



Welcome to Dissertation Defense of Jyotikrishna Dass

- Presentation will be Recorded
- Non-Committee attendees are requested
 - during presentation, to ensure
 - Audio is MUTE
 - Video is OFF
 - after presentation, join for Q&A and comments

Efficient and Scalable Machine Learning for Distributed Edge Intelligence

PhD Dissertation

Jyotikrishna Dass

Ph.D. Candidate

Computer Science and Engineering

email: dass.jyotikrishna@tamu.edu



TEXAS A&M
UNIVERSITY.

Artificial Intelligence is Everywhere

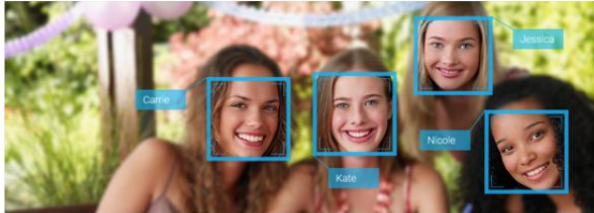


Figure: Face Recognition



Figure: Machine Translation

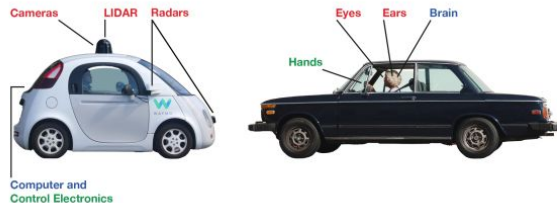


Figure: Self Driving Cars

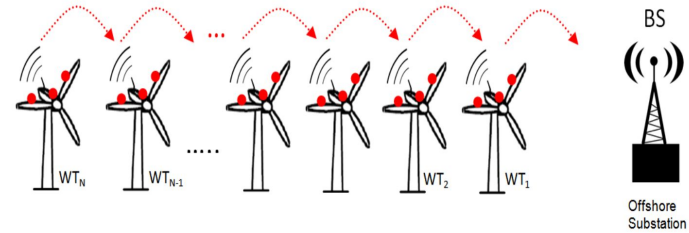


Figure: Speech Processing

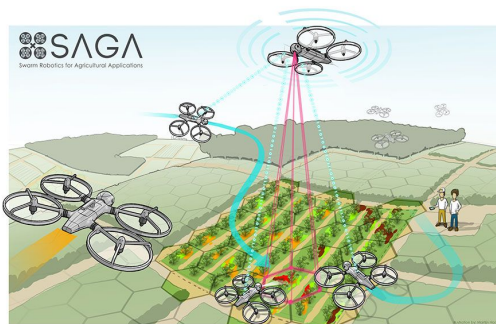
So are, Internet of Things



Intelligent Traffic



Smart Wind Farms

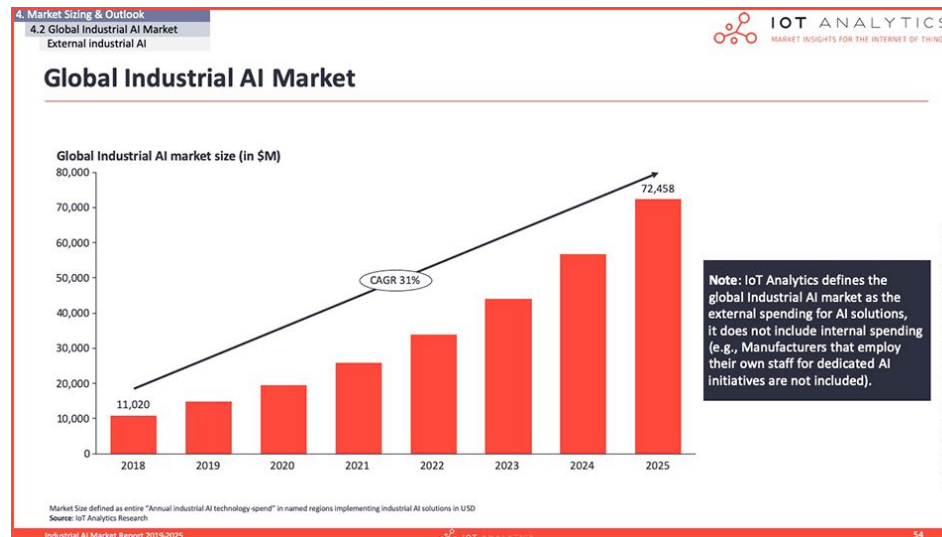
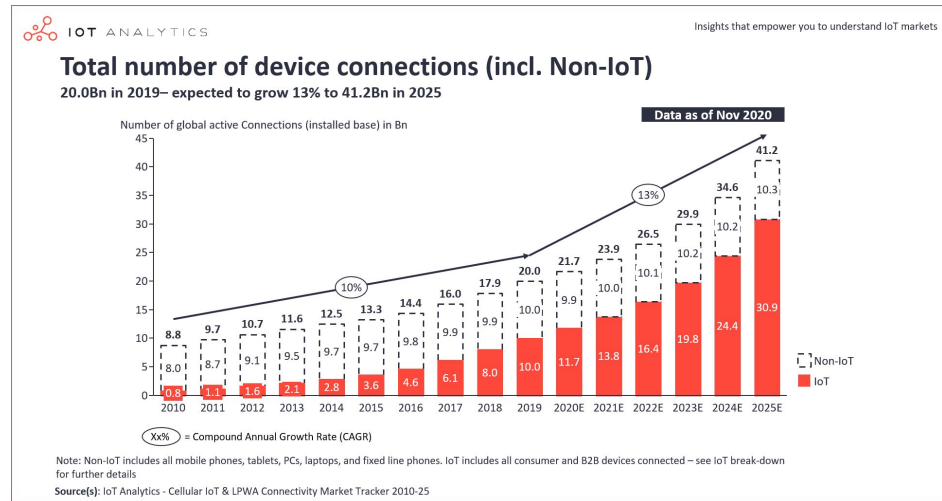


Precision Agriculture

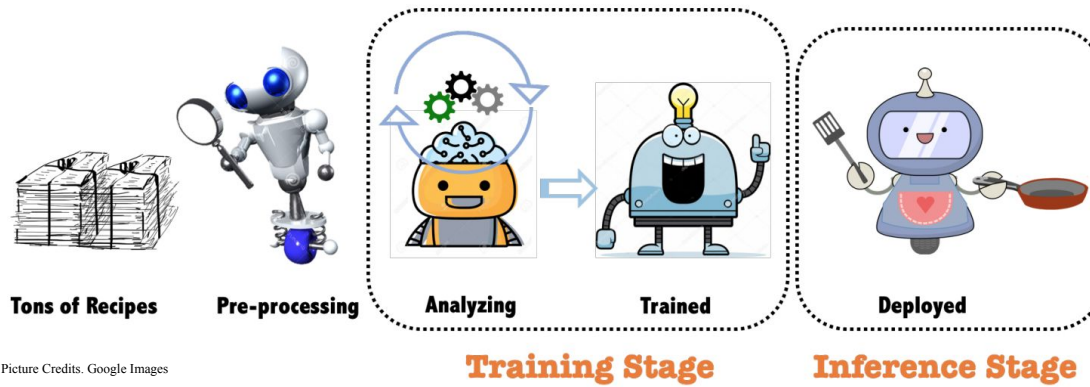


Disaster Relief

Market Trends: AI + IoT



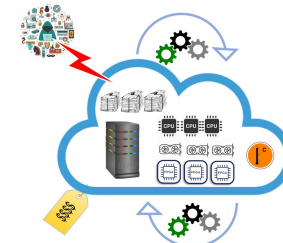
Machine Learning Process



ML Training, usually

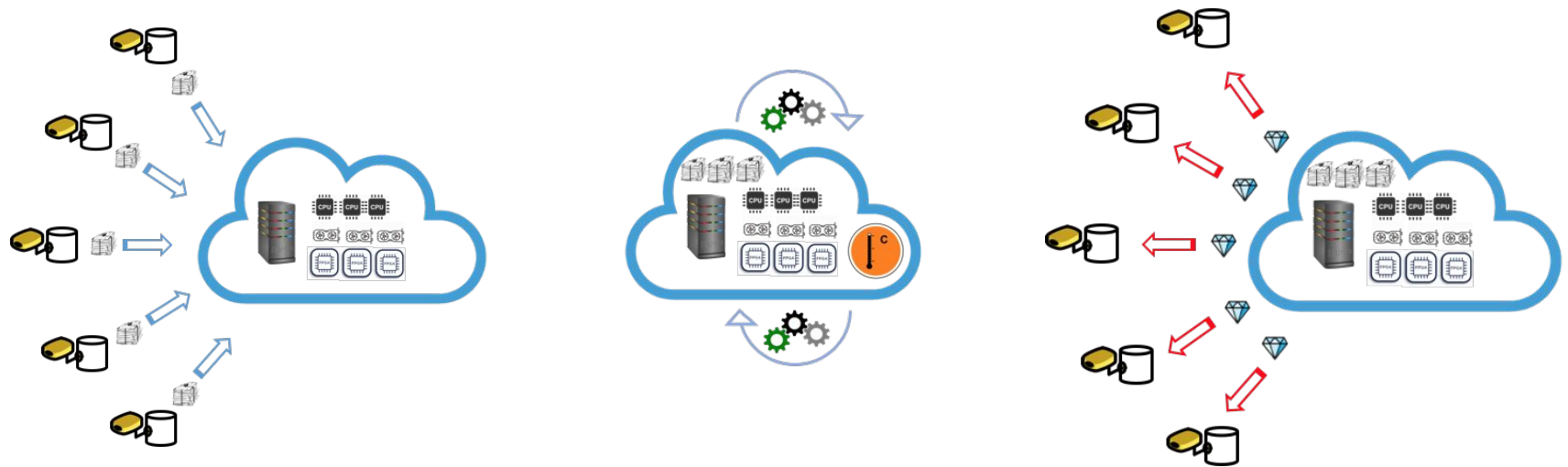
- Requires high memory to store large-scale training dataset (bigdata)
- Involves lots of iterations
- Has huge computational cost
- Is inherently sequential

Training is EXPENSIVE !

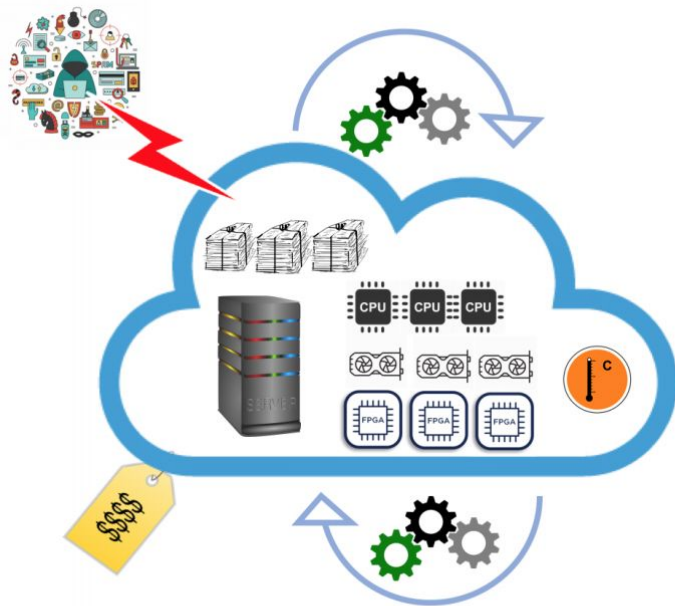


Centralized ML Framework

- 1 Collect data at Edge devices and transfer to HPC cloud server
- 2 Train ML model in the HPC cloud server using training algorithms
- 3 Deploy the trained model back at Edge devices for inference



Challenges of Centralized ML



At HPC cloud server

- Privacy at Risk
- High energy dissipation
- Expensive in cost \$\$\$\$
- High network latency

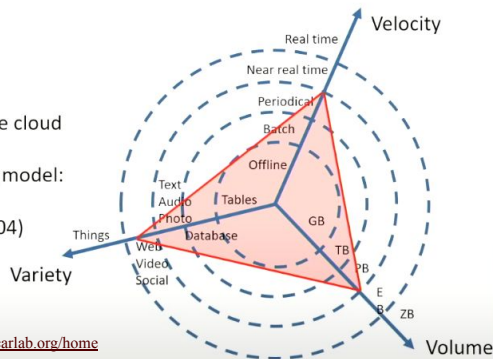
and, today there are Billions of connected small devices at Edge (IoT) having new requirements that can not be satisfied by Centralized framework

Paradigm Shift

Big Data Processing 1.0 (2005-2015)



- Push the data to the cloud
- Popular processing model:
MapReduce (OSDI'04)

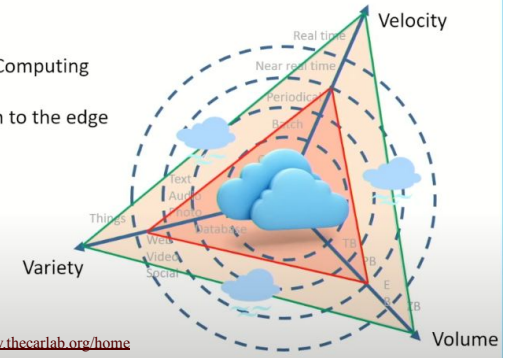


Picture Credits. <https://www.thecarlab.org/home>

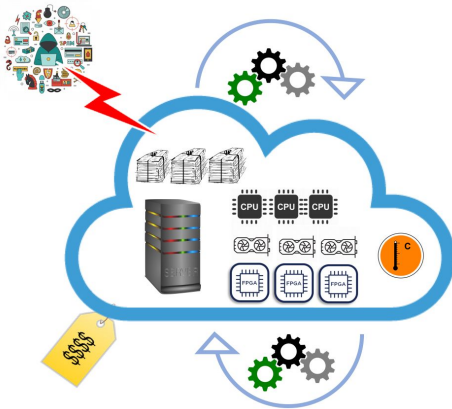
Big Data Processing 2.0 (2015-2025)



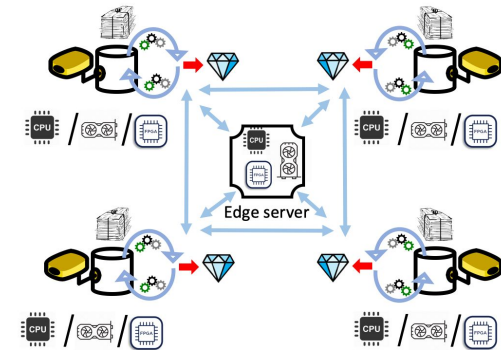
- Enabled by Edge Computing
- Push computation to the edge



Picture Credits. <https://www.thecarlab.org/home>

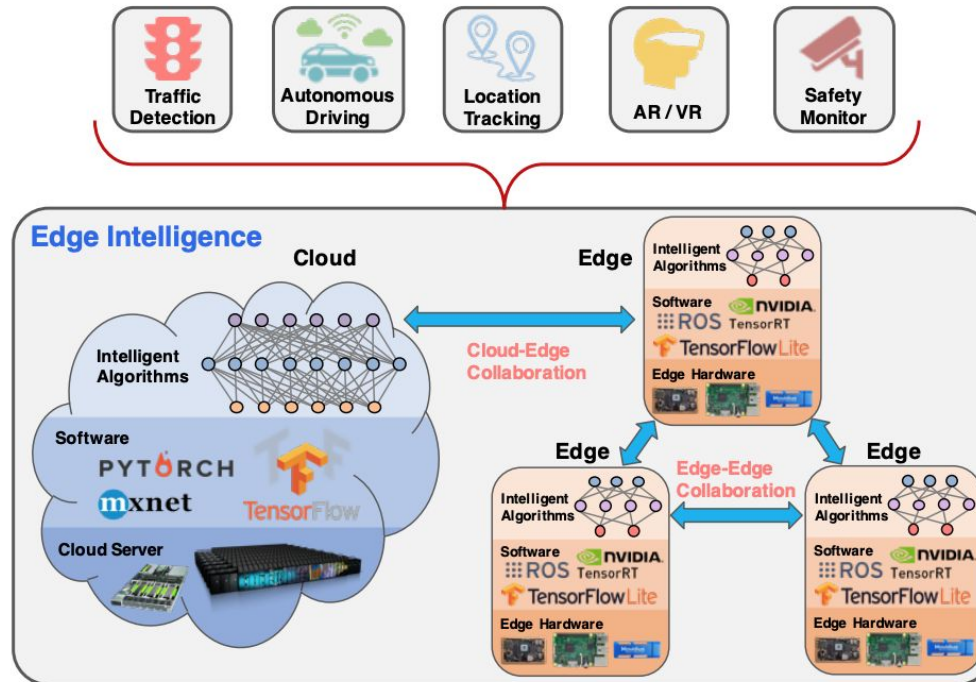


Cloud



Edge

Edge Intelligence

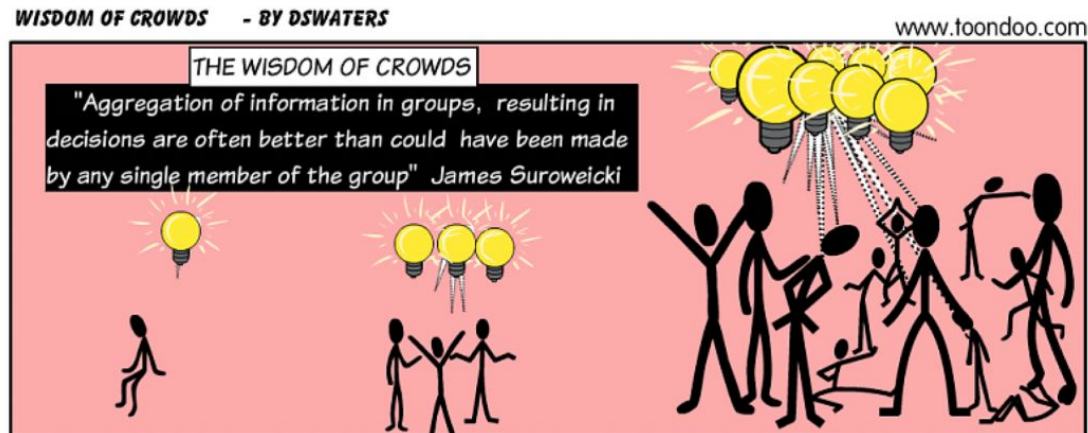


Zhang, Xingzhou, Yifan Wang, Sidi Lu, Liangkai Liu, and Weisong Shi. "Openel: An open framework for edge intelligence." In 2019 IEEE 39th International Conference on Distributed Computing Systems (ICDCS 2019)



Hypothesis

Decentralize Machine Learning via
Efficient and Scalable Algorithms and Architectures
for Distributed Edge Intelligence



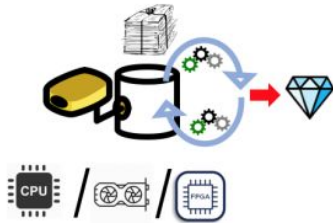
How to Decentralize ML?

