
Communication-Avoiding Methods for Regularized Least-Squares

— Aditya Devarakonda —

SIAM Annual

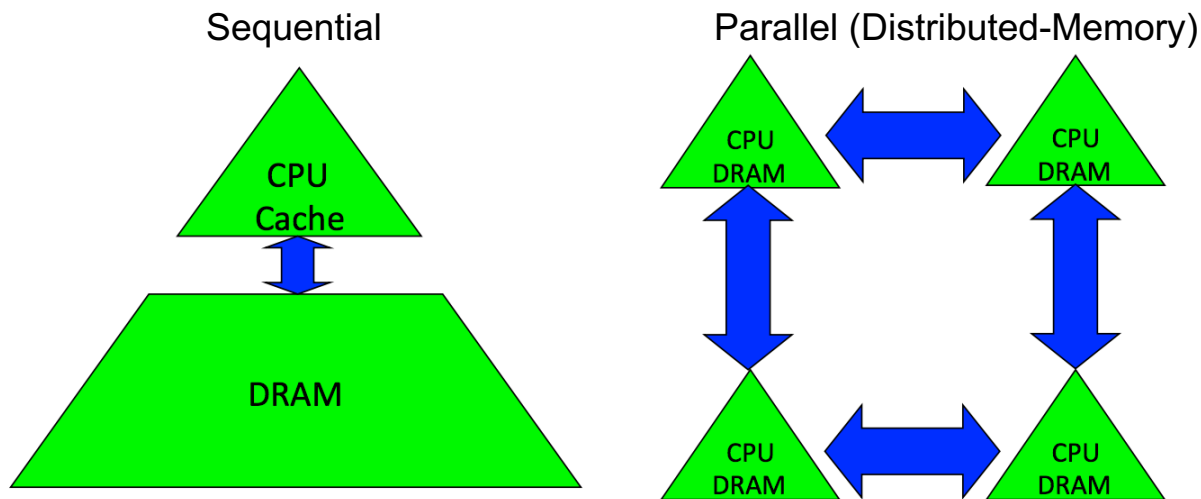
July 14th, 2017

Collaborators:

Jim Demmel (Adviser), Kimon Fountoulakis (Post-Doc), and Michael W. Mahoney (Co-adviser)

Definition

Communication is data movement.



Courtesy:
Demmel

Least-Squares (Linear Regression)

Many ways to solve.

Direct

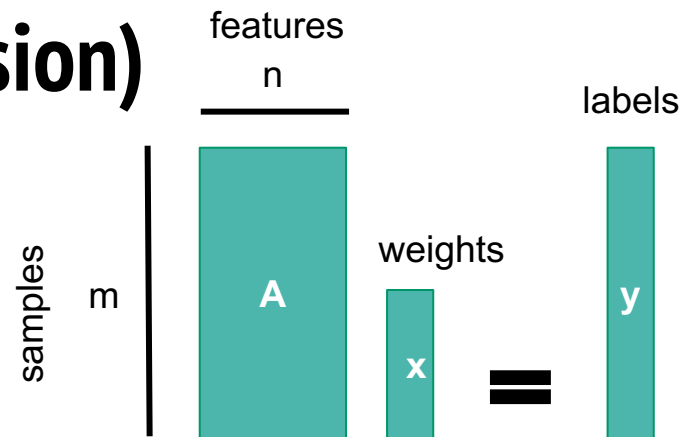
Explicitly solve normal equation.

Implicitly through matrix factorizations.

Iterative

Krylov methods (e.g. Conjugate Gradients).

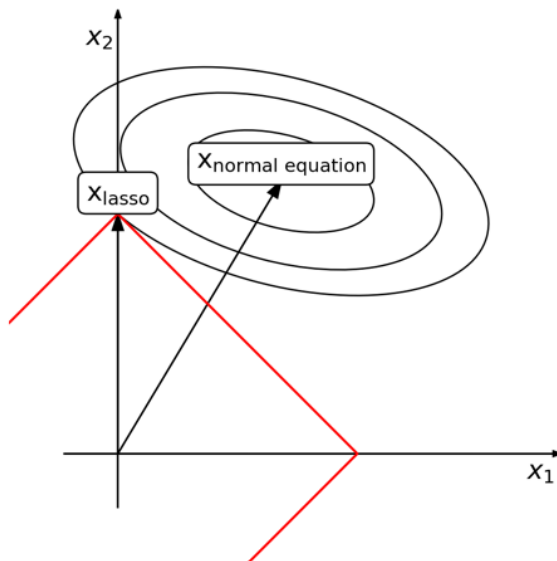
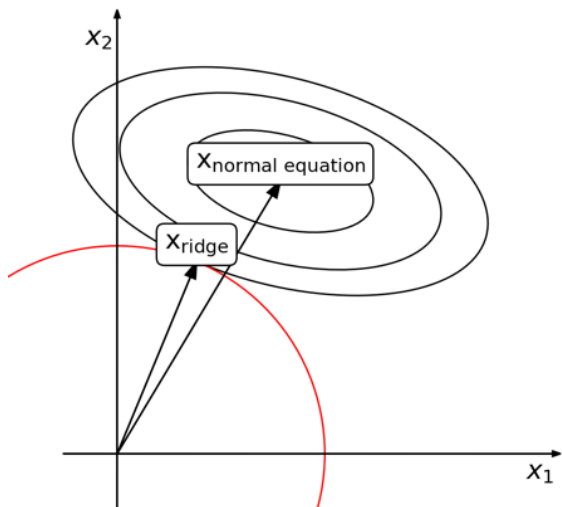
(Block) Coordinate Descent.



$$\text{Least-Squares: } \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - y\|_2^2$$

$$\text{Normal Equation: } x = (A^T A)^{-1} A^T y$$

Regularized Least-Squares (Regression)



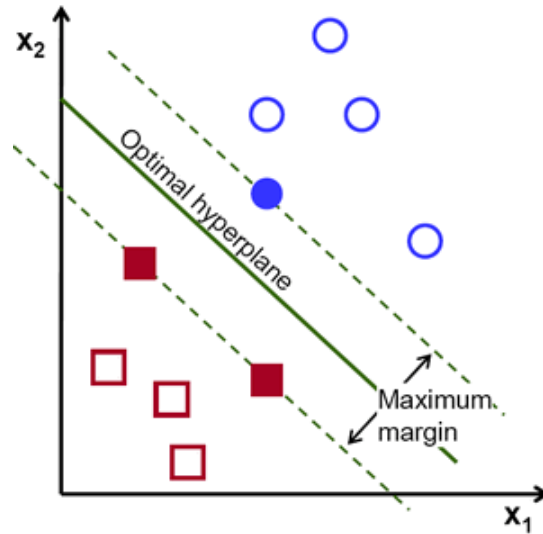
$$\text{Ridge: } \operatorname{argmin}_{x \in \mathbb{R}^n} \underbrace{\frac{1}{2} \|Ax - y\|_2^2}_{\text{Loss function.}} + \underbrace{\frac{\lambda}{2} \|x\|_2^2}_{\text{Regularization function.}}$$

$$\text{LASSO: } \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - y\|_2^2 + \boxed{\lambda \|x\|_1}$$

Loss function. Regularization function.

Binary Classification

Support Vector Machines



Ridge Regression with Block Coordinate Descent

Ridge solution:
(closed-form)

$$\begin{matrix} n \\ \left| \begin{matrix} x \end{matrix} \right. \end{matrix} = \left(\begin{matrix} \underbrace{A^T}_{m} & A \\ \lambda I_n \end{matrix} \right)^{-1} \begin{matrix} A^T \\ y \end{matrix}$$

λ = regularization parameter
 I_n = $n \times n$ Identity matrix
 y = labels

Similar to normal equation.

Ridge Regression with Block Coordinate Descent

Ridge solution:
(closed-form)

$$x = \left(A^T A + \lambda I_n \right)^{-1} A^T y$$

Block Coordinate Descent

λ = regularization parameter
 $I_b = b \times b$ Identity matrix
 r = residual

A' = subsampled columns.

