

Avoiding Communication in Logistic Regression

Aditya Devarakonda
Institute for Data Intensive Engineering and Science
Johns Hopkins University
Baltimore, Maryland
adi@jhu.edu

James Demmel
Department of EECS
Department of Mathematics
University of California, Berkeley
Berkeley, California
demmel@berkeley.edu

Abstract—Stochastic gradient descent (SGD) is one of the most widely used optimization methods for solving various machine learning problems. SGD solves an optimization problem by iteratively sampling a few data points from the input data, computing gradients for the selected data points, and updating the solution. However, in a parallel setting, SGD requires interprocess communication at every iteration. We introduce a new communication-avoiding technique for solving the logistic regression problem using SGD. This technique re-organizes the SGD computations into a form that communicates every s iterations instead of every iteration, where s is a tuning parameter. We prove theoretical flops, bandwidth, and latency upper bounds for SGD and its new communication-avoiding variant. Furthermore, we show experimental results that illustrate that the new Communication-Avoiding SGD (CA-SGD) method can achieve speedups of up to $4.97\times$ on a high-performance Infiniband cluster without altering the convergence behavior or accuracy.

Index Terms—Communication Avoidance, Logistic Regression, Stochastic Gradient Descent, Binary Classification.

I. INTRODUCTION

Optimization methods are at the core of many machine learning applications. For example, the areas of computer vision and natural language processing make use of large machine learning models that have been trained on vast amounts of data to enable computers to automatically classify images or translate speech to text. Much of the prediction power is derived from solving nonlinear optimization problems which often perform regression (linear and nonlinear) or classification (binary and multiclass). In order to solve these optimization problems, we often compute first- or second-order derivatives and incrementally update the solution until it converges. Such approaches to solving optimization problems [1]–[3] have been well-studied, however, with the rise of multi-core/multi-node processing these optimization methods must now be parallelized across multiple cores/nodes. As a result, studying and improving the parallel performance of these optimization methods is imperative and would impact many application areas.

In this paper, we will focus on the stochastic gradient descent (SGD) method [2] for solving binary classification problems using the logistic regression model. SGD solves the logistic regression problem by iteratively sampling a few data points from the input data, computing the gradient of the logistic loss function, and updating the solution. As a result, parallel variants of SGD require interprocess communication at

every iteration. On modern computing hardware, where communication cost often dominates computation cost, the running time of parallel SGD is often dominated by communication cost.

We model communication cost on a distributed-memory parallel cluster in terms of two costs: latency and bandwidth. Our goal is to show that the latency cost of SGD, which is often the dominant cost, can be improved by a tunable factor of s by trading a factor of sb additional bandwidth and computation, where b is a tunable batch size. *The main contributions of this paper are:*

- Derivation of a Communication-Avoiding SGD (CA-SGD) method for solving the logistic regression problem which reduces SGD latency cost by a tunable factor of s in exchange for a factor of sb additional bandwidth and computation.
- Theoretical analysis of the flops, bandwidth, and latency costs of SGD and CA-SGD under two input matrix partitioning schemes: 1D-block row and 1D-block column.
- Numerical experiments which illustrate that CA-SGD is numerically stable for very large values of s .
- Performance experiments which illustrate that CA-SGD can attain speedups of up to $4.97\times$ over SGD and can scale out to $4\times$ as many cores.

A. Logistic Regression

Logistic regression is a supervised learning model used to predict the probability of data points belonging to one of two classes (binary classification). This model is widely used in many applications like predicting disease risk, website click-through prediction, and fraud detection which often require classification of data in terms of two classes.

We will now briefly derive the optimization problem for logistic regression used in this paper. We begin by defining the logistic function:

$$\sigma(\theta) = \frac{e^\theta}{1 + e^\theta} \equiv \frac{1}{1 + e^{-\theta}}.$$

Now suppose we are given a dataset $A \in \mathbb{R}^{m \times n}$ with m data points (rows of A) and n features (columns of A) and a vector of labels (one label per data point), $y \in \mathbb{R}^m$ such that $y_i \in \{-1, +1\} \forall i = 1, \dots, m$. Given such a dataset, the goal is to compute a vector $x \in \mathbb{R}^n$ of weights for each feature that maximizes the probability of correctly classifying

the input data. Using the logistic function, we can model the probability for each data point as

$$P(y_i|a_i x) = \begin{cases} \sigma(a_i x) & y_i = +1 \\ 1 - \sigma(a_i x) & y_i = -1, \end{cases} \quad (1)$$

where a_i is the i -th data point (row) in A and x is the unknown vector of weights. In this paper we assume the labels are -1 and $+1$ because this label choice leads to fewer terms in the optimization problem. This results in a more concise derivation of communication-avoiding SGD (CA-SGD). Our results also hold for other mathematically equivalent formulations (i.e. labels that are 0 and $+1$, etc.) of the logistic regression problem.

Since $1 - \sigma(a_i x) = \sigma(-a_i x)$ by symmetry of the logistic function, (1) can be further simplified to,

$$P(y_i|a_i x) = \sigma(y_i a_i x).$$

From this the optimization problem can be defined as

$$\arg \max_x \prod_{i=1}^m \sigma(y_i a_i x) = \arg \max_x \prod_{i=1}^m \frac{1}{1 + \exp(-y_i a_i x)}, \quad (2)$$

which computes the weights, x , that maximizes the likelihood function. Finally, by taking the negative log of the likelihood, we can cast (2) into the form of empirical risk minimization,

$$\arg \min_x F(A, x, y), \quad (3)$$

$$\text{where } F(A, x, y) = \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i a_i x)).$$

Unlike linear regression, (3) does not have a closed-form solution and cannot be solved using direct methods (i.e. through matrix factorization). However, one approach to solving this problem is to update the solution iteratively using the gradient of (3) until the solution converges. The gradient with respect to x is given by

$$\nabla F(A, x, y) = \frac{1}{m} \sum_{i=1}^m \frac{-\tilde{a}_i^T}{1 + \exp(\tilde{a}_i x)}, \quad (4)$$

where $\tilde{a}_i = y_i a_i \forall i = 1, \dots, m$ (i.e. the rows of A scaled by their corresponding labels). In matrix form we will use the notation $\tilde{A} = A \circ y$, where \circ represents scaling the i -th row of A by the i -th element of y . For convenience we can rewrite (4) in matrix form as

$$\nabla F(A, x, y) = -\frac{1}{m} \tilde{A}^T \left(\vec{1} \oslash \left(\vec{1} + \exp(\tilde{A}x) \right) \right), \quad (5)$$

where \oslash is the elementwise division operation and $\exp(\cdot)$ is now the exponential function applied elementwise to the vector $\tilde{A}x$. For clarity we will use the notation $\text{sig}(\tilde{A}x) = \vec{1} \oslash \left(\vec{1} + \exp(\tilde{A}x) \right)$, which is the sigmoid function applied to the vector $\tilde{A}x$. We can then rewrite (5) as

$$\nabla F(A, x, y) = -\frac{1}{m} \tilde{A}^T \text{sig}(\tilde{A}x). \quad (6)$$

From (4) the solution vector, x_h , for iteration h can be obtained by,

$$x_h = x_{h-1} - \eta_h \nabla F(A, x_{h-1}, y), \quad (7)$$

where x_{h-1} is the solution vector from the previous iteration and η_h is the learning rate (or step size) at iteration h which determines by how much the solution moves in the $\nabla F(A, x_{h-1}, y)$ direction. This is the well-known gradient descent (GD) method for iteratively refining the solution, x_{h-1} , until it converges to the optimal solution. η_h is a tuning parameter which may drastically affect the convergence behavior of gradient descent.