Making training CHEAP



will REDUCE the need for HPC cloud server

Bring training on edge



Edge Device







- Has less storage
- Handles light computations
- Single low power CPU/GPU/FPGA
- Low cost \$\$



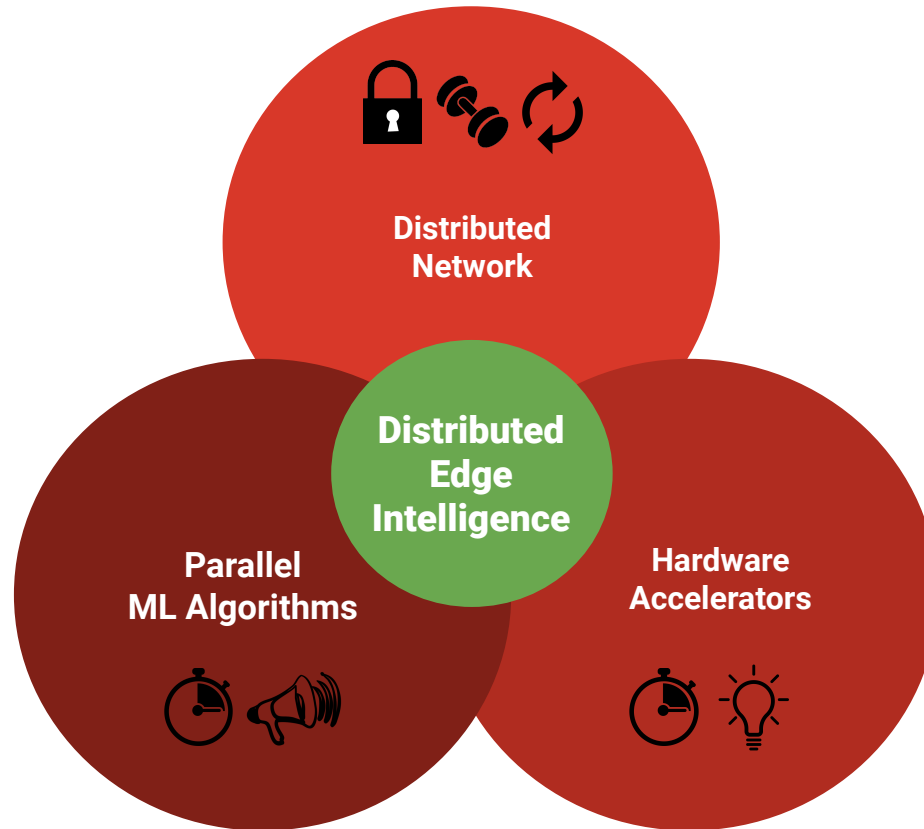
and AVOID centralization

Go for decentralized/distributed framework with parallel training

Requirements for Edge Intelligence

01	Protect Data Privacy		<ul style="list-style-type: none">• Keep data decentralized and local on devices• Design privacy-preserving ML algorithms
02	Reduce Latency		<ul style="list-style-type: none">• Efficient and scalable training algorithms• Cheap inference calculations to enable real-time analytics
03	Save Bandwidth		<ul style="list-style-type: none">• Communicate less data during training• Reduce synchronizations and idling during training
04	Energy Efficiency		<ul style="list-style-type: none">• Efficient computation and communication process• Build energy-efficient hardware accelerators for Green AI
05	Build Robust Models		<ul style="list-style-type: none">• Devise fault-tolerance for device failures or stragglers• Accurate and robust model to data perturbations
06	Streaming Data		<ul style="list-style-type: none">• Incremental federated learning to update the global model• Discard data after each update for memory and privacy

Research Focus



- | | | |
|----|----------------------|--|
| 01 | Protect Data Privacy | |
| 02 | Reduce Latency | |
| 03 | Save Bandwidth | |
| 04 | Energy Efficiency | |
| 05 | Build Robust Models | |
| 06 | Streaming Data | |







Research Contributions

1. Relaxed Synchronization for Parallel QP Problems
2. Householder Sketch for Machine Learning
3. Memory-efficient Framework for Distributed ML
4. Communication-efficient Framework for Scalable ML
5. Multiple FPGA-based System for Energy-efficient ML
6. Rapid Incremental Solver for Federated ML



Relaxed Synchronization for Parallel QP Problems



- 01 Protect Data Privacy 
- 02 Reduce Latency 
- 03 Save Bandwidth 
- 04 Energy Efficiency 
- 05 Build Robust Models 
- 06 Streaming Data 



Parallel QP Problems

Formulation

$$\min_{\mathbf{x}} \sum_{i=1}^p \left(\frac{1}{2} \mathbf{x}_i^T \mathbf{K}_i \mathbf{x}_i + \mathbf{c}_i^T \mathbf{x}_i \right)$$

subject to $\sum_{i=1}^p \mathbf{A}_i \mathbf{x}_i = \mathbf{b}$

where,

$\mathbf{x}_i, \mathbf{c}_i \in \mathbb{R}^{\frac{n}{p}}$, $\mathbf{b} \in \mathbb{R}^m$, $\mathbf{A}_i \in \mathbb{R}^{m \times \frac{n}{p}}$, and $\mathbf{A} = [\mathbf{A}_1, \dots, \mathbf{A}_p] \in \mathbb{R}^{m \times n}$

$\mathbf{K}_i \in \mathbb{R}^{\frac{n}{p} \times \frac{n}{p}}$ is symmetric positive definite matrix, and

$\mathbf{K} = \text{diag}(\mathbf{K}_1, \dots, \mathbf{K}_p) \in \mathbb{R}^{n \times n}$

n : Number of variables, m : Number of constraints,

p : Number of separable (parallel) sub-problems

Lagrangian \mathcal{L}_i for each sub-problem, $i = 1, \dots, p$

$$\mathcal{L}_i(\mathbf{x}_i, \boldsymbol{\beta}) = \left(\frac{1}{2} \mathbf{x}_i^T \mathbf{K}_i \mathbf{x}_i + \mathbf{c}_i^T \mathbf{x}_i \right) + \boldsymbol{\beta}^T (\mathbf{A}_i \mathbf{x}_i - \mathbf{b})$$

where, $\boldsymbol{\beta} \geq \mathbf{0}_m \in \mathbb{R}^m$ is a vector of Lagrangian dual variables

Parallel Dual Ascent

Dual sub-function: $g_i(\boldsymbol{\beta}) = \min_{\mathbf{x}_i} \mathcal{L}_i(\mathbf{x}_i, \boldsymbol{\beta})$

Dual problem: $\max_{\boldsymbol{\beta}} g(\boldsymbol{\beta})$



Parallel Dual Ascent steps

Gradient method - involves iterating through the following steps until convergence (error in β falls below stopping threshold)

Step 1: Minimization of Lagrangian for each sub-problem i

$$\begin{aligned}\mathbf{x}_i^{t+1} &= \arg \min_{\mathbf{x}_i} \mathcal{L}_i(\mathbf{x}_i, \beta^t) \\ &= -\left(\mathbf{K}_i\right)^{-1} \left(\mathbf{A}_i^T \beta^t + \mathbf{c}_i\right)\end{aligned}$$

solved in parallel using p machines which then broadcast local $\mathbf{A}_i \mathbf{x}_i^{t+1}$ during every iteration

Step 2: Dual variable update (Global Gathering of $\mathbf{A}_i \mathbf{x}_i^{t+1}$)

$$\beta^{t+1} = \beta^t + \eta \left(\sum_{i=1}^p \mathbf{A}_i \mathbf{x}_i^{t+1} - \mathbf{b} \right)$$

$\eta > 0$ is the step size, $\beta^{t=0} = \mathbf{0}_m$.

Synchronization is Necessary and Unavoidable in Step 2

- Results in idling
- Leads to waste of computing time



Relax Synchronization

We propose,

LSDA: Lazily Synchronized Dual Ascent

- Do not synchronize at every iteration
- Communicate data periodically
- Minimize frequency of communication

TSDA: Tightly Synchronized Dual Ascent

$$\beta^{t+1} = \beta^t + \eta \left(\sum_{i=1}^P \mathbf{A}_i \mathbf{x}_i^{t+1} - \mathbf{b} \right) = \beta^t + \eta \sum_{i=1}^P \left(\mathbf{A}_i \mathbf{x}_i^{t+1} - \frac{1}{p} \mathbf{b} \right)$$

LSDA: Lazily Synchronized Dual Ascent

$$\beta^{t+1} = \beta^t + \eta \sum_{i=1}^P \left(\mathbf{A}_i \mathbf{x}_i^{kP+1} - \frac{1}{p} \mathbf{b} \right), \quad kP \leq t < (k+1)P,$$

where, $k \in \mathbb{N}$, and $P \geq 1$ is the synchronization period

Synchronize across p parallel processes after every P iterations

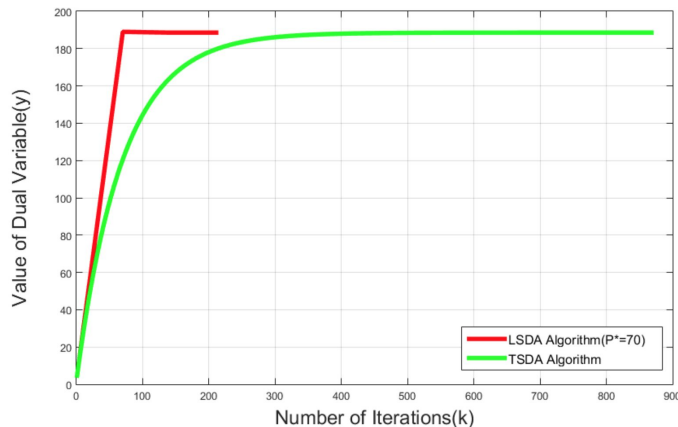


Results (1/2)

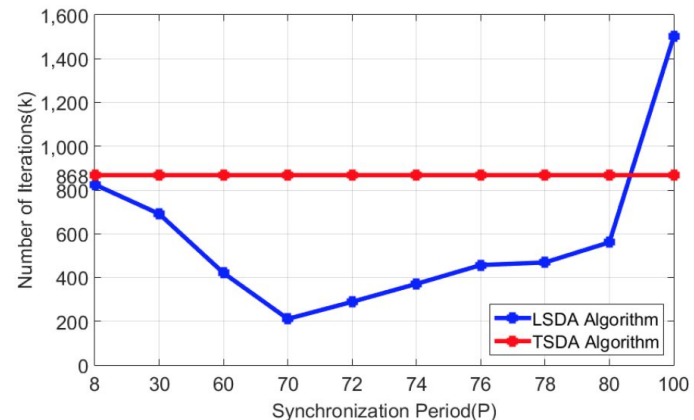
The data was synthetically generated with random values uniformly distributed over $[-1, 1]$. The problem specifics are as follows:

Experimental Setup

- 1) Number of instances in synthetic dataset, $d = 200,000$.
- 2) Step size, $\alpha = 0.27$.
- 3) Optimal Synchronization Period, $P^* = 70$.
- 4) Stopping threshold, $\epsilon = 10^{-5}$.
- 5) Cluster Size, $N = \{10, 20, 32, 40\}$.



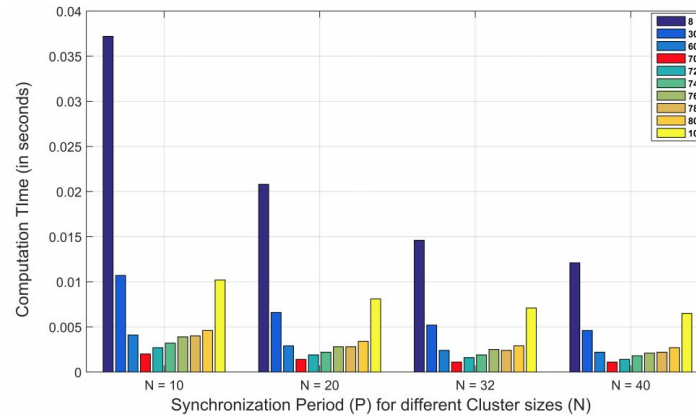
LSDA converges to the optimal solution of the dual variable significantly faster than the TSDA



Communication frequency is minimum at optimal synchronization period, $P=70$



Results (2/2)



Computation time is minimum at P=70
irrespective of number of parallel workers

Computing Performance Comparison between TSDA & LSDA algorithm

	TSDA algorithm	LSDA algorithm
No. of iteration	868	211
Sync. period	1	70
No. of Sync.	868 (=868/1) times	3 (=211/70) times
Comm. delay reduction	99.65%	
Speedup	160 times	



Research Contributions







✓ Relaxed Synchronization for Parallel QP Problems

[IEEE IPDPS'16] K. Lee, R. Bhattacharya, **J. Dass**, V. N. S. P. Sakuru, R. N. Mahapatra, "A Relaxed Synchronization Approach for Solving Parallel Quadratic Programming Problems with Guaranteed Convergence,"

- ❑ Householder Sketch for Machine Learning
- ❑ Memory-efficient Framework for Distributed ML
- ❑ Communication-efficient Framework for Scalable ML
- ❑ Multiple FPGA-based System for Energy-efficient ML
- ❑ Rapid Incremental Solver for Federated ML

Householder Sketch for Machine Learning



- | | | |
|----|----------------------|---|
| 01 | Protect Data Privacy |  |
| 02 | Reduce Latency |  |
| 03 | Save Bandwidth |  |
| 04 | Energy Efficiency |  |
| 05 | Build Robust Models |  |
| 06 | Streaming Data |  |

Sketching

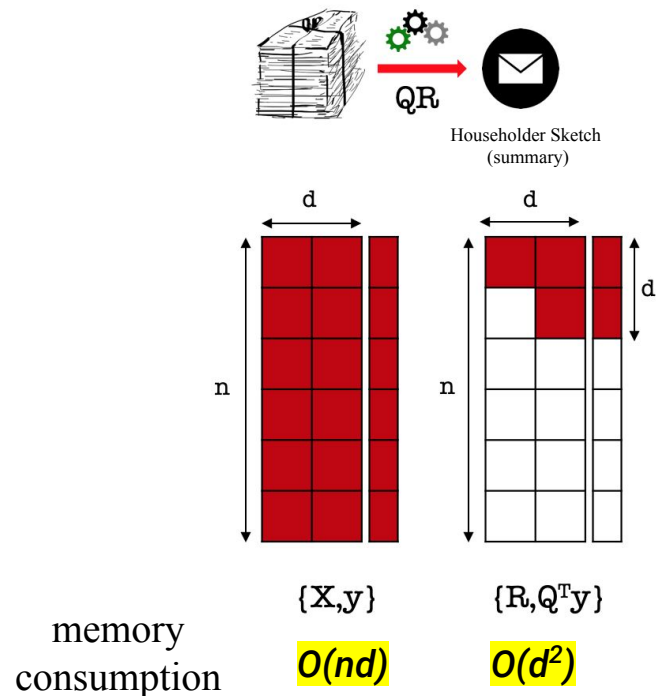
A compressed mapping of few or all data points (\mathbf{x}) in a data set to generate **data summary** called *Sketch* (\mathbf{s}) to preserve or approximate the covariance matrix, i.e.,

$$\mathbf{S}^T \mathbf{S} \approx \mathbf{X}^T \mathbf{X}$$



Householder Sketch

Theorem 3.2 (Householder Sketch). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the original data matrix, $\mathbf{y} \in \mathbb{R}^n$ be the corresponding output label or response vector, and $n \gg d$. Let $\mathbf{X} = \mathbf{QR}$ be Householder QR decomposition. Then, $(\mathbf{R}, \mathbf{Q}^T \mathbf{y})$ is a memory-efficient and theoretically accurate sketch of original data (\mathbf{X}, \mathbf{y}) such that $\mathbf{X}^T \mathbf{X} = \mathbf{R}^T \mathbf{R}$, and has memory footprint of $\left(\frac{d(d+3)}{2}\right)$ elements, computed in time $O(nd^2)$.



Householder-Sketch for LMS

Least-Mean-Squares

$$\min_{\mathbf{w}} f(\|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2) + g(\mathbf{w}).$$

$$\min_{\mathbf{w}} f(\|\mathbf{QR}\mathbf{w} - \mathbf{y}\|_2) + g(\mathbf{w}).$$

$$\|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2 = \|\mathbf{QR}\mathbf{w} - \mathbf{y}\|_2 = \|\mathbf{QR}\mathbf{w} - \mathbf{Q}\mathbf{Q}^T\mathbf{y}\|_2 = \|\mathbf{Q}\|_2 \|\mathbf{R}\mathbf{w} - \mathbf{Q}^T\mathbf{y}\|_2 = \|\mathbf{R}\mathbf{w} - \mathbf{Q}^T\mathbf{y}\|_2$$

(LMS)

Accurate Sketch $\mathbf{R}^T\mathbf{R} \approx \mathbf{X}^T\mathbf{X}$

(LMS-QR)

$$\mathbf{w} = \mathbf{X}^T\boldsymbol{\beta} = \mathbf{R}^T\mathbf{Q}^T\boldsymbol{\beta} = \mathbf{R}^T\bar{\boldsymbol{\beta}}$$

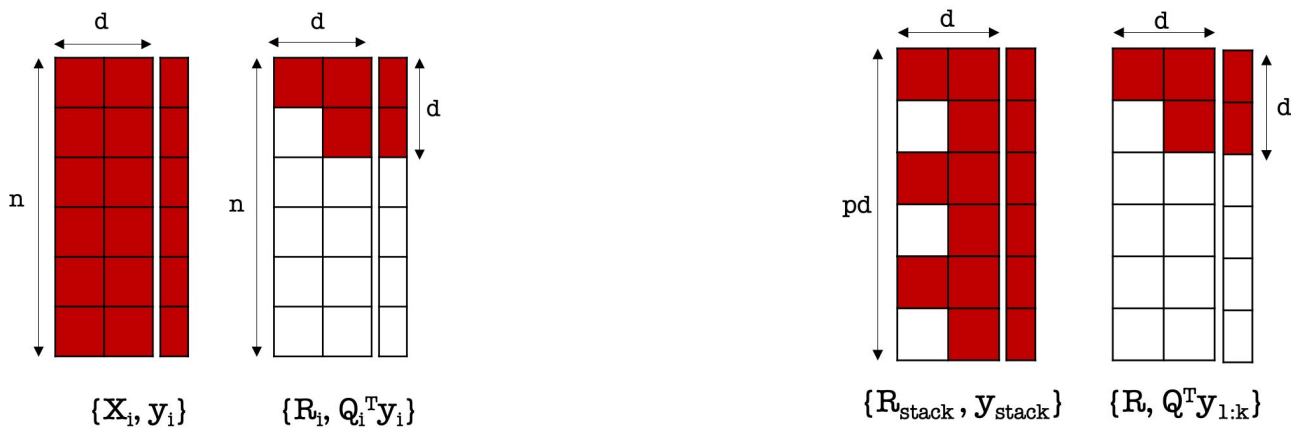
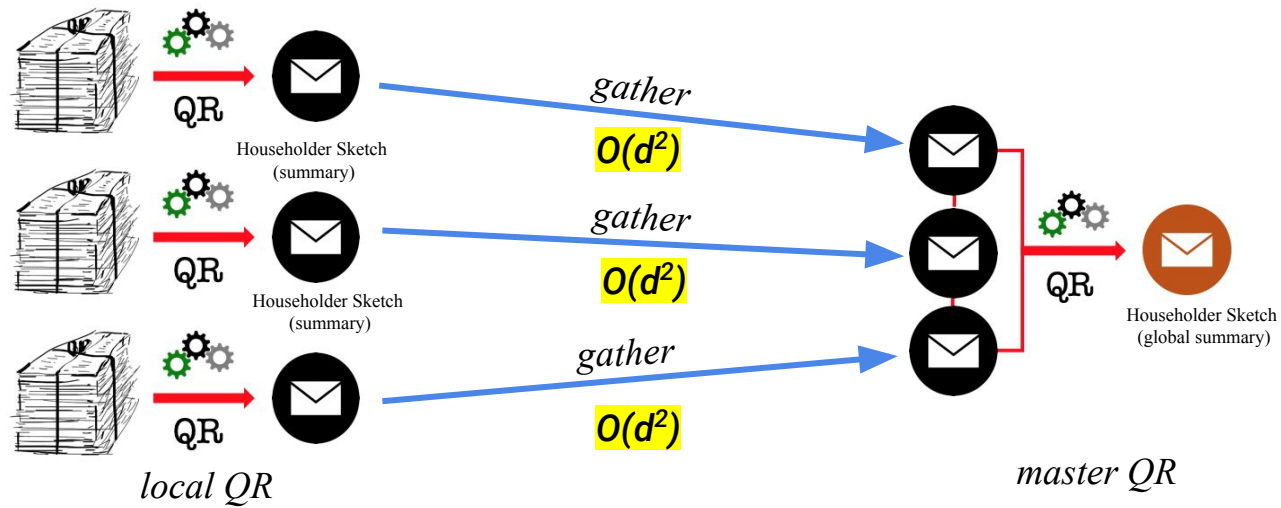


