

HogWild on Steroids: SGD via Migrating Threads



Peter M. Kogge
McCourtney Prof. of CSE
Univ. of Notre Dame
IBM Fellow (retired)



Observations

- Modern systems becoming heterogeneous
 - Accelerators for high compute intensive apps
 - Need complex “glue” code for multi-node scaling
- But many apps no longer good fits
 - Computation often memory B/W limited
 - Intra-node parallelism often coherency limited
 - Inter-node parallelism no longer regular or predictable
- Result: poor efficiency, poor scalability



Project Funding

- NSF SPX project joint with Vivek Sarkar
 - “Scalable Heterogeneous Migrating Threads for Post-Moore Computing”
- Goal: can we fix the heterogeneous problem of combining different ISAs/core types with “migrating threads?”
- Emu machine in CRNCH an excellent resource



This Study

- Can a modern **Machine Learning** problem
- Benefit from architectures with **Migrating Threads**
- And what is likely scalability

Today's talk: Very Preliminary
Paper Analysis



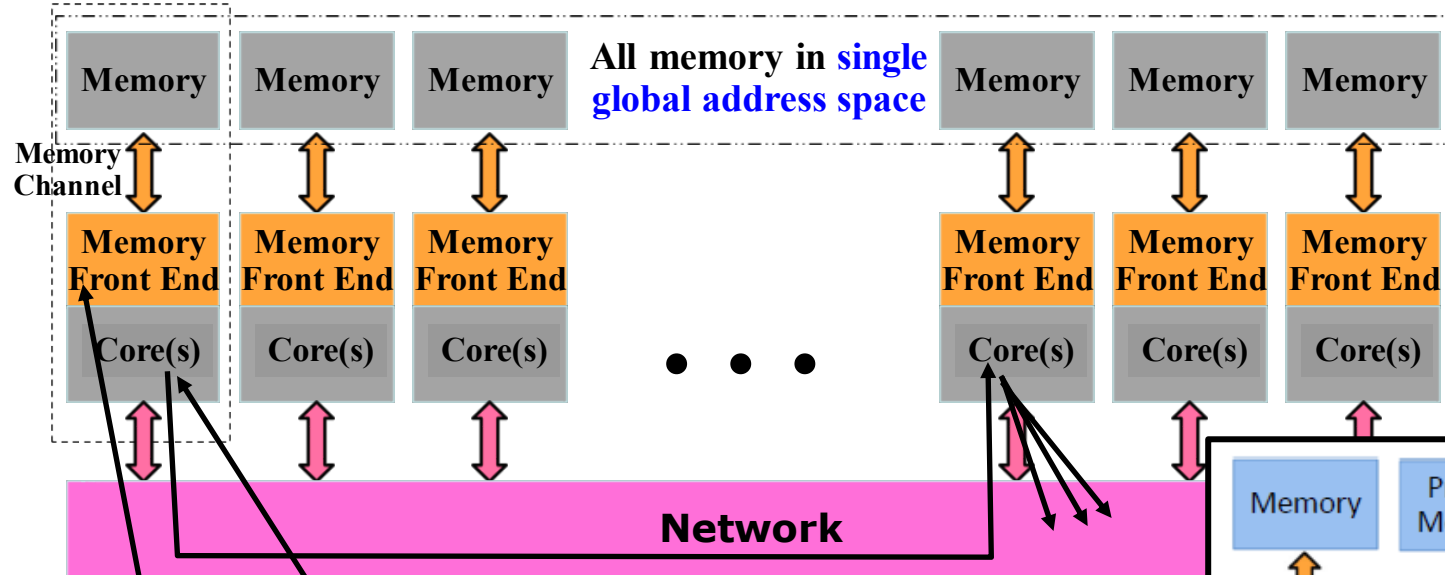
A Funny Thing Happened...

- For this problem
 - Sparsity can dominate computational complexity
 - Inter-node communication is still the problem
 - Heterogeneity probably not efficient match
- Efficient Multi/Migrating threads turn out to solve ***both*** the intra and inter node problems