Another method often used to solve the logistic regression problem is Stochastic Gradient Descent (SGD). Instead of using the entire matrix A to compute the gradient, $\nabla F(A, x_{h-1}, y)$, SGD computes the gradient for only a subset of the data points (rows) in A . Strictly speaking, SGD samples only 1 row of A at each iteration. However, we will generalize this to a tunable batch size of b rows sampled from the matrix A . The resulting update for SGD with batch size b becomes

$$x_h = x_{h-1} - \eta_h \nabla F(\mathbb{I}_h A, x_{h-1}, \mathbb{I}_h y), \quad (8)$$

where \mathbb{I}_h is a matrix in $\mathbb{R}^{b \times m}$ which corresponds to b rows sampled uniformly at random without replacement from the identity matrix, I_m . As a result, $\mathbb{I}_h A$ and $\mathbb{I}_h y$ simply select b rows of A and their corresponding b labels from y . Note that if $b = m$, SGD is equivalent to gradient descent except with the rows of A permuted every iteration. The resulting SGD algorithm is shown in Algorithm 1.

Algorithm 1 Stochastic Gradient Descent

Input: $A \in \mathbb{R}^{m \times n}$, $y \in \mathbb{R}^m$, η_0, b, H' .

Output: $x_H \in \mathbb{R}^n$

- 1: $x_0 = \vec{0}$, $\tilde{A} = A \circ y$
 - 2: **for** $h = 1, 2 \dots H'$ **do**
 - 3: choose $\{i_k \in [m] | k = 1, 2, \dots, b\}$ uniformly at random without replacement.
 - 4: $\mathbb{I}_h = [e_{i_1}, e_{i_2}, \dots, e_{i_k}, \dots, e_{i_b}]^T$ where $e_{i_k} \in \mathbb{R}^m$ is the k -th standard basis vector.
 - 5: $x_h = x_{h-1} + \frac{\eta_h}{m} \tilde{A}^T \mathbb{I}_h^T \text{sig}(\mathbb{I}_h \tilde{A} x_{h-1})$
 - 6: **end for**
 - 7: **return** x_H
-

Figure 1 compares the convergence behavior of GD and SGD for the best η_h setting on the a6a dataset from the LIBSVM [4] repository. We perform tuning of η_h offline and show the best setting for GD and SGD, respectively. Note that many strategies exist for finding optimal, static learning rates and recent results have also illustrated that adaptive learning rates work well for convex optimization methods. In this paper, we focus on introducing the communication-avoiding (CA) derivation and studying its numerical and performance characteristics. We leave the effects of various learning rate strategies on the CA technique for future work.

In Figure 1, we observe that GD ($\eta_h = 1$) and SGD ($\eta_h = 10$) on the a6a dataset. We can observe that SGD

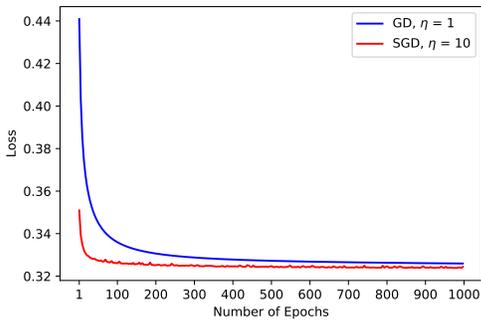


Fig. 1. Comparison of Gradient Descent (GD, blue) and Stochastic Gradient Descent (SGD, red) loss over 1000 epochs. We select the best learning rate for GD and SGD through offline tuning.

converges faster than GD over the 1000 epochs. This suggests that SGD is a better choice of algorithm for applications of logistic regression. Furthermore, the fast initial convergence of SGD suggests that it is a much better algorithm if a low-accuracy solution is sufficient.

When $b \ll m$, each iteration of SGD requires less computation than GD (by a factor of $\frac{m}{b}$) while one epoch of SGD performs the same amount of computation as one iteration (one epoch) of GD. In the distributed-memory parallel setting with A distributed across several processors each iteration of GD and SGD requires communication. Since SGD requires $\frac{m}{b}$ iterations to match GD, SGD requires a factor of $\frac{m}{b}$ more rounds of communication. On modern parallel hardware where communication is often the dominant cost, SGD requires orders of magnitude more communication than GD. This paper focuses on reducing the communication bottleneck in SGD without altering the convergence rate and behavior up to floating-point error.

II. RELATED WORK

Many techniques exist in literature which attempt to reduce the communication bottleneck in machine learning. For example, HOGWILD! [5] uses an asynchronous SGD method for the shared-memory setting where each thread computes gradients and updates the solution vector without synchronization. Due to the lack of synchronization, a thread may overwrite (and undo) the progress another thread has made. Convergence of HOGWILD! is not guaranteed, but will converge with high probability if the solution updates are sufficiently sparse and if there is bounded delay. In HOGWILD!, the latency bottleneck is reduced at the expense of convergence rate.

CoCoA [6], [7] is a general framework for reducing the synchronization cost of solving various machine learning problems in the distributed-memory setting. CoCoA reduces the synchronization cost by performing coordinate ascent on only the locally stored rows of A . After a tunable number of local iterations, the solution from each processor are sum-reduced (or averaged). If too many local iterations are performed, then global convergence will be slow. As with HOGWILD!, the reduction in latency (by deferring communication) comes at the

expense of convergence rate. In contrast, our approach does not alter the convergence rate of SGD. Instead, we introduce a tunable communication-avoiding parameter, s , that trades off additional computation and bandwidth in order to reduce latency by a factor of s . This means that if latency is the dominant cost in SGD, then we can reduce it by a factor of s and attain s -fold speedup with our new CA-SGD method.

Our technique is closely related to the one introduced in s -step and communication-avoiding Krylov (CA-Krylov) methods [8]–[13]. The s -step and CA-Krylov methods work showed that the recurrence relations in Krylov methods can be unrolled by a tunable factor of s and the remaining computation rearranged to avoid synchronization cost in the distributed-memory parallel setting. While the new methods have been shown to be faster, they suffered from numerical instability which subsequent work addressed by introducing techniques to improve numerical stability [9], [14]–[16].

The same recurrence unrolling technique has been shown to be effective for primal/dual coordinate and block coordinate descent methods and quasi-Newton’s method for solving ridge regression, LASSO, and SVM [17]–[21]. This paper extends prior results by illustrating that the technique works for solving the logistic regression problem using SGD where the loss function is nonlinear instead of linear (linear/ridge regression) or piecewise linear (LASSO, SVM).

Unlike CA-Krylov methods, our CA-SGD method does not exhibit any numerical instability even for very large values of s . This allows CA-SGD to simply select the value of s that balances the additional computation and bandwidth with the reduction in latency. The prior work on primal and dual block coordinate descent focused on piecewise linear problems. The piecewise linearity ensures that the distributive property can be applied in order to simplify the communication-avoiding derivation. However, this is not true for the logistic regression problem which requires computation of $\text{sig}(\tilde{A}x)$ for the gradient. We will show in Section III that this issue can be managed by making use of additional memory and communication properties of matrix-vector multiply (Lemmas IV.1 and IV.2) for the 1D-block column partitioned algorithm. The 1D-block row partitioned algorithm, however, requires an entirely new approach in order to reduce latency by a factor of s . This new approach for the 1D-block row partitioned case is described in Section IV.

