel hit different detector channels at different views. To update each voxel, cache lines must read and write data for the voxel following the same sinusoidal trace in memory. Cache lines, however, has a high cache miss rate as no cache lines can efficiently operate on data following a sinusoidal path in memory. In addition, each voxel has a unique trace in Sinogram and different voxel's traces often overlap at some intersections in Sinogram. Spatially close voxels have highly overlapped traces and more intersections. Spatially separated voxels have fewer overlapped traces with few or no intersections. Figure 1(j) shows an intersection colored in white between two voxel traces. To read and write data for voxel updates, locks or atomic operations must be applied to data at the intersections.

2.2 MBIR Formulation

Unfortunately, the solution to Equation (5) cannot be directly computed by taking a linear inverse, given that the data error e is unknown and cannot be measured. In addition, the size of reconstruction for super-resolution, N , is often much bigger than the size of measurements, M , leading to undetermined solutions to Equation (5). For example, Equations (1-4) have 4 measurements but 7 voxels to be computed. Therefore, voxels are undetermined with many possible solutions. For reconstruction x to be a determined solution, the size of measurements, M , must be at least as large as the size of reconstruction, N . To have a larger M size, the

CT detector resolution must be higher as M is proportional to the detector resolution.

To address the above problems without spending millions of dollars on a scanner for a higher detector resolution, the super-resolution MBIR is the solution to the following cost function:

$$x = \operatorname{argmin}_x \left\{ \frac{1}{2} \|y - Ax\|_{\Lambda}^2 + R(x) \right\} \quad (6)$$

where Λ is a diagonal weighting matrix that models the CT system statistical noise. $R(x)$ is a prior model to estimate image details when measurements y are insufficient. In Equation (6), $\frac{1}{2} \|y - Ax\|_{\Lambda}^2$ is the forward model to fit reconstruction x from measurements y . The prior model, $R(x)$, is to penalize irregular image spatial property, suppress noise and estimate image details based on prior information. In this paper, we use a strictly convex Generalized Markov Random Field (GMRF) prior model to estimate each voxel's image details based on its neighboring voxels [23, 30] and Equation (6) has a unique global optimum given that both the forward and the prior models are strictly convex. Equation (6) can be viewed as a supervised machine learning problem that trains on a high variance dataset with limited size, whose forward model is the mean-square-error (MSE) of a weighted linear regression model and the prior model is a regularizer to prevent data over-fitting.

2.3 Related Work

To compute Equation (6) efficiently, there are two approaches: the volume partition and measurement partition methods. The volume partition methods include Grouped Coordinate Descent [6, 31], Parallel Iterative Coordinate Descent [10, 16, 20], and the Non-Uniform Parallel Super-Voxel (NU-PSV) algorithm [27]. To achieve parallelism, the volume partition methods split the reconstruction volume, x , into smaller voxel groups and assign each core with a voxel group to be updated in parallel. The scalability in such methods, however, is limited by the number of voxels in a reconstruction. In addition, the volume partition methods often have a high communication overhead with little speedup [6, 31]. Furthermore, the volume partition method cannot reduce the memory footprint for MBIR and the attainable image resolution is therefore limited.

The measurement partition methods, such as Parallel Mini-Batch Gradient Descent [15] and Ordered-Subset methods [1, 5], split measurements, y , and the system matrix into subsets and distribute them among cores. Each core reconstructs a private volume using its assigned subset, and then the private reconstruction volumes from all cores are merged into a consensus reconstruction. The measurement partition methods, however, are approximation methods and a convergence to the global optimum is not guaranteed. If measurements y are uniformly sampled among subsets, the methods converge close to the global optimum. If y are partitioned non-uniformly, the converged solution is far from the global optimum. In addition to the convergence issue, the amount of parallelism for the measurement partition methods is very limited as the number of parallel cores is at most the number of subsets, and subsets are usually few (≤ 40) as too many subsets slow down convergence [1].

Recently, the Consensus Equilibrium Method (CE) shows that it can partition measurements in any order and always converge to the global optimum with a proof in [4]. The paper in [4], however,

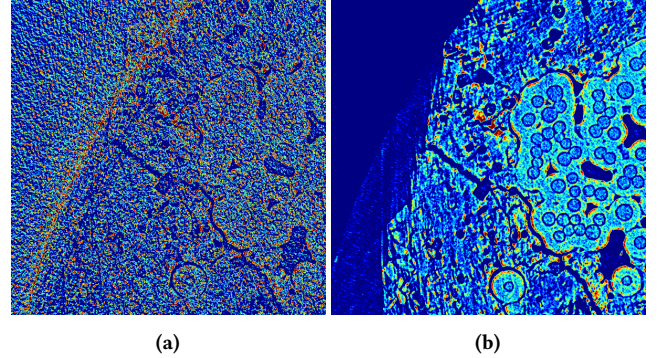


Figure 3: (a) An LR image with significant noise reconstructed from 1/4 of the uniformly sampled measurements. (b) A consensus HR image with much less noise and clearer details.

only discusses CE mathematical theory and does not have any application to massive parallelism, memory footprint reduction or image reconstruction. The Multi-Agent Consensus Equilibrium (MACE) algorithm in [22] uses the CE Method for CT image reconstruction and partitions measurements into *uniform sparse-view subsets*, which are subsets of interleaved view angles. Taking Figure 1 as an example with 2 uniform sparse-view subsets, measurements for Figure 1(a) and (c) are in a subset, and measurements for Figure 1(b) and (d) are in the other subset. The MACE algorithm, however, is for single slice 2D reconstruction, and does not have a distributed node implementation or experiments to measure the reduction in the memory footprint.

In comparison, AC-MBIR utilizes both volume partition and measurement partition for a fully 3D reconstruction. AC-MBIR not only scales to a large supercomputer with little communication overhead, but also guarantees a convergence to the global optimum. In addition, AC-MBIR proposes a new measurement partition method, grouped-view partition, that optimizes memory footprint reduction for super-resolution CT reconstruction.

2.4 The State of The Art Implementation

This paper compares AC-MBIR with the NU-PSV algorithm, which is currently the fastest MBIR implementation and scales to hundreds of thousands of cores, and was selected as a Gordon Bell Prize finalist [27]. We do not compare AC-MBIR with Mini-Batch Gradient Descent or Ordered-Subset (except convergence experiments in Figure 7) because their reconstructions are approximated and incomparable with AC-MBIR's exact solution. AC-MBIR is not compared with the MACE algorithm as MACE is for 2D single-slice reconstruction but AC-MBIR is for fully 3D reconstruction.

NU-PSV achieves its scalability by dividing a full volume x into sub-volumes, where a sub-volume is a consecutive group of slices. Nodes update different sub-volumes simultaneously, and different cores in a node update spatially separated super-voxels (SVs) in a sub-volume in parallel. Figure 2(a) shows an example volume divided into 2 sub-volumes and two SVs in a sub-volume, and each SV is a group of neighboring voxels in the shape of a $3 \times 2 \times 2$ rectangular cuboid. Since both spatially separated SVs and different

sub-volumes have little shared data to access and update in Sinogram, as explained in Section 2.1, NU-PSV has a low communication overhead by minimizing shared data among nodes and cores.

To achieve a high per-core performance, each core updates voxels in its assigned SV in sequence using the Coordinate Descent method [26, 27]. Note that measurements for an SV are closely grouped together in the shape of a sinusoidal band in the Sinogram, as shown in Figure 2(b), and each voxel in the SV follows a sinusoidal trace within the band. Example voxels and their corresponding traces in Sinogram are colored in red in Figure 2(a) and (b). Different voxels in an SV follow different, but highly overlapped, traces in the same band with many shared data in common. Therefore, measurements for voxels in an SV are close in memory and updating voxels in an SV in sequence is cache efficient with up to 187 times speedup on 20 cores [26].