Thesis: Use “Migrating Threads” as Glue

Nodelet: New unit of parallelism



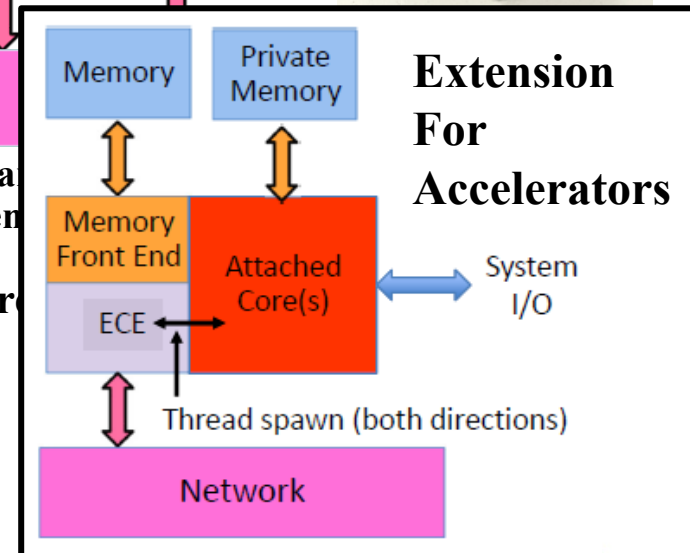
Smart Memory Controllers that also do atomics

Until they make a non-local reference
And then moved to correct nodelet

And they are independent

Threads execute here
Spawned threads include fixed function r

**Really good for large sparse,
irregular data sets
& memory bound apps**



Machine Learning

- Given set E of “training samples”
- Find some model vector x^*
- That minimizes $\hat{f}(x) = \sum_{e \in E} f_e(x)$
 - Where $f_e(x)$ is defined by e ’th sample



Stochastic Gradient Descent

Iteratively improve an estimate of x^*

- Randomly choose a new sample
- Update estimate: $x^{(k+1)} = x^{(k)} - \alpha \nabla \hat{f}(x^{(k)})$
- Where $\nabla \hat{f}(x^{(k)})$ is multidimensional gradient



Naïve Parallel Algorithm

- Obvious Algorithm
 - Break sample set into “mini-batches”
 - Process each independently
 - Combine updates
- If updates done asynchronously, results bad
- If updates done in single critical section, poor scalability



Hogwild Algorithms

- Observation: if samples *sparse*, updates can be done *incrementally*, without locking
- **HogWild!**: use shared memory multi-threaded platform
 - One thread/mini-batch
 - Atomic updates to elements of shared model vector
 - Some scalability but suffers from coherency traffic
- **BuckWild**: use short precision
- **HogWild++**: Keep multiple local model vectors & perform cyclic updates
 - Much better scalability