Solution to sub-problem.

$$\Delta x_h = \left(\begin{matrix} m \\ A'^T \end{matrix} A' + \lambda I_b \right)^{-1} A'^T r$$

Block Coordinate Descent in Parallel

$$b \mid \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \vdots \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \Delta x_h = \left(\begin{array}{c|c} \text{---} & \text{---} \\ \text{---} & \text{---} \\ \text{---} & \text{---} \\ \vdots & \vdots \\ \text{---} & \text{---} \\ \text{---} & \text{---} \\ \text{---} & \text{---} \end{array} A'^T \mid \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \vdots \\ \text{---} \\ \text{---} \\ \text{---} \end{array} A' + \lambda_b \mid \begin{array}{c|c} \text{---} & \text{---} \\ \text{---} & \text{---} \\ \text{---} & \text{---} \\ \vdots & \vdots \\ \text{---} & \text{---} \\ \text{---} & \text{---} \\ \text{---} & \text{---} \end{array} A'^T \mid \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \vdots \\ \text{---} \\ \text{---} \\ \text{---} \end{array} r_h \right)^{-1}$$

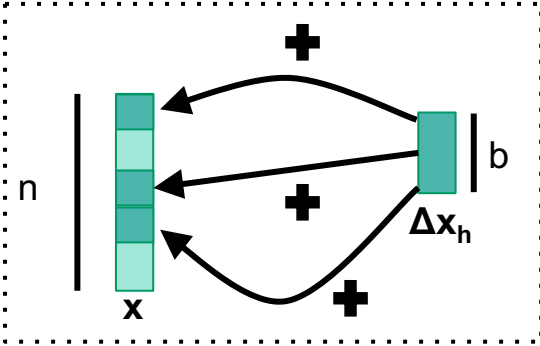
MPI_Allreduce

Block Coordinate Descent in Parallel

$$b \mid \Delta x_h = \left(\begin{array}{c|c} A'^T & A' \\ \hline \vdots & \vdots \end{array} + \lambda_b \right)^{-1} \begin{array}{c} r \\ h \\ \vdots \\ \vdots \end{array}$$

MPI_Allreduce

Update solution

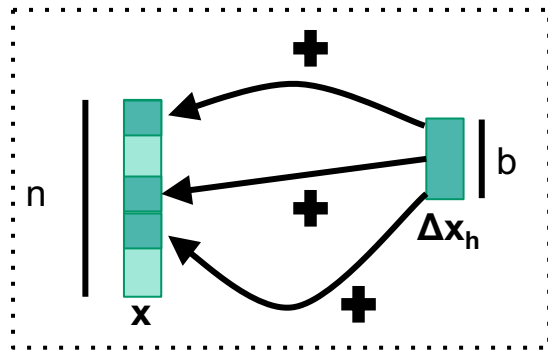


Block Coordinate Descent in Parallel

$$b \mid \begin{array}{|c|} \hline \Delta x_h \\ \hline \end{array} = \left(\begin{array}{|c|} \hline A'^T \\ \hline \vdots \\ \hline \end{array} \mid \begin{array}{|c|} \hline A' \\ \hline \vdots \\ \hline \end{array} + \lambda I_b \right)^{-1} \begin{array}{|c|} \hline r \\ \hline \vdots \\ \hline \end{array} \mid \begin{array}{|c|} \hline h \\ \hline \end{array}$$

MPI_Allreduce

Update solution



Select new coordinates.
Repeat until convergence.

Communication-Avoiding Block Coordinate Descent

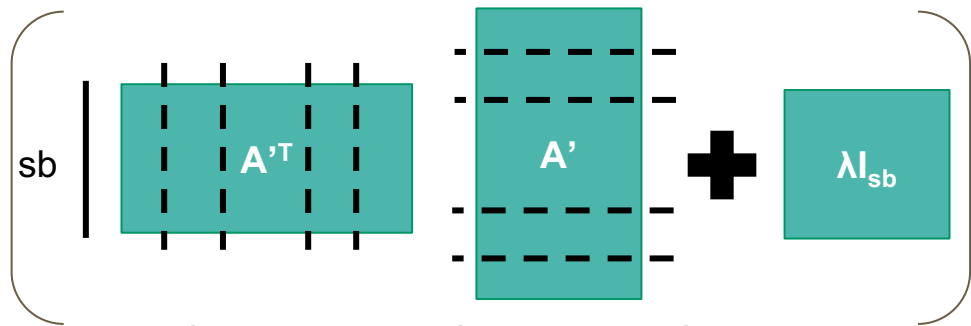
$$b \begin{vmatrix} \Delta x_h \end{vmatrix} = \left(\begin{vmatrix} A'^T & A' \end{vmatrix} + \lambda I_b \right)^{-1} \begin{vmatrix} r_h \end{vmatrix}$$

Block Coordinate Descent

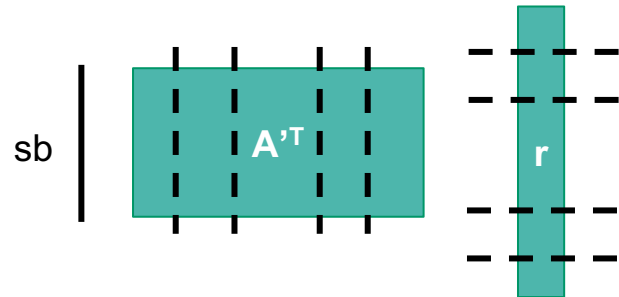
Communication-Avoiding Block Coordinate Descent

$$b \begin{vmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{vmatrix} \Delta x_h = \left(\begin{matrix} \text{---} & \text{---} & \text{---} & \text{---} \\ \text{---} & A'^T & \text{---} & \text{---} \\ \text{---} & \text{---} & A' & \text{---} \\ \text{---} & \text{---} & \text{---} & \text{---} \end{matrix} + \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \lambda_{sb} \right)^{-1} \begin{matrix} \text{---} & \text{---} & \text{---} & \text{---} \\ \text{---} & A'^T & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} & \text{---} \end{matrix} \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} r_h$$

$s = \text{CA parameter}$
 $I_{sb} = sb \times sb \text{ Identity}$
 $r = \text{residual}$



Compute larger Gram matrix, \mathbf{G} .
(Do not invert!)

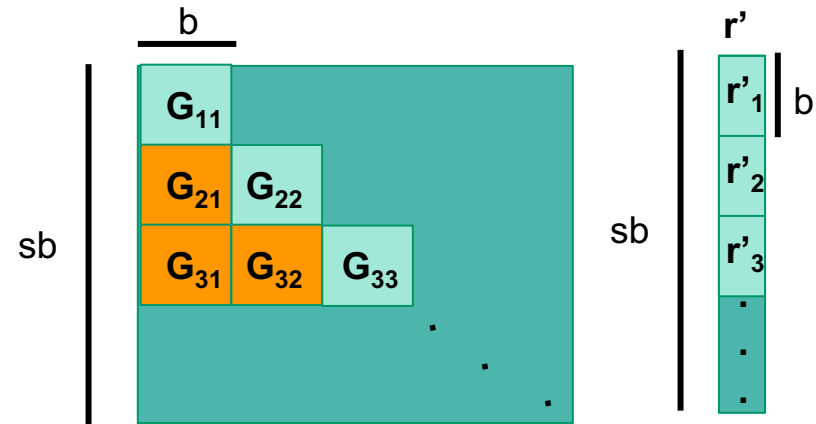


Larger residual contribution.
 r is from this iteration 12

Redundantly on all processors.

Dimensions exaggerated for clarity.

Communication-Avoiding Block Coordinate Descent



Partition G and r' .

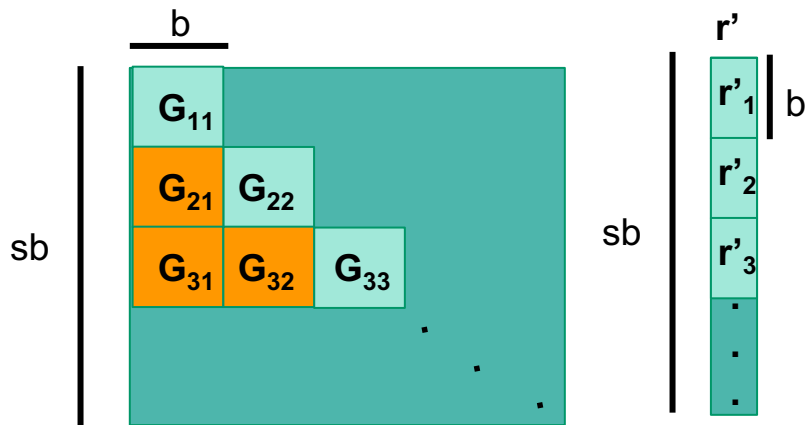
(G is symmetric)

Redundantly on all processors.

No communication for inner iterations.