III. DERIVATION

The Stochastic Gradient Descent (SGD) method is defined by the solution update,

$$x_h = x_{h-1} - \eta_h \nabla F(\mathbb{I}_h \tilde{A}, x_h, y),$$

where b is the batch size and $\mathbb{I}_h \in \mathbb{R}^{b \times m}$ is a matrix that contains b rows sampled uniformly at random without replacement from the m -dimensional identity matrix, I_m . We will use the following form in our derivation for the communication avoiding variant,

$$x_h = x_{h-1} + \frac{\eta_h}{m} \tilde{A}^T \mathbb{I}_h^T \text{sig}(\tilde{A}x_{h-1}). \quad (9)$$

Algorithm 2 Communication-Avoiding SGD

Input: $A \in \mathbb{R}^{m \times n}, y \in \mathbb{R}^m, \eta_0, b, s, H'$.

Output: $x_H \in \mathbb{R}^n$

```
1:  $x_0 = \vec{0}, \tilde{A} = A \circ y$ 
2: for  $h = 0, 2, \dots, \frac{H'}{s}$  do
3:   for  $j = 1, 2, \dots, s$  do
4:     choose  $\{i_k \in [m] | k = 1, 2, \dots, b\}$  uniformly at
       random without replacement.
5:      $\mathbb{I}_{sh+j} = [e_{i_1}, e_{i_2}, \dots, e_{i_k}, \dots, e_{i_b}]^T$  where  $e_{i_k}$  is the
        $k$ -th standard basis vector.
6:   end for
7:   Let  $Y = \begin{bmatrix} \mathbb{I}_{sh+1} \\ \mathbb{I}_{sh+2} \\ \vdots \\ \mathbb{I}_{sh+s} \end{bmatrix} \tilde{A}$ 
8:    $G = YY^T$ 
9:    $r = Yx_{sh}$ 
10:  for  $j = 1, 2, \dots, s$  do
11:    Update  $x_{sh+j}$  according to eq. (11).
12:  end for
13: end for
14: return  $x_{H'}$ 
```

By unrolling the recurrence, we can write x_{h+1} in terms of x_{h-1} ,

$$\begin{aligned} x_{h+1} &= x_{h-1} + \frac{\eta_h}{m} \tilde{A}^T \mathbb{I}_h^T \text{sig} \left(\mathbb{I}_h \tilde{A} x_{h-1} \right) \\ &\quad + \frac{\eta_{h+1}}{m} \tilde{A}^T \mathbb{I}_{h+1}^T \\ &\quad \text{sig} \left(\mathbb{I}_{h+1} \tilde{A} x_{h-1} + \frac{\eta_h}{m} \mathbb{I}_{h+1} \tilde{A} \tilde{A}^T \mathbb{I}_h^T \text{sig} \left(\mathbb{I}_h \tilde{A} x_{h-1} \right) \right). \end{aligned} \quad (10)$$

Note that the term $\text{sig} \left(\mathbb{I}_h \tilde{A} x_{h-1} \right)$ is used twice, once in (9) and then again to correct the gradient for x_{h+1} in (10). Since $\text{sig} \left(\mathbb{I}_h \tilde{A} x_{h-1} \right)$ has already been computed from (9), it can be reused in future solution updates. As a result, the recurrence unrolled solution updates still require only one $\text{sig}(\cdot)$ computation per solution update. This is important since the exponential operation is more expensive than typical arithmetic operations.

For convenience we will change the loop iteration counter from h to $sh + j$ where $0 \leq h < H'$ is the outer iteration counter (where communication occurs) and $1 \leq j \leq s$ is the inner iteration counter (where a sequence of s solution vectors are computed). By induction we can show that

$$\begin{aligned} x_{sh+j} &= x_{sh} + \sum_{i=1}^{j-1} \frac{\eta_{sh+i}}{m} \tilde{A}^T \mathbb{I}_{sh+i}^T \text{sig} \left(\mathbb{I}_{sh+i} \tilde{A} x_{sh+i} \right) \\ &\quad + \frac{\eta_{sh+j}}{m} \tilde{A}^T \mathbb{I}_{sh+j}^T \text{sig} \left(\mathbb{I}_{sh+j} \tilde{A} x_{sh} \right. \\ &\quad \left. + \sum_{i=1}^{j-1} \frac{\eta_{sh+i}}{m} \mathbb{I}_{sh+j} \tilde{A} \tilde{A}^T \mathbb{I}_{sh+i}^T \text{sig} \left(\mathbb{I}_{sh+i} \tilde{A} x_{sh+i} \right) \right). \end{aligned} \quad (11)$$

Note that we omit the expansion of x_{sh+i} in (11) for clarity and will show how it is handled in Section IV. The resulting CA-SGD algorithm is shown in Algorithm 2.

IV. ANALYSIS OF ALGORITHMS

Note that CA-SGD requires the matrix-matrix multiplications $\mathbb{I}_{sh+j} \tilde{A} \tilde{A}^T \mathbb{I}_{sh+i}^T$ in addition to the matrix-vector multiplications using $\tilde{A}^T \mathbb{I}_{sh+i}^T \text{sig}(\cdot)$ and $\mathbb{I}_{sh+i} \tilde{A} x_{sh+i}$. Unlike prior work on ridge regression, LASSO, and SVM [18]–[21], the nonlinear vector operation $\text{sig}(\cdot)$ prevents simplification of (11). We will show in this section that CA-SGD can reduce the latency cost despite the nonlinearity under two data partitioning schemes (1D-block column and 1D-block row). Due to the nonlinearity in (11) we will rely on the communication properties of the distributed matrix-vector products $v = Ax$ and $w = A^T v$ in order to avoid communication in the x_h update step.

Lemma IV.1. *Given a matrix $A \in \mathbb{R}^{m \times n}$ stored in 1D-block column format across p processors and a vector $x \in \mathbb{R}^n$ distributed across the p processors, the matrix vector product $v = Ax$ with v replicated on all processors requires $O(m)$ words moved and $O(\log p)$ messages.*

Proof: Computing Ax requires that each row of A be multiplied by x . With the given partitioning, each processor can multiply the portion of row elements stored locally with the corresponding locally stored elements of x . Each processor produces a partial vector $v^{(i)} \in \mathbb{R}^m \forall i \in \{1, 2, \dots, p\}$ s.t. $v = \sum_{i=1}^p v^{(i)}$. Computing $\sum_{i=1}^p v^{(i)}$ and replicating v on all processors requires one all-reduce with summation which costs $O(m)$ words moved and $O(\log p)$ messages. Note that the MPI implementation makes runtime decisions about the optimal routing algorithm based on message-size and number of processors [22]. This proof selects bandwidth and latency bounds with the lowest latency cost. ■

Lemma IV.2. *Given a matrix $A \in \mathbb{R}^{m \times n}$ stored in 1D-block column format across p processors and a vector $v \in \mathbb{R}^m$ replicated on all p processors, the matrix vector product $w = A^T v$ with w distributed across the p processors does not require communication.*

Proof: A similar cost analysis to Lemma IV.1 proves this lemma. ■

From Lemmas IV.1 and IV.2 we can see that in (11) computing $\mathbb{I}_{sh+i} \tilde{A} x_{sh+i}$ requires communication, whereas apply the sigmoid function and computing $\tilde{A}^T \mathbb{I}_{sh+j}^T \text{sig} \left(\mathbb{I}_{sh+i} \tilde{A} x_{sh+i} \right)$ does not.

We will now analyze the computation and communication costs of SGD and CA-SGD. We assume that A is sparse with nonzeros distributed uniformly between the rows and that the vectors y and x are dense. We will use the notation fmn to refer to the number of nonzeros in A , where $0 < f \leq 1$. This allows us to bound the number of nonzeros of $\mathbb{I}_h \tilde{A}$ by $\text{nnz}(\mathbb{I}_h \tilde{A}) = fbn$. Furthermore, we will also assume that the nonzeros are distributed uniformly between the processors. Note that logistic regression requires an $\exp(\cdot)$ operation

and element-wise division on b -dimensional vectors. This operation requires more floating-point operations to compute than typical arithmetic operations. We will model this by introducing the parameter ω to represent the cost of a single $\text{sig}(\cdot)$ operation. The elementwise $\text{sig}(\cdot)$ operation on a b -dimensional vector, as a result, costs ωb flops.