Generic HogWild++

```

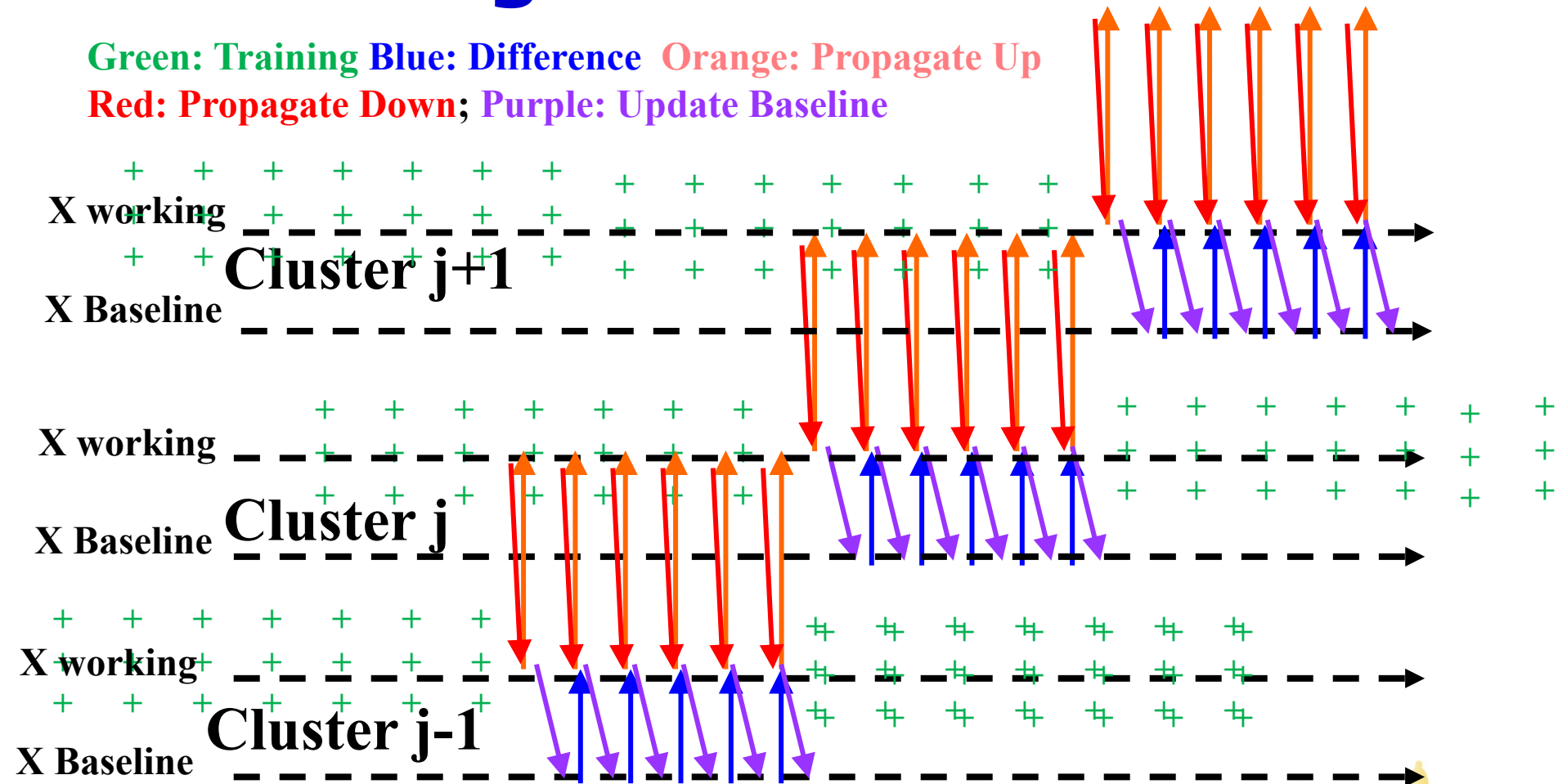
1: repeat
2:   for all threads in cluster j do
3:     Pick a random training sample e
4:      $atomic(\hat{x}_j^{(\tau+1)} = \hat{x}_j^{(\tau)} - \eta_0 \gamma^t \nabla f_e(\hat{x}^{(\tau)}))$ 
5:      $atomic(\tau = \tau + 1)$ 
6:     if Token is received then
7:        $\Delta x_j^{(k)} = \hat{x}_j^{(\tau)} - x_j^{(k)}$  {Compute local changes}
8:        $atomic(x_j^{(k+1)} = \lambda \hat{x}_{j+1}^{(\tau')} + (1-\lambda)x_j^{(k)} + \beta \gamma^t \Delta x_j^{(k)})$ 
9:       {Now update neighbor's  $\hat{x}$ }
10:       $atomic(\hat{x}_{j+1}^{(\tau')} = \hat{x}_{j+1}^{(\tau')} + \beta \gamma^t \Delta x_j^{(k)})$ 
11:       $\hat{x}_j^{(0)} = x_j^{(k+1)}$  {Update local model}
12:       $k = k + 1; \tau = 0$ 
13:      Pass token after at least  $\tau_0$  iterations
14:    end if
15:  end for
16: until  $\tau \geq m$ 
  
```

- Threads divided into clusters
- Each cluster has own model vector
- Local updates incrementally
- Token passed cyclically around
- Cluster j exchanges updates with cluster j+1



HogWild++ Traffic

Green: Training Blue: Difference Orange: Propagate Up
Red: Propagate Down; Purple: Update Baseline



Support Vector Machine (SVM)

- Each sample (z_e, y_e)
 - a vector z_e of F features
 - Member of one of two classes ($y_e = +1$ or -1)
- Sign of $x^T z_e$ determines class
- $f_e(x) = \max(0, 1 - y_e x^T z_e)$ $\hat{f}(x) = (\sum_e \max(0, 1 - y_e \boxed{x^T z_e})) + \lambda \|x\|_2^2$ $O(F)$
- $\nabla \hat{f}(x)[i] = (\hat{f}(x + h b_i) - \hat{f}(x)) / h$
- $\nabla f_e(x)[i] = (\max(0, g_e - y_e h z_e[i]) - \max(0, g_e)) / h$
- If z_e is sparse: $\hat{f}(x) = \sum_e (\max(0, 1 - y_e \boxed{\sum_{i \in Z_e} x[i] z_e[i]})) + \lambda \sum_{i=1}^n x[i]^2$
 - **Huge reduction** if $F' \ll F$ $O(F')$ where $F' = \# \text{ non-zeros}$