Distributed Householder Sketches

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \cdot \\ \cdot \\ \mathbf{X}_p \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_1 \mathbf{R}_1 \\ \mathbf{Q}_2 \mathbf{R}_2 \\ \cdot \\ \cdot \\ \mathbf{Q}_p \mathbf{R}_p \end{pmatrix} = \text{diag}(\mathbf{Q}_1, \dots, \mathbf{Q}_p) \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{R}_2 \\ \cdot \\ \cdot \\ \mathbf{R}_p \end{pmatrix}, \quad \mathbf{R}_{stack} = \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{R}_2 \\ \cdot \\ \cdot \\ \mathbf{R}_p \end{pmatrix} = \mathbf{Q}_M \mathbf{R}_M.$$

$$\mathbf{X} = \text{diag}(\mathbf{Q}_1, \dots, \mathbf{Q}_p) \mathbf{R}_{stack} = \underbrace{\text{diag}(\mathbf{Q}_1, \dots, \mathbf{Q}_p) \mathbf{Q}_M}_{\mathbf{Q}} \mathbf{R}_M$$





memory consumption

$O(nd)$

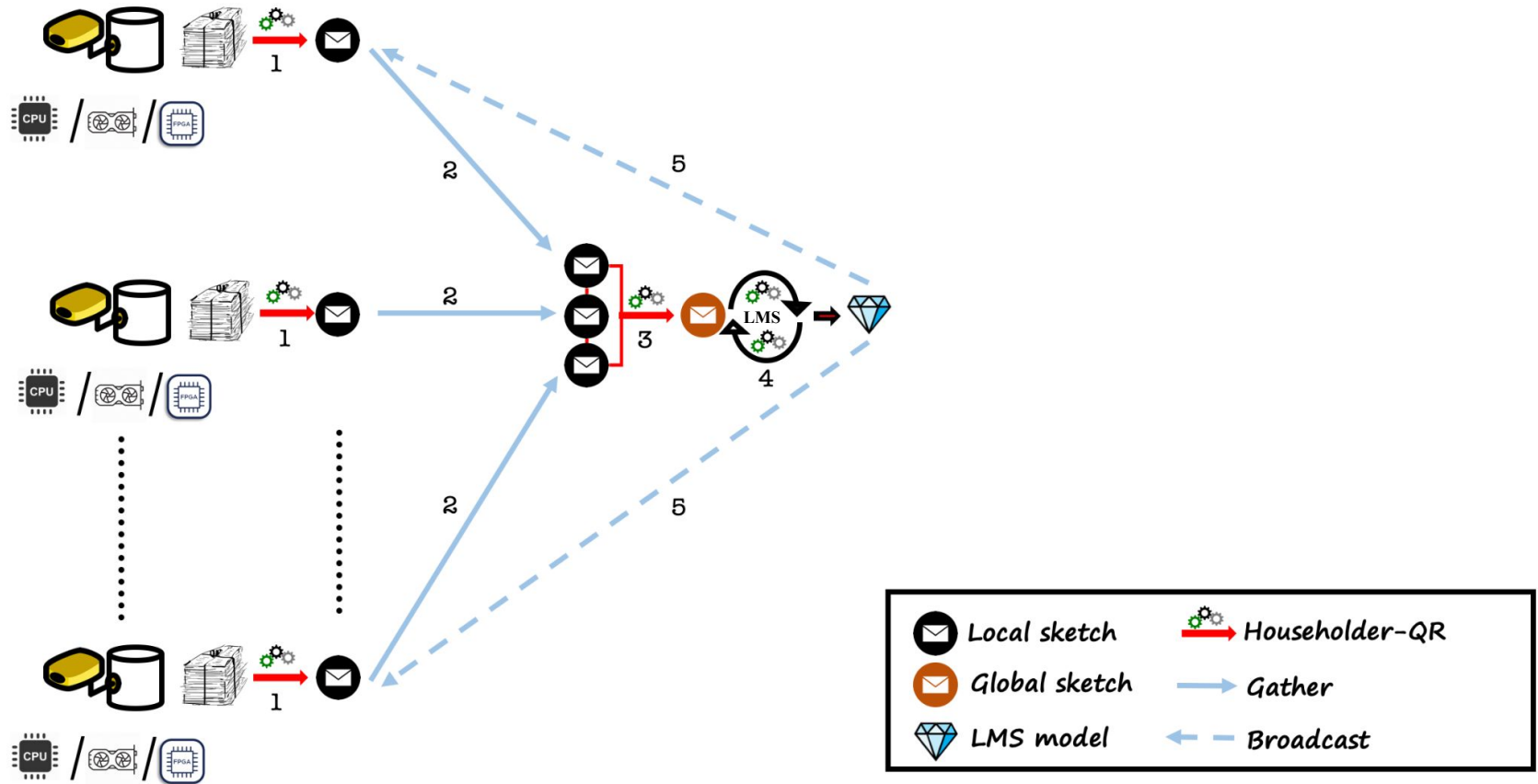
$O(d^2)$

$O(pd^2)$

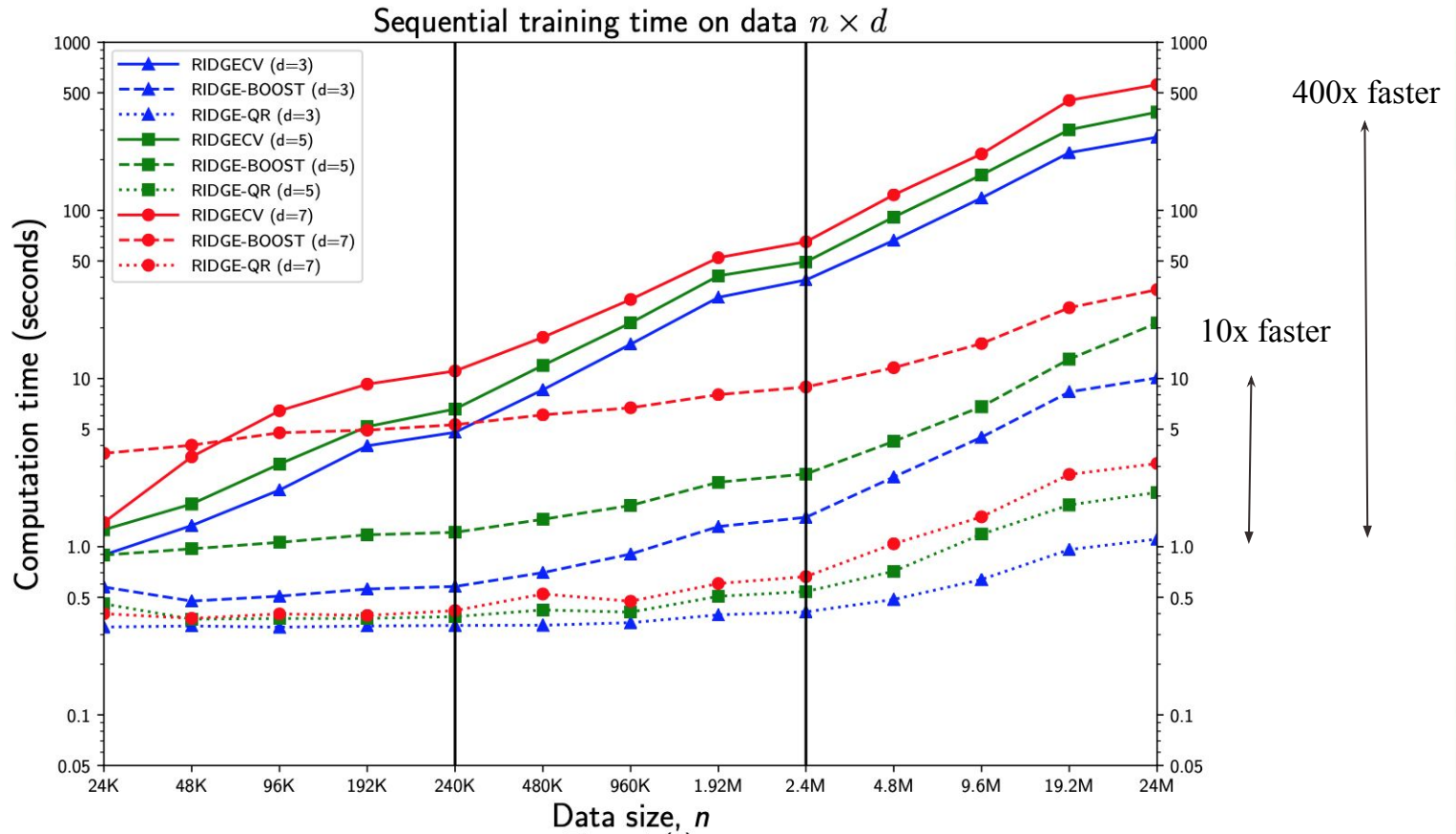
$O(d^2)$



Distributed Framework for LMS



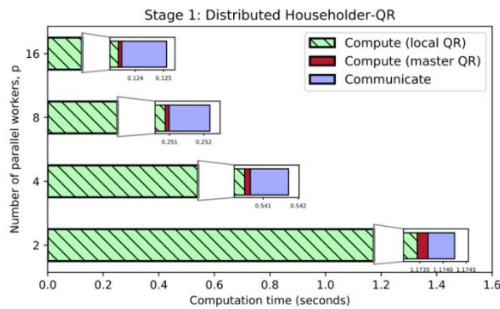
Results (1/2)



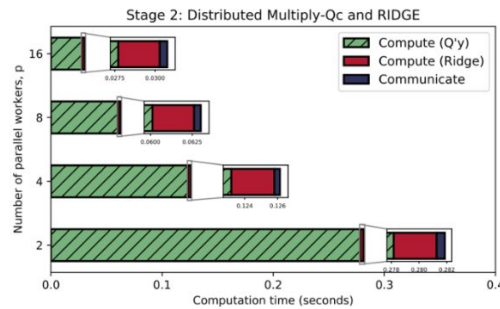
Sequential Training Time (RIDGE-QR vs others)



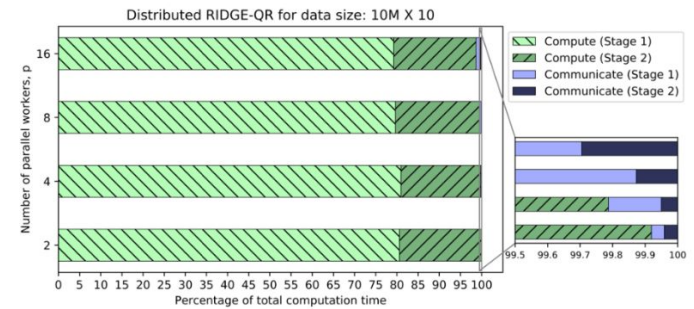
Results (2/2)



(a)



(b)



(c)

Breakdown of DISTRIBUTED RIDGE-QR (10M x 10 dataset) training time with zoomed insets depicting communication time (a): Stage 1: DISTRIBUTED HOUSEHOLDER-QR, (b): Stage 2: DISTRIBUTED MULTIPLY-QC and RIDGE, (c): Combined percentage



Research Contributions

✓ Relaxed Synchronization for Parallel QP Problems

[IEEE IPDPS'16] K. Lee, R. Bhattacharya, **J. Dass**, V. N. S. P. Sakuru, and R. N. Mahapatra, "A Relaxed Synchronization Approach for Solving Parallel Quadratic Programming Problems with Guaranteed Convergence,"

✓ Householder Sketch for Machine Learning







[ICML'21] **J. Dass**, and R. N. Mahapatra, "Householder Sketch for Accurate and Accelerated Least-Mean-Squares Solvers,"

- ❑ Memory-efficient Framework for Distributed ML
- ❑ Communication-efficient Framework for Scalable ML
- ❑ Multiple FPGA-based System for Energy-efficient ML
- ❑ Rapid Incremental Solver for Federated ML



Memory-efficient Framework for Distributed ML



- 01 Protect Data Privacy 
- 02 Reduce Latency 
- 03 Save Bandwidth 
- 04 Energy Efficiency 
- 05 Build Robust Models 
- 06 Streaming Data 

Distributed Machine Learning

K is kernel matrix which is generally non-separable

$$\text{(SVM)} \quad \min_{\alpha} \frac{1}{2} \alpha^T (\mathbf{X}\mathbf{X}^T + \rho\mathbf{I})\alpha + \mathbf{e}^T \alpha$$

subject to, $-\mathbf{I}_n \alpha \leq \mathbf{0}_n$

By using Householder-QR, and $\hat{\alpha} = \mathbf{Q}^T \alpha$

$$\text{(QRSVM)} \quad \min_{\hat{\alpha}} \frac{1}{2} \hat{\alpha}^T (\mathbf{R}\mathbf{R}^T + \rho\mathbf{I})\hat{\alpha} + \hat{\mathbf{e}}^T \hat{\alpha}$$

subject to, $-\mathbf{Q}\hat{\alpha} \leq \mathbf{0}_n$

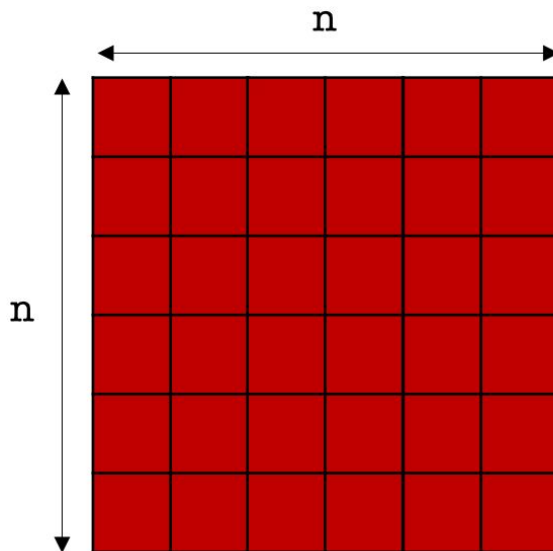


Memory-efficient Distributed ML

K is kernel matrix which is generally non-separable

$$\min_{\alpha} \frac{1}{2} \alpha^T (\mathbf{X}\mathbf{X}^T + \rho\mathbf{I}) \alpha + \mathbf{e}^T \alpha$$

subject to, $-\mathbf{I}_n \alpha \leq \mathbf{0}_n$



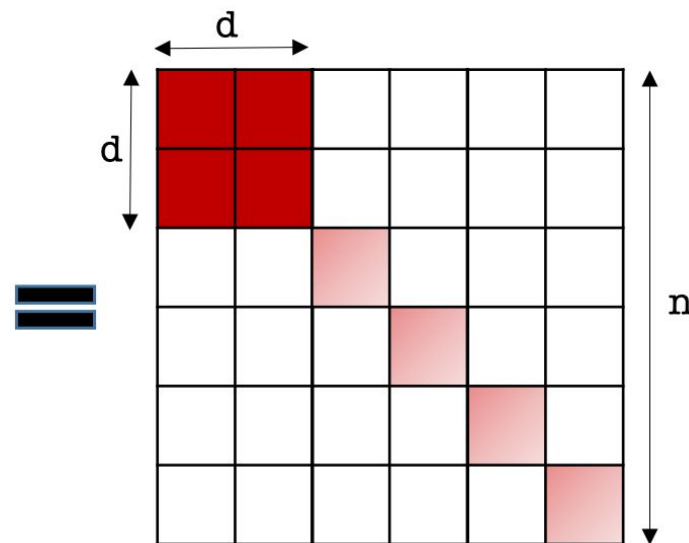
$\mathbf{X}\mathbf{X}^T + \rho\mathbf{I}$
Non - Separable

memory consumption

Dense $O(n^2)$

$$\min_{\hat{\alpha}} \frac{1}{2} \hat{\alpha}^T (\mathbf{R}\mathbf{R}^T + \rho\mathbf{I}) \hat{\alpha} + \hat{\mathbf{e}}^T \hat{\alpha}$$

subject to, $-\mathbf{Q}\hat{\alpha} \leq \mathbf{0}_n$



$\mathbf{R}\mathbf{R}^T + \rho\mathbf{I}$
Block-Separable

Sparse $O(d^2)$



Parallel Dual Ascent

Step 1: Minimization of Lagrangian - In Parallel

At core $i = \{1, \dots, p\}$

$$\hat{\alpha}_i^{t+1} = \mathbf{F}_i^{-1}(-\hat{\beta}_i^t + \hat{\mathbf{e}}_i)$$

where,

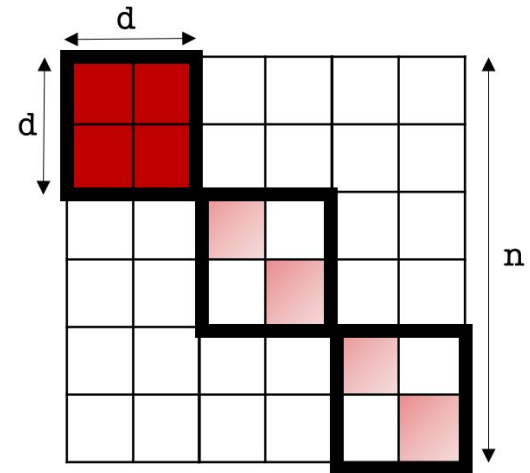
$$\mathbf{F}_i^{-1} = \begin{cases} \mathbf{F}_1^{-1} & \text{if } i = 1 \\ -2C & \text{if } i = 2, \dots, p \end{cases}$$

Step 2: Dual variable update - In Parallel

At core i

$$\hat{\beta}_i^{t+1} = \hat{\beta}_i^t + \eta^*(-\hat{\alpha}_i^{t+1})$$

where, η^* is the optimal step size, and $\hat{\beta}^t = \mathbf{Q}^T \beta^t$



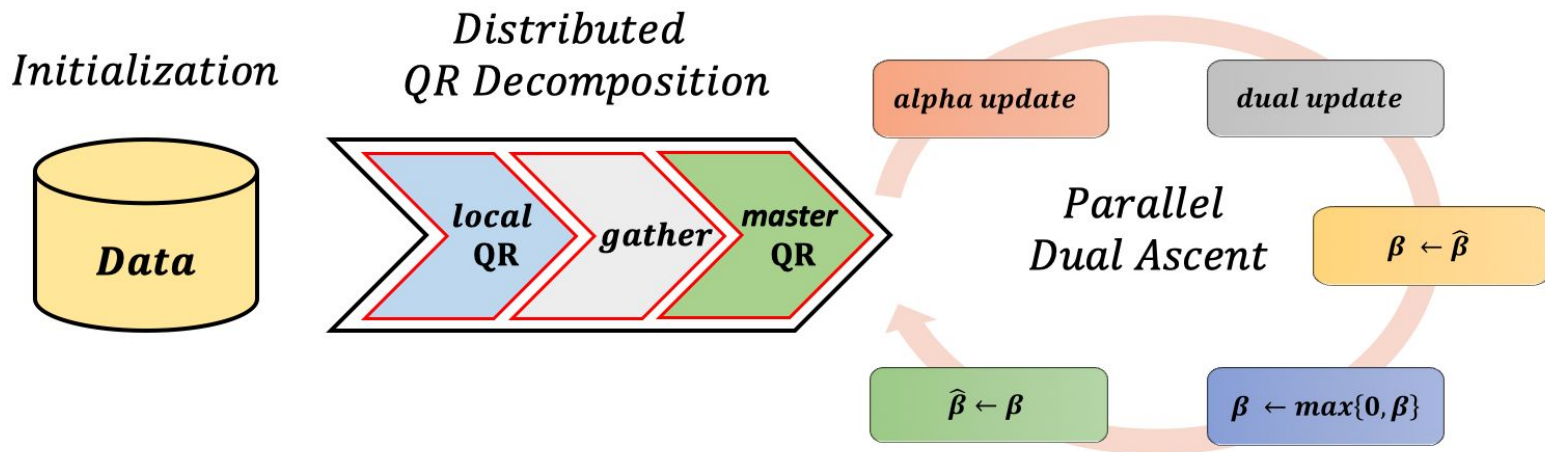
$$\mathbf{F} = \text{diag}(\mathbf{F}_1, \mathbf{F}_2, \mathbf{F}_3)$$

$$(\mathbf{R}\mathbf{R}^T + \rho\mathbf{I})$$

Householder-QR based QRSVM renders SVM separable into independent sub-problems!