Furthermore, NU-PSV reorganizes the sinusoidal memory layout for an SV into blocks of memory, called *super-voxel buffer* (SV buffer), that reads and writes data from memory in a completely regular pattern. To reorganize the memory layout for a good hardware prefetching and SIMD performance, NU-PSV introduces redundant measurements to the SV buffers and pads zero-value elements to the system A matrix [20, 26, 27]. Figure 2(b) shows two private SV buffers and outlines an example block in an SV buffer with a yellow boundary. Figure 2(c) shows the memory layout for the example block and features a completely regular access pattern, where red cells in the figure are memory locations for a voxel's measurements.

Although NU-PSV has the above benefits, its resolution is limited by the memory footprint and its scalability is limited by the number of slices. As each sub-volume must have at least 1 slice, the number of parallel nodes is at most the same as the number of slices. NU-PSV cannot increase its scalability further because it has depleted all sources of parallelism from dividing reconstruction x .

3 METHODS

3.1 Reducing Memory Footprint and Ensuring Convergence

To reduce the memory footprint, the system matrix A and measurements y should be partitioned and stored in distributed memory without corrupting their information. In this section, we will apply the CE Method to super-resolution MBIR, so that the system matrix and the measurements can be partitioned in any order, including any non-uniform sampling, and the final reconstruction from a partitioned system matrix and measurements is exactly the same as the solution to unpartitioned measurements. With the CE Method, the memory footprint for MBIR reduces significantly as a memory only needs to hold data for a partitioned subset.

To use the CE Method, we first partition measurements y into G subsets in any order, and the i^{th} measurement subset is denoted as y^i . Similarly, the system matrix A is partitioned in the same order, and the system matrix rows related to y^i is denoted as A^i . For each subset y^i , MBIR produces a low-resolution (LR) reconstruction, denoted as v^i . More precisely, an MBIR LR reconstruction from the i^{th} subset is the solution to a sub-problem, $f^i(v^i)$, defined below:

$$v^i \leftarrow \operatorname{argmin}_{v^i} f^i(v^i) = \operatorname{argmin}_{v^i} \left\{ \frac{1}{2} \|y^i - A^i v^i\|_{\Lambda^i}^2 + \frac{R(v^i)}{G} \right\}, \quad (7)$$

where Equation (7) is the same as Equation (6), except that v^i is not a reconstruction from all measurements, but from the i^{th} subset. Then LR reconstructions from all subsets are merged into a consensus high-resolution (HR) reconstruction. Figure 3(a) shows an LR image reconstructed from 1/4 of the uniformly sampled measurements. Figure 3(b) shows a consensus HR solution after merging all LR reconstructions. Figure 3(a) has much image noise and lacks image details while Figure 3(b) suppresses noise and recovers image details.

A crucial question is how to merge LR reconstructions into a consensus HR solution, x . A popular method is Parallel Mini-Batch Gradient Descent with its pseudo-code in Algorithm 1 [15], where α is a fixed learning rate and ∇f^i is the gradient of f^i with respect to v^i . In step 4, v^i is updated along the direction of ∇f^i . After all LR reconstructions are updated, the consensus solution, x , is updated in step 6 by averaging all LR reconstructions. In step 8, the averaged LR reconstruction is used as the input to update v^i in the next iteration. If each subset only holds a single measurement, then Algorithm 1 is the Parallel Stochastic Gradient Descent [32]. If the learning rate, α , is replaced with a relaxed learning rate that diminishes across iterations, then Algorithm 1 is the Ordered-Subset algorithm [1]. To converge to an approximated solution close to the global optimum, Mini-Batch Gradient Descent requires that measurements and the system matrix are partitioned uniformly and the gradient, ∇f^i , needs to be approximately the same across different subsets. Otherwise, Mini-Batch Gradient Descent converges to a solution far from the actual global optimum.

Algorithm 1 Parallel Mini-Batch Gradient Descent Algorithm

INPUT: Uniformly partitioned measurements y and matrix A .
LOCAL: G : the number of subsets. α : A fixed learning rate. ∇f^i : A gradient of f^i . v_i : the i^{th} LR reconstruction.
OUTPUT: x : the consensus HR solution.

- 1: Initialize x, v^i
- 2: **while** Not Converged **do**
- 3: **for** each measurement subset i from 1 to G **do in parallel**
- 4: $v^i \leftarrow v^i - \alpha \nabla f^i(v^i)$
- 5: **end for**
- 6: $x \leftarrow \text{Average}(v^1, \dots, v^G)$
- 7: **for** each measurement subset i from 1 to G **do in parallel**
- 8: $v^i \leftarrow x$
- 9: **end for**
- 10: **end while**

To ensure a convergence to the global optimum of Equation (6) and enable non-uniform subsets, we introduce the CE Method's proximal function to MBIR, denoted as $F_i(x)$, whose goal is to find a balance between an LR reconstruction, v^i , and the consensus HR solution, x :

$$v^i \leftarrow \operatorname{argmin}_{v^i} F^i(x) = \operatorname{argmin}_{v^i} \left\{ f^i(v^i) + \frac{\|v^i - x\|^2}{2\sigma^2} \right\}, \quad (8)$$

where the proximal map function input is x and output is v_i . The goal of $f^i(v^i)$ is to find v^i that fits the measurements for the subset,

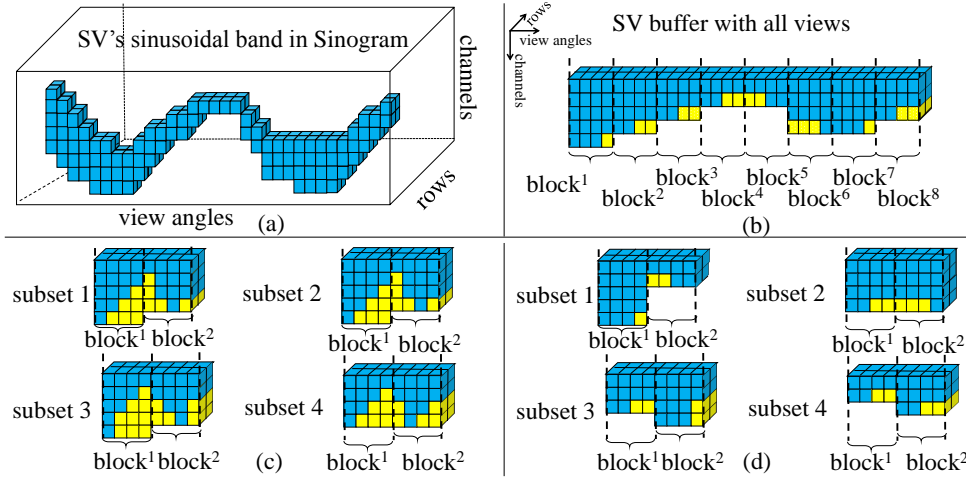


Figure 4: (a) A zoom-in on an SV's sinusoidal band in Sinogram. (b) An SV buffer copied from the band in Sinogram and every four columns is a block. Yellow cells in the SV buffer are redundant measurements. (c) 4 measurement subsets uniformly sampled across view angles. The yellow cell redundant measurements are significantly more than those in (b). (d) 4 grouped-view partitioned subsets. Notice that the redundant measurements are a lot fewer than (c).

and the goal of the penalizing term, $\frac{\|v^i - x\|^2}{2\sigma^2}$, is to constrain v^i to be close to the consensus HR solution, x . σ is a constant scalar that controls the strength of the penalizing term for fastest convergence. The optimal value for σ is found through convergence experiments and σ value increases when the number of subsets, G , increases. In each iteration, every subset evaluates its proximal map function, F^i , and produces an LR reconstruction. Then the consensus solution, x , is updated by fusing all LR reconstructions, where the fusing operation is more than a simple averaging, as we will discuss below. Iterations repeat until $x = v^i$ for $i = 1, 2, \dots, N$ and the final consensus solution x is exactly the same as the solution to Equation (6).