Communication-Avoiding Block Coordinate Descent

Solution to 1st sub-problem.



$$\Delta x_1 = \left(G_{11} \right)^{-1} r'_1$$

Inner iteration 1

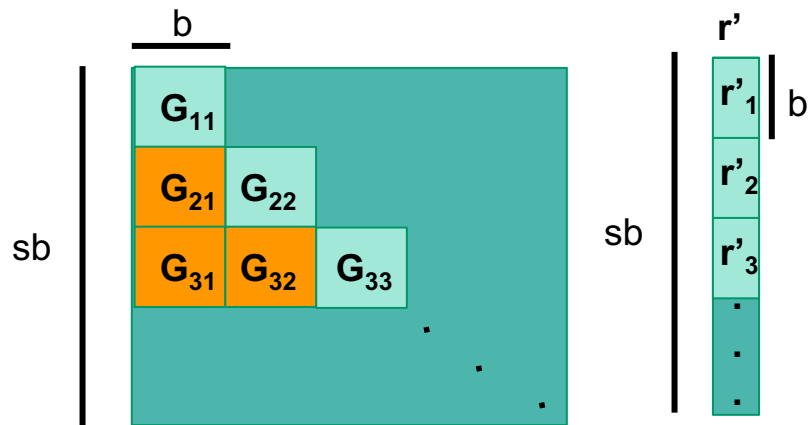
r'_2, r'_3, \dots become stale as soon as Δx_1 is computed.

Orange blocks are the updates

Redundantly on all processors.

No communication for inner iterations.

Communication-Avoiding Block Coordinate Descent



$$\Delta x_1 = \begin{pmatrix} G_{11} \end{pmatrix}^{-1} r'_1$$

Inner iteration 1

$$\Delta x_2 = \begin{pmatrix} G_{22} \end{pmatrix}^{-1} r'_2 - G_{21} \Delta x_1$$

Inner iteration 2

Solution to 2nd sub-problem.

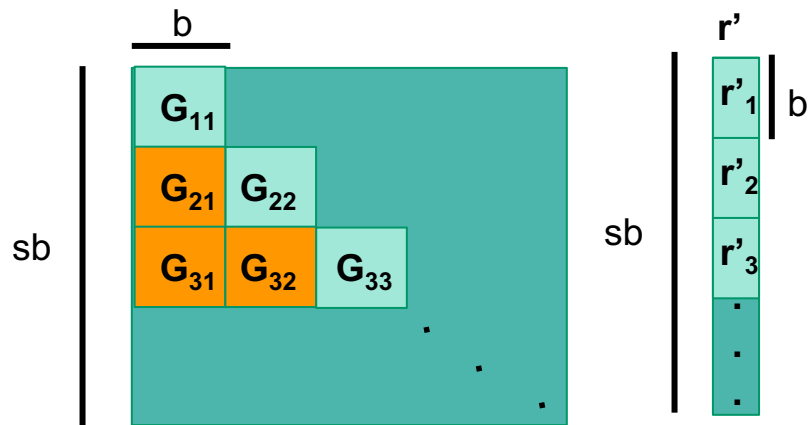
r'_2, r'_3, \dots become stale as soon as Δx_1 is computed.

Orange blocks are the updates

Redundantly on all processors.

No communication for inner iterations.

Communication-Avoiding Block Coordinate Descent



$$\begin{aligned} \Delta x_1 &= \begin{bmatrix} G_{11} \end{bmatrix}^{-1} r'_1 && \text{Inner iteration 1} \\ \dots & & & \dots \\ \Delta x_2 &= \begin{bmatrix} G_{22} \end{bmatrix}^{-1} r'_2 - G_{21} \Delta x_1 && \text{Inner iteration 2} \\ \dots & & & \dots \\ \Delta x_3 &= \begin{bmatrix} G_{33} \end{bmatrix}^{-1} r'_3 - G_{31} \Delta x_1 - G_{32} \Delta x_2 && \text{Inner iteration 3} \end{aligned}$$

Solution to 3rd sub-problem.

Subtraction terms are needed for CA and non-CA solutions to agree.

Algorithm 1 Block Coordinate Descent (BCD) Algorithm

- 1: **Input:** $X \in \mathbb{R}^{d \times n}$, $y \in \mathbb{R}^n$, $H > 1$, $w_0 \in \mathbb{R}^d$, $b \in \mathbb{Z}_+$ s.t. $b \leq d$
- 2: **for** $h = 1, 2, \dots, H$ **do**
- 3: choose $\{i_m \in [d] | m = 1, 2, \dots, b\}$ uniformly at random without replacement
- 4: $\mathbb{I}_h = [e_{i_1}, e_{i_2}, \dots, e_{i_b}]$
- 5: $\Gamma_h = \frac{1}{n} \mathbb{I}_h^T X X^T \mathbb{I}_h + \lambda \mathbb{I}_h^T \mathbb{I}_h$
- 6: $\Delta w_h = \Gamma_h^{-1} (-\lambda \mathbb{I}_h^T w_{h-1} - \frac{1}{n} \mathbb{I}_h^T X z_{h-1} + \frac{1}{n} \mathbb{I}_h^T X y)$
- 7: $w_h = w_{h-1} + \mathbb{I}_h \Delta w_h$
- 8: $z_h = z_{h-1} + X^T \mathbb{I}_h \Delta w_h$
- 9: **Output** w_H

Communication
(every iteration)

Algorithm 2 Communication-Avoiding Block Coordinate Descent (CA-BCD) Algorithm

- 1: **Input:** $X \in \mathbb{R}^{d \times n}$, $y \in \mathbb{R}^n$, $H > 1$, $w_0 \in \mathbb{R}^d$, $b \in \mathbb{Z}_+$ s.t. $b \leq d$
- 2: **for** $k = 0, 1, \dots, \frac{H}{s}$ **do**
- 3: **for** $j = 1, 2, \dots, s$ **do**
- 4: choose $\{i_m \in [d] | m = 1, 2, \dots, b\}$ uniformly at random without replacement
- 5: $\mathbb{I}_{sk+j} = [e_{i_1}, e_{i_2}, \dots, e_{i_b}]$
- 6: let $Y = [\mathbb{I}_{sk+1}, \mathbb{I}_{sk+2}, \dots, \mathbb{I}_{sk+s}]^T X$.
- 7: compute the Gram matrix, $G = \frac{1}{n} Y Y^T + \lambda I$.
- 8: **for** $j = 1, 2, \dots, s$ **do**
- 9: Γ_{sk+j} are the $b \times b$ diagonal blocks of G .
- 10: $\Delta w_{sk+j} = \Gamma_{sk+j}^{-1} \left(-\lambda \mathbb{I}_{sk+j}^T w_{sk} - \lambda \sum_{t=1}^{j-1} \left(\mathbb{I}_{sk+j}^T \mathbb{I}_{sk+t} \Delta w_{sk+t} \right) - \frac{1}{n} \mathbb{I}_{sk+j}^T X z_{sk} \right. \\ \left. - \frac{1}{n} \sum_{t=1}^{j-1} \left(\mathbb{I}_{sk+j}^T X X^T \mathbb{I}_{sk+t} \Delta w_{sk+t} \right) + \frac{1}{n} \mathbb{I}_{sk+j}^T X y \right)$
- 11: $w_{sk+j} = w_{sk+j-1} + \mathbb{I}_{sk+j} \Delta w_{sk+j}$
- 12: $z_{sk+j} = z_{sk+j-1} + X^T \mathbb{I}_{sk+j} \Delta w_{sk+j}$
- 13: **Output** w_H

Communication
(every outer iteration)

No communication

Theoretical Bounds

No free lunch: Reduce latency, but increase flops and bandwidth.

Suppose we perform **H iterations**.