Workflow



Results

Hardware

- Ada Supercomputing Cluster at TAMU
- Intel Xeon E5-2670 v2 (Ivy Bridge-EP), 10-core, 2.5GHz
- 64 GB/node and 16 cores/node
- Message-Passing Interface (MPI), InfiniBand interconnect

Dataset	n	d	Description
a9a	32560	123	predict annual income
covtype	464810	54	predict forest cover type

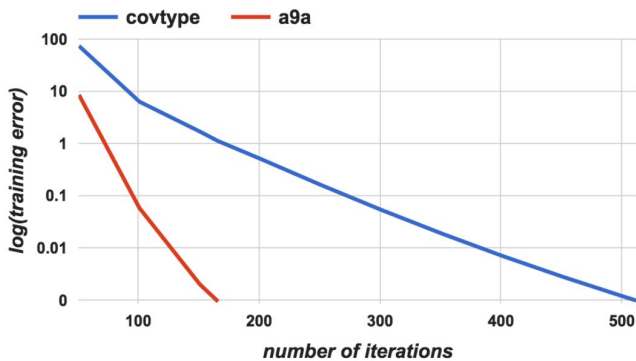
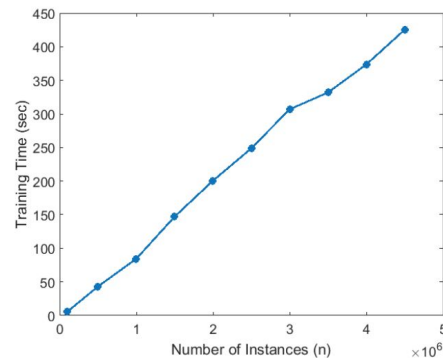
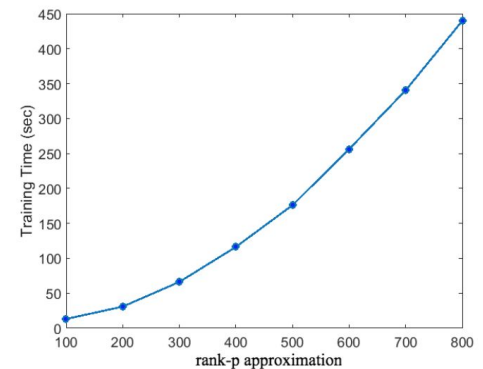


Figure: a9a: $t=166$, covtype: $t=512$, threshold= 10^{-3}

Convergence



Linear scaling with #samples



Quadratic scaling with rank



Research Contributions

✓ Relaxed Synchronization for Parallel QP Problems

[IEEE IPDPS'16] K. Lee, R. Bhattacharya, **J. Dass**, V. N. S. P. Sakuru, and R. N. Mahapatra, "A Relaxed Synchronization Approach for Solving Parallel Quadratic Programming Problems with Guaranteed Convergence"

✓ Householder Sketch for Machine Learning

[ICML'21] **J. Dass**, and R. N. Mahapatra, "Householder Sketch for Accurate and Accelerated Least-Mean-Squares Solvers"

✓ Memory-efficient Framework for Distributed ML

[CDCS'17] **J. Dass**, V. N. S. P. Sakuru, V. Sarin and R. N. Mahapatra, "Distributed QR Decomposition Framework for Training Support Vector Machines"

- ❑ Communication-efficient Framework for Scalable ML
- ❑ Multiple FPGA-based System for Energy-efficient ML
- ❑ Rapid Incremental Solver for Federated ML

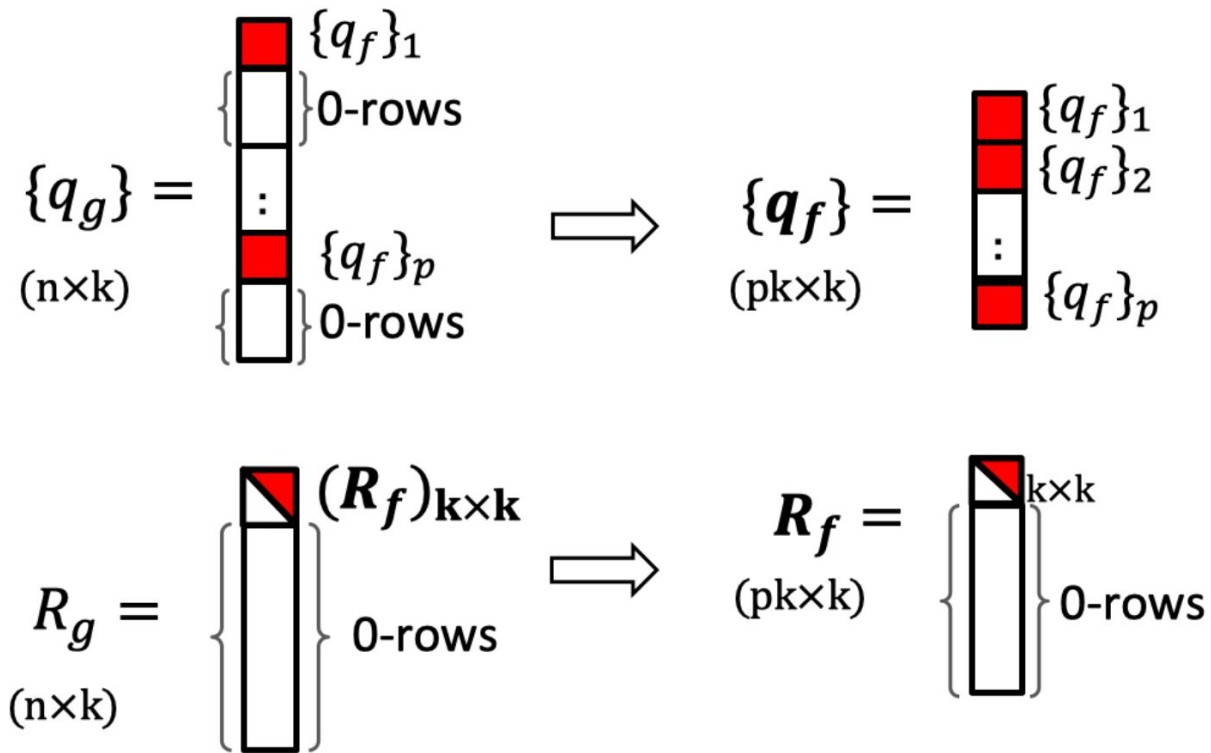


Communication-efficient Framework for Scalable ML



- | | | |
|----|----------------------|--|
| 01 | Protect Data Privacy | |
| 02 | Reduce Latency | |
| 03 | Save Bandwidth | |
| 04 | Energy Efficiency | |
| 05 | Build Robust Models | |
| 06 | Streaming Data | |

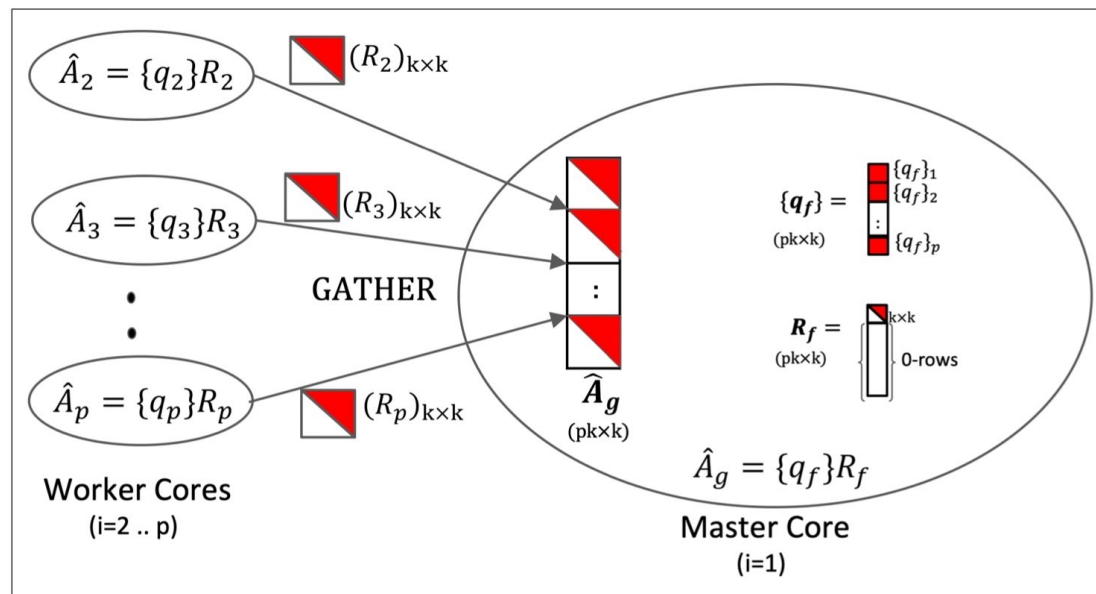
Improved Master QR



At master core: Memory improvement and representation of QR factors



Communication-Efficient Workflow

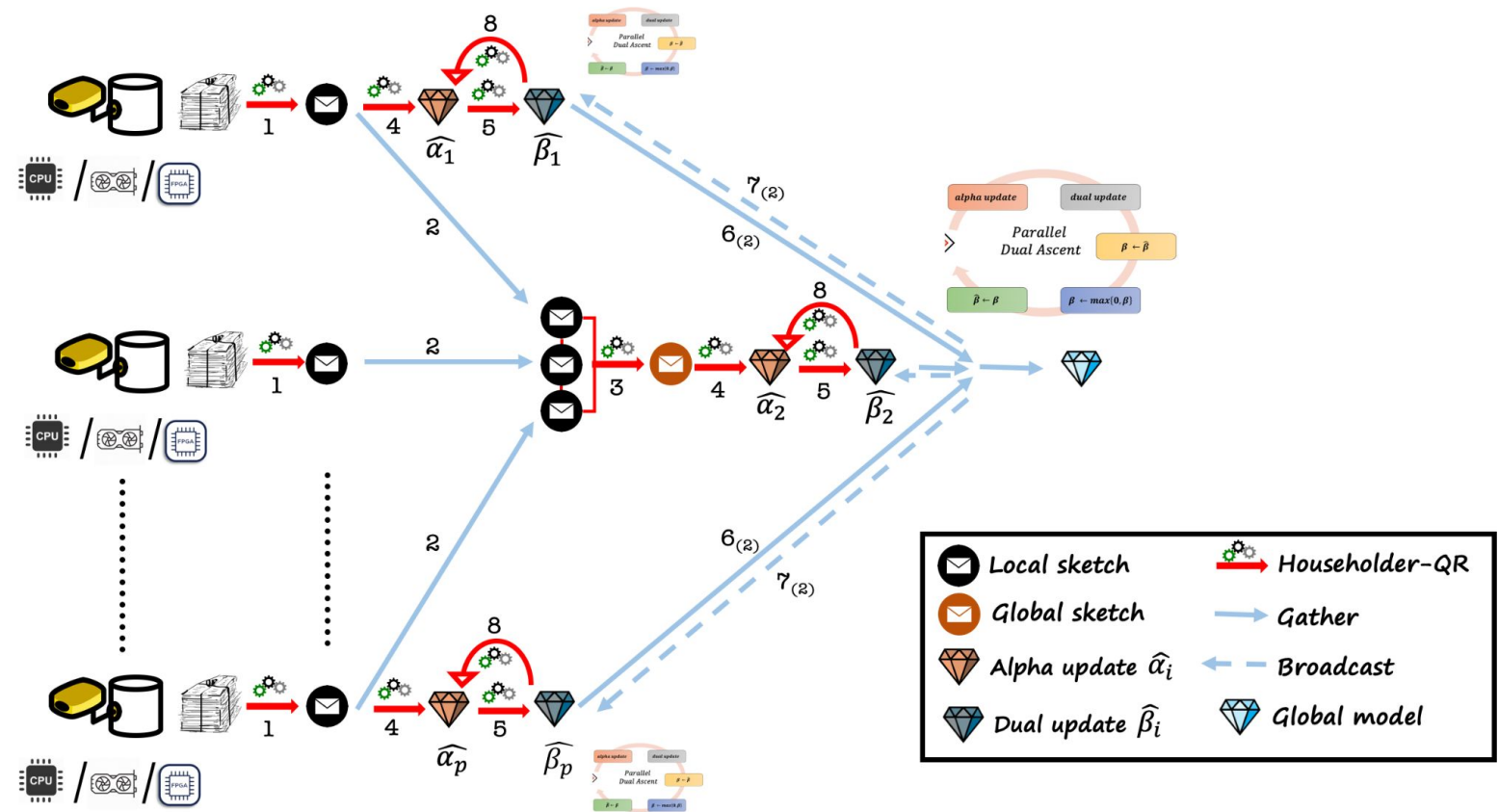


Communication Efficient \Rightarrow High Scalability

- During every iteration of parallel dual ascent, each **Gather** and **Scatter** involves communicating **negligibly small** $O(k)$ data volume per core with Master core, where, $k \ll \frac{n}{p}$
- Scalable with number of cores (p) and handle larger datasets than previous implementation



Communication-efficient Framework for Scalable ML



Results (1/2)

TABLE : Benchmark Dataset Description

Dataset	#training samples (n)	#features (d)	k-rank approx.
covtype.binary	464,810	54	64
webspam(unigram)	350,000	254	128
SUSY	5,000,000	18	128

TABLE : Scalability of distributed-QRSVM.

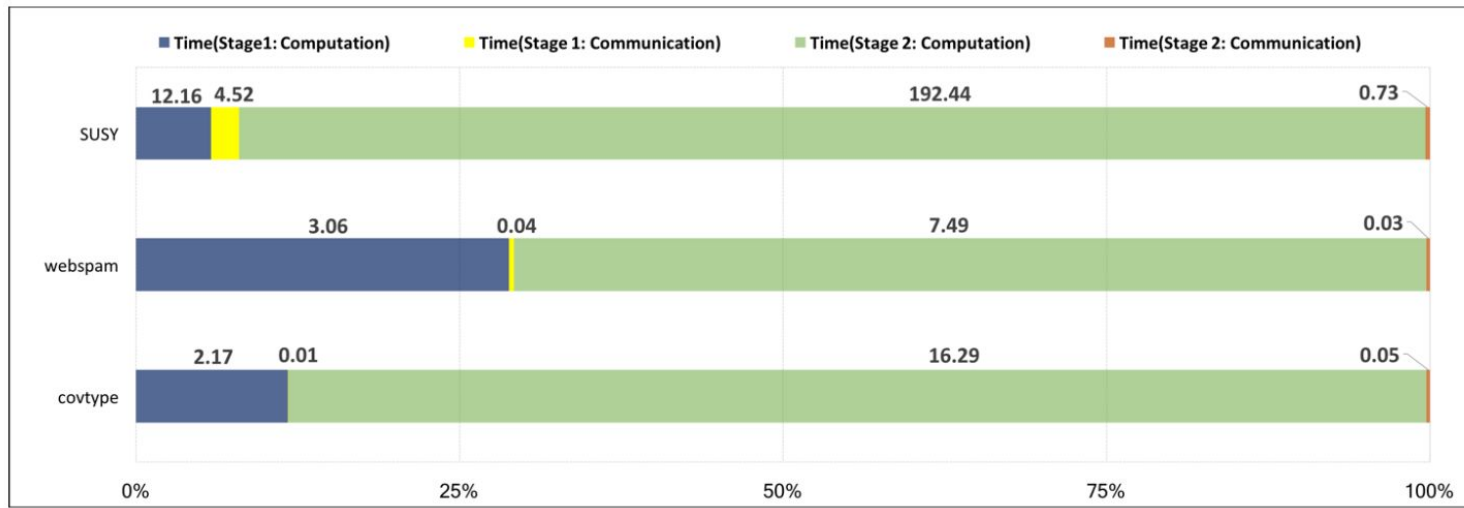
Training time (in seconds) and Speedup, S_p wrt sequential-QRSVM (S_p is shown in parenthesis).

Dataset	sequential	p=2	p=4	p=8	p=16	p=32	p=64
covtype	268 (1x)	132 (2x)	64 (4x)	33 (8x)	18 (15x)	10 (27x)	6 (45x)
webspam	258 (1x)	120 (2x)	58 (4x)	32 (8x)	19 (14x)	11 (23x)	9 (29x)
SUSY	28,614 (1x)	11,284 (2x)	4,405 (7x)	1,686 (17x)	804 (36x)	380 (75x)	210 (136x)

**Near linear scalability on large datasets
with increasing parallel workers**



Results (2/2)



Overall training time, T_{train} analysis (in seconds) for benchmarks covtype ($p = 16$), webspam ($p = 32$) and SUSY ($p = 64$).

Communication overhead is Negligible.



