

# FPGA-based Distributed Edge Training of SVM

# Overview

#### Goal

- To design and implement distributed training of machine learning, here, Support Vector Machines (SVM), on multiple FPGA system
- To reduce network communication while achieving fast training , memory-efficiency, and energy savings

#### Motivation

- > SVM training in computationally expensive with high memory requirement for *kernel* matrix
- > Traditional SVM training accelerators are based on inherently sequential algorithms using a single FPGA board
- > Need to distribute and accelerate training on edge where the data is generated and stored across multiple devices

# Background

#### Support Vector Machines

> Supervised machine learning algorithm for

classification and regression problems quadratic Mathematically. programming а problem which solves for maximal separating hyperplane as a classifier



Fig. Kernel SVM: Learning hyperplane in higher feature dimension

### Distributed QRSVM

- QR decomposition-based distributed SVM
- Memory-efficient + negligible Communication
- Comprises of 3 stages:
  - 1. Initialization
  - 2. Distributed QR decomposition (formulation)
  - 3. Parallel Dual Ascent (solver)



Fig. Process Flow for Distributed QRSVM algorithm

training dataset,  $\mathcal{D} = \{(x_i, y_i), i = 1, ..., n\}$ class label vector,  $y = \{y_i \in \{-1, 1\}, i = 1...n\}$ 

$\min_{\hat{\alpha}} \frac{1}{2} \hat{\alpha}^{T} \Big( R R^{T} + \frac{1}{2C} I_{n} \Big) \hat{\alpha} + (\hat{e})^{T} \hat{\alpha}$
subject to $-Q\hat{lpha}\leq oldsymbol{0}_n$

<b>Step 1:</b> Minimization of Lagrangian - In Parallel								
At edge unit, <i>i</i>								
$\hat{\alpha}_{k}^{k+1} - F^{-1}(-\hat{\beta}_{k}^{k} + \hat{\alpha}_{k})$								
$\alpha_i = r_i (-\rho_i + e_i)$								
where.								
$\left( E_{i}^{-1}  \text{if } i = 1 \right)$								
$F_{i}^{-1} = \begin{cases} r_{1} & r_{-1} \end{cases}$								
(-2C) if $i = 2p$								

# Algorithmic Design

## Distributed QR Decomposition

Algorithm 1 Distributed QR decon	nposition	
1: $k \leftarrow rank$		
2: <b>for</b> each edge <i>i</i> <b>do</b>		
3: $\hat{A}_i \leftarrow \text{local partitioned}$	data	
4: Parallel Compute $\{q_i\}, R_i$	$\leftarrow \hat{A_i} \rightarrow Algor$	rithm 2
5: GATHER $(R_i)_{k \times k}$ at Master	r unit	
6: <b>end for</b>		
7: $\hat{A}_g \leftarrow$ gathered or stacked	$(R_i)_{k \times k}$	
8: Compute $\{q_f\}, R_f \leftarrow \hat{A}_q$ at Ma	ster unit ⊳Algor	rithm 2
9: Use $(R_f)_{k \times k}$		
<b>Algorithm 2</b> $\{q_i\}, R_i \leftarrow A_i$ , via Ho	useholder algorithm	
Algorithm 2 $\{q_i\}, R_i \leftarrow A_i$ , via Ho	ouseholder algorithm	
1: $Q_{\hat{n}\times k}, (\hat{A}_i)_{\hat{n}\times k}$	⊳ $\hat{n}$ : samples per	• edge
2: <b>for</b> $j \leftarrow 1$ to $k$ <b>do</b>		
3: $q_{ij} \leftarrow A_i(j:\hat{n},j)$	()	
4: $q_{ij}(1) \leftarrow q_{ij}(1) + sign(q_{ij}(1))$	$))  imes \ q_{ij}\ _2  ightarrow  ext{scalar u}$	ıpdate
5: $q_{ij} \leftarrow \frac{q_{ij}}{\ q_{ii}\ _2}$	⊳vector normaliz	ation
6: $\hat{A}_i(j:\hat{n},j:k) \leftarrow \hat{A}_i(j:\hat{n},j)$	$j:k) - 2q_{ij} < q_{ij}, \hat{A}_i(j)$	$: \hat{n}, j :$
( <i>k</i> ) >	⊳ Algor	rithm 4
7: $R_i = \hat{A}_i(j:\hat{n}, j:k)$		
8: end for		
9: $\{q_i\} \leftarrow [q_{i1}, q_{i2}, \ldots, q_{ik}]$	⊳set of k-refle	ectors



Algorithm 4 Computing Step 6 in Algorithm 2					
1: $\triangleright \hat{A}_i(i:\hat{n},i:k) \leftarrow \hat{A}_i(i:\hat{n},i:k) - 2a_{ii} < k$	$(a_{ij}, \hat{A}_{i}(i:\hat{n}, i:k)) >$				
2: for $m \leftarrow j$ to k do					
3: $a_{BRAM} \leftarrow A(j:\hat{n},m)$	⊳Load into BRAM				
4: Compute $sum = \langle q_{ij}, a_{BRAM} \rangle$	⊳ Inner Product				
5: $a_{BRAM} \leftarrow a_{BRAM} - 2 \times sum \times q_{ij}$	⊳ saxpy				
6: $\overline{A(j:n,m)} \leftarrow a_{BRAM}$	⊳Write to DDR				
7: end for					