Algorithm 2 Consensus Equilibrium MBIR

INPUT: Measurements y and matrix A partitioned in any order.

LOCAL: w^i : a temporary copy for v^i . u^i : input change to proximal map function, $F^i(x)$. Other variables are the same as in Algorithm 1.

OUTPUT: x : the consensus HR solution.

```

1: Initialize  $x, v^i, w^i$ , and  $u^i$ 
2: while  $x \neq v^i, i = 1, 2, \dots, G$  do
3:   for each subset  $i$  from 1 to  $G$  do in parallel
4:      $v^i \leftarrow \operatorname{argmin}_{v^i} F^i(x + u^i)$ 
5:      $(w')^i \leftarrow w^i$ 
6:      $w^i \leftarrow (2v^i - x - u^i)$ 
7:      $w^i \leftarrow \rho w^i + (1 - \rho)(w')^i, 0 \leq \rho \leq 1$ 
8:   end for
9:    $x \leftarrow \operatorname{Average}(w^1, \dots, w^G)$ 
10:  for each subset  $i$  from 1 to  $G$  do in parallel
11:     $u^i \leftarrow x - w^i$ 
12:  end for
13: end while
    
```

Algorithm 2 summarizes how to apply the CE Method to MBIR and what the fusing operation does. For each subset, step 4 adds the proximal function input with u^i and produces an LR reconstruction. Step 6 adjusts the LR reconstruction and stores the new value in w^i . Step 7 damps the adjustment in step 6, where ρ is a constant scalar between 0 and 1 to control the degree of damping for fast convergence. The optimal value for ρ is chosen through convergence experiments and in general, more parallelism requires a smaller ρ for fastest convergence. Step 9 updates the consensus solution, x , by averaging w^i across all subsets. If $x \neq v^i$, the proximal function input change, u^i , is updated and iterations repeat until a convergence is reached, and a convergence proof is provided in Appendix A.

Note that Algorithm 2 allows the measurements and the system matrix to be partitioned in any order. In the special case of a uniform sparse-view partition, Algorithm 2 is the same as the MACE algorithm [22]. To show convergence rate comparison between Algorithm 1 and Algorithm 2, experiments in Figure 7 show that AC-MBIR has a much faster convergence than Mini-Batch Gradient Descent for both the uniform sparse-view and the non-uniform partitions. In addition, Mini-Batch Gradient Descent with few sparse-view subsets converges closer to the global optimum than the same algorithm with non-uniform partition. In contrast, AC-MBIR always converges to the global optimum for both the uniform and non-uniform partitions.

3.2 More Memory Footprint Reduction

The uniform sparse-view partition is widely used but does not optimize memory reduction. This section proposes a non-uniform grouped-view partition that further reduces memory footprint.

Figure 4(a) zooms in on an SV's 3D sinusoidal band in Sinogram and the band curves up and down in Sinogram, as explained in Section 2.1, with a varying width at different view angles [27]. In addition, neighboring view angles have similar widths and distant

view angles have dissimilar widths. Figure 4(b) copies data from the sinusoidal band to an SV buffer column by column (in channel direction) and every 4 view angles is a block. Each block is padded with redundant data to make the block a rectangular cuboid. For each redundant data, its corresponding system matrix row is also padded with zeros so that computations on the redundant data do not change the final reconstruction. In the example of Figure 4(b), redundant data are colored in yellow and the essential data for reconstruction are colored in blue. In addition, Figure 4(b) has 8 blocks and neighboring blocks are separated by short vertical dash lines. Columns in the same block have a constant width, although the width varies across different blocks. To have a good data vectorization, cache lines operate on data in the SV buffer block by block with a completely regular data access pattern.

Figure 4(c) uniformly partitions the SV buffer into 4 subsets across view-angles. Subset 1, for example, takes the first column of each block in the SV buffer and groups these columns into 2 new blocks in the shape of rectangular cuboids. Subset 2 takes the second column of each block in the SV buffer and group columns into 2 new blocks. Note that Figure 4(c) has 31% more redundant data than in the SV buffer. This phenomenon is due to the fact that uniformly partitioned subsets consist of data from distant view angles and have widely varying widths in each block. Therefore, significant amounts of redundant data are needed to pad each block into rectangular cuboids and the benefit of memory footprint reduction diminishes. To make matters even worse, the uniform sparse-view partition has more redundant data with more subsets.

To reduce redundant data, Figure 4(d) partitions the SV buffer into 4 grouped-view subsets. Each block in a subset consists of data from neighboring view angles, and different blocks in the same subset consist of data from distant view angles. In the example of Figure 4(d), subset 1 consists of SV buffer's first and fifth blocks. Subset 2 consists of SV buffer's second and sixth blocks. With the grouped-view partition design, Figure 4(d) has the same amount of redundant data as the unpartitioned SV buffer in Figure 4(b), and significantly fewer redundant data than the uniform sparse-view partition in Figure 4(c). Synchrotron imaging experiments in Table 1 shows that the redundant data for 16 uniform sparse-view subsets increase the memory footprint and the computations on data by 2.1 times, compared to the unpartitioned baseline. 16 grouped-view subsets, however, reduce the uniform sparse-view's memory footprint by 33% and the amount of computations by 38%. Table 2 further shows that the grouped-view partition has a 3.9 times speedup over the uniform sparse-view partition at 8640 nodes (587520 cores).

3.3 Improving Scalability

The CE Method has a unique source of parallelism across measurement subsets as different LR reconstructions can be updated in parallel. To improve scalability, the AC-MBIR algorithm is a new implementation that improves on the NU-PSV algorithm with an orthogonal level of parallelism across subsets, and can scale to the entire supercomputer even for a dataset with few slices.

The AC-MBIR algorithm divides the consensus HR solution, x , into sub-volumes with a unique color for each sub-volume. Example Figure 5(a) shows a consensus solution with 4 different

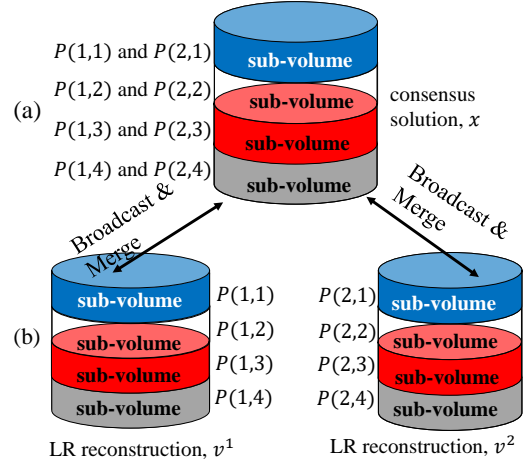


Figure 5: (a) Shows a consensus HR solution with 4 sub-volumes. (b) Shows 2 LR reconstructions with 4 sub-volumes. Nodes in the same cabinet update the same color sub-volumes in v^1 and v^2 as well as the same color sub-volume in the consensus solution.

color sub-volumes. Similarly, each LR reconstruction (2 subsets) in Figure 5(b) is also divided into 4 sub-volumes, with different sub-volumes across subsets colored the same if they merge to the same color sub-volume in the consensus HR solution. Then the AC-MBIR algorithm divides nodes in a supercomputer into clusters and cabinets, where the number of clusters equals to the number of measurement subsets and the number of cabinets is equal to the number of sub-volumes. Each node, denoted as $P(i, j)$, has dual identities with i as the conceptual cluster it belongs to and j as the conceptual cabinet it belongs to. Example Figure 5(b) shows 8 nodes divided into 2 clusters with $P(1, 1), P(1, 2), P(1, 3), P(1, 4)$ in a cluster. At the same time, nodes are also divided into 4 cabinets with $P(1, 1)$ and $P(2, 1)$ in a cabinet.