Theorem IV.3. *Given a matrix $A \in \mathbb{R}^{m \times n}$ stored in 1D-block column layout on p processors, labels $y \in \mathbb{R}^m$ replicated on all processors, and $x \in \mathbb{R}^n$ partitioned across p processors, H iterations of SGD (Alg. 2) with batch size b requires $O\left(H \frac{fbn}{p} + H \frac{n}{p} + H\omega b\right)$ flops, $O(Hb)$ words moved, and $O(H \log p)$ messages sent.*

Proof: Each iteration of SGD requires computation of $\mathbb{I}_h \tilde{A} x_{h-1}$ which costs $\frac{fbn}{p}$ flops and produces a b -dimensional vector on each processor which must be summed. The all-reduce with summation requires b words moved and $\log p$ messages. Computing $\text{sig}(\mathbb{I}_h \tilde{A} x_{h-1})$ costs ωb flops and no communication. The matrix-vector product $\frac{n_h}{m} \tilde{A}^T \mathbb{I}_h^T \text{sig}(\mathbb{I}_h \tilde{A} x_{h-1})$ costs $\frac{fbn}{p}$ multiplications and does not require communication (from Lemma IV.2). Finally, updating x_h costs $\frac{n}{p}$ flops and does not require any communication. Multiplying each cost by H gives the results of this proof. ■

We will now show that our new CA-SGD algorithm can asymptotically reduce the communication cost.

Theorem IV.4. *Given a matrix $A \in \mathbb{R}^{m \times n}$ stored in 1D-block column layout on p processors, labels $y \in \mathbb{R}^m$ replicated on all processors, and $x \in \mathbb{R}^n$ partitioned across p processors, H iterations of CA-SGD (Alg. 3) with batch size b requires $O\left(H \frac{f^2 s b^2 n}{p} + H \frac{n}{p} + H s b^2 + H\omega b\right)$ flops, $O(H s b^2)$ words moved, and $O\left(\frac{H}{s} \log p\right)$ messages sent.*

Proof: Each iteration of CA-SGD begins by computing

the matrix vector product $\begin{bmatrix} \mathbb{I}_h \\ \mathbb{I}_{h+1} \\ \dots \\ \mathbb{I}_{h+s} \end{bmatrix} \tilde{A} x_{h-1}$. This computation requires $O\left(\frac{f s b n}{p}\right)$ flops. In addition to the matrix vector prod-

uct, the Gram matrix $\begin{bmatrix} \mathbb{I}_h \\ \mathbb{I}_{h+1} \\ \dots \\ \mathbb{I}_{h+s} \end{bmatrix} \tilde{A} \tilde{A}^T \begin{bmatrix} \mathbb{I}_h^T & \mathbb{I}_{h+1}^T & \dots & \mathbb{I}_{h+s}^T \end{bmatrix}$

must be computed. This costs $O\left(\frac{f^2 s^2 b^2 n}{p}\right)$ when computing pair-wise inner products¹. There are $s^2 b^2$ possible inner products and each costs $O(f^2 n)$ flops. Once the partial matrix-vector products and Gram matrices are computed, an all-reduce with summation is required to combine the partial products. This communication requires $s^2 b^2 + s b$ words moved and $\log p$ messages. Then, we can compute s gradient vectors each of which requires a b -dimensional elementwise $\text{sig}(\cdot)$ operation. This costs $\omega s b$ flops and no communication. In order to complete the gradient computation, a matrix vector

product with blocks of \tilde{A}^T are required which costs $O\left(\frac{f s b n}{p}\right)$ flops and no communication (from Lemma IV.2). Note that after the first gradient is computed, the subsequent $s - 1$ gradients require additional computation in order to correct for missed solution updates. This additional computation requires $s^2 b^2$ flops (there are $\frac{(s-2)(s-1)}{2}$ total matrix vector products). Once all gradients have been computed, the solution vector x_{sh+s} can be computed by taking a sum over all s gradients which requires $O\left(\frac{sn}{p}\right)$ flops. Unlike SGD, each iteration of CA-SGD computes s gradients. Therefore, $\frac{H}{s}$ outer iterations of CA-SGD are required to perform the equivalent of H SGD iterations. Multiplying the costs by $\frac{H}{s}$ gives the results of this proof. ■

We will now assume that A is stored in 1D-block row layout, y is distributed across the processors, and x is replicated on all processors. Note that under this partitioning scheme choosing b rows of A uniformly at random may cause load imbalance², so assume that each processor will chose an equal number of local rows (i.e. $b \geq p$ s.t. $b/p \in \mathbb{Z}^+$). This variant of SGD interpolates between sequential SGD when $p = 1$ and GD when $p = m$. Note that this variation has not been discussed in prior work [18]–[21].

Theorem IV.5. *Given a matrix $A \in \mathbb{R}^{m \times n}$ stored in 1D-block row layout on p processors, labels $y \in \mathbb{R}^m$ distributed across all processors, and $x \in \mathbb{R}^n$ replicated on all processors, H iterations of SGD (Alg. 2) with batch size b requires $O\left(H \frac{f b n}{p} + H n + H \frac{\omega b}{p}\right)$ flops, $O(H n)$ words moved, and $O(H \log p)$ messages sent.*

Proof: Each iteration of SGD requires computation of $\mathbb{I}_h \tilde{A} x_{h-1}$ which costs $O\left(\frac{f b n}{p}\right)$ flops since each processor selects $\frac{b}{p}$ rows from locally stored data. Computing $\text{sig}(\mathbb{I}_h \tilde{A} x_{h-1})$ requires each processor to perform the $\text{sig}(\cdot)$ operation on a $\frac{b}{p}$ -dimensional vector which costs $\frac{\omega b}{p}$ flops. The matrix-vector product $\frac{n_h}{m} \tilde{A}^T \mathbb{I}_h^T \text{sig}(\mathbb{I}_h \tilde{A} x_{h-1})$ costs $O\left(\frac{f b n}{p}\right)$ flops and requires an all-reduce with summation. The all-reduce communicates n words using $\log p$ messages. Finally updating x_h costs n flops and no communication since all processors have a copy of x_{h-1} and a copy of the gradient. Multiplying by the number of iterations H provides the results of this proof. ■

From this proof we can observe that SGD with 1D-block row layout also requires one round of communication for each iteration. We will now prove the computation and communication cost of CA-SGD with 1D-block row layout.

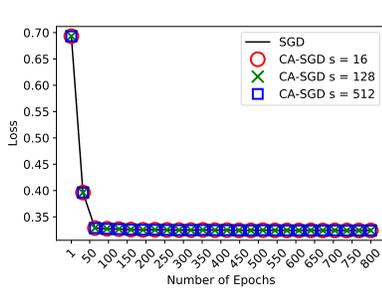
Theorem IV.6. *Given a matrix $A \in \mathbb{R}^{m \times n}$ stored in 1D-block row layout on p processors, labels $y \in \mathbb{R}^m$ distributed across all processors, and $x \in \mathbb{R}^n$ replicated on all processors, H iterations of CA-SGD (Alg. 3) with batch size b requires $O\left(H \frac{f^2 s b^2 n}{p} + H n + H \frac{s b^2}{p} + H \frac{\omega b}{p}\right)$ flops, $O(H f b n + H b + H \frac{n}{s})$ words moved, and $O\left(\frac{H}{s} \log p\right)$ messages sent.*

²Some processors may have more than the average number of rows chosen from locally stored data while other may have no rows selected.