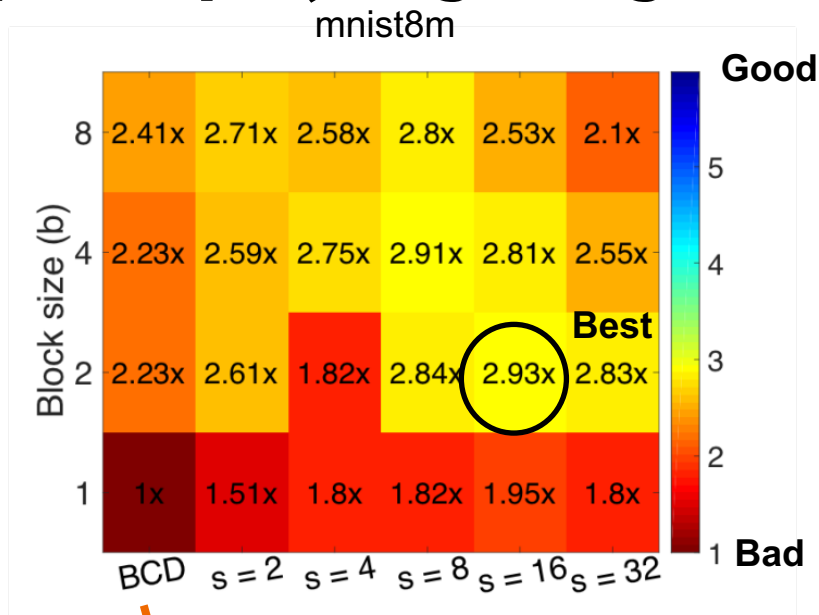
	Flops	Bandwidth	Latency
BCD	$O\left(\frac{Hb^2n}{P} + Hb^3\right)$	$O(Hb^2)$	$O(H \log P)$
CA-BCD	$O\left(\frac{Hsb^2n}{P} + Hb^3\right)$	$O(Hsb^2)$	$O\left(\frac{H}{s} \log P\right)$

Similar bounds for **sparse**.

So, **Latency-** or **Synchronization-Avoiding**.

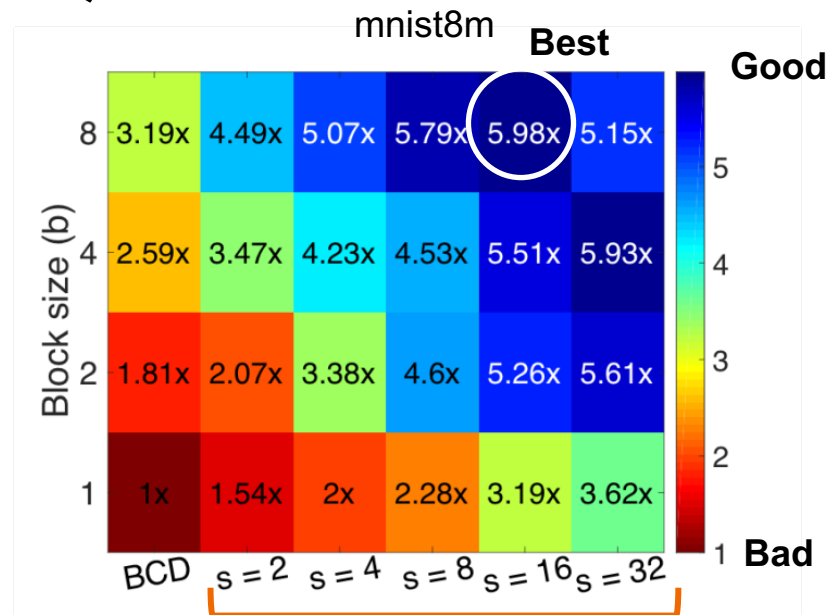
Speedups (Ridge Regression)

8M samples x 768 features



64 nodes of Edison

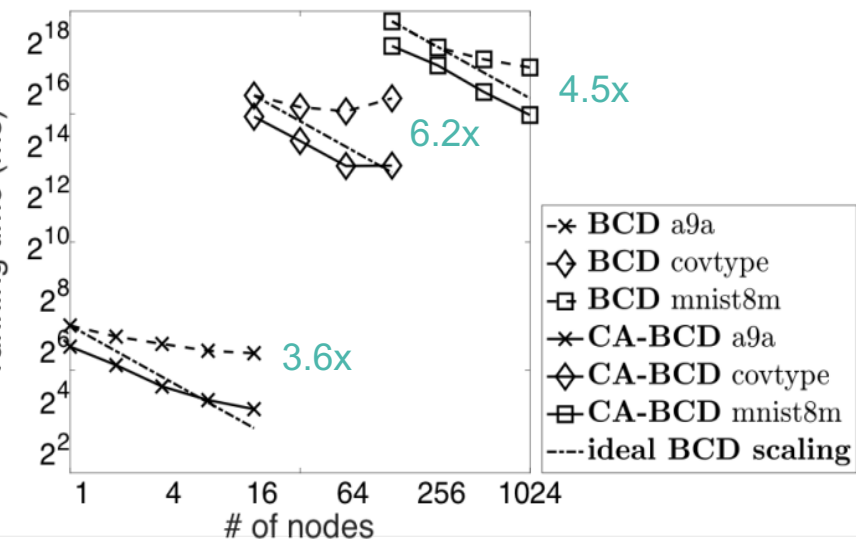
Standard Algorithm



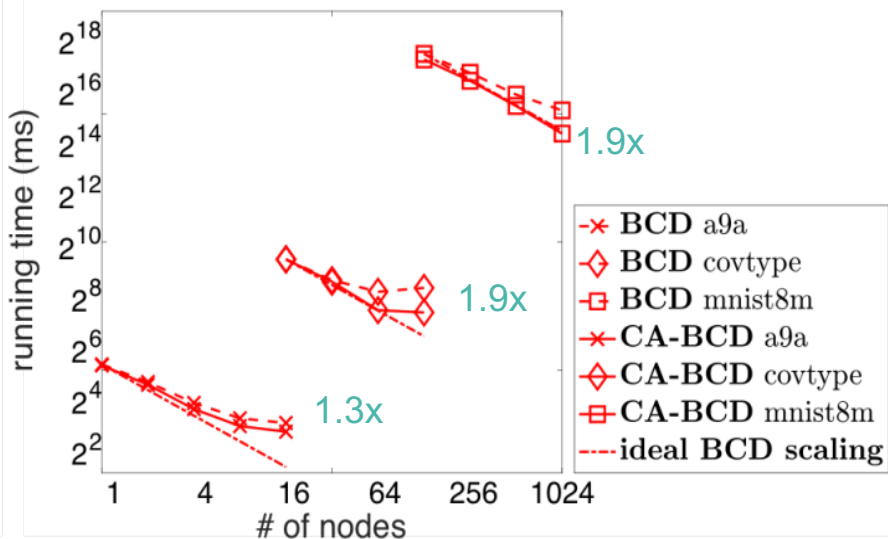
1K nodes of Edison

Communication-Avoiding

Strong Scaling (Ridge Regression)



Blocksize = 1

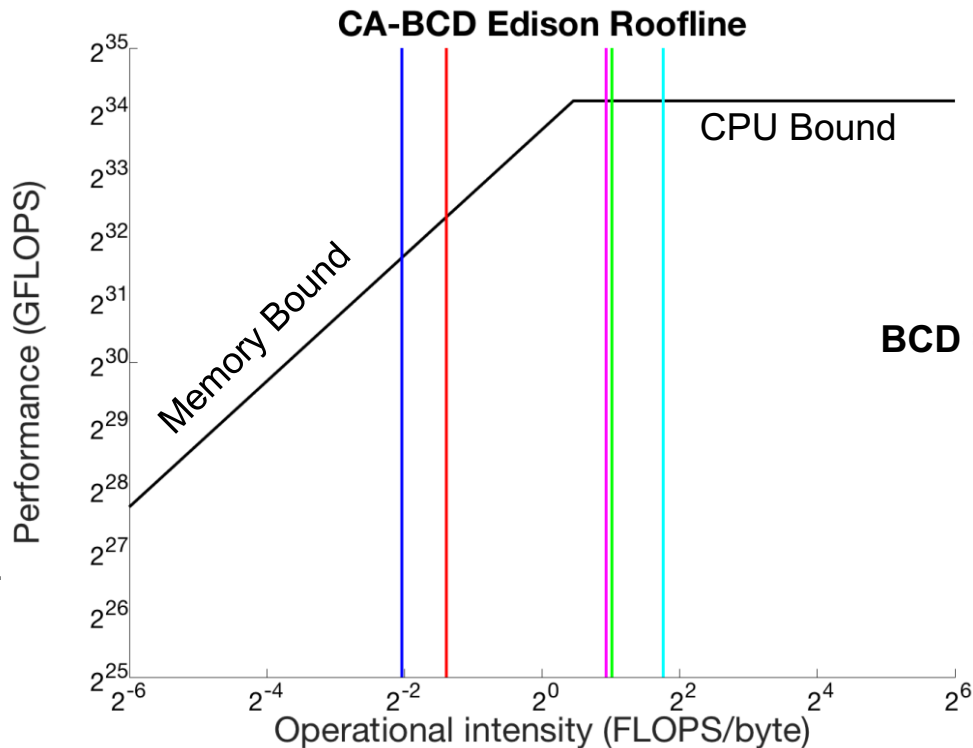


Blocksize = 8

Strong scaling = fixed problem size as P increases.