To map the node configuration with computations for image reconstruction, AC-MBIR lets each node $P(i, j)$ update the j^{th} sub-volume of the i^{th} LR reconstruction, and nodes in the same cluster share measurements and system matrix rows from the same subset, but update different color sub-volumes simultaneously. Nodes in the same cabinet have different measurements and system matrix rows, but update and merge the same color sub-volumes simultaneously across different subsets. After merging, nodes in the same cabinet broadcast the merged sub-volume in x to the same color sub-volumes in the LR reconstructions. To fully utilize parallel cores and SIMD computations, cores of a node are mapped to computations for a SV update in a sub-volume, and SIMD vectors are mapped to computations for a voxel update in a manner similar to the NU-PSV algorithm. See the pseudo code in Algorithm 3 for details. The fusing operation across subsets is achieved through MPI all reduce operations among nodes in the same cabinet. In step 14, neighboring sub-volumes exchange sub-volumes' boundary voxels and use the boundary voxels in evaluating the prior model in Equation (7).

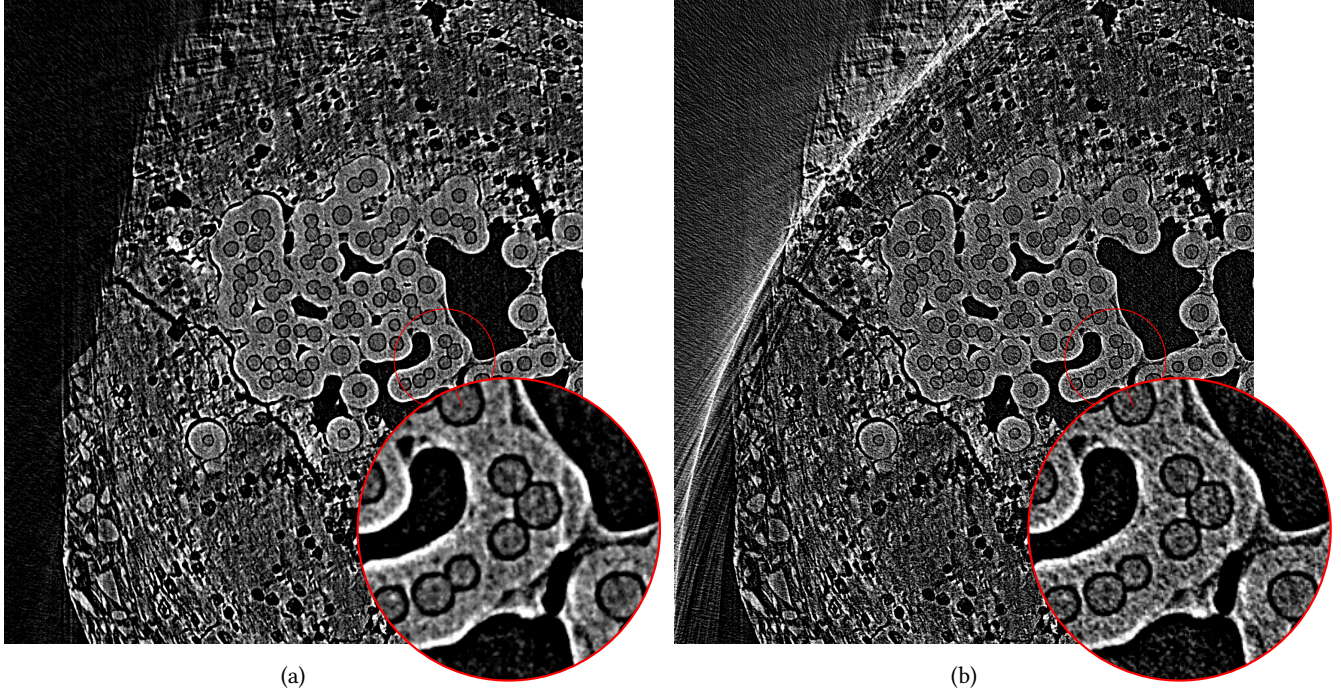


Figure 6: (a) A zoom-in image of a slice reconstructed by AC-MBIR with a 4.9 times super-resolution. (b) The same slice reconstructed by interpolation method with more noise and artifacts.

Algorithm 3 AC-MBIR

INPUT: The partitioned CT measurements and system matrix.
LOCAL: G : the number of subsets. v^i : the i^{th} LR reconstruction.
 $[v^i]^j$: the j^{th} sub-volume of v^i . $[v^i]^{j,k}$: the k^{th} SV of $[v^i]^j$.
Similar notations apply to w^i, u^i .
OUTPUT: x : the consensus HR solution.

- 1: Initialize x, v^i, w^i , and u^i
- 2: **while** $x \neq v^i, i = 1, 2, \dots, G$ **do**
- 3: **for** each subset i **do in parallel across clusters**
- 4: **for** each sub-volume $[v^i]^j \in v^i$, **do in parallel across nodes**
- 5: **for** each SV $[v^i]^{j,k} \in [v^i]^j$, **do in parallel across cores**
- 6: $[v^i]^{j,k} \leftarrow \operatorname{argmin}_{[v^i]^{j,k}} F^i([x]^{j,k} + [u^i]^{j,k})$
- 7: $[(w^i)^j]^{j,k} \leftarrow [w^i]^{j,k}$
- 8: $[w^i]^{j,k} \leftarrow (2[v^i]^{j,k} - [x]^{j,k} - [u^i]^{j,k})$
- 9: $[w^i]^{j,k} \leftarrow \rho[w^i]^{j,k} + (1 - \rho)[(w^i)^j]^{j,k}$
- 10: $[x]^{j,k} \leftarrow \operatorname{Average}([w^1]^{j,k}, \dots, [w^G]^{j,k})$,
through MPI all reduce among nodes in a cabinet
- 11: $[u^i]^{j,k} \leftarrow [x]^{j,k} - [w^i]^{j,k}$
- 12: **end for**
- 13: **end for**
- 14: Neighboring sub-volumes exchange boundary voxels
through MPI send/receive operations.
- 15: **end for**
- 16: **end while**

Note that AC-MBIR is a unique asynchronous consensus optimization for reducing communication overhead. The common approach for consensus optimization selects a master node to update the consensus solution while the remaining worker nodes update LR reconstructions [30]. This approach, however, has a large global communication overhead as worker nodes' updates are bounded by how fast the master node updates and vice versa. AC-MBIR, however, has asynchronous updates and much smaller communication overhead because each node is a worker node and the master node at the same time. For example in Figure 5(b), $P(1, 1)$ and $P(2, 1)$ are in the same cabinet updating the blue sub-volumes of the LR reconstructions, but also update, merge and broadcast the blue sub-volume in the consensus solution. Therefore, AC-MBIR reduces communication overhead by restricting inter-node communications to nodes in the same cabinet and avoids global communications. Experimental results in Section 4.1 show that the asynchronous updates keep AC-MBIR's communication overhead at a low level, with 16% of its runtime spent on communication overhead at 272 nodes (18496 cores) and 13% of its runtime spent on communication overhead at 1080 nodes (73440 cores).