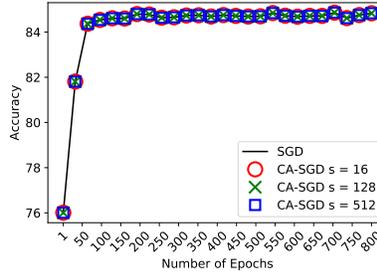
¹Note that the diagonal blocks of the Gram matrix are not required.

TABLE I
PROPERTIES OF THE LIBSVM DATASETS FOR NUMERICAL EXPERIMENTS

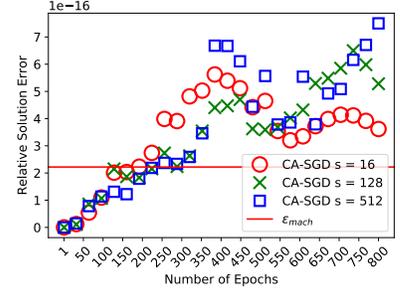
Name	m	n	$nnz(A)$	$nnz(A)/(mn)$	σ_{\max}	$\sigma_{\min} > 0$
a6a	11,220	123	155,608	0.1137	47.9015	0.9972
Mushrooms	8,124	112	170,604	0.1875	289.8993	1.2841
w7a	24,692	300	288,148	0.0389	9.8112	0.6846



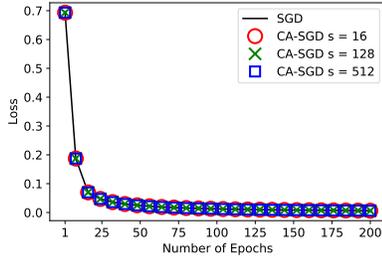
(a) a6a Loss vs. Epochs



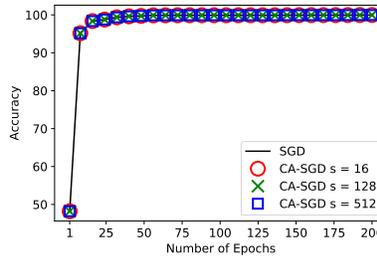
(b) Training Accuracy vs. Epochs



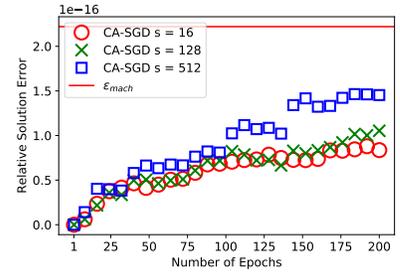
(c) Relative Solution Error vs. Epochs



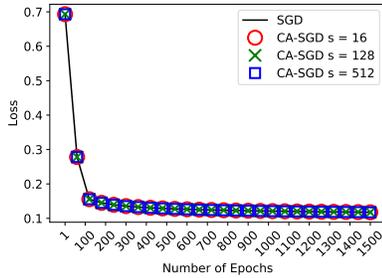
(d) Mushrooms Loss vs. Epochs



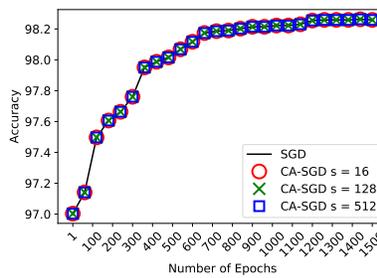
(e) Training Accuracy vs. Epochs



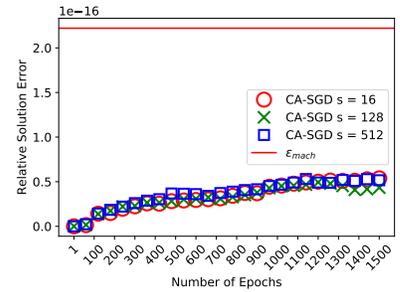
(f) Relative Solution Error vs. Epochs



(g) w7a Loss vs. Epochs



(h) Training Accuracy vs. Epochs



(i) Relative Solution Error vs. Epochs

Fig. 2. Comparison of SGD and CA-SGD convergence behavior on the a6a, mushrooms, and w7a (see Table I) for various values of s . The loss function, relative solution error, and training accuracy are reported over 100 epochs of training.

Proof: Each iteration of CA-SGD requires computation

$$\text{of } Yx_{sh} \text{ where } Y = \begin{bmatrix} \mathbb{I}_{sh+1} \\ \mathbb{I}_{sh+2} \\ \dots \\ \mathbb{I}_{sh+s} \end{bmatrix} \tilde{A} \text{ which costs } O\left(\frac{f s b n}{p}\right) \text{ flops}$$

and no communication. The resulting sb -dimensional vector is partitioned across p processors. We must also compute the Gram matrix $G = YY^T$. Since Y is 1D-block row partitioned across p processors with each processor storing $\frac{sb}{p}$ rows, computing G requires communication. By using an all-gather routine [22], each processor can obtain the necessary rows of Y from other processors in order to compute G .

In addition to communicating Y we also communicate the vector Yx_{sh} . This costs $O\left(\frac{f^2 s^2 b^2 n}{p}\right)$ flops and communicates $O((s-1)fbn + sb)$ words using $\log p$ messages³. Each processor computes blocks of G such that G is stored in 1D-block cyclic row layout with each processor storing $\frac{b}{p}$ consecutive rows in each $b \times b$ upper triangular block of G . We can now compute all s of the $\text{sig}(\cdot)$ vectors shown in (11) which requires $\frac{s^2 b^2 + \omega s b}{p}$ flops on each processor (the resulting b -

³If the message size for the all-gather is large, then many MPI implementations will switch to a 1D-ring routing algorithm. CA-SGD reduces the latency cost by a factor of s in both situations.

dimensional vectors are partitioned across all processors). We can now compute the gradient by multiplying the local $\frac{sb}{p}$ columns of \tilde{A}^T with the local $\frac{sb}{p}$ elements of the $\text{sig}(\cdot)$ vectors. This costs $\frac{fsbn}{p}$ flops and no communication. The resulting s gradient vectors on each processor are n -dimensional and must be sum-reduced in order to update the solution vector. The all-reduce communicates n words⁴ using $\log p$ messages. Once all processors have a copy of the s gradient vectors, they can be summed to obtain the update to the solution vector, which costs sn flops. Combining the costs results in $O\left(\frac{f^2s^2b^2n}{p} + sn + \frac{s^2b^2}{p} + \frac{\omega sb}{p}\right)$ flops and communicates $O(fsbn + sb + n)$ words in $O(\log p)$ messages per outer iteration. Multiplying the per iteration costs by $\frac{H}{s}$ gives the results of this proof. ■

These proofs show that CA-SGD reduces the latency cost by a tunable factor of s at the expense of additional bandwidth and computation cost. If latency is the dominant cost, then CA-SGD can attain s -fold speedup over SGD. This section primarily focuses on 1D layouts of A , however, 2D layouts of A might yield better performance for large, (nearly) square matrices. In this setting we can combine a divide-and-conquer local SGD algorithm with our (1D-column layout) CA-SGD method. We leave the analysis, implementation, and performance comparison of this 2D layout variant for future work.

V. EXPERIMENTAL RESULTS

Prior work on CA-Krylov methods [8], [9], [14]–[16] illustrated that applying the CA-technique can result in numerical instability due to the additional computation and rearrangement of the solution updates. This held true for even modest values of s and required development of residual replacement and orthogonal basis functions. However, unlike prior work, we will show that CA-SGD is numerically stable for very large values of s . Then we will show the practical performance tradeoffs and scaling properties of CA-SGD with an MPI implementation targeting a high-performance Infiniband cluster.