Theoretical Peak Attainable

Intel "Ivy Bridge": ~19 Gflops
64 GB DDR3 1866 MHz: 14GB/s



BCD = Block Coordinate Descent

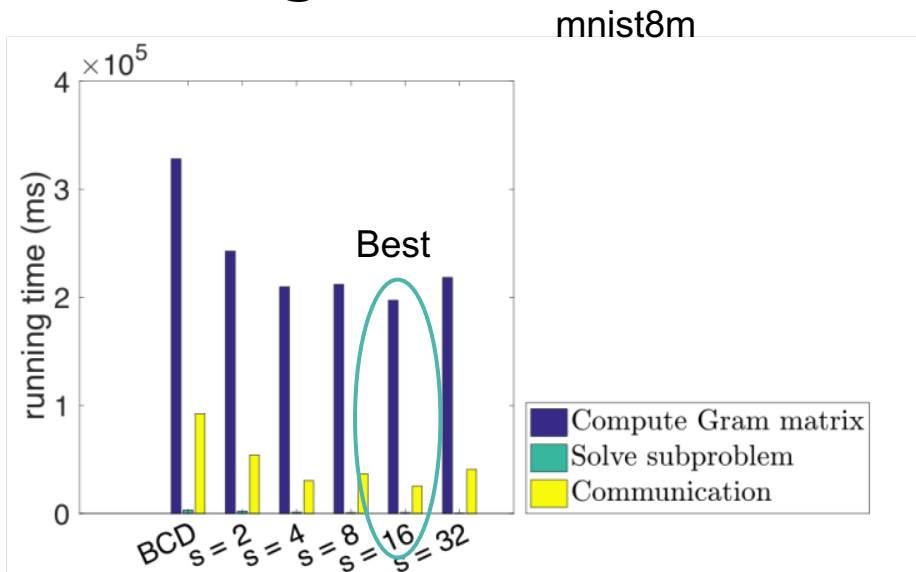
- roofline
- BCD $b = 1$ (dot-product)
- CA-BCD $b = 1, s = 64$
- BCD $b = 8$
- CA-BCD $b = 8, s = 8$
- $A^T A$

Dataset

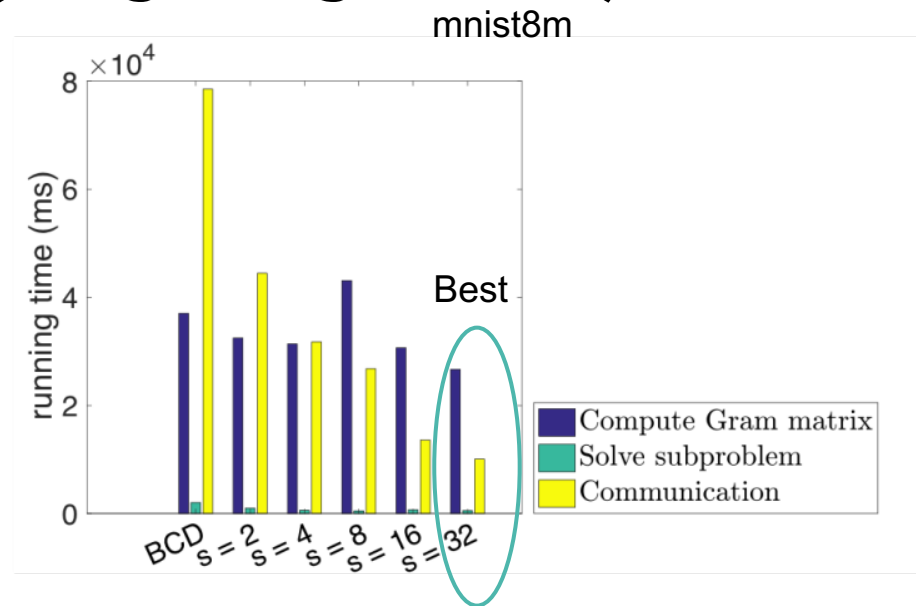
a9a: 32k x 123 matrix.

Sparsity: 11.28%

Running Time Breakdown (Ridge Regression)



64 nodes of Edison

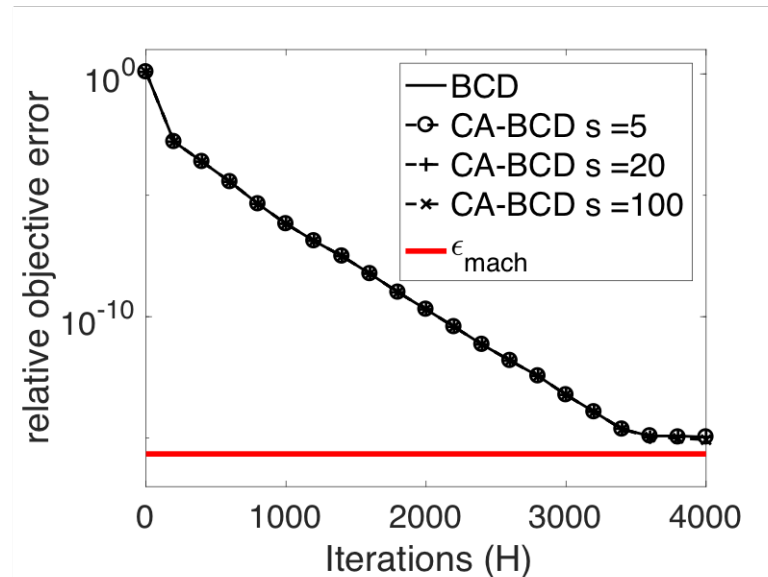
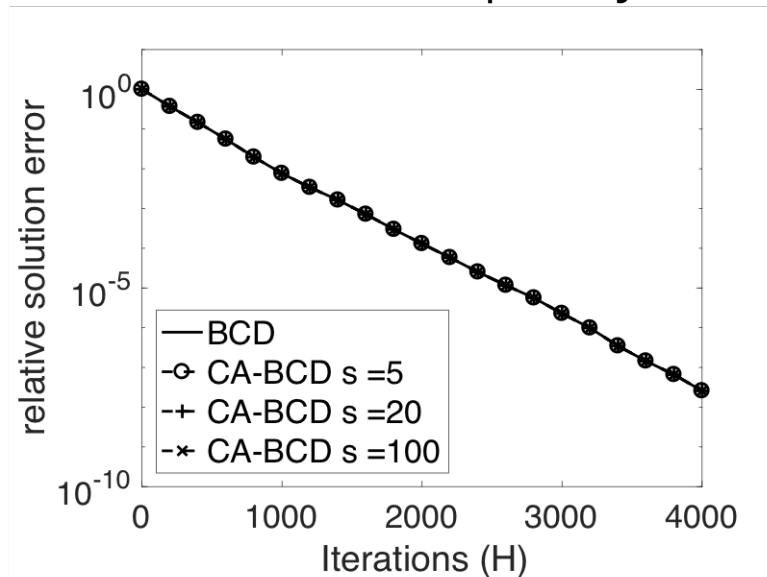


1K nodes of Edison

Blocksize = 1

Numerical Stability (Ridge Regression)

a9a dataset. 32k samples by 123 features. 11.28% non-zeros.



Blocksize = 16, $\sigma_{\min} = 4.884e-6$
 $\lambda = 1000 * \sigma_{\min}$

Not Just Ridge Regression

Applies to **proximal methods** (e.g. LASSO, Group LASSO, elastic-net).

Support Vector Machines (binary classification).

And **associated Dual problems**.

Kernel Ridge and SVM (current work).

more generally to **Generalized Linear Models?** (future work).

Other optimization methods: **Stochastic Gradient Descent?** (future work).

Summary and Future Work

Large speedups when **latency dominates**.

Provably communication-avoiding.

CA-technique applies to **non-linear optimization**.

How far can we go (e.g. Logistic regression)?

Speedups on **other platforms and frameworks**?

Example: Cloud + Spark is latency dominated.

Expect greater speedups!

Problem	MPI Speedup
Ridge Regression	Up to 6.1x
Proximal Least-Squares	Up to 5.1x
SVM	Similar expected

Questions?