4 RESULTS

All performance numbers are evaluated on two datasets, a Ceramic Matrix Composite (CMC) dataset, imaged under synchrotron at Lawrence Berkeley National Laboratory, and an Iron Hydroxide (FeOOH) dataset, imaged under X-ray microscope at Air Force Research Lab. The CMC is a composite material for the next-generation

| Grouped-view subsets | $G = 1$ (NU-PSV) | 2 | 4 | 8 | 16 |
|---------------------------------------|---------------------|-------|-------|-------|-------|
| Nodes | 108 | 216 | 432 | 864 | 1728 |
| Memory footprint per node (Gigabytes) | 58.0 | 31.4 | 18.3 | 11.2 | 8.5 |
| Memory footprint reduction | 1.0x | 1.8x | 3.2x | 5.2x | 6.8x |
| Run time (secs) | 1681.0 | 779.8 | 342.6 | 258.4 | 215.7 |
| Computation increase | 1.0x | 1.0x | 1.0x | 1.1x | 1.3x |
| Strong scaling efficiency | 100% | 108% | 123% | 89% | 63% |
| Bandwidth (petabytes/sec) | <0.1 | <0.1 | <0.1 | <0.1 | <0.1 |

(a) AC-MBIR with grouped-view subsets

| Uniform-view subsets | $G = 1$ (NU-PSV) | 2 | 4 | 8 | 16 |
|---------------------------------------|---------------------|-------|-------|-------|-------|
| Nodes | 108 | 216 | 432 | 864 | 1728 |
| Memory footprint per node (Gigabytes) | 58.0 | 33.4 | 21.3 | 15.0 | 11.4 |
| Memory footprint reduction | 1.0x | 1.7x | 2.7x | 3.9x | 5.1x |
| Run time (secs) | 1681.0 | 807.1 | 385.7 | 347.1 | 347.5 |
| Computation increase | 1.0x | 1.0x | 1.2x | 1.5x | 2.1x |
| Strong scaling efficiency | 100% | 104% | 131% | 91% | 63% |
| Bandwidth (petabytes/sec) | <0.1 | <0.1 | <0.1 | <0.1 | <0.1 |

(b) AC-MBIR with uniform sparse-view subsets

Table 1: CMC dataset experiments for AC-MBIR with 108 sub-volumes but with different numbers of subsets, G . The left table is for grouped-view subsets and the right table is for uniform sparse-view subsets. More subsets lead to a rapidly decreasing memory footprint per node and an increased number of operations. Note that the grouped-view partition with all numbers of subsets has less memory footprint, shorter runtimes, and fewer computations than the uniform sparse-view partition.

aviation turbine-engine blades and space shuttles engine components [2]. The FeOOH is a material under development for nanorods.

The computing platform is the Cori Knights Landing supercomputer at the National Energy Research Scientific Computing Center (NERSC) with 68 cores per node. More information on the hardware can be found in the artifact description appendix. All data are stored and loaded into memory as single-precision floating point except for the system matrix. The system matrix is pre-computed as single-precision floating point, but stored in the NERSC data filesystem in a sparse matrix format and quantized as 4-bit unsigned characters for minimal memory footprint. When performing reconstructions, the quantized system matrix is loaded into memory and converted back to single-precision floating point [3, 20].

All runtimes are based on the entire program execution minus the IO time, and are measured at least three times with the average taken. The memory footprint is measured as the allocated memory for the program. Communication overhead is measured through NERSC IPM tool. Both the peak computing speed in petaflops and the average bandwidth in petabytes per second are measured through Intel SDE tool. To compute the peak performance and the average bandwidth, we count the number of FLOPS (72% single precision, 28% double) as well as the total bytes read and written on 4 nodes, and then extrapolate the performance to 8640 nodes.

For fastest convergence, we determine the optimal values for ρ and σ empirically. We use $\sigma = 10$ for 2 subsets, $\sigma = 20$ for 4 subsets, $\sigma = 45$ for 8 subsets, and $\rho = 0.8$ for all experiments. To decouple the convergence rate from runtimes, all programs in tables 1, 2, and 3 are run with 16 full iterations, where a full iteration is a cycle through N voxel updates. We pick 16 full iterations because the Normalized Root Mean Square Error (NRMSE) at 16 iterations is measured to be less than 2%, compared with a fully converged reconstruction.

4.1 Ceramic Matrix Composite (CMC) Dataset

The CMC dataset has a measurement size M of $2560 \times 2160 \times 1024$ (2560 channels, 2160 slices, and 1024 view angles), and the highest

image resolution limited by the CT detector is 1609^2 voxels per slice. To demonstrate AC-MBIR’s super-resolution performance, we reconstruct at an image resolution 4.9 times higher than the hardware limit, with 3600^2 voxels per slice and 2160 slices. A zoom-in picture of an example reconstruction is shown in Figure 6. The left image is a region reconstructed by AC-MBIR, and the right image is the same region reconstructed by the FBP interpolation method, whose reconstruction has more image noise and artifacts than AC-MBIR’s. For example, the upper left corner of the FBP interpolation image has many braiding artifacts, while AC-MBIR image has no braiding artifacts. Therefore, material scientists can use AC-MBIR to better examine the micro-structure at super-resolution.

Table 1 summarizes AC-MBIR and NU-PSV’s performance with different numbers of subsets, G , but with the same number of sub-volumes (108 sub-volumes). The left table is for grouped-view subsets and the right table is for uniform sparse-view subsets. Nodes are configured such that the number of clusters equals to the number of subsets G , and the number of cabinets equals to 108. Row 1 of the tables shows different numbers of subsets and AC-MBIR with $G = 1$ is identical to the NU-PSV algorithm.

Row 2 is the number of nodes (68 cores per node), computed as the number of cabinets multiplied by the number of clusters, and AC-MBIR with 16 subsets has 16 times more scalability than NU-PSV. Row 3 is the memory footprint per node and row 4 is AC-MBIR’s memory footprint reduction over NU-PSV. Both the grouped-view and the uniform sparse-view partitions significantly reduce the memory footprint by 6.8 times and 5.1 times respectively, although the grouped-view partition has 25% more memory footprint reduction. In addition, both the grouped-view and the uniform sparse-view partitions reduce the memory footprint non-linearly with more subsets, as the number of redundant data in the subsets increases non-linearly with more subsets.

Row 5 is the program runtimes and NU-PSV’s runtime is 1681 seconds. In contrast, the grouped-view partition with 16 subsets has a smaller runtime of 215.7 seconds by reducing the memory pressure. The uniform sparse-view partition also has a much smaller

| Number of sub-volumes | 68 | 135 | 270 | 540 | 1080 | 2160 |
|---------------------------------------|-------|-------|-------|------|------|------|
| Nodes | 272 | 540 | 1080 | 2160 | 4320 | 8640 |
| Memory footprint per node (Gigabytes) | 20.7 | 17.4 | 15.7 | 14.9 | 14.5 | 14.3 |
| Memory footprint reduction | 2.8x | 3.3x | 3.7x | 3.9x | 4.0x | 4.1x |
| Run time (secs) | 558.0 | 271.7 | 147.2 | 84.4 | 51.4 | 40.0 |
| Computation increase | 1.0x | 1.0x | 1.0x | 1.0x | 1.0x | 1.0x |
| Strong scaling efficiency | 120% | 124% | 114% | 100% | 82% | 53% |
| Bandwidth (petabytes/sec) | <0.1 | <0.1 | 0.1 | 0.4 | 1.4 | 3.5 |

(a) AC-MBIR with grouped-view subsets

| Number of sub-volumes | 68 | 135 | 270 | 540 | 1080 | 2160 |
|---------------------------------------|-------|-------|-------|-------|-------|-------|
| Nodes | 272 | 540 | 1080 | 2160 | 4320 | 8640 |
| Memory footprint per node (Gigabytes) | 23.7 | 20.4 | 18.7 | 17.9 | 17.5 | 17.3 |
| Memory footprint reduction | 2.4x | 2.8x | 3.1x | 3.2x | 3.3x | 3.4x |
| Run time (secs) | 654.6 | 326.4 | 174.3 | 159.2 | 154.0 | 154.0 |
| Computation increase | 1.2x | 1.2x | 1.2x | 1.2x | 1.2x | 1.2x |
| Strong scaling efficiency | 122% | 124% | 116% | 63% | 33% | 16% |
| Bandwidth (petabytes/sec) | <0.1 | <0.1 | 0.1 | 0.3 | 0.5 | 1.1 |

(b) AC-MBIR with uniform sparse-view subsets

Table 2: CMC dataset experiments for AC-MBIR with 4 subsets but with different numbers of sub-volumes. The left table is for grouped-view subsets and the right table is for uniform sparse-view subsets. With more sub-volumes, the memory footprint per node decreases slowly and the number of operations stays the same. The grouped-view subsets with 1080 sub-volumes have a 17% smaller memory footprint, a 3.9 times faster runtime, and a 3.3 times higher strong scaling efficiency than the uniform sparse-view subsets.

runtime than NU-PSV with 347.5 seconds at 16 subsets. The runtime for the uniform sparse-view partition, however, is larger than that for the grouped-view partition because the uniform sparse-view partition has a larger memory footprint and more redundant data.