Research Contributions

✓ Relaxed Synchronization for Parallel QP Problems

[IEEE IPDPS'16] K. Lee, R. Bhattacharya, **J. Dass**, V. N. S. P. Sakuru, and R. N. Mahapatra, "A Relaxed Synchronization Approach for Solving Parallel Quadratic Programming Problems with Guaranteed Convergence"

✓ Householder Sketch for Machine Learning

[ICML'21] **J. Dass**, and R. N. Mahapatra, "Householder Sketch for Accurate and Accelerated Least-Mean-Squares Solvers"

✓ Memory-efficient Framework for Distributed ML

[IEEE ICDCS'17] **J. Dass**, V. N. S. P. Sakuru, V. Sarin and R. N. Mahapatra, "Distributed QR Decomposition Framework for Training Support Vector Machines"

✓ Communication-efficient Framework for Scalable ML

[IEEE TPDS'18] **J. Dass**, V. Sarin and R. N. Mahapatra "Fast and Communication-Efficient Algorithm for Distributed Support Vector Machine Training"

❑ Multiple FPGA-based System for Energy-efficient ML

❑ Rapid Incremental Solver for Federated ML

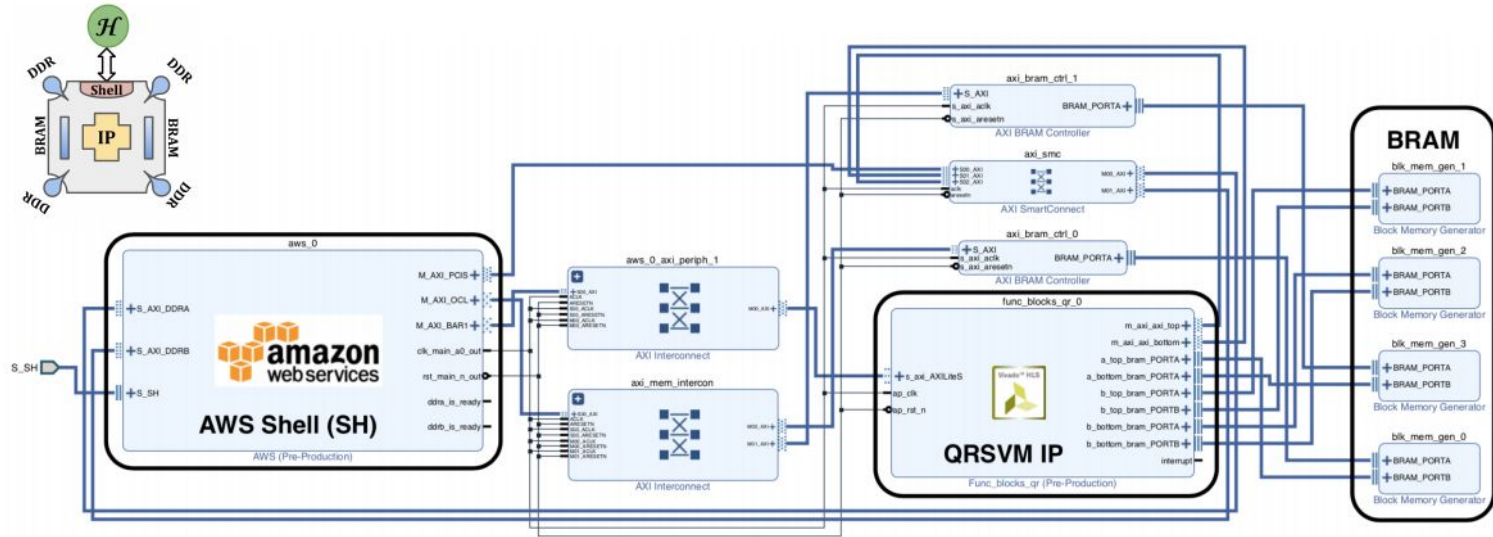


Multiple FPGA-based System for Energy-efficient ML



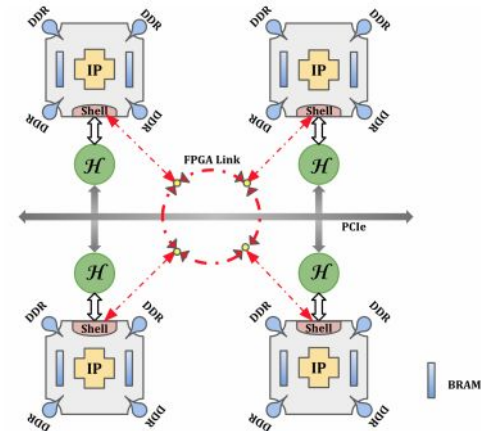
- | | | |
|----|----------------------|--|
| 01 | Protect Data Privacy | |
| 02 | Reduce Latency | |
| 03 | Save Bandwidth | |
| 04 | Energy Efficiency | |
| 05 | Build Robust Models | |
| 06 | Streaming Data | |

Multiple FPGA-based System



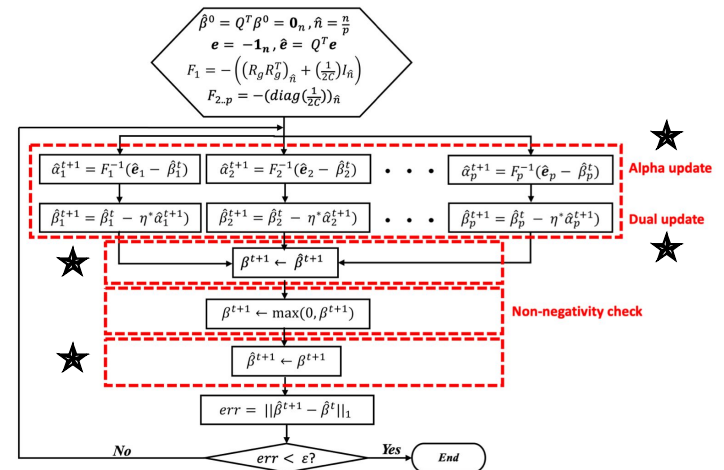
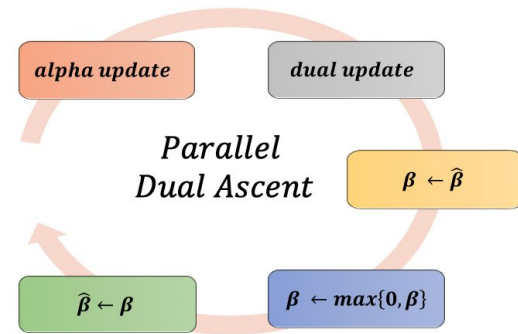
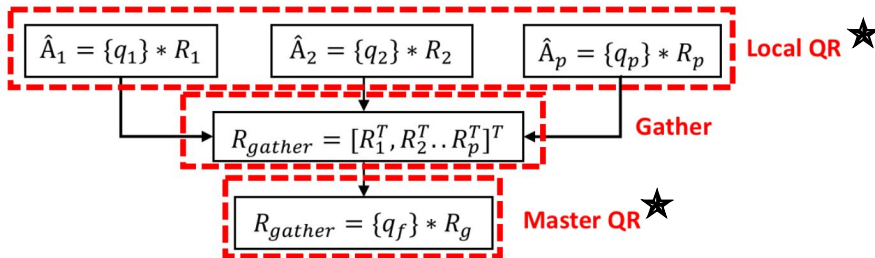
Multiple FPGA Network

- Each edge device comprises of FPGA IP logic + Host processor
- Illustration for $p = 4$ edges in a network
- A single QRSVM IP is synthesized per FPGA device to operate at clock frequency of 125 MHz with 39 Watts of power
- Computational workload is entirely with the FPGA while communication is through Host processor (PCIe)
- Each FPGA handles maximum of $256K$ samples



Computational Workflow

Distributed QR Decomposition



Inner Product and SAXPY

Algorithm 6.1 $\{q_i\}, R_i \leftarrow \hat{A}_i$, via Householder algorithm

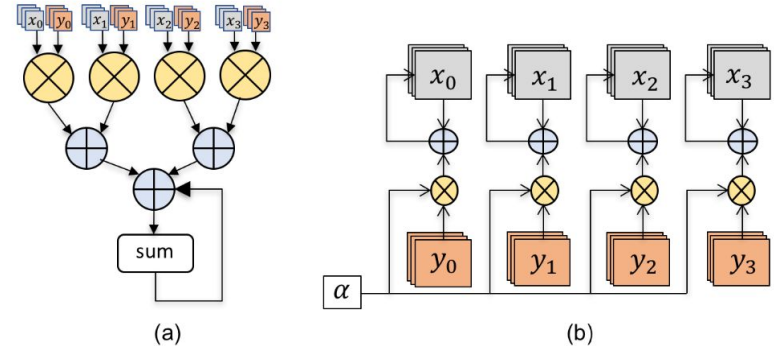
Data: Matrix \hat{A}_i
 $Q_{\hat{n} \times k}, (\hat{A}_i)_{\hat{n} \times k}$ $\triangleright \hat{n}$: samples per compute unit
1 **for** $j \leftarrow 1$ **to** k **do**
2 $q_{ij} \leftarrow \hat{A}_i(j : \hat{n}, j)$
3 $q_{ij}(1) \leftarrow q_{ij}(1) + \text{sign}(q_{ij}(1)) \times \|q_{ij}\|_2$ \triangleright scalar update
4 $q_{ij} \leftarrow \frac{q_{ij}}{\|q_{ij}\|_2}$ \triangleright vector normalization
5 $\hat{A}_i(j : \hat{n}, j : k) \leftarrow \hat{A}_i(j : \hat{n}, j : k) - 2q_{ij} \langle q_{ij}, \hat{A}_i(j : \hat{n}, j : k) \rangle$ \triangleright Algorithm 6.2
6 $R_i = \hat{A}_i(j : \hat{n}, j : k)$
7 **end**
8 $\{q_i\} \leftarrow \{q_{i1}, q_{i2}, \dots, q_{ik}\}$ \triangleright set of k-reflectors

Algorithm 6.2 Computing Step 6 in Algorithm 6.1

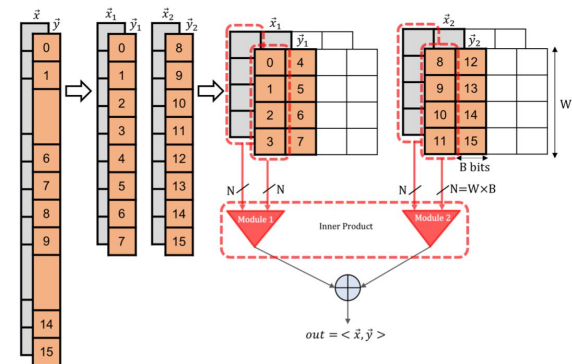
Data: Matrix \hat{A}_i
 $\hat{A}_i(j : \hat{n}, j : k) \leftarrow \hat{A}_i(j : \hat{n}, j : k) - 2q_{ij} \langle q_{ij}, \hat{A}_i(j : \hat{n}, j : k) \rangle$
1 **for** $m \leftarrow j$ **to** k **do**
2 $a_{BRAM} \leftarrow \hat{A}_i(j : \hat{n}, m)$ \triangleright Load into BRAM
3 Compute $sum = \langle q_{ij}, a_{BRAM} \rangle$ \triangleright Inner Product
4 $a_{BRAM} \leftarrow a_{BRAM} - 2 \times sum \times q_{ij}$ \triangleright saxpy
5 $\hat{A}_i(j : \hat{n}, m) \leftarrow a_{BRAM}$ \triangleright Write to DDR
6 **end**

$$\beta = Q\hat{\beta} = \text{diag}(Q_1, Q_2, \dots, Q_p) \times Q_p \times \hat{\beta}$$

$$\hat{\beta}_{i[j:\hat{n}]} \leftarrow \hat{\beta}_{i[j:\hat{n}]} - 2q_{ij} \langle q_{ij}, \hat{\beta}_{i[j:\hat{n}]} \rangle \quad j \leftarrow k \text{ to } 1$$



SIMD hardware kernels for (a) inner product, $\langle x, y \rangle$ (b) saxpy, $x = x + \alpha y$. The vectorized kernels process $W = 4$ elements in each pass with pipeline stages, $D=2$ (inner product)



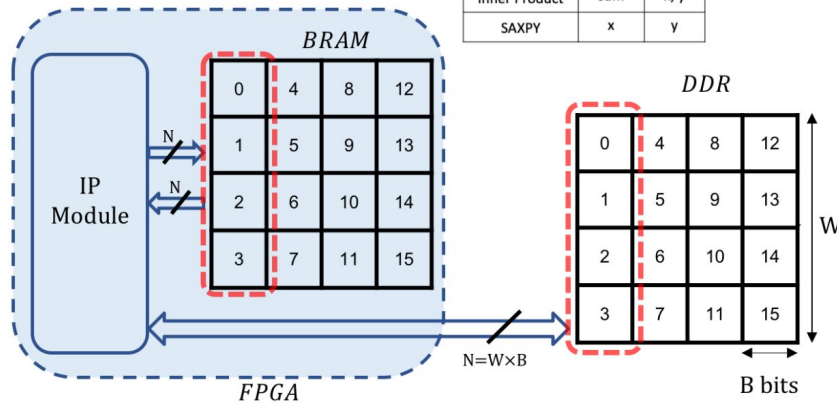
Doubling Throughput for Inner Product



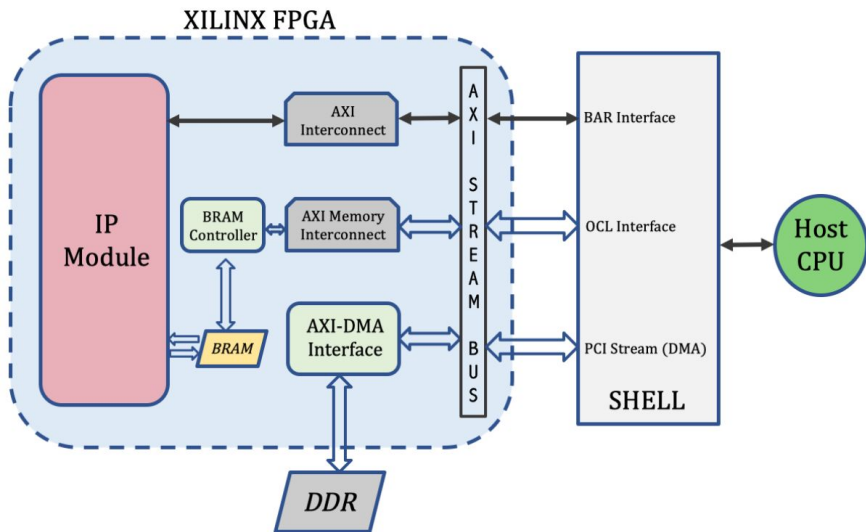
Data Layout + Memory Interface

Memory Assignment

Kernel	BRAM	DDR
Inner Product	sum	x, y
SAXPY	x	y



Data layout in column-major order and memory interface for on-chip Block RAM (Full-duplex) and off-chip DDR (Half-duplex) with the IP. **4MB to synthesize 2x2MB BRAMs on each FPGA. Can store 256K data samples**



FPGA Block Diagram with memory interfacing. Maximum bus width supported by this interface is $N = 1024$ bits. Support for double-precision ($B=64$ bits) floating point numbers without any loss of accuracy when compared to software implementation, the **maximum parallel compute units**, $W = \lfloor N/B \rfloor = 16$



Results (1/3)

Benchmark	Application	#samples (n)	#features (d)	k-rank
MNIST	Image	60,000	780	128
Skin	Health	200,000	3	64
Webspam	Email	350,000	254	128
Covtype	Geography	464,810	54	64
SUSY	Physics	2,000,000	18	128

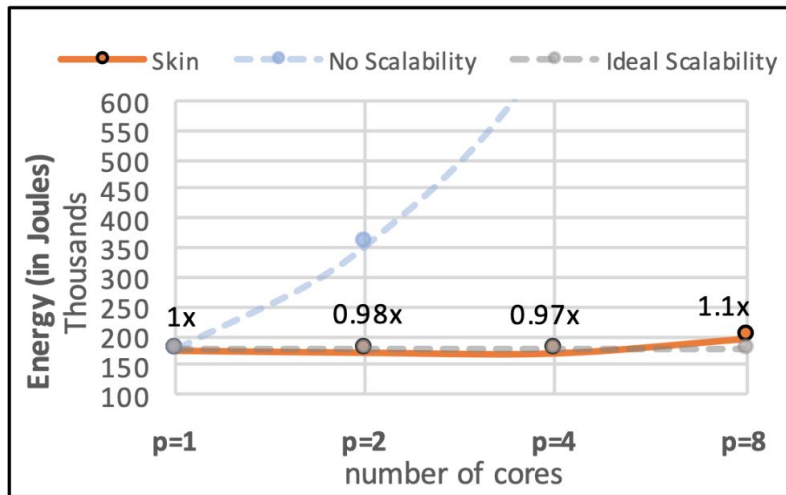
Datasets

Resource	BRAM	DSP	FF	LUT
Used	1405	1221	545248	449113
Available	2160	6840	2363536	1181768
Utilization	65%	18%	23%	38%

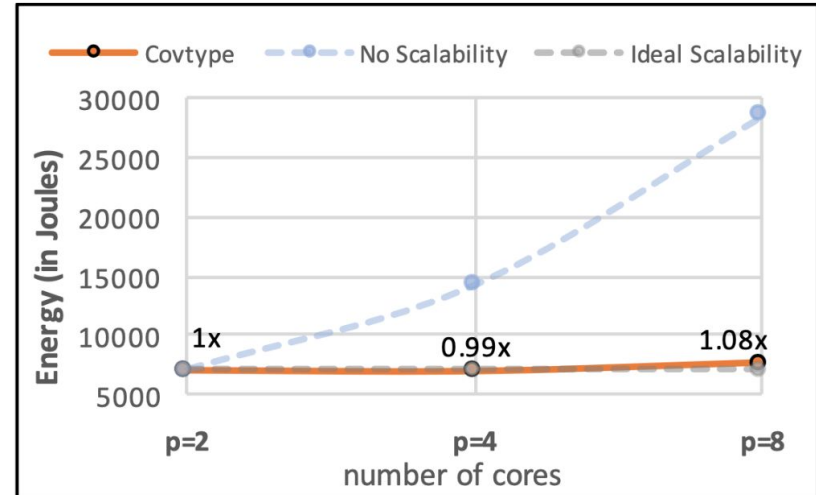
Utilization for FPGA Xilinx Virtex xcvu9p-flgb2104-2-i



Results (2/3)



(a)



(b)

Energy consumption Analysis of Multi-FPGA system under strong scaling scenario
Benchmarks: (a) Skin and (b) Covtype

Nearly CONSTANT (ideal) energy consumption across #FPGA units
Fully-Parallel Implementation



Results (3/3)

TABLE 4
Comparison With Embedded Edge Processor (ARM Cortex A15) and Cloud Processor (Broadwell) Platforms

#units	MNIST						Skin					
p	T_p^{FPGA}	T_p^{ARM}	T_p^{Broad}	E_p^{FPGA}	E_p^{ARM}	E_p^{Broad}	T_p^{FPGA}	T_p^{ARM}	T_p^{Broad}	E_p^{FPGA}	E_p^{ARM}	E_p^{Broad}
1	10.92	31.06	18.78	0.43	0.44	2.72	4,536	21,773	7,167	177	305	1,039
2	5.84	20.12	8.25	0.45	0.56	2.40	2,228	6,121	3,093	174	172	897
4	3.58	12.9	4.35	0.56	0.72	2.52	1,108	3,044	1,607	173	170	932

#units	Webspam						Covtype					
p	T_p^{FPGA}	T_p^{ARM}	T_p^{Broad}	E_p^{FPGA}	E_p^{ARM}	E_p^{Broad}	T_p^{FPGA}	T_p^{ARM}	T_p^{Broad}	E_p^{FPGA}	E_p^{ARM}	E_p^{Broad}
1	-	895	236.54	-	12.5	34.30	-	1079	292.63	-	15	42.43
2	76.14	477	133.99	5.9	13.4	38.86	91.45	520	160.08	7.1	14.5	46.42
4	39.36	254	65.92	6.1	14.2	38.23	45.36	251	77.19	7	14	44.77

#units	#samples	SUSY					
p	n	T_p^{FPGA}	T_p^{ARM}	T_p^{Broad}	E_p^{FPGA}	E_p^{ARM}	E_p^{Broad}
1	250K	108.08	2452	171.29	4.2	34.3	24.84
2	500K	131.02	3131	232.01	10.2	87.7	67.28
4	1M	176.18	4189	319.03	27.5	234	185

In a multiple compute system, #units, p , corresponds to #FPGA units (QRSVM IP cores), #ARM processors, and #Broadwell processors. Training time (in s), T_p^{FPGA} , T_p^{ARM} , and T_p^{Broad} . Energy consumption (in kJ), E_p^{FPGA} , E_p^{ARM} , and E_p^{Broad}

Broadwell at 145W
ARM at 14W
FPGA at 39W

Multi-FPGA implementation is upto
1.7x faster and **6x lower energy** than the cloud processor (Broadwell)
3x-24x faster and **2x-8x lower energy** than edge processor (ARM)



Research Contributions

✓ Relaxed Synchronization for Parallel QP Problems

[IEEE IPDPS'16] K. Lee, R. Bhattacharya, **J. Dass**, V. N. S. P. Sakuru, and R. N. Mahapatra, "A Relaxed Synchronization Approach for Solving Parallel Quadratic Programming Problems with Guaranteed Convergence"

✓ Householder Sketch for Machine Learning

[ICML'21] **J. Dass**, and R. N. Mahapatra, "Householder Sketch for Accurate and Accelerated Least-Mean-Squares Solvers"