Reported Results

Data set	S: Training Samples	Sparsity	Per Sample		Max Speedup	Best Configuration			τ_0	$\log_{64}(S)$
			F: Features	F': Non-Zeros		Cores	Cores/Cluster	Clusters		
news20	16,000	0.0336%	1,355,191	455	9.5	49	4	10	16	3.4
covtype	464,810	22.12%	54	12	30	40	1	40	16	1
webspam	280,000	33.52%	254	85	40	40	1	40	16	1.4
rcv1	677,399	0.155%	47,236	73	38	40	1	40	16	2.6
epsilon	400,000	100%	2,000	2,000					16	1.9

- **Webspam**: small feature set F, good speedup
- **News20**: large feature set, poor speedup



Issues to Consider

- Sparsity of samples: reduce training costs
 - But we know a priori which are non zero
- Sparsity of changes to handle during token passing updates
 - Not known a priori at either j or $j+1$
- Reducing inter-cluster traffic during token update
 - Avoid multiple large vector transfers
- Performing individual updates atomically



Single Cluster SVM- Training Only

Algorithm 2 Single Cluster SVM Epoch Computation

SY: vector of pairs (index into IZ, y_e) for sample e

IZ: vector of pairs (feature index,value) for samples

X: vector of current model values

$\lambda, \beta, \eta\gamma^t, h$: coefficients for update

```

1: for  $\tau = 1$  to  $|E|$  by 1 do
2:   {Multiple threads can execute asynchronously}
3:    $e = 2 * \text{random\_int}(1, |E|)$  {Select sample}
4:    $iz2 = iz = SY[e]$  { Index to 1st non-zero}
5:    $y_e = SY[e + 1]$ ;  $iend = SY[e + 2]$ ;  $g = 0$ 
6:   {Compute  $y_e \sum_{i \in Z_e} x[i]z_e[i]$ }
7:   repeat
8:      $i = IZ[iz]$ ;  $zei = IZ[iz + 1]$ 
9:      $g = g + X[i] * zei$ 
10:  until  $(iz ++ 2) \geq iend$ 
11:   $g = 1 - y_e g$ ;  $gmax1 = \max(0, g)$ ;  $yh = y_e * h$ 
12:  {Compute each  $\nabla f_s(x)[i]$ }
13:  repeat
14:     $i = IZ[iz2]$ ;  $zei = IZ[iz2 + 1]$ 
15:     $dx2i = 2h * zei$ 
16:     $gmax2 = \max(0, g - yh * zei - dx2i)$ 
17:     $gi = \eta_0 \gamma^t (gmax2 - gmax1) / h$ 
18:     $\text{atomic}(X[i] += gi)$ 
19:  until  $(iz2 ++ 2) \geq iend$ 
20: end for
    
```

$O(F')$

$O(F')$

Features/Sample	No Cache			64B Cache Lines		
	10	100	1000	10	100	1000
Transfers/Sample	89	809	8009	49	409	4009
Flops/Transfer	1.04	1.12	1.12	2.02	2.22	2.35
Flops/Byte Accessed	0.13	0.14	0.14	0.03	0.03	0.04

- Modern architectures have flops/byte of 4-8



Single Token Update – Full Vector

Algorithm 3 Single Token Update - Full Vector

F : Number of features (zero or non-zero)

$\hat{x}_{j+1}^{(\tau')}$ is on the next cluster

algorithm.3

```
1: {Iterate over each feature in model vector}
2: for  $i = 1$  to  $F$  by 1 do
3:   {Line 9:  $\Delta x_j^{(k)} = \hat{x}_j^{(\tau)} - x_j^{(k)}$ }
4:    $x1 = x_j^{(k)}[i]$  {Local read}
5:    $x2 = \hat{x}_j^{(\tau)}[i]$  {Local read}
6:    $dx = \beta \gamma^t (x2 - x1)$ 
7:   {Line 14:  $x_j^{(k+1)} = \lambda \hat{x}_{j+1}^{(\tau')} + (1 - \lambda) x_j^{(k)} + \beta \gamma^t \Delta x_j^{(k)}$ }
8:    $x3 = \hat{x}_{j+1}^{(\tau')}[i]$  {Read from cluster  $j + 1$ }
9:    $x4 = \lambda x3 + (1 - \lambda) x1 + dx$ 
10:   $x_j^{(k+1)}[i] = x4$  {Local store}
11:  {Line 16:  $\hat{x}_{j+1}^{(\tau')} = \hat{x}_{j+1}^{(\tau')} + \beta \gamma^t \Delta x_j^{(k)}$ }
12:   $\text{atomic}(\hat{x}_{j+1}^{(\tau')}[i] += x4)$  {Remote update}
13:  {Line 17:  $\hat{x}_j^{(0)} = x_j^{(k+1)}$ }
14:   $\hat{x}_j^{(0)}[i] = x4$  {Local store}
15:  {The increment of  $i$  should be an AMO.}
16: end for
```