Row 6 is the increase in floating point operations, compared to those for NU-PSV. Floating point operations increase non-linearly with more subsets as more subsets incur more redundant data and operations on the data. With 16 subsets, the grouped-view partition has 38% less redundant data than the uniform sparse-view partition. Row 7 is the parallel strong scaling efficiency compared to NU-PSV ($G = 1$) and row 8 is the average bandwidth. The grouped-view and the uniform sparse-view subsets have similar strong scaling efficiencies with a super-linear speedup up to 4 subsets. From the strong scaling efficiency numbers, we can conclude that the grouped-view partition does not have a better strong scaling efficiency over the uniform sparse-view partition. The grouped-view partition, however, has smaller runtimes by reducing memory footprint and computations.

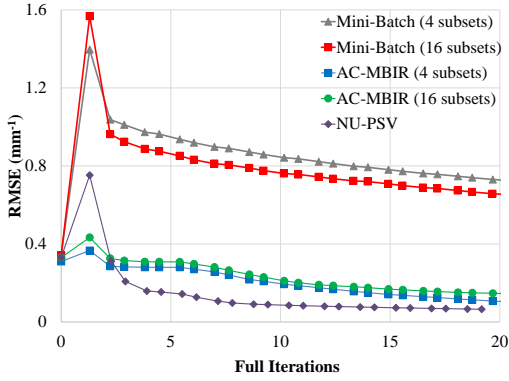
Table 2 summarizes AC-MBIR’s performance with the same number of view-subsets (4 view-subsets) but with different numbers of sub-volumes. With the same number of subsets, the memory allocation for the system matrix does not change across different numbers of sub-volumes. However, more sub-volumes reduce the number of slices and measurements per sub-volume. The memory footprint in row 3 thereby decreases slowly with more sub-volumes. In row 4, the grouped-view partition reduces more memory footprint than the sparse-view partition, as the former has fewer redundant data than the latter. In addition, the number of operations does not increase with more sub-volumes as the amount of redundant data remains the same for the same number of subsets.

Row 5 shows the runtimes and the grouped-view partition with 2160 sub-volumes (90% of Cori supercomputer with 587520 cores) has a 3.9 times speedup over the uniform sparse-view partition, a 46.5 times speedup over the NU-PSV algorithm, and has a measured peak computing performance of 1.5 petaflops at 8640 nodes. If not

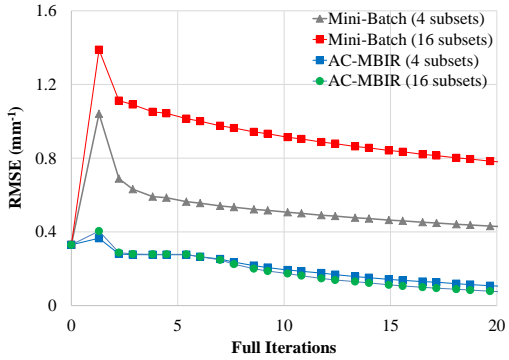
counting the MPI Allreduce overhead in step 10 of Algorithm 3, the peak performance number is 5.6 petaflops. The strong scaling efficiencies in row 7 show that the grouped-view partition has a 100% strong scaling efficiency at 2160 nodes, compared to NU-PSV. AC-MBIR has a high strong scaling efficiency for two reasons: (1) a smaller memory footprint leads to more data kept in cache, and (2) a modest communication overhead. The IPM tool shows that grouped-view subsets has a low communication overhead with 16% of its runtime spent on communication overhead at 272 nodes, 13% communication overhead at 1080 nodes, and an increased communication overhead of 73% at 4320 nodes. Without AC-MBIR’s asynchronous updates, the communication overhead is more than 50% at few numbers of nodes.

Row 8 is AC-MBIR’s average bandwidth and the grouped-view partition achieves a bandwidth of 3.5 petabytes per second at 8640 nodes, which is more than 35 times higher than that for NU-PSV. AC-MBIR has a good bandwidth because AC-MBIR’s runtimes are bounded by memory and most of the runtimes are spent on reorganizing sinusoidal data for high cache locality, regular data access pattern, and a low memory footprint. In contrast, AC-MBIR’s floating point operations are much fewer than its memory operations, leading to an impressive bandwidth of 3.5 petabytes per second, but a modest peak performance of 1.5 petaflops.

Convergence evaluations are another interesting topic. Figure 7 compares the convergence rate for the first 20 full iterations for AC-MBIR, NU-PSV and Mini-Batch Gradient Descent on an example CMC dataset sub-volume with 20 slices, with Figure 7(a) for the grouped-view partition and Figure 7(b) for the uniform sparse-view partition. For the grouped-view partition, the Mini-Batch Gradient Descent never converges to the global optimum even with more iterations, and converges to a solution far from the global optimum, whereas AC-MBIR converges to the global optimum rapidly both with 4 and 16 subsets. Compared to NU-PSV, AC-MBIR has a similar convergence rate at 20 full iterations. In addition, we can observe that all algorithms has a RMSE spike in early iterations



(a) Convergence for grouped-view subsets.



(b) Convergence for uniform sparse-view subsets.

Figure 7: (a) The CMC dataset convergence plot for Mini-Batch Gradient Descent and AC-MBIR with 4 and 16 grouped-view subsets. (b) The convergence plot for the two algorithms with 4 and 16 uniform sparse-view subsets.

from a convergence overshoot caused by the excessive updates. For the uniform sparse-view partition, Mini-Batch Gradient Descent with 4 subsets converges to a solution close to the global optimum, although its convergence rate is still much slower than AC-MBIR.

4.2 Iron Hydroxide (FeOOH) Dataset

The FeOOH dataset is much smaller than the CMC dataset, with a measurement size M of $1024 \times 1024 \times 900$ (1024 channels, 1024 slices and 900 view angles). The highest image resolution limited by the CT detector is 960^2 voxels per slice. We reconstruct at an image resolution 1.14 times higher than the hardware limit, with 1024^2 voxels per slice and 1024 slices. Table 3 summaries AC-MBIR and NU-PSV’s performance. The second row is NU-PSV’s runtimes and NU-PSV only scales to 1024 nodes, as each node needs to operate on a sub-volume with at least 1 slice. The third row is AC-MBIR’s runtimes with 4 grouped-view subsets and AC-MBIR has a 4 times higher scalability than NU-PSV. The fourth row is AC-MBIR’s speedup over NU-PSV with a 3.7 times speedup at 4096 nodes (278528 cores). Since NU-PSV cannot scale to 4096 nodes, the speedup at 4096 nodes is AC-MBIR’s runtime compared with NU-PSV’s best runtime at 1024 nodes. The fifth row is AC-MBIR’s

| Nodes | 16 | 64 | 256 | 1024 | 4096 |
|----------------------------|--------|-------|-------|------|------|
| NU-PSV runtimes (secs) | 1776.4 | 457.5 | 110.8 | 33.1 | NA |
| AC-MBIR runtimes (secs) | 556.2 | 142.2 | 47.1 | 18.0 | 10.1 |
| Speedup | 3.2x | 3.2x | 2.4x | 1.8x | 3.3x |
| Memory footprint reduction | 0.9x | 1.9x | 3.0x | 3.5x | 3.7x |

Table 3: Performance comparison for the FeOOH dataset between AC-MBIR (with 4 grouped-view subsets) and NU-PSV.

memory footprint reduction over NU-PSV and AC-MBIR’s memory footprint at 4096 nodes is 3.7 times smaller than that for NU-PSV.