Thanks!

Backup Slides

Running Time Breakdown

Strong Scaling (Ridge Regression)

Numerical Stability

Proof of communication-avoidance (Ridge Regression)

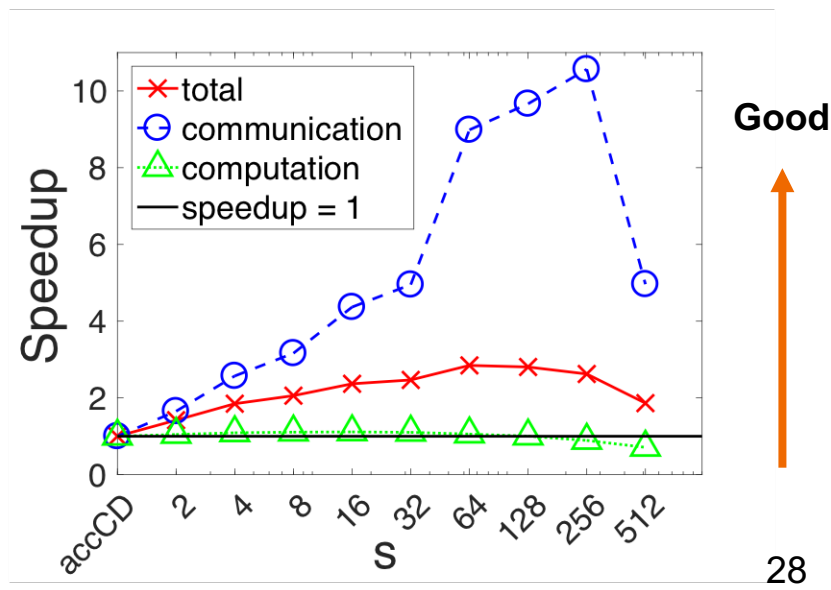
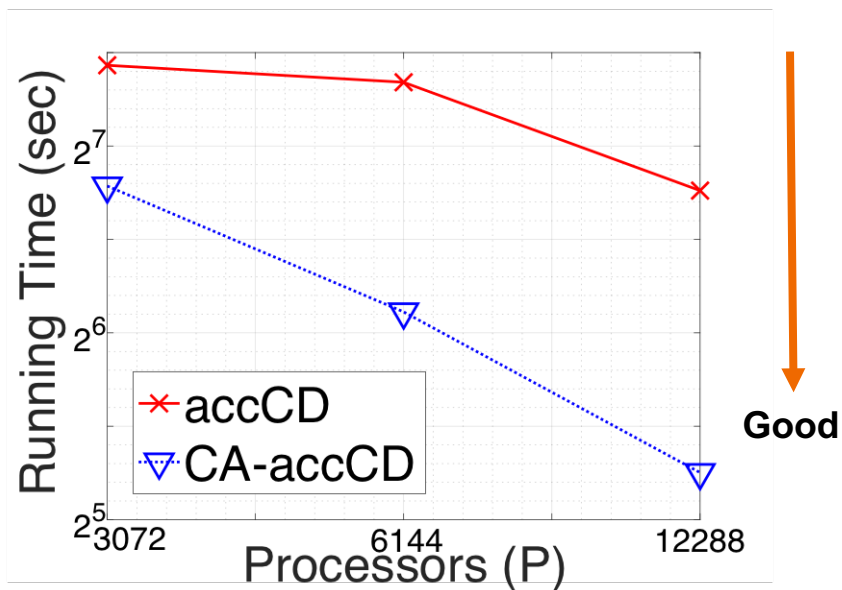
CALASSOS: Scalable Proximal Methods

accCD = accelerated
Coordinate Descent



CA-technique
applies to accelerated
methods.

Strong scaling and speedups on **Url dataset (2M by 3M)**.



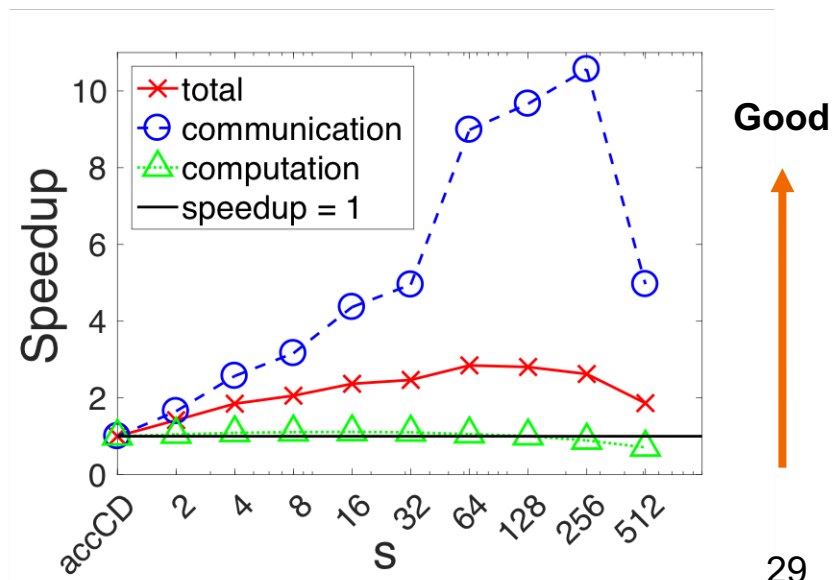
CALASSOS: Scalable Proximal Methods

CA-method perform s^2 more flops.

But still get **computation speedup**.

Due to BLAS-3 calls instead of BLAS-1 in non-CA.

BLAS-3 = Cache-efficient computation + higher flops rate.



CA-SVM: Preliminary Results

Based on Dual Coordinate Descent for Linear SVM (Hsieh, et. al.)

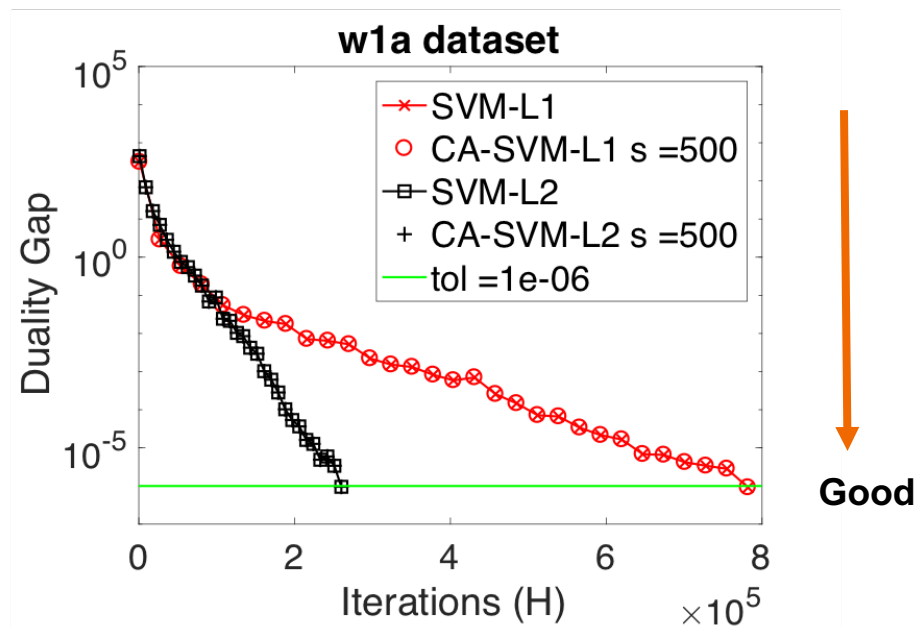
SVM-L1 = Hinge loss = $\max(0, 1 - A^T_i x y_i)$

SVM-L2 = (Hinge loss)²

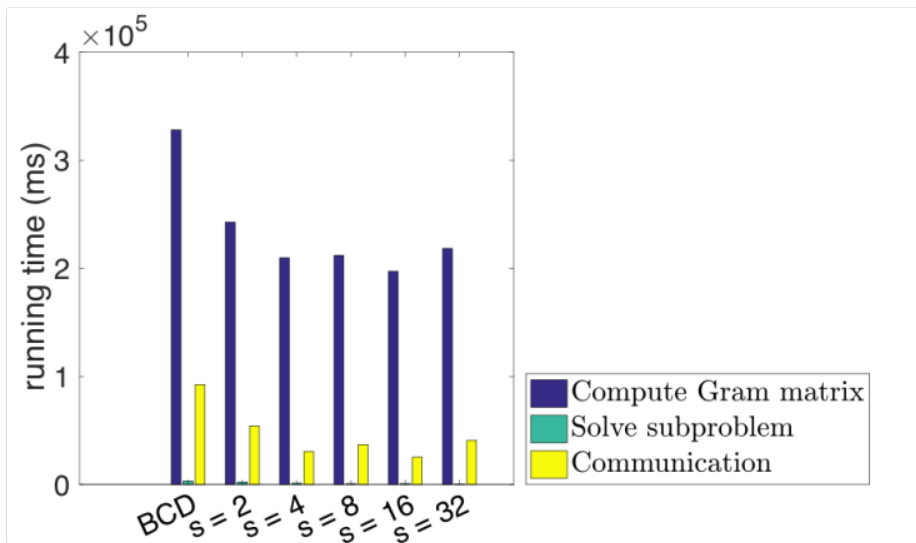
easier so, converges quickly.