**All these read
full vector of
F components,
even if most are 0**

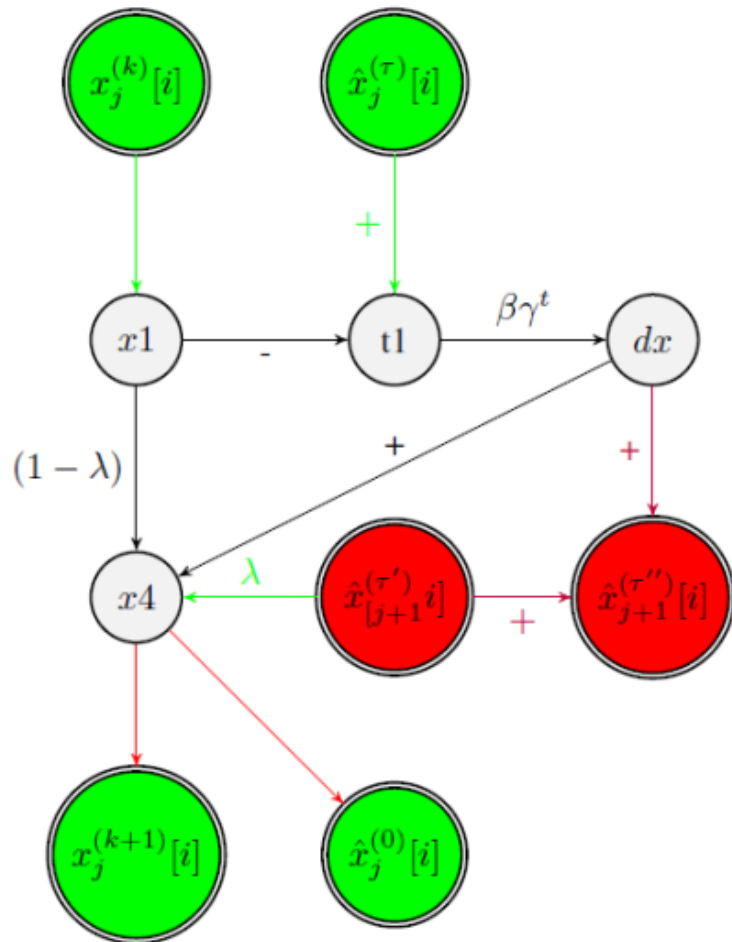


Options for Tracking Dynamic Changes

- Keep list of changes
 - Add each time local model vector is updated
 - Cost of removing duplicates likely large
- Keep bit vector of changes
 - But still requires $O(F)$ reads to scan
- Keep tree of bit vectors
 - Four deep of 64b words covers up to 16M samples
 - Requires up to 4 operations to record change
 - Must be done atomically if multi-threaded
- Also remember changes in cluster $j+1$ likely different from changes in cluster j



Complexity of Inter-Cluster Updates



Double circle: memory location

Green circle: local memory; Red circle: remote memory

Green line: read; Red line: write; Purple line: AMO

- Complex flow of data
 - in both directions
- Migrating Threads to rescue:
 - Separate cluster j and cluster $j+1$ processing
 - Spawn thread from cluster j to do processing at cluster $j+1$
 - Have cluster $j+1$ issue remote atomic adds for updates back to cluster j