5 DISCUSSION

AC-MBIR has a significant memory footprint reduction, a smaller runtime and more scalability over NU-PSV. Its parallelism across data subsets and sub-volumes, however, has advantages and disadvantages. The data subsets parallelism has better memory footprint reduction over the sub-volume parallelism, but has worse strong scaling efficiency and a slower runtime. The sub-volume parallelism has a better strong scaling efficiency and a faster runtime, but performs worse for memory footprint reduction. To achieve the best performance, node configurations between the two parallelisms have to be determined empirically as there is no theory that can accurately pre-determine the node configuration. One way for node configuration is to follow the ablation studies in Table 1 to find the number of subsets with best performance, while fixating on the same number of sub-volumes. Then, nodes are configured by following the studies in Table 2 to find the number of sub-volumes with best performance, while fixating on the same number of subsets.

6 IMPLICATIONS

AC-MBIR can be applied to supervised machine learning problems with the following form:

$$x = \operatorname{argmin}_x \{G(x) + R(x)\}, \quad (9)$$

where x is the output; $G(x)$ is a statistical model that estimates x from a limited size noisy training dataset, denoted as y , and $R(x)$ is a regularizer to prevent data over-fitting. If $G(x)$ is a weighted linear regression model, then Equation (9) has the same form as Equation (6), where M is the dataset size; A is a sparse system matrix; Λ is a weighting matrix; and $R(x)$ is an L_0 norm, L_1 norm or other regularizer.

If the output x is viewed as an image or a volume, then Equation (9) with a weighted linear regression training model is a cost function for image reconstruction problems, such as CT, PET, MRI, electron microscopy, synchrotron, neutron, proton, and ultrasound imaging. If x is a matrix of measurements or an object under test, then Equation (9) with a weighted linear regression training model is a cost function for compressive sensing problems and predicts the state of the physical world from noisy sensor measurements. Such applications include autonomous navigation, depth sensors, digital holography, geophysical sensing, radar, radio astronomy, crystallography, machine learning techniques such as the least absolute shrinkage and selection operator.

AC-MBIR partitions Equation (9) into sub-problems and merge the sub-problem solutions into a consensus one with a guaranteed convergence if Equation (9) is convex. This algorithm has two benefits: (1) reducing memory pressure by distributing data among sub-problems, and (2) mapping to a supercomputer for fast computations. To solve Equation (9), AC-MBIR non-uniformly partitions the dataset y into G data subsets and data in the same subset are highly similar and different subsets are least similar. Each subset has a new training model, $g_i(x)$, such that $G(x) = \sum_{i=1}^G g_i(x)$. Each subset solves the proximal map function formulated in Equation (8) and the individual solutions from all subsets are merged asynchronously into the consensus solution by following Algorithm 3.

A MATHEMATICAL PROOF FOR ALGORITHM 2

The consensus-equilibrium Method [4, 22] provides a parallel method to solve the following MBIR optimization problem:

$$x = \operatorname{argmin}_x \left\{ \frac{1}{2} \|y - Ax\|_{\Lambda}^2 + R(x) \right\} = \operatorname{argmin}_x \sum_{i=1}^G f^i(x), \quad (10)$$

where f^i has the same definition as in (7), and represents a sub-problem for the i^{th} measurement subset. G is the total number of subsets. Theorem 1 of [4] shows that finding the global minimum to Equation (10), x , is exactly the same as finding the average of $\{w^1, \dots, w^G\}$, denoted as \bar{w} , for the following system of equations:

$$\begin{bmatrix} F^1(w^1) \\ \vdots \\ F^N(w^G) \end{bmatrix} = \frac{1}{G} \begin{bmatrix} w^1 + w^2 + \dots + w^G \\ \vdots \\ w^1 + w^2 + \dots + w^G \end{bmatrix} = \begin{bmatrix} \bar{w} \\ \vdots \\ \bar{w} \end{bmatrix}, \quad (11)$$

where F^i is the proximal map function for the i^{th} measurement subset, as defined before in Equation (8). Then, to prove that the output of Algorithm 2 is the same as the solution to Equation (10), we need to demonstrate that the output of Algorithm 2 equals to \bar{w} .

If we let F represents the stacked function for $[F^1, \dots, F^G]^T$, w as the list, $[w^1, \dots, w^G]^T$, and function D stacks G copies of \bar{w} , then we can rewrite Equation (11) with a more compact short-hand notation:

$$F(w) = D(w), \quad (12)$$

After multiplying both sides of Equation (12) with a constant 2 and subtract them with w , we get

$$(2F - I)w = (2D - I)w \quad (13)$$

where I is the identity function. Note that for any w , applying function D twice to w yields the same result as applying D only once to w . So, $D^2w = Dw$, and consequently, $(2D - I)^2w = w$. Based on this special property, we apply function $(2D - I)$ to both sides of equation (13) and we obtain

$$(2D - I)(2F - I)w = w \quad (14)$$

To find w that satisfies Equation (14), w can be computed iteratively by finding the fixed-point to Equation (14). In every iteration, w is updated as below:

$$w \leftarrow \rho w + (1 - \rho)w', \quad (15)$$

where ρ is a convergence parameter chosen between 0 and 1 and w' stores a copy of w from the previous iteration. Note that Equation (15) is exactly the same as step 7 in Algorithm 2. After updates in Equation (15) are converged to a fixed-point vector, the final output, \bar{w} , is computed as the average for all elements in w , namely $\{w^1, \dots, w^G\}$.

ACKNOWLEDGMENTS

The authors would like to thank the National Energy Research Scientific Computing Center (NERSC) for providing supercomputing resources under contract No. DE-AC02-05CH11231. This research was supported by the National Science Foundation (NSF) under Award CCF-1763896. Additional support was provided by the DHS ALERT Center for Excellence under Grant Award 2013-ST-061-ED0001.