A. Numerical Experiments

We will now show how the convergence behavior of SGD compares to CA-SGD as s is varied. The datasets used in the experiments are binary classification problems obtained from the LIBSVM repository [4]. Table I summarizes properties of the datasets tested in this section. The SGD and CA-SGD methods have been implemented in Python using NumPy for linear algebra subroutines. Figure 2 illustrates the loss function convergence, relative solution error, and training accuracy of SGD and CA-SGD. From Figures 2a, 2d, and 2g we can observe that CA-SGD converges at the same rate as SGD and to the same final loss value for all datasets and values of s up to 512. Similarly, Figures 2b, 2e, and 2h show that CA-SGD attains the same training accuracy as SGD for all

⁴Note that we perform local summations on the s gradient vectors to obtain one partially summed vector that is subsequently sum-reduced across processors. However, doing so means we can only obtain the final solution vector x_{sh+s} and not the intermediate solutions $x_{sh+j} \forall j = 1, 2, \dots, s-1$.

values of s . The loss convergence and training accuracy figures experimentally validate that CA-SGD is simply a mathematical reformulation and computes the same sequence of partial solutions (up to floating-point error) as SGD.

The next set of experiments aim to quantify the floating-point error in the partial solutions computed by CA-SGD when compared to SGD. To show this, we will plot the relative solution error of CA-SGD with respect to SGD. For these experiments relative solution error is defined as, $\|x_h - x'_h\|_2 / \|x_h\|_2$, where h is the epoch of training, x_h is the SGD solution vector, and x'_h is the CA-SGD solution vector. We will also plot the machine precision, ϵ_{mach} , of the target computer as a reference line. Figures 2c, 2f, and 2i illustrate the results of this experiment. For the mushrooms and w7a datasets, we observe that the relative solution error is below machine precision over all epochs of training, which means there is negligible accumulation of floating-point error from CA-SGD for all values of s that were tested. The relative solution error for the a6a dataset is greater than machine precision but only by a constant factor and is still accurate up to 15-digits. For all datasets we observe that as s increases, the relative solution error also increases. This is to be expected since large values of s require computation of larger Gram matrices and additional matrix-vector products. However, despite the additional computation we see that CA-SGD is numerically stable.

For other datasets that have similar singular value spread, CA-SGD is likely to remain numerically stable. Furthermore, if techniques like data normalization and regularization are incorporated into the logistic regression model, the datasets become more well-conditioned. As a result, we do not expect CA-SGD to exhibit numerical instability for most practical applications and desired accuracies. In addition, numerical analysis of CA-SGD would be helpful in provide bounds on error accumulation.

B. Performance Experiments

To show the practical tradeoffs between SGD and CA-SGD, we implement both algorithms in C++ with MPI for parallel processing. The input matrix is stored in CSR 3-array format and partitioned in 1D-block column layout. Since SGD and CA-SGD only operate on b and sb rows of A , respectively, we reimplemented a subset of sparse BLAS-1 and BLAS-2 so that they operate only on the sampled rows. This avoids the overhead of explicitly copying sampled rows of A into a buffer at every iteration. In addition, we implement a new sparse BLAS-3 kernel to compute the Gram matrix required by CA-SGD. Since the Gram matrix is symmetric, we only compute and store the upper-triangular portion.

The experiments were performed on a high-performance cluster provided by the Maryland Advanced Research Computing Center. The compute nodes are dual-socket 2.5GHz Intel Haswell 12-core processors (24 cores per node) which are interconnected by an Infiniband network using a fat-tree topology. The code is built with the Intel 18.0.3 C++ compiler and OpenMPI 3.1 [23]. We experimented with hybrid OpenMP

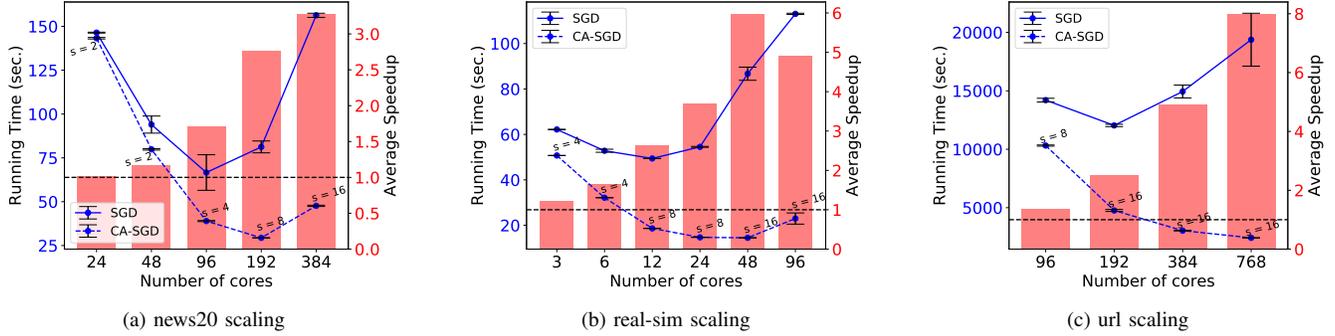


Fig. 3. Strong scaling comparison between SGD ($b = 1$) and CA-SGD on the news20, real-sim and url datasets (see Table II). We report the mean running times (blue points) and the standard deviation (error bars) over 5 trials. CA-SGD results are annotated with the value of s that achieved the fastest running times.

TABLE II
PROPERTIES OF THE LIBSVM DATASETS FOR PERFORMANCE EXPERIMENTS

Name	m	n	$nnz(A)$
news20	19,996	1,355,191	1,674,113
real-sim	72,309	20,958	3,709,083
url	2,396,130	3,231,961	277,058,644

and MPI configurations but found that flat MPI performed best. Table II summarizes the LIBSVM [4] datasets used for the performance experiments. All matrices and vectors are stored in double-precision format. The datasets were chosen to illustrate the SGD vs CA-SGD methods at various machine scales (real-sim being small scale, news20 being medium scale, and url being large scale).

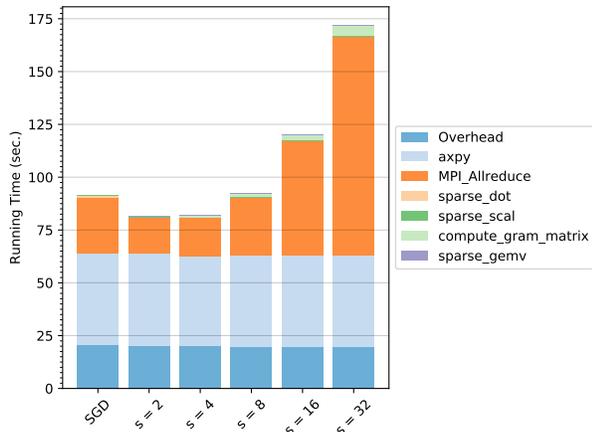
1) *Scaling*: This benchmark is intended to show how CA-SGD and SGD behave as the number of cores is varied. As the number of cores increase, the SGD running time becomes more latency dominant. Since CA-SGD reduces latency cost by s , we expect to see performance improvements over SGD. Figures 3a-3c illustrate the strong scaling (left y-axis in blue) and speedups (right y-axis in red) for the datasets in Table II. Each plot in Figure 3 shows the SGD (solid blue) and CA-SGD (dashed blue) running times with batch size of 1 for each dataset. Both methods were trained for 100 epochs and we report the mean running time and standard deviation (error bars) over 5 trials. The CA-SGD running times are annotated with the value of s that achieved the best performance.