Migrating Thread Pseudocode

Algorithm 4 Notional Migrating Thread Training Loop

```

1: repeat
2:   {Process samples using Algorithm 2}
3:   {Multiple threads can execute asynchronously}
4:    $e = 2 * \text{random\_int}(1, |E|)$  {Select sample}
5:    $ix2 = ix = SY[e]$ ;
6:    $y_e = SY[e + 1]$ ;  $iend = SY[e + 2]$ ;  $g = 0$ 
7:   repeat
8:      $g = g + X[IZ[ix]] * IZ[ix + 1]$ 
9:   until  $(ix + 2) \geq iend$ 
10:   $gmax1 = \max(0, g - 1 - y_e g)$ ;  $yh = y_e * h$ 
11:  repeat
12:     $i = IZ[ix2]$ ;  $zei = IZ[ix2 + 1]$ ;  $dx2i = 2h * zei$ 
13:     $gmax2 = \max(0, g - yh * zei - dx2i)$ 
14:     $\text{atomic}(X[ix] += \eta_0 \gamma^t (gmax2 - gmax1) / h)$ 
15:     $\text{set\_change\_bit}_j(i)$ 
16:  until  $(ix2 + 2) \geq iend$ 
17:  {Test for token arrival}
18:  if  $\text{atomic}(\tau += 1) > \tau_0$  then
19:    if  $\text{CAS}(\text{token}_j = 1, -1)$  then
20:      Break {1st thread to see it}
21:    else if  $\text{token} < 0$  then
22:      Quit execution {Later threads}
23:    end if
24:  end if
25: until  $\tau > m_j$ 
26: {Token has arrived - drop to Token Processing}
  
```

Algorithm 5 Notional Migrating Thread Token Update

```

1: {Following code is executed on cluster j}
2:  $last\_i = 1$ 
3: for  $i = \text{find\_next\_change\_bit}_j()$  do
4:    $\text{reset\_change\_bit}_j(i)$ 
5:    $x1 = x_j[i]$  {Local read}
6:    $dx = \beta \gamma^t (\hat{x}_j[i] - x1)$ 
7:    $x4 = (1 - \lambda)x1 + dx$  {Partial update}
8:    $x_j[i] = x4$  {Local store of partial}
9:    $\hat{x}_j[i] = x4$  {Local store}
10:  spawn remote_update( $i, dx, last\_i$ )
11:   $last\_i = i$ 
12: end for
13: Restart sample training loop
14:
15: {Following code is executed on cluster j + 1}
16: procedure remote_update( $i, dx, last\_i$ )
17:    $\text{atomic}((x3 = x_{j+1}[i]) += dx)$  {Update local x}
18:    $\text{remote\_atomic}(x_j[i] += \lambda x3)$  {Complete  $x_j[i]$  update}
19:    $\text{remote\_atomic}(\hat{x}_j[i] += \lambda x3)$  {Complete  $\hat{x}_j[i]$  update}
20:    $\text{set\_change\_bit}_{j+1}(i)$ 
21:   for  $n = i - 1$  down to  $last\_i + 1$  do
22:     if  $\text{change\_bit}_{j+1}[n] == 1$  then
23:        $x3 = \hat{x}_{j+1}^{(\tau')}$  {Now a local read}
24:        $\text{remote\_atomic}(\hat{x}_j[n] += x3)$ 
25:        $\text{set\_remote\_change\_bit}_j(n)$ 
26:     end if
27:   end for
28: quit
  
```

Migrating Thread

Float atomic

Integer atomic

Key Model Parameters

- **F**: number of features per sample
- **F'**: average non-zero features per sample
- **F''**: average # of changes in x in one cluster between tokens
- **F'''**: average # of changes in cluster $j+1$ that were not changes in cluster j