✓ Memory-efficient Framework for Distributed ML

[IEEE ICDCS'17] **J. Dass**, V. N. S. P. Sakuru, V. Sarin and R. N. Mahapatra, "Distributed QR Decomposition Framework for Training Support Vector Machines"

✓ Communication-efficient Framework for Scalable ML

[IEEE TPDS'18] **J. Dass**, V. Sarin and R. N. Mahapatra "Fast and Communication-Efficient Algorithm for Distributed Support Vector Machine Training"

✓ Multiple FPGA-based System for Energy-efficient ML







[ACM FPGA'19 | IEEE TC'20] **J. Dass**, Y. Narawane, R. N. Mahapatra, and V. Sarin, "Distributed Training of Support Vector Machine on a Multiple-FPGA System"

❑ Rapid Incremental Solver for Federated ML



Rapid Incremental Solver for Federated ML



- 01 Protect Data Privacy 
- 02 Reduce Latency 
- 03 Save Bandwidth 
- 04 Energy Efficiency 
- 05 Build Robust Models 
- 06 Streaming Data 

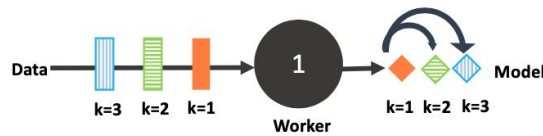


Incremental Federated ML Requirements

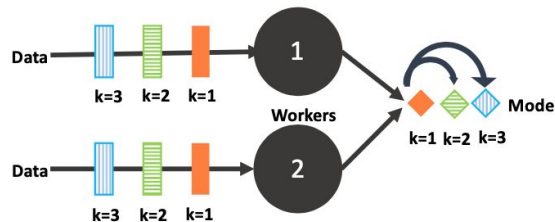
1. Data should **never be saved** on a centralized server **nor shared** among peers
2. Data samples must **not be stored** between successive model updates, i.e., update using **current data only**
3. **No retraining** from scratch is allowed
4. **Accurately** solve for global model in **collaboration** with other workers
5. **Robustness** and **Fault-tolerant** to straggling/offline workers



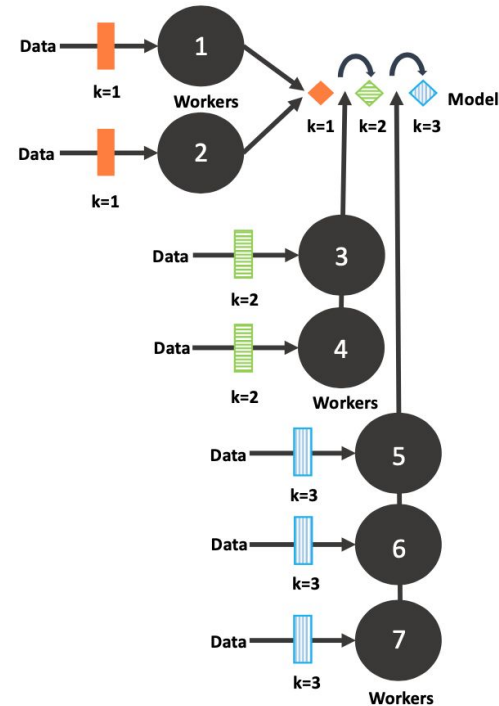
RIVER Setups



(a) Stream



(b) Tributary

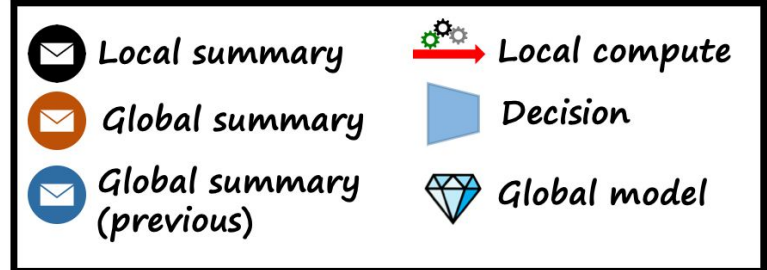
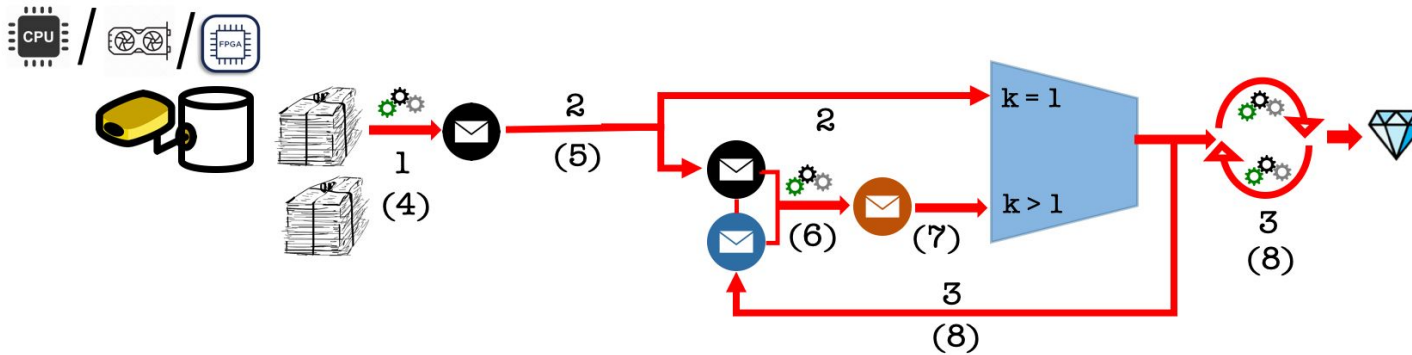


(c) Basin

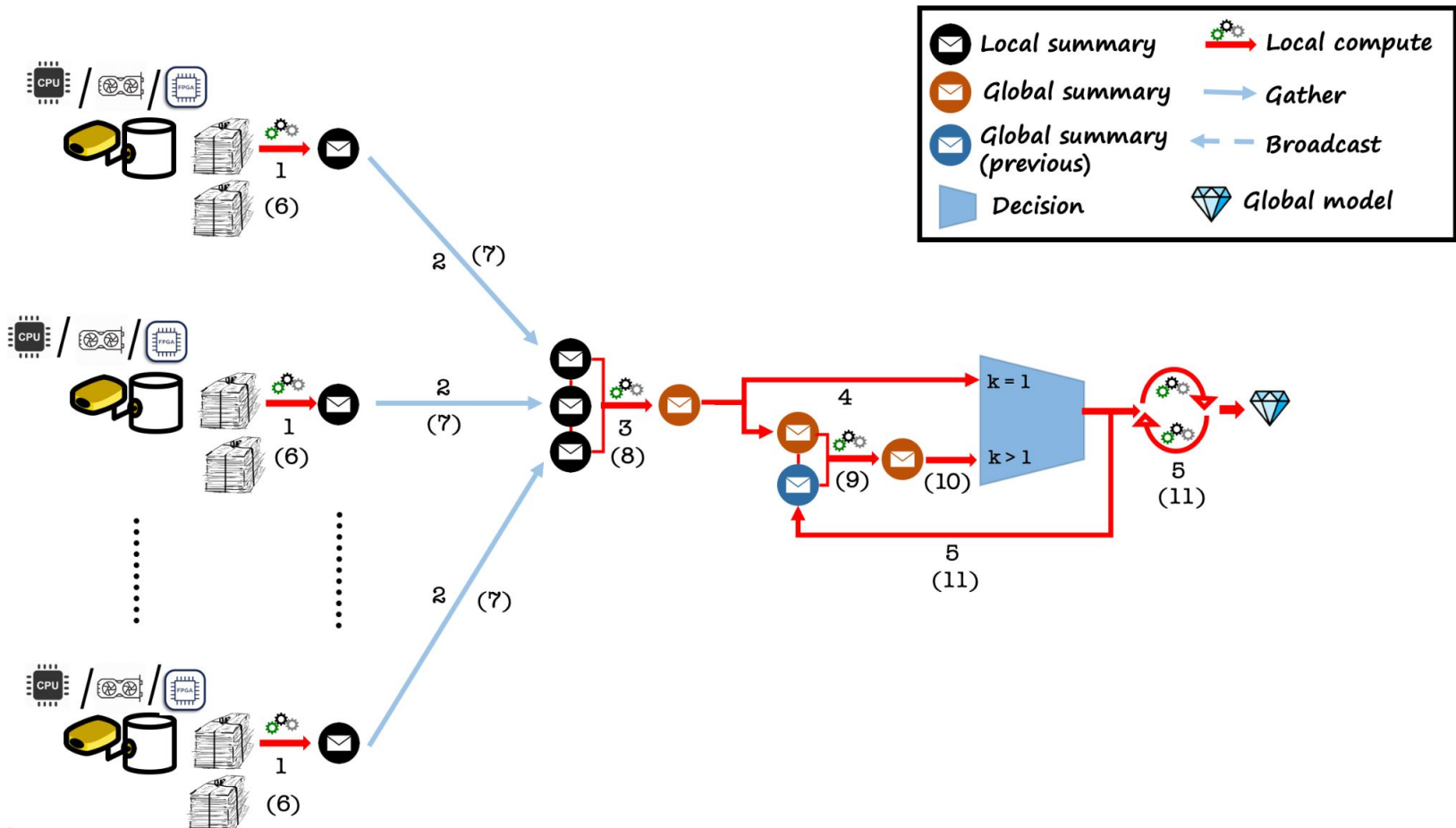
Various setups, round $k \in [3]$ **(a) Stream:** only one worker, where data is generated in every round **(b) Tributary:** fixed number of workers, where data is generated by each worker in every round **(c) Basin:** dynamic number of workers, where in each round different groups of workers participate in the network with their data



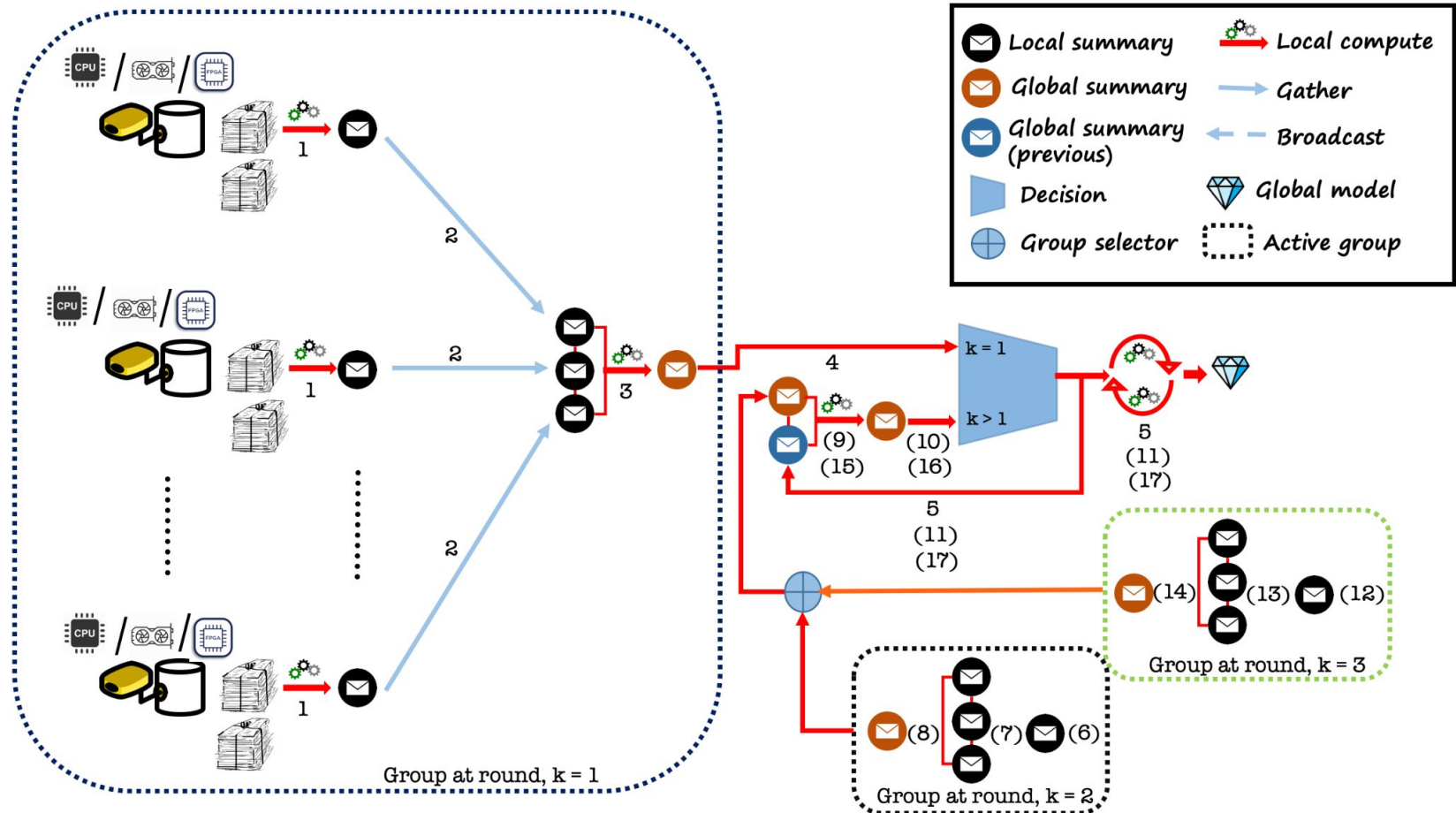
RIVER-Stream



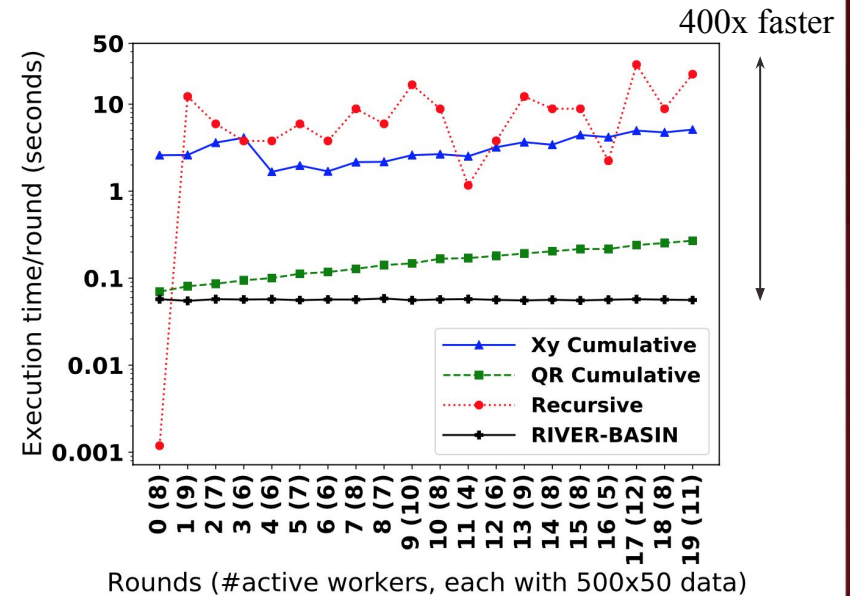
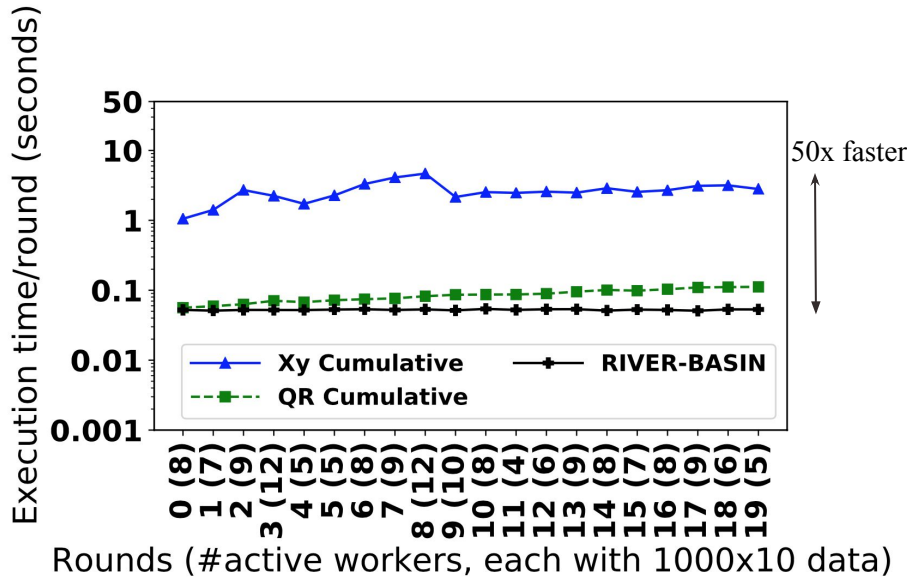
RIVER-Tributary



RIVER-Basin



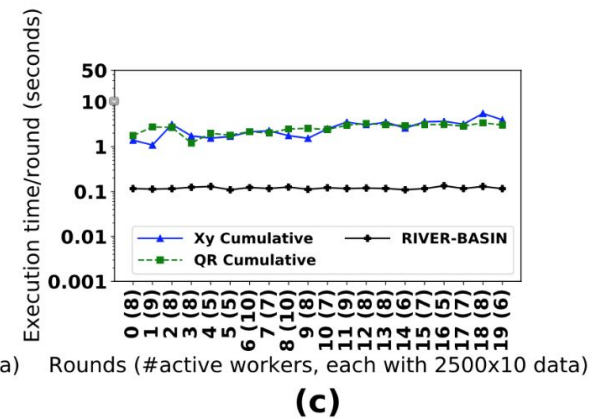
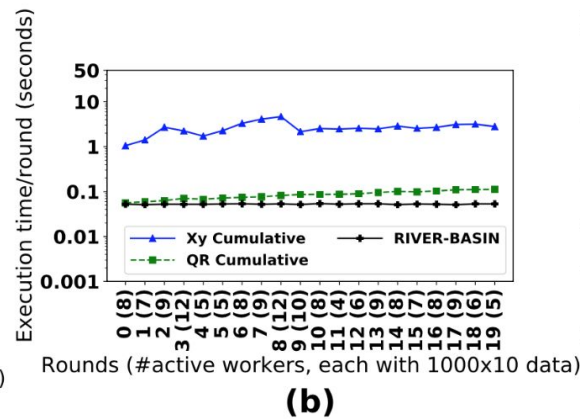
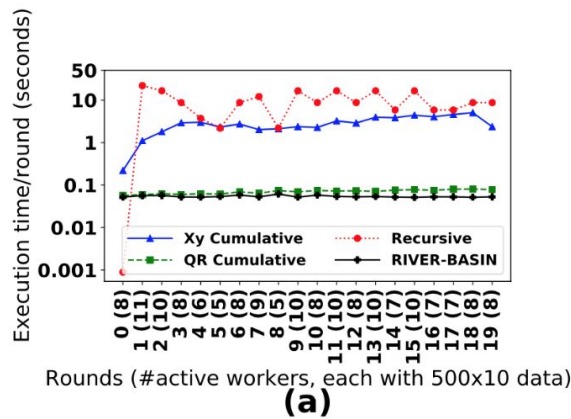
Results (1/5)



RIVER-BASIN execution time per streaming round
 (addition of new data samples or workers)
 is **CONSTANT**, i.e., computations depends **ONLY on the current data** and **FASTER**