REFERENCES

- [1] S. Ahn and J. A. Fessler. 2003. Globally Convergent Image Reconstruction for Emission Tomography Using Relaxed Ordered Subsets Algorithms. *IEEE Transactions on Medical Imaging* 22, 5 (May 2003), 613–626. <https://doi.org/10.1109/TMI.2003.812251>
- [2] GE Aviation. 2016. Space Age Ceramics Are Aviation's New Cup Of Tea. <https://www.ge.com/reports/space-age-cmc-aviations-new-cup-of-tea/>. (2016).
- [3] T. Balke, S. Majeed, G. T. Buzzard, S. Poveromo, P. Howard, M. A. Groeber, J. McClure, and C. A. Bouman. 2018. Separable Models for cone-beam MBIR Reconstruction. In *Computational Imaging XVI, Burlingame, California, USA, 28 Jan 2018 - 1 Feb 2018*.
- [4] G. Buzzard, S. Chan, S. Sreehari, and C. Bouman. 2018. Plug-and-Play Unplugged: Optimization-Free Reconstruction Using Consensus Equilibrium. *SIAM Journal on Imaging Sciences* 11, 3 (2018), 2001–2020. <https://doi.org/10.1137/17M1122451> arXiv:<https://doi.org/10.1137/17M1122451>
- [5] H. Erdogan and J. Fessler. 1999. Ordered Subsets Algorithms for Transmission Tomography. *Physics in Medicine & Biology* 44(11) (1999).
- [6] J. A. Fessler, E. P. Ficaro, N. H. Clinthorne, and K. Lange. 1997. Grouped-Coordinate Ascent Algorithms for Penalized-Likelihood Transmission Image Reconstruction. *IEEE Transactions on Medical Imaging* 16(2) (1997).
- [7] J. A. Fessler and D. Kim. 2011. Axial Block Coordinate Descent (ABCD) Algorithm for X-ray CT Image Reconstruction. In *11th International Meeting on Fully Three-Dimensional Image Reconstruction in Radiology and Nuclear Medicine*.
- [8] M. C. Fonseca, B. H. S. Araujo, C. S. B. Dias, N. L. Archilha, D. P. A. Neto, E. Cavalheiro, H. Westfahl, A. J. R. da Silva, and K. G. Franchini. 2018. High-Resolution Synchrotron-Based X-ray Microtomography as a Tool to Unveil The Three-Dimensional Neuronal Architecture of The Brain. *Scientific Reports* 8, 1 (Aug 2018), 12074.
- [9] C. Fournier, F. Jolivet, L. Denis, N. Verrier, E. Thiebaut, C. Allier, and T. Fournel. 2017. Pixel Super-Resolution in Digital Holography by Regularized Reconstruction. *Applied Optics* 56, 1 (Jan 2017), 69–77. <https://doi.org/10.1364/AO.56.000069>
- [10] S. Ha and K. Mueller. 2018. A GPU-Accelerated Multivoxel Update Scheme for Iterative Coordinate Descent (ICD) Optimization in Statistical Iterative CT Reconstruction (SIR). *IEEE Transactions on Computational Imaging* 4, 3 (Sep. 2018), 355–365. <https://doi.org/10.1109/TCL.2018.2833622>
- [11] C. Jiang, Q. Zhang, R. Fan, and Z. Hu. 2018. Super-Resolution CT Image Reconstruction Based on Dictionary Learning and Sparse Representation. *Scientific Reports* 8(1) (2018).
- [12] P. Jin, C. A. Bouman, and K. D. Sauer. 2013. A Method for Simultaneous Image Reconstruction and Beam Hardening Correction. In *2013 IEEE Nuclear Science Symposium and Medical Imaging Conference (NSS/MIC)*. 1–5.
- [13] E. A. Kazerooni. 2001. High-resolution CT of the lungs. *American Journal of Roentgenology* 177, 3 (Sep 2001), 501–519.
- [14] A. Korn, M. Fenchel, B. Bender, S. Danz, T. K. Hauser, D. Ketelsen, T. Flohr, C. D. Claussen, M. Heuschmid, U. Ernemann, and H. Brodoefel. 2012. Iterative Reconstruction in Head CT: Image Quality of Routine and Low-Dose Protocols in Comparison with Standard Filtered Back-Projection. *American Journal of Neuroradiology* 33, 2 (Feb 2012), 218–224.
- [15] M. Li, T. Zhang, Y. Chen, and A. J. Smola. 2014. Efficient Mini-batch Training for Stochastic Optimization. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '14)*. ACM, New York, NY, USA, 661–670. <https://doi.org/10.1145/2623330.2623612>
- [16] X. Li, Y. Liang, W. Zhang, T. Liu, H. Li, G. Luo, and M. Jiang. 2018. cuMBIR: An Efficient Framework for Low-dose X-ray CT Image Reconstruction on GPUs.

- In *Proceedings of the 2018 International Conference on Supercomputing (ICS '18)*. ACM, New York, NY, USA, 184–194. <https://doi.org/10.1145/3205289.3205309>
- [17] P. Milanfar. 2010. *Super-Resolution Imaging*. CRC Press.
- [18] K. A. Mohan, S. V. Venkatakrisnan, J. W. Gibbs, E. B. Gulsoy, X. Xiao, M. D. Graef, P. W. Voorhees, and C. A. Bouman. 2015. TIMBER: A Method for Time-Space Reconstruction from Interlaced Views. *IEEE Transactions on Computational Imaging* 1, 2 (June 2015), 96–111.
- [19] Z. Nadir, M. S. Brown, M. L. Comer, and C. A. Bouman. 2017. A Model-Based Iterative Reconstruction Approach to Tunable Diode Laser Absorption Tomography. *IEEE Transactions on Computational Imaging* 3, 4 (Dec 2017), 876–890.
- [20] A. Sabne, X. Wang, S. J. Kisner, C. A. Bouman, A. Raghunathan, and S. P. Midkiff. 2017. Model-Based Iterative CT Image Reconstruction on GPUs. In *Proceedings of the 22Nd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '17)*. ACM, New York, NY, USA, 207–220. <https://doi.org/10.1145/3018743.3018765>
- [21] S. Sreehari, S. V. Venkatakrisnan, K. L. Bouman, J. P. Simmons, L. F. Drummy, and C. A. Bouman. 2016. Multi-Resolution Data Fusion for Super-Resolution Electron Microscopy. *CoRR abs/1612.00874* (2016). arXiv:1612.00874 <http://arxiv.org/abs/1612.00874>
- [22] V. Sridhar, G. T. Buzzard, and C. A. Bouman. 2018. Distributed Framework for Fast Iterative CT Reconstruction from View-subsets. In *Computational Imaging XVI, Burlingame, California, USA, 28 Jan 2018 - 1 Feb 2018*.
- [23] J. B. Thibault, K. D. Sauer, C. A. Bouman, and J. Hsieh. 2007. A Three-Dimensional Statistical Approach to Improved Image Quality for Multi-Slice Helical CT. *Medical Physics* 34(11) (2007).
- [24] S. V. Venkatakrisnan, L. F. Drummy, M. Jackson, M. D. Graef, J. Simmons, and C. A. Bouman. 2013. High Angle Annular Dark Field- Scanning Transmission Electron Microscope (HAADF-STEM) Tomography. *IEEE Transactions on Image Processing* 22(1) (2013).
- [25] X. Wang, K. A. Mohan, S. J. Kisner, C. A. Bouman, and S. P. Midkiff. 2016. Fast Voxel Line Update for Time-Space Image Reconstruction. In *The 41st IEEE International Conference on Acoustics, Speech and Signal Processing*.
- [26] X. Wang, A. Sabne, S. J. Kisner, A. Raghunathan, C. A. Bouman, and S. P. Midkiff. 2016. High Performance Model Based Image Reconstruction. *SIGPLAN Notice* 51, 8, Article 2 (Feb. 2016), 12 pages. <https://doi.org/10.1145/3016078.2851163>
- [27] X. Wang, A. Sabne, P. Sakdhnagool, S. J. Kisner, C. A. Bouman, and S. P. Midkiff. 2017. Massively Parallel 3D Image Reconstruction. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC '17)*. ACM, New York, NY, USA, Article 3, 12 pages. <https://doi.org/10.1145/3126908.3126911>
- [28] R. Yan, S. V. Venkatakrisnan, J. Liu, C. A. Bouman, and W. Jiang. 2019. MBIR: A Cryo-ET 3D Reconstruction Method that Effectively Minimizes Missing Wedge Artifacts and Restores Missing Information. *Journal of Structural Biology* (Mar 2019).
- [29] Z. Yu, J.-B. Thibault, C. A. Bouman, K. D. Sauer, and J. Hsieh. 2011. Fast Model-Based X-Ray CT Reconstruction using Spatially Nonhomogeneous ICD Optimization. *IEEE Transactions on Image Processing* 20(1) (2011).
- [30] R. Zhang and J. T. Kwok. 2014. Asynchronous Distributed ADMM for Consensus Optimization. In *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32 (ICML '14)*. JMLR.org, II–1701–II–1709.
- [31] J. Zheng, S. S. Saquib, K. Sauer, and C. A. Bouman. 2000. Parallelizable Bayesian Tomography Algorithms with Rapid, Guaranteed Convergence. *IEEE Transactions on Image Processing* 9(10) (2000).
- [32] M. Zinkevich, M. Weimer, L. Li, and A. J. Smola. 2010. Parallelized Stochastic Gradient Descent. In *Advances in Neural Information Processing Systems 23*. Curran Associates, Inc., 2595–2603. <http://papers.nips.cc/paper/4006-parallelized-stochastic-gradient-descent.pdf>