Numerically stable (**unlike CA-Krylov**)!

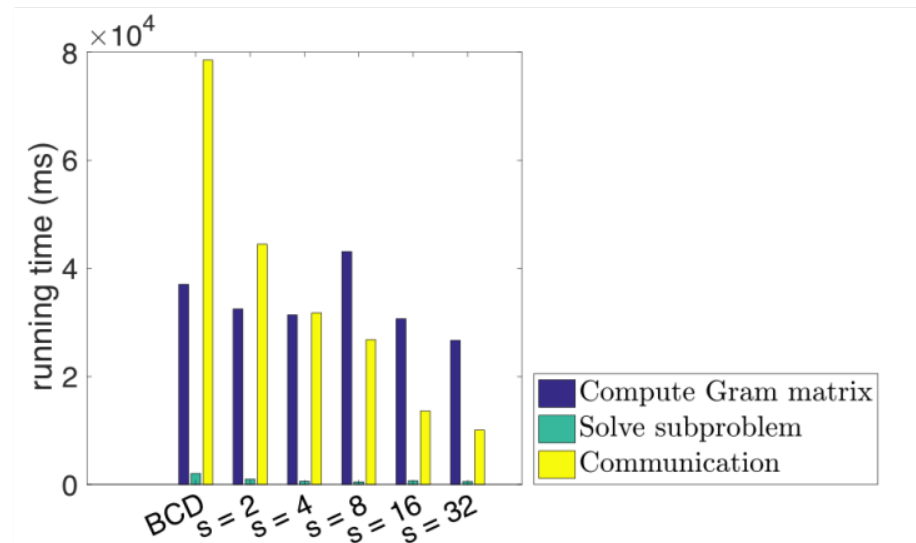
Ditto for Ridge and Proximal.



Running Time Breakdown (Ridge Regression)



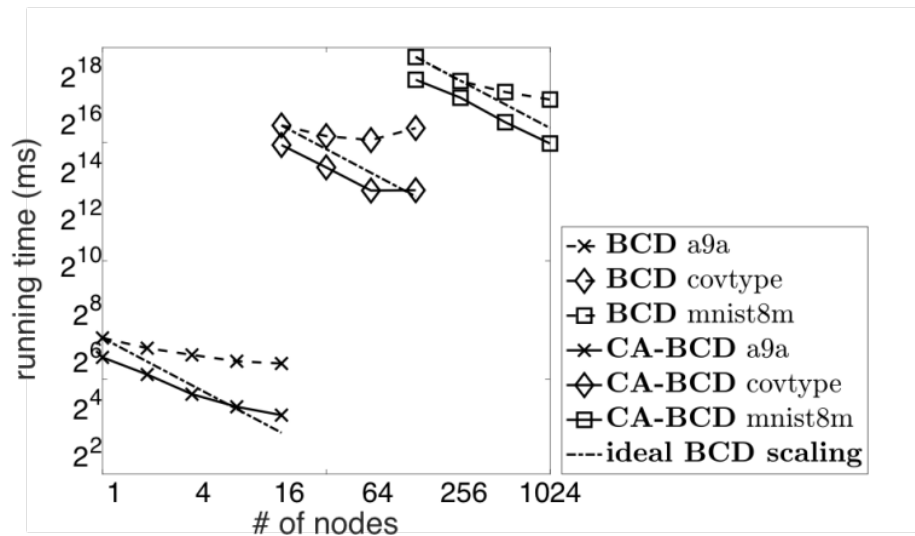
64 nodes of Edison



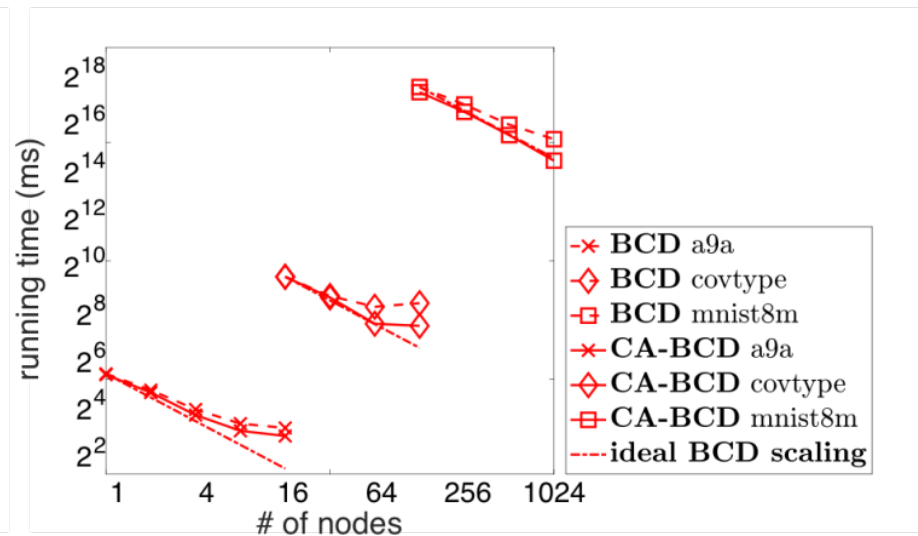
1K nodes of Edison

Blocksize = 1
(Coordinate Descent)

Strong Scaling (Ridge Regression)



Blocksize = 1



Blocksize = 8

Strong scaling = fixed problem size as \mathbf{P} increases.

Provably Communication-Avoiding

No free lunch: Reduce latency, but increase flops and bandwidth.

Suppose we perform **H iterations**.

	Flops	Bandwidth	Latency
BCD	$O\left(\frac{Hb^2n}{P} + Hb^3\right)$	$O(Hb^2)$	$O(H \log P)$
CA-BCD	$O\left(\frac{Hsb^2n}{P} + Hb^3\right)$	$O(Hsb^2)$	$O\left(\frac{H}{s} \log P\right)$

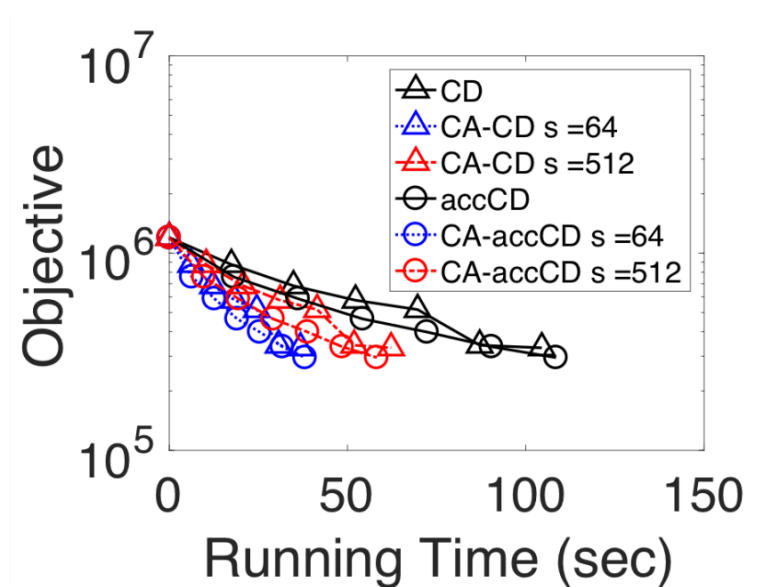
Similar bounds for **sparse**.

CALASSOS: Scalable Proximal Methods

Re-organized Accelerated BCD (Fercoq and Richtarik) for LASSO.