Traffic Counts

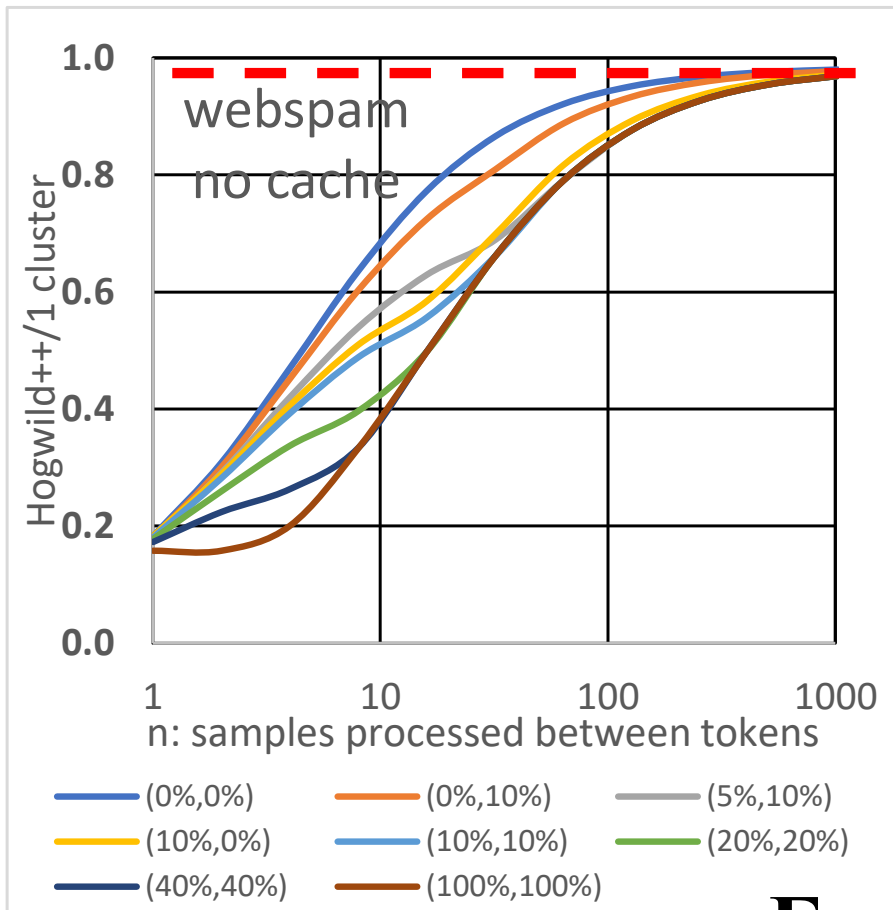
	Training	Local Update	Remote Update	Total Transfers	
				No Cache	Cache
Local Random Reads	$n + 3nF'$	$2F''(n)$	$2F'''(n)$	$n + 3nF' + 2F''(n) + 2F'''(n)$	same
Local Sequential Reads	$n + 2nF'$			$n + 2nF'$	0
Local Random Writes		$2F''(n)$		$2F''(n)$	$4F''(n)$
Local Integer Atomics	$2n$	$F''(n)$		$6n + 3F''(n)$	$4n + 2F''(n)$
Local Float Atomics	nF'		$F''(n)$	$3nF' + 3F''(n) + 3F'''(n)$	$3nF' + 2F''(n) + 2F'''(n)$
Remote Float Atomics			$2F''(n) + F'''(n)$	$6F''(n) + 3F'''(n)$	$2nF' + 4F''(n) + 2F'''(n)$
<i>random_int()</i>	n			$6n$	$6n$
<i>Find_next_change_bit()</i>		$F''(n)$		$\log_{64}(F)F''(n)$	$\log_{64}(F)F''(n)$
Local (re)set_change_bits	n	$F''(n)$	$F''(n)$	$3\log_{64}(F)(n + 2F''(n))$	$2\log_{64}(F)(2n + F''(n))$
Remote (re)set_change_bits			$F'''(n)$	$3\log_{64}(F)F'''(n)$	$2\log_{64}(F)F'''(n)$
Full Thread Spawns		$F''(n)$			
Migrations from Cluster		$F''(n)$			
Float Adds	$n + 4nF'$	$2F''(n)$	$3F''(n) + F'''(n)$	$n + 4nF' + 5F''(n) + F'''(n)$	
Float Multiplies	$2n + 4nF'$	$2F''(n)$	$F''(n)$	$2n + 4nF' + 3F''(n)$	

F' : number of samples processed between tokens. F' : ave. non-zeros per sample.
 F'' : Ave. number of changes in cluster j during the processing of n samples.
 F''' : Ave. number of changes in cluster $j + 1$ during the processing of n samples by cluster j that are not changes by j .

- Transfer: movement of data to/from memory
- **No cache**: current Emu machine
- **Cache**: cache added to memory controller
- No cache AMOs: 3 transfers
- Cache AMOs: 2 transfers