In Figure 3a we can observe that at small scale ($p = 24$ and $p = 48$) the computational cost dominates with $s = 2$ resulting in the best CA-SGD running times. Since the latency cost increases with the number of cores, we can see that CA-SGD can use larger values of s at larger core counts. This eventually leads CA-SGD ($p = 192$ and $s = 8$) to attain an average speedup of $2.27\times$ over SGD ($p = 96$). Finally, at $p = 384$ we see that CA-SGD performance degrades despite increasing s . This is because the additional computation and bandwidth costs dominate the reduction in latency cost. Figure 3b shows the scaling results for the smaller real-sim dataset. This dataset has fewer columns per core which means that

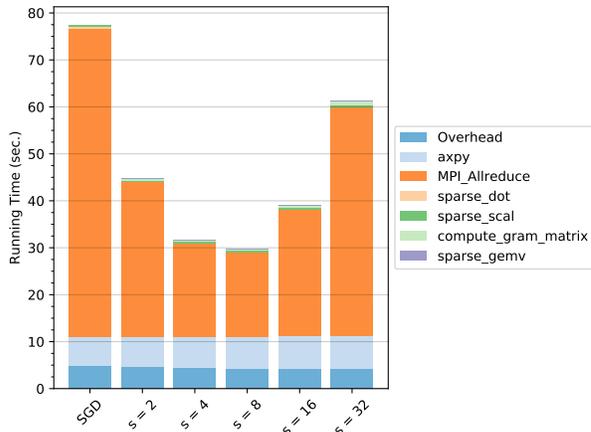
the latency cost is more dominant at smaller scales. This is evidenced by the fact that we can start at $s = 4$ for this dataset. As the number of cores increases, we observe that s becomes larger and the speedup from CA-SGD increases. However, at $p = 96$ we see the additional computation and bandwidth costs begin to dominate and performance of CA-SGD degrades. For the real-sim dataset, CA-SGD ($p = 48$ and $s = 16$) achieves an average speedup of $3.41\times$ over SGD ($p = 12$). Figure 3c illustrates scaling results for the larger url dataset. Due to the size of this dataset, SGD and CA-SGD can scale to larger numbers of cores (potentially higher latency costs). For this dataset, CA-SGD ($p = 768$ and $s = 16$) achieves an average speedup of $4.97\times$ over SGD ($p = 192$) and scales to $4\times$ as many cores. The scaling results in this section suggest that CA-SGD can achieve large average speedups of up to $4.97\times$ over SGD and scale out further on a parallel cluster. These experiments further validate the theoretical analysis and illustrate that reducing latency at the expense of bandwidth and computation can lead to significant performance improvements.

2) *Running time breakdown*: This benchmark is intended to show a breakdown of how much time is spend on computational kernels and communication routines in the SGD and CA-SGD algorithms. We will compare SGD vs. CA-SGD and at two different core counts. At smaller core counts, latency is less dominant so the benefits of CA-SGD will be less pronounced. However, once we transition to large core counts, the latency reduction of CA-SGD should result in larger speedups. We report the running time breakdown of SGD vs CA-SGD with $b = 1$ for several values of s on the news20 dataset. We obtained the running time breakdown by using the Tuning and Analysis Utilities (TAU) to instrument our code [24]. Since TAU generates profiles for each MPI process, we show the average over all MPI processes. Some operations such as the row sampling, scalar operations, loop overheads, and memory management are grouped into overhead.

Figures 4a and 4b illustrate the running time breakdown of SGD and CA-SGD at $p = 48$ and $p = 192$, respectively. Note that the MPI_Allreduce times include bandwidth and latency costs. Compute Gram Matrix, sparse dot, sparse scal,



(a) news20 Running Time Breakdown ($p = 48$)



(b) news20 Running Time Breakdown ($p = 192$)

Fig. 4. Running time breakdown of SGD and CA-SGD (with $b = 1$) on the news20 dataset (see Table II) for $p = 48$ and $p = 192$. Overhead includes the time spent on row sampling, scalar operations, loop overhead, and memory management. SGD does not require the sparse gemv nor the compute Gram matrix operations. CA-SGD replaces the sparse dot with sparse gemv computations. Note that MPI_Allreduce in the legend is a combination of the bandwidth and latency costs analyzed and the remaining legend items except Overhead correspond to the flops cost analyzed in Section IV.

sparse gemv, and aaxy correspond to the flops cost analyzed in Section IV. At small scale (Fig. 4a), the computation cost (aaxy) dominates the communication cost (MPI_Allreduce). As s increases we begin to see a reduction in MPI_Allreduce times due to a reduction in latency cost. However, starting at $s = 8$ the additional computation and bandwidth costs of CA-SGD dominate and cause the running times to grow. In the best case (at $s = 2$) CA-SGD achieves a communication speedup of $1.58\times$ and an overall speedup of $1.12\times$ over SGD.

At large scale (Fig. 4b), the communication time is dominated by latency due to synchronization with $4\times$ more cores. Since latency dominates, CA-SGD achieves speedups for a wider range of values for s . At $s = 8$ CA-SGD attains a communication speedup of $6.5\times$ and overall speedup of $2.6\times$ over SGD. In both figures we see cases where CA-SGD is much slower than SGD. In those cases, the additional bandwidth cost of CA-SGD is the bottleneck and not the additional computation. This suggests that if a candidate parallel cluster is bandwidth-limited, then the maximum values of s and the speedups attained by CA-SGD will be limited.

3) *Batch size vs s* : This benchmark will explore how setting $b > 1$ affects the speedups CA-SGD can obtain over SGD. From the previous results with $b = 1$ we see that the additional bandwidth and computation cost introduced by $s > 1$ does degrade CA-SGD performance when s is too high (e.g. $s > 16$ for the url dataset). The performance results thus far compare SGD, which samples a single row every iteration, to its CA-SGD variant. As expected, SGD is latency dominated which results in large speedups for CA-SGD. If the batch size is increased then the sb additional bandwidth and computation costs will likely dominate. We should expect that s must be decreased in order to compensate for increasing the batch size, b . Figure 5 illustrates a speedup heatmap comparing SGD (with $s = 1$ and $b \geq 1$, bottom-left corner) and CA-SGD (with

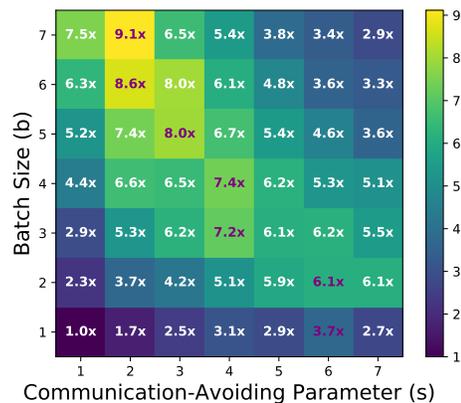


Fig. 5. Comparison of SGD vs CA-SGD for various batch sizes and values of s on the url dataset with $p = 384$. All speedups are relative to SGD ($s = 1$ and $b = 1$). As s and b increase, the bandwidth cost for CA-SGD increases by a factor of sb . Speedups relative to SGD with $b > 1$ can be obtained by dividing CA-SGD ($s > 1$, $b > 1$) speedups by the corresponding SGD ($s = 1$, $b > 1$) speedup.

$s > 1$ and $b \geq 1$) on the url dataset with $p = 384$. Speedups are relative to SGD with $s = 1, b = 1$ (bottom left corner) which means that speedups greater than $1\times$ in column $s = 1$ are due exclusively to increasing batch size. Speedups greater than $1\times$ in row $b = 1$ are due exclusively to communication-avoidance. For data points with $s > 1$ and $b > 1$, speedups greater than $1\times$ are due to a combination of larger batch sizes and communication-avoidance. Note that for $b > 1$, the speedups are due to using BLAS-2 (with increasing matrix sizes) instead of BLAS-1 functions. The heatmap illustrates that CA-SGD is most effective for small batch sizes where latency dominates.