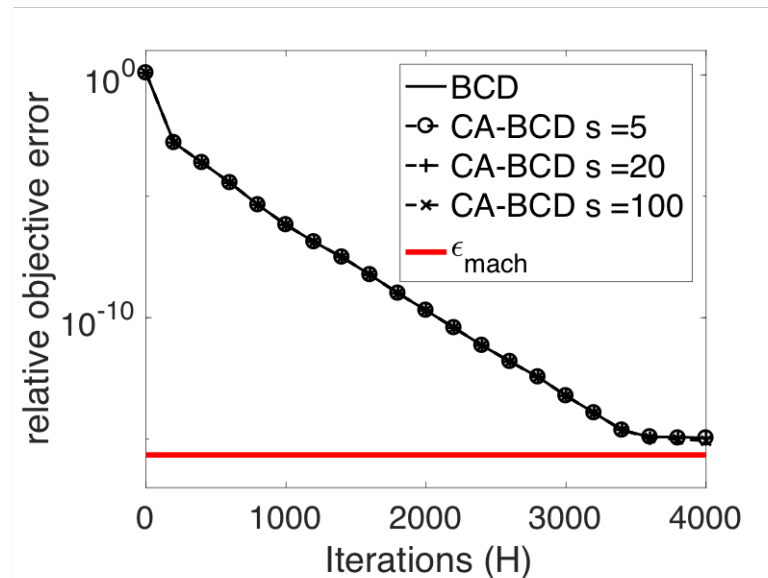
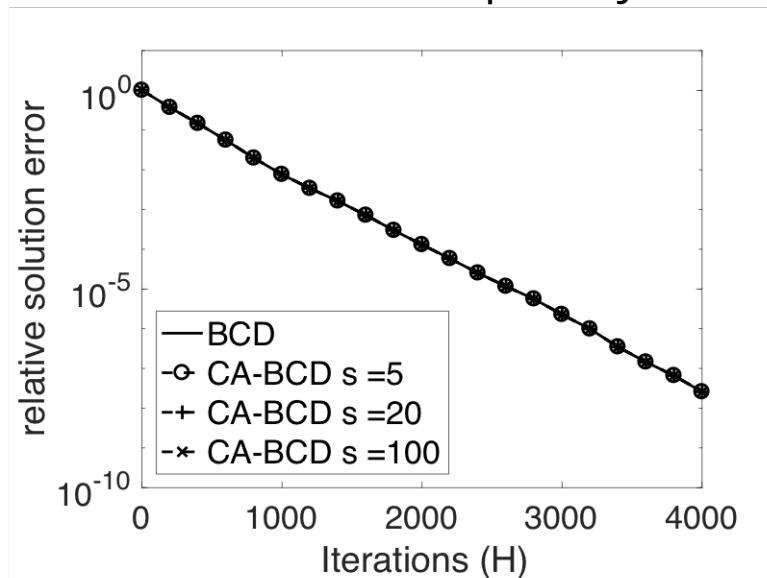
CA vs. non-CA on 1K nodes of Cray XC30.

Choose **s** carefully.



Numerical Stability (Ridge Regression)

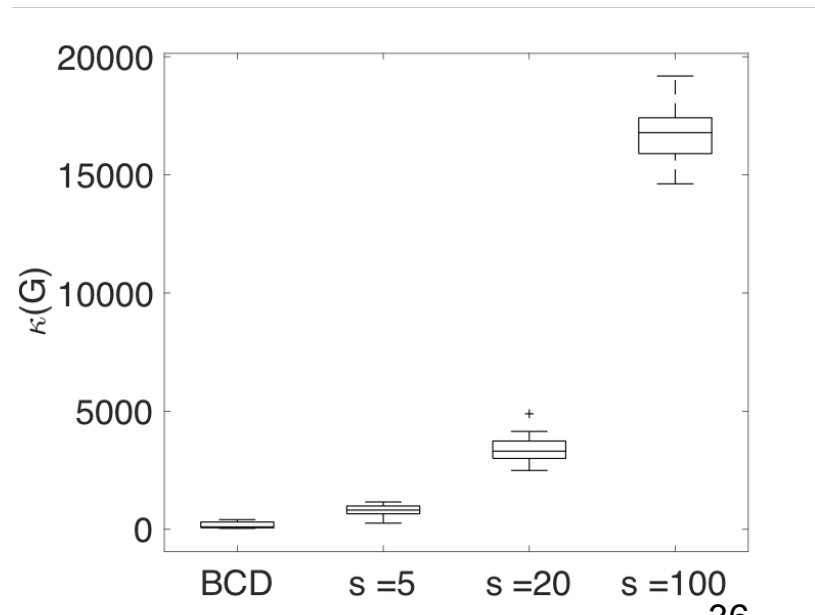
a9a dataset. 32k samples by 123 features. 11% non-zeros.



Numerical Stability (Ridge Regression)

We extract blocks from \mathbf{G} .

So, $\text{cond}(\mathbf{G})$ not a problem.



Numerical Stability (Ridge Regression)

CA-Krylov: numerical stability issues

Fix: Orthogonal polynomials and residual replacement strategies.

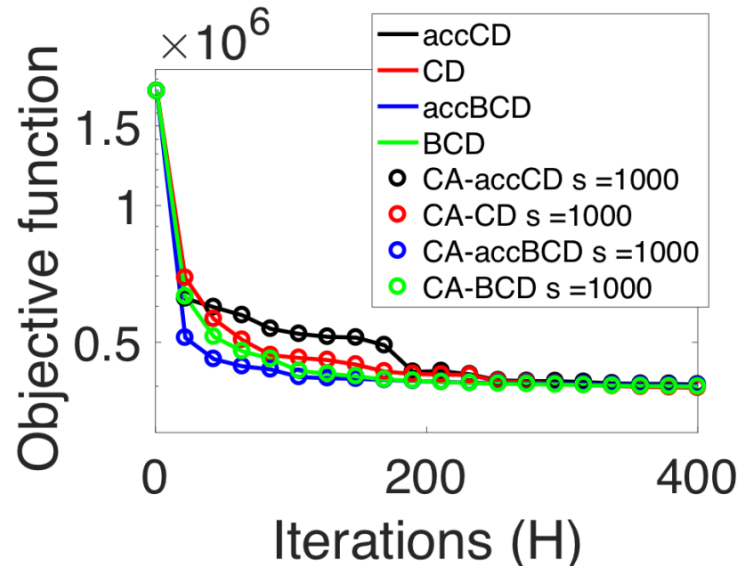
Dissertations by Hoemmen and Carson.

CA-Machine Learning: Magnitude of numerical error is small $\sim O(10^{-15})$.

If we want **low-accuracy**, then CA-ML **essentially stable**.

CALASSOS: Scalable Proximal Methods

Re-organized Accelerated BCD (Fercoq and Richtarik) for LASSO.



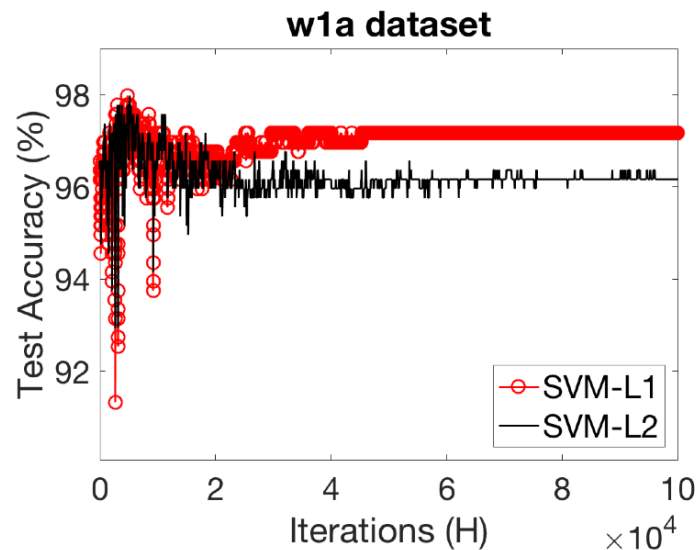
CA-SVM: Preliminary Results

Dual Coordinate Descent for Linear SVM (Hsieh, et. al.)

Once again, non-linearity in “inner loop”.

CA technique applies.

Similar speedups to LASSO expected.



Modeling Communication

Hardware Parameters

Algorithm Parameters

Running Time = Computation time + Communication time

(time per flop) **x** (# of flops)

(time per word) **x** (# of words) + (time per message) **x** (# of messages)