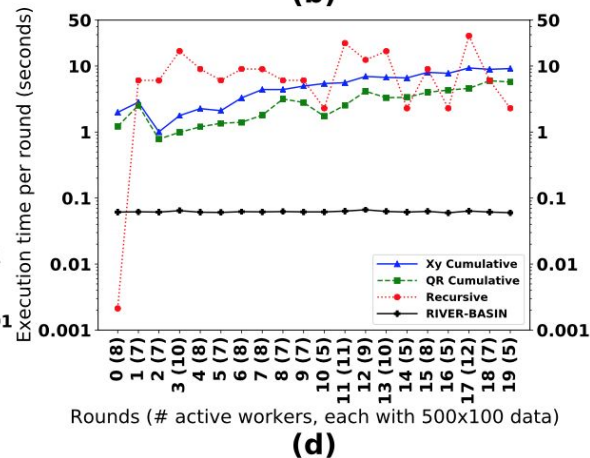
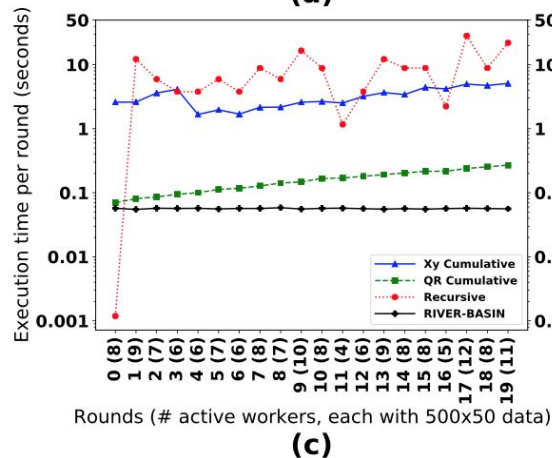
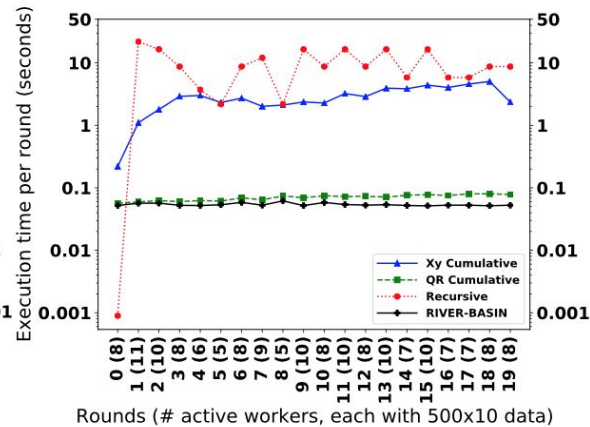
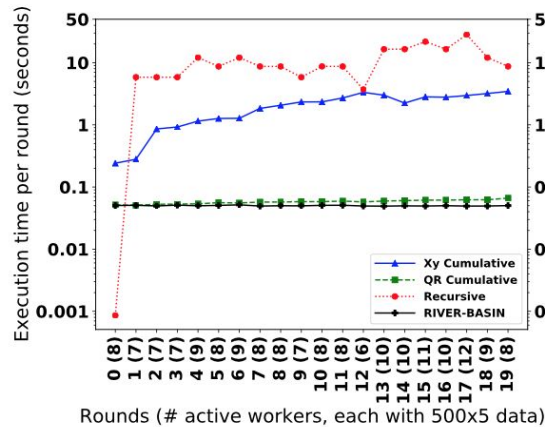
Results (2/5)



BASIN Scalability across varying streaming batch size, n
 (a) 500×10 (b) 1000×10 (c) 2500×10



Results (3/5)

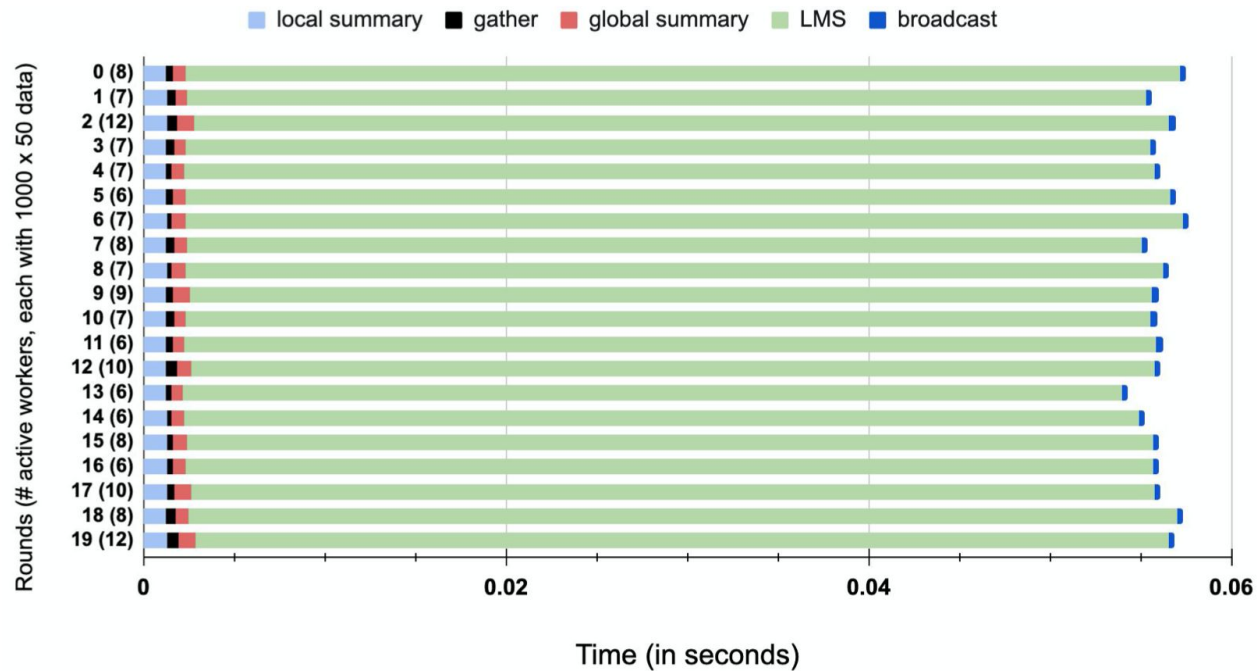


BASIN Scalability across varying feature dimension, d

(a) 500x5 (b) 500x10 (c) 500x50 (d) 500x100



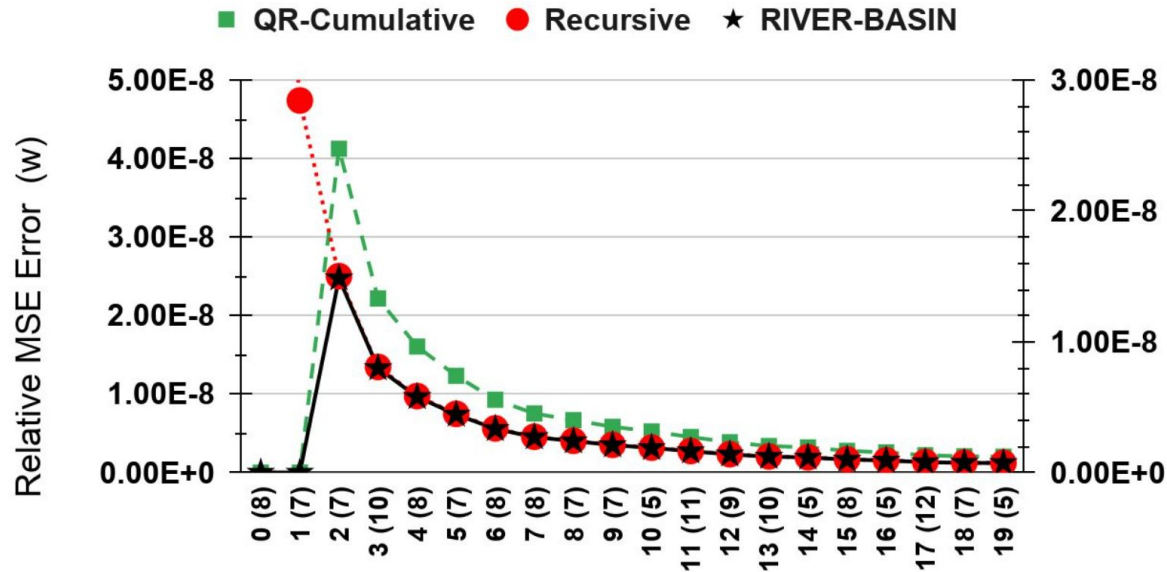
Results (4/5)



RIVER- BASIN: Timing breakdown analysis



Results (5/5)



Rounds (# active workers, each with 500x100 data)

BASIN model error relative to Xy-Cumulative

Model is accurately learnt in each round



Research Contributions

✓ Relaxed Synchronization for Parallel QP Problems

[IEEE IPDPS'16] K. Lee, R. Bhattacharya, **J. Dass**, V. N. S. P. Sakuru, and R. N. Mahapatra, "A Relaxed Synchronization Approach for Solving Parallel Quadratic Programming Problems with Guaranteed Convergence"

✓ Householder Sketch for Machine Learning

[ICML'21] **J. Dass**, and R. N. Mahapatra, "Householder Sketch for Accurate and Accelerated Least-Mean-Squares Solvers"

✓ Memory-efficient Framework for Distributed ML

[IEEE ICDCS'17] **J. Dass**, V. N. S. P. Sakuru, V. Sarin and R. N. Mahapatra, "Distributed QR Decomposition Framework for Training Support Vector Machines"

✓ Communication-efficient Framework for Scalable ML

[IEEE TPDS'18] **J. Dass**, V. Sarin and R. N. Mahapatra "Fast and Communication-Efficient Algorithm for Distributed Support Vector Machine Training"

✓ Multiple FPGA-based System for Energy-efficient ML

[ACM FPGA'19 | IEEE TC'20] **J. Dass**, Y. Narawane, R. N. Mahapatra, and V. Sarin, "Distributed Training of Support Vector Machine on a Multiple-FPGA System"







✓ Rapid Incremental Solver for Federated ML

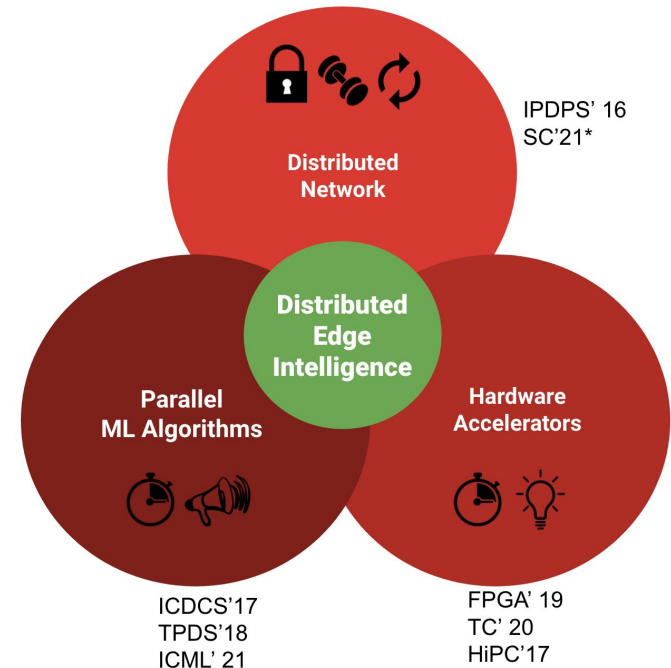
[Under Review, ACM SC'21] **J. Dass**, N. Purwosumarto, R. N. Mahapatra, and X. Hu, "Rapid Incremental Solver for Federated Regression"



Conclusions

Distributed Edge Intelligence Requirements

01	Protect Data Privacy	
02	Reduce Latency	
03	Save Bandwidth	
04	Energy Efficiency	
05	Build Robust Models	
06	Streaming Data	



We proposed,

- **LSDA** for Relax Synchronization (IPDPS)
- **LMS-QR** for Householder Sketch (ICML)
- **QRSVM** for Memory-efficient Distributed Machine Learning (ICDCS)
- **QRSVM.v2** for Communication-efficient Scalable Machine Learning (TPDS)
- **Multi-FPGA** system for Energy-efficient Machine Learning (TC)
- **RIVER** for Incremental learning under Federated ML setups (SC, under Review)

Future Directions

- **Secure Multi-Party Decentralized Machine Learning**
 - Cryptographic methods
 - Differential Privacy
 - Safeguards against malicious workers
 - Privacy-preserving sketch computations
- **Codesigned AutoML Systems for Distributed Edge Intelligence**
 - automatically generate ML model-Accelerator codesigned pair
 - Hardware-aware Neural Architecture Search (HW-NAS)
 - Incorporate design parameters/constraints from distributed computing framework to generate optimal pairs
- **Systems for Lifelong Multi-Agent Learning**
 - Learning from little data, retain the acquired knowledge, share knowledge with other agents and apply to learning in new settings
 - Continuously update model: drone swarms, connected vehicles+users



ACKNOWLEDGMENTS

Dissertation Committee

- Dr. Rabi Mahapatra (advisor)
- Dr. Eun Jung Kim
- Dr. Xia Hu
- Dr. Raktim Bhattacharya

Collaborators

- Dr. Vivek Sarin, *TAMU*
- Dr. Kooktae Lee, *New Mexico Tech*
- Dr. Dharanidhar Dang, *UCSD*
- VNS Prithvi Sakuru, *Amazon*
- Yashwardhan Narawane, *NVIDIA*
- Rengang “Angelo” Yang, *TAMU*
- Nathan Purwosumarto, *TAMU*

CODES Lab Members

- Biplab Patra, *Microsoft*
- Akash Sahoo, *Samsung*
- Karl Ott, *PhD continuing*
- Dr. Ying Yung “Jerry” Yiu, *Intel*
- and everyone

Special Thanks

- Dr. Aakash Tyagi, *TAMU*
- Dr. Joseph Hurley, *TAMU (retired)*
- To all my students, PTs, and graders

Family, Friends, and Teachers



References

1. **J. Dass**, N. Purwosumarto, R. N. Mahapatra and X. Hu, RIVER: Rapid Incremental Solver for Federated Regression, in ACM/IEEE International Conference for High Performance Computing, Networking, Storage, and Analysis (SC 2021), Under Review
2. **J. Dass**, R. N. Mahapatra, Householder Sketch for Accurate and Accelerated Least-Mean-Squares Solvers, in International Conference on Machine Learning (ICML 2021)
3. **J. Dass**, Y Narawane, R. N. Mahapatra and V. Sarin, Distributed Training of Support Vector Machine on a Multiple- FPGA System, in IEEE Transactions on Computers (TC 2020),
4. **J. Dass**, Y Narawane, R. N. Mahapatra and V. Sarin, FPGA-based Distributed Edge Training of SVM, in ACM/SIGDA 27th International Symposium on Field Programmable Gate Arrays (FPGA 2019), Seaside, CA.
5. **J. Dass**, V. Sarin and R. N. Mahapatra, Fast and Communication-Efficient Algorithm for Distributed Support Vector Machine Training, in IEEE Transactions on Parallel and Distributed Systems (TPDS 2018)
6. D. Dang, **J. Dass** and R. Mahapatra, ConvLight: A Convolutional Accelerator with Memristor Integrated Photonic Computing, in IEEE 24th International Conference on High Performance Computing (HiPC 2017), Jaipur,
7. **J. Dass**, V. N. S. P. Sakuru, V. Sarin and R. N. Mahapatra, Distributed QR Decomposition Framework for Training Support Vector Machines, in IEEE 37th International Conference on Distributed Computing Systems (ICDCS 2017), Atlanta, GA,
8. K. Lee, R. Bhattacharya, **J. Dass**, V. N. S. P. Sakuru and R. N. Mahapatra, A Relaxed Synchronization Approach for Solving Parallel Quadratic Programming Problems with Guaranteed Convergence, in IEEE International Parallel and Distributed Processing Symposium (IPDPS 2016), Chicago, IL





Thank You !









TEXAS A&M
UNIVERSITY.

APPENDIX



Requirements for Edge Intelligence

01	Keep Data Local		<ul style="list-style-type: none">• Keep data decentralized and local on devices• Design privacy-preserving ML algorithms
02	Reduce Latency		<ul style="list-style-type: none">• Efficient and scalable training algorithms• Cheap inference calculations to enable real-time analytics
03	Communicate Less		<ul style="list-style-type: none">• Communicate less data during training• Reduce synchronizations and idling during training
04	Energy Efficiency		<ul style="list-style-type: none">• Efficient computation and communication process• Build energy-efficient hardware accelerators for Green AI
05	Build Robust Models		<ul style="list-style-type: none">• Devise fault-tolerance for device failures or stragglers• Accurate and robust model to data perturbations
06	Streaming Data		<ul style="list-style-type: none">• Incremental federated learning to update the global model• Discard data after each update for memory and privacy

Use Case

HOUSTON'S SMART CITY VISION

Transportation

Scroll right to see more



Intelligent Transportation System

A powerhouse of smart devices and real-time data helps manage Houston traffic on a comprehensive scale.

[More on this project](#)

Dockless Bikes & Scooters

More mobility options means less traffic congestion, convenient travel, better air quality and more.

[More on this project](#)

Parking Pay-By-Plate

A modern take on parking meters promises more seamless parking, easier enforcement and cost savings for the city.

[More on this project](#)

Parking Analytics

We're managing our city's parking inventory in real time to boost parking satisfaction while driving big efficiencies.

[More on this project](#)

Rice/Kinder Foundation Study

Research study addressing our environmentally-friendly vehicle replacement to modern electric vehicles.

[More on this project](#)

Autonomous Transit Circulator

Autonomous circulator system fast-tracks travel time, helps with traffic congestion and improves public transit accessibility.

[More on this project](#)

Public Safety

Scroll right to see more



Downtown Sandbox

This collaborative hub enables ongoing testing and installation of the latest tech to keep our city safe.

[More on this project](#)

Project Edison

This platform alerts authorities and the public during active emergencies, enabling safe passage — and safer communities.

[More on this project](#)

Firefighting Drones

Firefighting in Houston is safer than ever with smart drones that assess conditions and identify hazards.

[More on this project](#)

Resiliency & Sustainability

Scroll right to see more



Flood Detection Sensors

A low-cost, high-impact solution from Houston's brightest minds uses sensors to alert drivers of high-water conditions.

[More on this project](#)

Sanitary Sewer Outflow Monitors

A powerful network of monitors is keeping our citywide sewers in check and our city streets cleaner.

[More on this project](#)

Mobile Air Quality Sensors

Sensor-enabled vehicles monitor air pollutants throughout our city to embrace green environment initiatives.

[More on this project](#)

Water Innovation Hub

The nation's first large-scale water technology demonstration hub is testing and showcasing innovative solutions to water and wastewater challenges.

[More on this project](#)

Bridge To Clean Air

Cutting-edge filtration systems and air monitors work together to purify polluted air protecting neighboring communities.

[More on this project](#)

Roadway Flood Warning System

High water sensors and advanced traffic signaling will keep our city safe and informed during flood threats.