In these experiments, we focused on square and rectangular sparse matrices stored in CSR format. However, in some

situations the input data may be dense. In the dense case, CA-SGD becomes more compute-bound due to the additional elements present in the matrix. As a result, CA-SGD speedups over SGD will be more modest than for sparse matrices. Given that CA-SGD achieves better speedups for latency-dominated/distributed environments, it is well placed to attain large speedups over SGD in cloud environments and when using programming models like Spark/MapReduce (due to the higher latencies in those settings). CA-SGD is unlikely to attain large speedups in shared-memory environments where inter-core latencies are orders of magnitude lower than inter-node latencies. However, exploring and quantifying the performance difference between CA-SGD and SGD on shared-memory would be interesting. We intend to study the performance evaluation of CA-SGD on the various hardware and programming environments in future work.

VI. CONCLUSION AND FUTURE WORK

In this paper, we derived a communication-avoiding variant of SGD for solving the logistic regression problem. We proved theoretical bounds on the computation and communication costs which showed that CA-SGD reduces latency costs by a tunable factors of s . We showed that CA-SGD is numerically stable and achieves speedups of up to $4.97\times$ over SGD on a high-performance Infiniband cluster. When latency is the dominant cost CA-SGD can achieve large speedups despite the additional bandwidth and computation costs. However, as the computation and bandwidth costs increase (by increasing batch size), the speedups decrease. This suggests that CA-SGD might perform even better on cloud/commodity resources. Implementing and benchmarking CA-SGD on cloud resources and programming models would be very interesting.

1) *Implications for neural networks:* Backpropagation in neural networks introduces a set of nested recurrence relations with nonlinear activation functions at each layer and hidden unit. Since logistic regression can be interpreted as a single-layer neural network, we can likely apply our technique to feedforward neural networks (FNN) with nonlinear activation functions. The extension to convolution layers should also be straightforward given that convolutions are linear operations. However, for practical applications to CNNs we need to assess whether the s -step derivation can be applied to batch normalization and pooling layers. While we believe that the s -step technique can be extended to FNNs and CNNs, hand deriving CA-variants for each individual FNN/CNN model is unscalable. Therefore, it is critical to develop tools and techniques that can help automate the CA-derivation process. Finally, as models get wider and deeper, they become more compute and bandwidth bound. This suggests that our approach is most impactful when large models are scaled out to a latency-bound setting.

ACKNOWLEDGMENT

We would like to thank the anonymous reviewers for their helpful feedback. Computational resources were provided by

the Maryland Advanced Research Computing Center. AD is supported by the Gordon and Betty Moore Foundation.

REFERENCES

- [1] S. J. Wright, "Coordinate descent algorithms," *Mathematical Programming*, vol. 151, no. 1, pp. 3–34, 2015.
- [2] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proceedings of COMPSTAT*. Springer, 2010, pp. 177–186.
- [3] Y. Nesterov, "Efficiency of coordinate descent methods on huge-scale optimization problems," *SIAM Journal on Optimization*, vol. 22, no. 2, pp. 341–362, 2012.
- [4] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM transactions on intelligent systems and technology (TIST)*, vol. 2, no. 3, p. 27, 2011.
- [5] B. Recht, C. Re, S. Wright, and F. Niu, "Hogwild: A lock-free approach to parallelizing stochastic gradient descent," in *Advances in Neural Information Processing Systems 24*, 2011, pp. 693–701.
- [6] V. Smith, S. Forte, C. Ma, M. Takáč, M. I. Jordan, and M. Jaggi, "CoCoA: A general framework for communication-efficient distributed optimization," *Journal of Machine Learning Research*, vol. 18, no. 230, pp. 1–49, 2018.
- [7] C. Ma, J. Konečný, M. Jaggi, V. Smith, M. I. Jordan, P. Richtárik, and M. Takáč, "Distributed optimization with arbitrary local solvers," *Optimization Methods and Software*, vol. 32, no. 4, pp. 813–848, 2017.
- [8] M. F. Hoemmen, "Communication-avoiding Krylov subspace methods," Ph.D. dissertation, EECS Department, University of California, Berkeley, Apr 2010.
- [9] E. Carson, "Communication-avoiding Krylov subspace methods in theory and practice," Ph.D. dissertation, EECS Department, University of California, Berkeley, Aug 2015.
- [10] A. T. Chronopoulos and C. W. Gear, "On the efficient implementation of preconditioned s -step conjugate gradient methods on multiprocessors with memory hierarchy," *Parallel computing*, vol. 11, no. 1, pp. 37–53, 1989.
- [11] A. Chronopoulos, "A class of parallel iterative methods implemented on multiprocessors," Ph.D. dissertation, Chicago, IL, USA, 1987.
- [12] M. Mohiyuddin, M. Hoemmen, J. Demmel, and K. Yelick, "Minimizing communication in sparse matrix solvers," in *Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*, Nov 2009, pp. 1–12.
- [13] J. Demmel, M. Hoemmen, M. Mohiyuddin, and K. Yelick, "Avoiding communication in sparse matrix computations," in *IEEE International Symposium on Parallel and Distributed Processing*, April 2008, pp. 1–12.
- [14] E. Carson, N. Knight, and J. Demmel, "Avoiding communication in nonsymmetric Lanczos-based Krylov subspace methods," *SIAM Journal on Scientific Computing*, vol. 35, no. 5, pp. S42–S61, 2013.
- [15] E. Carson and J. Demmel, "A residual replacement strategy for improving the maximum attainable accuracy of s -step Krylov subspace methods," *SIAM Journal on Matrix Analysis and Applications*, vol. 35, no. 1, pp. 22–43, 2014.
- [16] —, "Accuracy of the s -step Lanczos method for the symmetric eigenproblem in finite precision," *SIAM Journal on Matrix Analysis and Applications*, vol. 36, no. 2, pp. 793–819, 2015.
- [17] Z. A. Zhu, W. Chen, G. Wang, C. Zhu, and Z. Chen, "P-packSVM: Parallel primal gradient descent kernel SVM," in *Proceedings of the 9th IEEE International Conference on Data Mining*, 2009, pp. 677–686.
- [18] S. Soori, A. Devarakonda, Z. Blanco, J. Demmel, M. Gurbuzbalaban, and M. M. Dehnavi, "Reducing communication in proximal Newton methods for sparse least squares problems," in *Proceedings of the International Conference on Parallel Processing*. ACM, 2018, pp. 22:1–22:10.
- [19] A. Devarakonda, "Avoiding communication in first order methods for optimization," Ph.D. dissertation, EECS Department, University of California, Berkeley, Jul 2018.
- [20] A. Devarakonda, K. Fountoulakis, J. Demmel, and M. W. Mahoney, "Avoiding communication in primal and dual block coordinate descent methods," *SIAM Journal on Scientific Computing*, vol. 41, no. 1, pp. C1–C27, 2019.
- [21] —, "Avoiding synchronization in first-order methods for sparse convex optimization," in *IEEE International Parallel and Distributed Processing Symposium*. IEEE, 2018, pp. 409–418.

- [22] R. Thakur, R. Rabenseifner, and W. Gropp, "Optimization of collective communication operations in MPICH," *Int. J. High Perform. Comput. Appl.*, vol. 19, no. 1, pp. 49–66, Feb. 2005.
- [23] R. L. Graham, T. S. Woodall, and J. M. Squyres, "Open MPI: A flexible high performance MPI," in *Proceedings of the 6th International Conference on Parallel Processing and Applied Mathematics*, 2005, p. 228–239.
- [24] S. S. Shende and A. D. Malony, "The tau parallel performance system," *Int. J. High Perform. Comput. Appl.*, vol. 20, no. 2, pp. 287–311, May 2006.