Jyotikrishna Dass, Yashwardhan Narawane, Rabi Mahapatra, Vivek Sarin Texas A&M University



### **COMPUTER SCIENCE & ENGINEERING** TEXAS A&M UNIVERSITY

# **Experimental Results and Discussions**

#### Hardware Platform

Proof-of-Concept implemented on Amazon AWS F1 instance with  $p=\{1,2,4,8\}$  16nm Xilinx Virtex Ultrascale+ VU9P FPGA units forming a multiple FPGA network

### Training Time

units	#iterations	M	<b>MNIST</b> ( $C = 1, \gamma = 2^{-6}, \eta^* = 0.9$ )				
р	t	$T_{QR}$	$T_{DA}$	$T_p^{FPGA}$	comp	comm	
1	181	3.42	7.49	10.92	99.9%	0.1%	
2	181	1.76	4.07	5.84	99.8%	0.2%	
4	182	0.93	2.51	3.58	96%	4%	
8	182	0.46	2.14	2.61	99.6%	0.4%	
units	#iterations		<b>Skin</b> ( <i>C</i> = 1, $\gamma = 2^{-8}$ , $\eta^* = 0.9$ )				
р	t	$T_{QR}$	$T_{DA}$	$T_p^{FPGA}$	comp	comm	
1	67441	2.80	4533	4536	99.9%	0.1%	
2	64424	1.46	2226	2228	99.9%	0.1%	
4	59761	0.74	1107	1108	99.9%	0.1%	
8	54744	0.38	625	626	99.9%	0.1%	
units	#iterations	W	ebspan	$\mathbf{n} (C = 1, \gamma)$	$= 1, \eta^* =$	: 0.9)	
р	t	$T_{QR}$	$T_{DA}$	$T_p^{FPGA}$	comp	comm	
1	-	-	-	-	-	-	
2	566	9.80	66.20	76.14	99.8%	0.2%	
4	564	4.88	34.40	39.36	99.8%	0.2%	
8	569	2.60	17.92	20.59	99.7%	0.3%	
units	#iterations	C	<b>Covtype</b> ( $C = 1, \gamma = 2^3, \eta^* = 0.9$ )				
р	t	$T_{QR}$	$T_{DA}$	$T_p^{FPGA}$	comp	comm	
1	-	-	-	-	-	-	
2	1125	3.35	88.02	91.45	99.9%	0.1%	
4	1080	1.70	43.58	45.36	99.8%	0.2%	
8	1068	0.80	23.80	24.75	99.4%	0.6%	
units	#samples	S	USY (C =	$= 1, \gamma = 2^{-1}$	$^{-3}, \eta^* = 0$	.9)	
p	n	$T_{QR}$	$T_{DA}$	$T_p^{FPGA}$	comp	comm	
1	250K	14.01	94.04	108.08	99.9%	0.1%	
2	500K	14.04	116.8	131.02	99.8%	0.2%	
4	1M	14.07	162.1	176.18	99.9%	0.1%	
8	2M	14.14	285.47	299.63	99.9%	0.1%	

#### Parallel Dual Ascent is computationally dominant than Distributed QR decomposition Near negligible communication overhead <1%</p>

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





- seems to be overkill









# Future Work

- □ To design and implement distributed training for Deep Learning models for resource constrained devices with limited memory and low power
- □ To explore online/incremental learning capabilities for machine learning models at edge

# References

- Burges, Christopher JC. "A tutorial on support vector machines for pattern recognition." Data mining and knowledge discovery 2, no. 2 (1998): 121-167. Papadonikolakis, Markos, and Christos-Savvas Bouganis. "A scalable FPGA architecture for non-linear SVM training." In 2008 International Conference on Field-Programmable Technology, pp. 337-340. IEEE, 2008.
- Cadambi, Srihari, Igor Durdanovic, Venkata Jakkula, Murugan Sankaradass, Eric Cosatto, Srimat Chakradhar, and Hans Peter Graf. "A massively parallel FPGA-based coprocessor for support vector machines." In 2009 17th IEEE Symposium on Field Programmable Custom Computing Machines, pp. 115-122. IEEE, 2009.
- Rabieah, Mudhar Bin, and Christos-Savvas Bouganis. "FPGA based nonlinear support vector machine training using an ensemble learning." In 2015 25th International Conference on Field Programmable Logic and Applications (FPL), pp. 1-4. IEEE, 2015. Si, Si, Cho-Jui Hsieh, and Inderjit S. Dhillon. "Memory efficient kernel approximation." The Journal of Machine Learning Research 18, no. 1 (2017): 682-713.

Dass, Jyotikrishna, Vivek Sarin, and Rabi N. Mahapatra. "Fast and Communication-Efficient Algorithm for Distributed Support Vector Machine Training." IEEE Transactions on Parallel and Distributed Systems (2018).