[More on this project](#)

Lead

A pov buildi prote issue

[More](#)

<http://houstontx.gov/smartcity/>



Gaps in Parallel ML Training

- Expensive **communication** (synchronization) overheads
- **Memory cost** is nearly **quadratic** with sample size
- **Computational cost** per iteration is **quadratic** with sample size
- Some parallel techniques are limited to small workload, thereby, can **not scale**

For efficient parallel ML training

- Reduce the communication (synchronization) overhead via asynchronous or relaxed synchronous techniques
- Incorporate memory-efficient representation of data by inducing sparsity
- Design highly separable and independent sub-problems amenable for parallelization
- Build algorithms than can scale linearly with number of parallel machines



Opportunities

Untapped Private Data

- Available public data with big companies is just a tip of iceberg
- Utilize untapped private data in individual devices for ML with privacy protected



Idle Processing Power

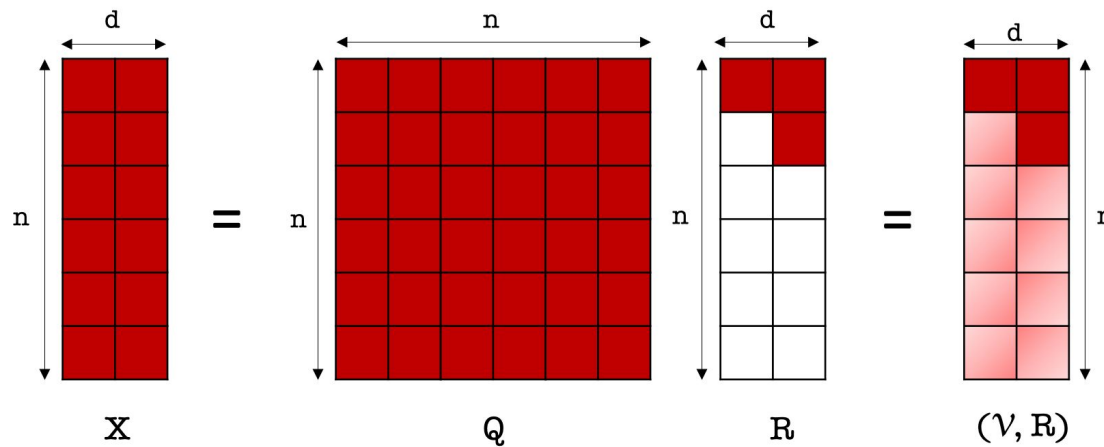
- Billions of connected devices with underutilized processing capacity
- Leverage their idle processing power to build Decentralized processor for ML !



Householder-QR

Theorem 3.1 (Householder-QR [57]). *Let matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ with $n > d$. Householder QR decomposition of \mathbf{X} generates set of d Householder matrices \mathcal{H} and an $n \times d$ upper trapezoidal matrix \mathbf{R} . The Householder matrices are stored as a set of d Householder reflectors \mathcal{V} . Total memory footprint of above factors is nd elements with time complexity of $O(nd^2)$ for $n \gg d$.*

$$\mathbf{X} = \mathbf{QR}, \text{ where, } \mathbf{Q}^T \mathbf{Q} = \mathbf{Q} \mathbf{Q}^T = \mathbf{I}$$



memory
consumption $O(nd)$

$O(n^2)$

$O(d^2)$

$O(nd)$



Distributed Householder Sketches

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \cdot \\ \cdot \\ \mathbf{X}_p \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_1 \mathbf{R}_1 \\ \mathbf{Q}_2 \mathbf{R}_2 \\ \cdot \\ \cdot \\ \mathbf{Q}_p \mathbf{R}_p \end{pmatrix} = \text{diag}(\mathbf{Q}_1, \dots, \mathbf{Q}_p) \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{R}_2 \\ \cdot \\ \cdot \\ \mathbf{R}_p \end{pmatrix}, \quad \mathbf{R}_{stack} = \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{R}_2 \\ \cdot \\ \cdot \\ \mathbf{R}_p \end{pmatrix} = \mathbf{Q}_M \mathbf{R}_M.$$

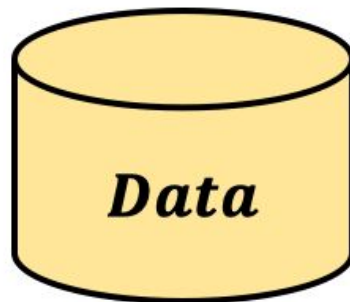
$$\mathbf{X} = \text{diag}(\mathbf{Q}_1, \dots, \mathbf{Q}_p) \mathbf{R}_{stack} = \text{diag}(\mathbf{Q}_1, \dots, \mathbf{Q}_p) \underbrace{\mathbf{Q}_M}_{\mathbf{Q}} \mathbf{R}_M$$

$$\mathbf{Q}^T \mathbf{c} = \mathbf{Q}_M^T \text{diag}(\mathbf{Q}_1^T, \dots, \mathbf{Q}_p^T) \mathbf{c} = \mathbf{Q}_M^T \begin{pmatrix} \mathbf{Q}_1^T \mathbf{c}_1 \\ \mathbf{Q}_2^T \mathbf{c}_2 \\ \cdot \\ \cdot \\ \mathbf{Q}_p^T \mathbf{c}_p \end{pmatrix}$$



Workflow

Initialization



*Distributed
QR Decomposition*



Properties

Lemma

(Stability) The dual variable for LSDA algorithm is stable if and only if

$$\rho(\mathbf{A}(\mathbf{P})) < 1. \quad (8)$$

where $\mathbf{A}(\mathbf{P}) := I - P \sum_{i=1}^N \alpha_i (A_i Q_i^{-1} A_i^T)$ and the symbol $\rho(\cdot)$ denotes the spectral radius of the given matrix (i.e., the largest magnitude of the eigenvalue).

Proposition

(Convergence) Consider the QP problem that is separable. If the condition (8) holds, then the dual variables y_{LSDA} for LSDA and y_{TSDA} for TSDA converge to the same fixed-point value $y^ := \lim_{k \rightarrow \infty} y_{\text{TSDA}}^k = \lim_{t \rightarrow \infty} y_{\text{LSDA}}^{tP}$.*

Theorem

(Optimality) For the given parallel QP problem with LSDA technique, the optimal synchronization period P^ is obtained by*

$$P^* = \max_{P \in \mathbb{N}} \operatorname{argmin} \max\{|1 - \underline{\lambda}(\beta)P|, |1 - \bar{\lambda}(\beta)P|\} \quad (9)$$

where $\beta := \sum_{i=1}^N \alpha_i A_i Q_i^{-1} A_i^T$, $\underline{\lambda}(\cdot)$ and $\bar{\lambda}(\cdot)$ denote the smallest and the largest eigenvalues of the square matrix, respectively.



Memory-efficient Distributed ML

We propose,

- 1 QRSVM: QR decomposition framework for **memory-efficient** modeling and training of SVM
- 2 Optimal step size calculation for **fast convergence**
- 3 Distributed QRSVM for decomposing SVM into parallel equivalent sub-problems that are trained **in parallel**

Lagrangian \mathcal{L} of QRSVM

$$\mathcal{L}(\hat{\alpha}, \beta) = \frac{1}{2} \hat{\alpha}^T \left(\mathbf{R}\mathbf{R}^T + \frac{1}{2C} \mathbf{I}_n \right) \hat{\alpha} + (\hat{\mathbf{e}})^T \hat{\alpha} + \beta^T (-\mathbf{Q}\hat{\alpha})$$

where, $\beta \geq \mathbf{0}_n$ is the Lagrangian dual variable.

Dual Decomposition via projected sub-gradient

Dual function: $g(\beta) = \min_{\hat{\alpha}} \mathcal{L}(\hat{\alpha}, \beta)$

Dual Problem: $\max_{\beta} g(\beta)$



Results (2/2)

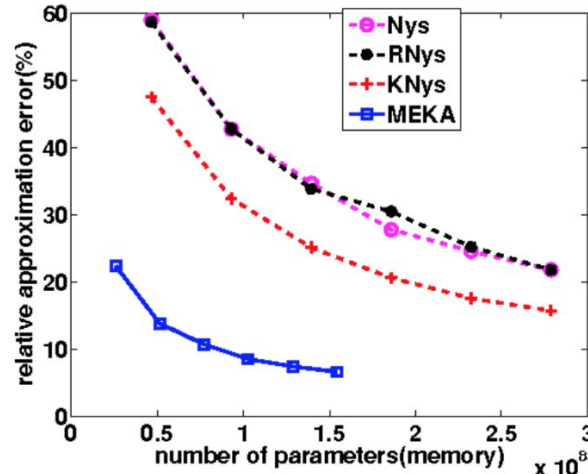


Figure 4.7: Low-rank Gaussian kernel approximation results using MEKA and other methods. [5]

Parameters	a9a	covtype
rank, k	40	64
C	2^{-1}	2^{-1}
γ	2^{-3}	2^3
approx. K_{error}	0.51	0.58
#cores, p	16	16
stopping threshold	10^{-3}	10^{-3}
optimal step size, η^*	1.9	1.9
#iterations, t	166	512

Table 4.2: Distributed-QRSVM: Parameter values. [2]

Distributed QRSVM: Timing Discussions for $p = 16$

Stage 1: Distributed QR

- Computation: $T(p_{localQR}) + T(p_{masterQR})$
- Communication: $T(c_{gatherR})$

Stage2: Parallel Dual Ascent

- Computation: $T(p_{pda})$
- Communication: $T(c_{pda})$
Gather+Scatter

Time details	a9a (in ms)	covtype (in s)
$T(p_{meka})$	460	2.1
$T(p_{localQR})$	24	1.89
$T(p_{masterQR})$	4	0.02
$T(c_{gatherR})$	0.5	0.04
$T(p_{pda})$	1628.1	120.18
$T(c_{pda})$	17.1	0.36
$T(train)$	1674.2	122.50

Optimal step size

Lemma

For any $\eta > 0$, the optimal step size η^* can be computed using

$$\eta^* = P^* \eta, \quad P^* \in \mathbb{N}$$

where,

$$P^* = \begin{cases} 1 & \text{if } 0 < \bar{\lambda}^{-1} < 2 \\ \lfloor \bar{\lambda}^{-1} \rfloor & \text{if } \bar{\lambda}^{-1} \geq 2 \end{cases}$$

and $\bar{\lambda} = (\lambda_{\max}(\mathbf{M}) + \lambda_{\min}(\mathbf{M}))/2$

$$\mathbf{M} := \eta \left(\mathbf{R}\mathbf{R}^T + \frac{1}{2c} \mathbf{I}_n \right)^{-1},$$

$\eta > 0$ is step size

$\lambda_{\min}(\cdot)$ and $\lambda_{\max}(\cdot)$ eigenvalues of matrix \mathbf{M}



(a) a9a



(b) covtype

Figure 4.9: Optimal step size: For both datasets, optimal step size is observed to be 1.9. [2]



In prior work,

Observe,

$$\{q_g\} = \begin{matrix} \text{red} & \{q_f\}_1 \\ \text{white} & \text{0-rows} \\ \vdots & \\ \text{red} & \{q_f\}_p \\ \text{white} & \text{0-rows} \end{matrix} \quad (n \times k)$$

$$R_g = \begin{matrix} \text{red} & (R_f)_{k \times k} \\ \text{white} & \text{0-rows} \end{matrix} \quad (n \times k)$$

Figure: At Master core

$$\beta = Q \hat{\beta} \text{ using } \{q_g\}$$

- Gather $\hat{\beta}_i$ at Master core
- Scatter $(Q_g \hat{\beta})_i$ to each core i

$$\hat{\beta} = Q^T \beta \text{ using } \{q_g\}$$

- Gather $(Q_i^T \beta)_i$ at Master core
- Scatter $\hat{\beta}_i$ to each core i

Communication Overhead in distributed QRSVM

During every iteration of parallel dual ascent, each **Gather** and **Scatter** involves communicating **huge** $O(\frac{n}{p})$ data volume per core with Master core \Rightarrow Poor Scalability

Results (3/3)

TABLE : Comparing *dis*-QRSVM with PSVM, P-packSVM and [12] on T_{train} (in seconds) for *covtype* dataset.

Algorithm	p=2	p=4	p=8	p=16	p=32	p=64
PSVM	8,562	4,396	2,352	1,270	635	341
P-packSVM	-	-	2,019	1,022	295	110
<i>dis</i> -QRSVM [12]	390	309	271	261	256	454
<i>dis</i> -QRSVM	132	64	33	18	10	6

- data unavailable

TABLE : Comparing proposed *dis*-QRSVM with [12] on Stage 2 computation time, T_{update} and communication time, T_{2g+2s} (in seconds) for *covtype* dataset.

Time	p=2	p=4	p=8	p=16	p=32	p=64
T_{update} [12]	379	304	269	259	252	448
T_{update} (proposed)	122	59	30	16	8	5
T_{2g+2s} [12]	1.63	1.81	0.34	0.05	1.19	3.02
T_{2g+2s} (proposed)	0.02	0.02	0.14	0.05	0.03	0.11

- *dis*-QRSVM trains upto $70\times$, upto $60\times$, and upto $75\times$ faster than PSVM, P-packSVM, and earlier [ICDCS'17] implementation, respectively
- *dis*-QRSM is communication-efficient and scales better on PDA with $p = \{2, 4, 8, 16, 32, 64\}$ than [ICDCS'17]
- *dis*-QRSVM can handle larger workloads than [ICDCS'17]

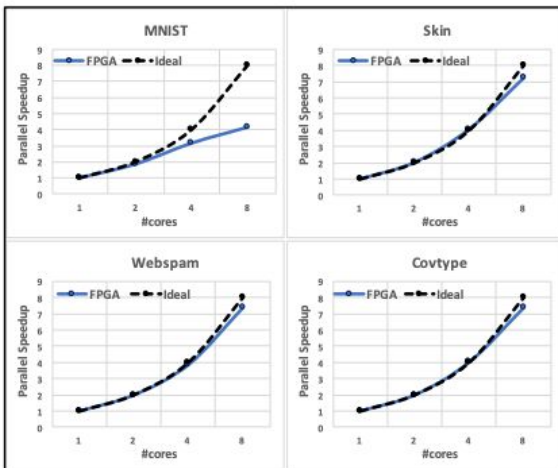


Motivation

- Edge devices are getting efficient in processing data
- Make them capable for accelerating ML training
- Help reduce latency for critical applications
- Lower the dependency on HPC server grade CPUs, GPUs
- Provide energy-efficient modeling and training
- First-of-its-kind work for training SVM in a multiple FPGA environment

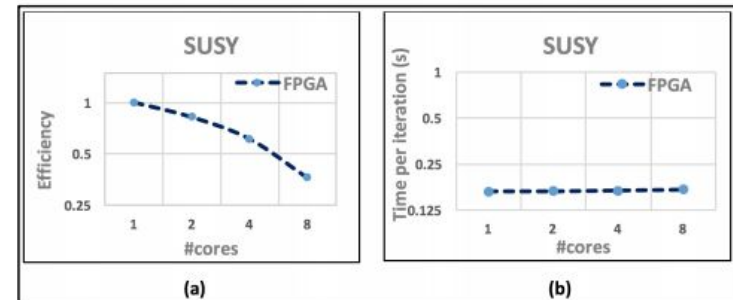
Results (1/3)

Strong Scaling Analysis



- Achieves near **linear parallel speedup** for larger datasets Skin, Webspam, Covtype
- For small dataset MNIST, going beyond $p = 4$ seems to be overkill

Weak Scaling Analysis



- Workload per FPGA fixed at 250K samples
- (a) T_{QR} is constant while $T_{DA} \uparrow$ with #iterations
- (b) (T_p^{FPGA}/t) is constant as desired

Results (2/3)

Energy Analysis

- Under strong scaling, the proposed FPGA design follows the **ideal energy consumption** trend that is **constant** across #FPGA units.
- Validates **fully parallel** implementation

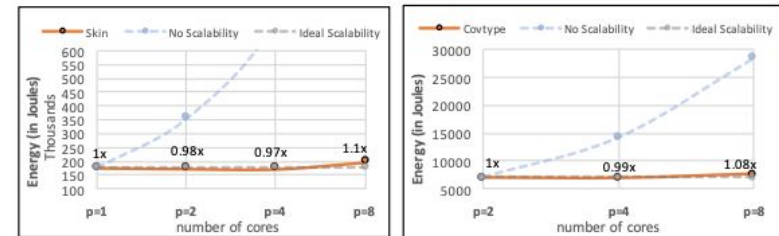


Fig. Energy consumption under Strong Scaling
Skin and Covtype

- Under weak scaling, the ideal energy consumption trend is linear while no scalability trend is quadratic.
- The proposed design is closer to being linear than quadratic. Aberration at $p=8$ due to large #iterations for fine tuning model with increasing overall problem size
- $(Energy/p)$ is nearly **constant** as expected with uniform workload per device

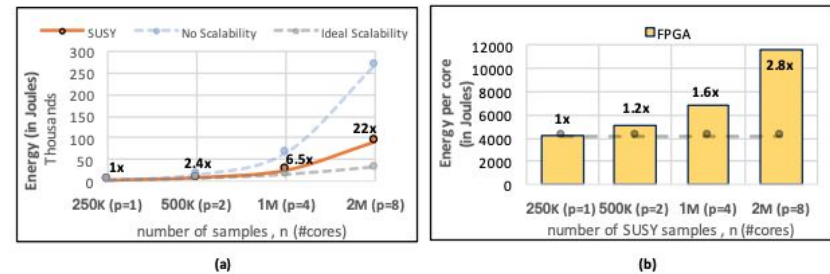
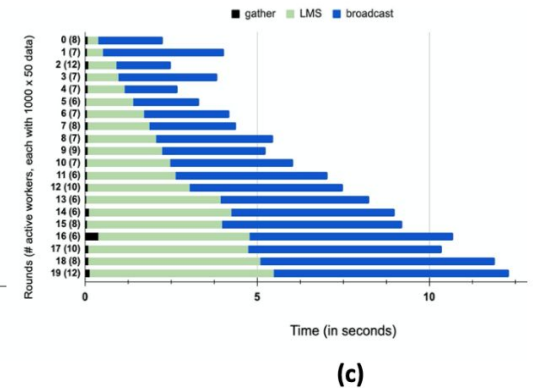
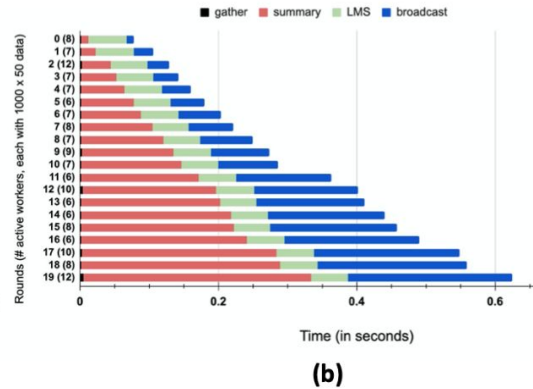
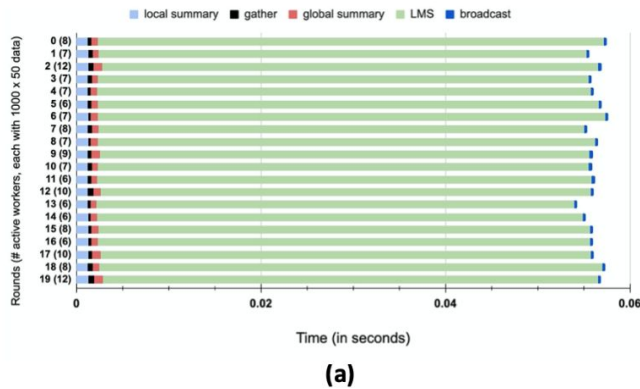


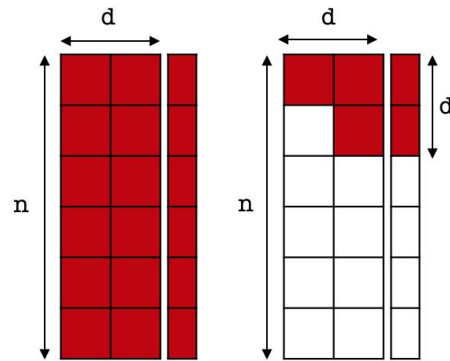
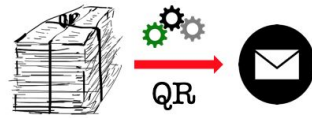
Fig. (a) Total Energy (b) Energy per core
consumption under Weak Scaling for SUSY

Results (4/5)



BASIN: Timing breakdown analysis (a) RIVER (b) QR-Cumulative (c) Xy-Cumulative





memory
consumption

$\{X,y\}$
 $O(nd)$

$\{R,Q^T y\}$
 $O(d^2)$