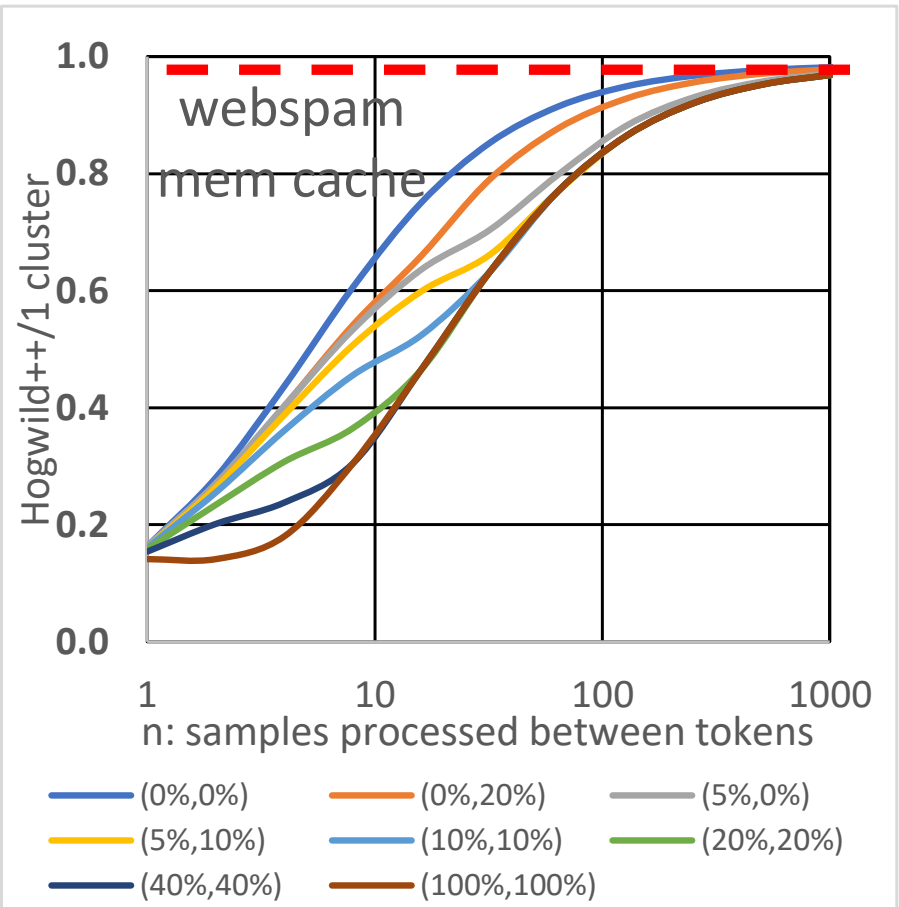
Webspam Dataset



No Cache

F = 254

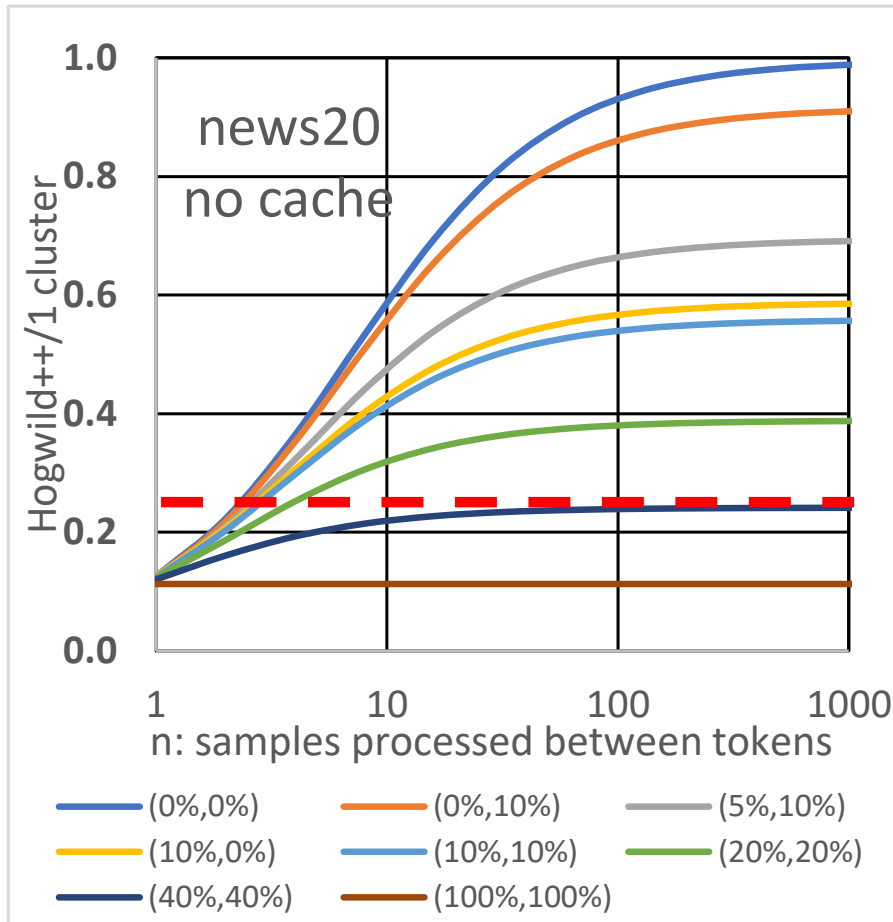
F' = 85



Mem-Cache



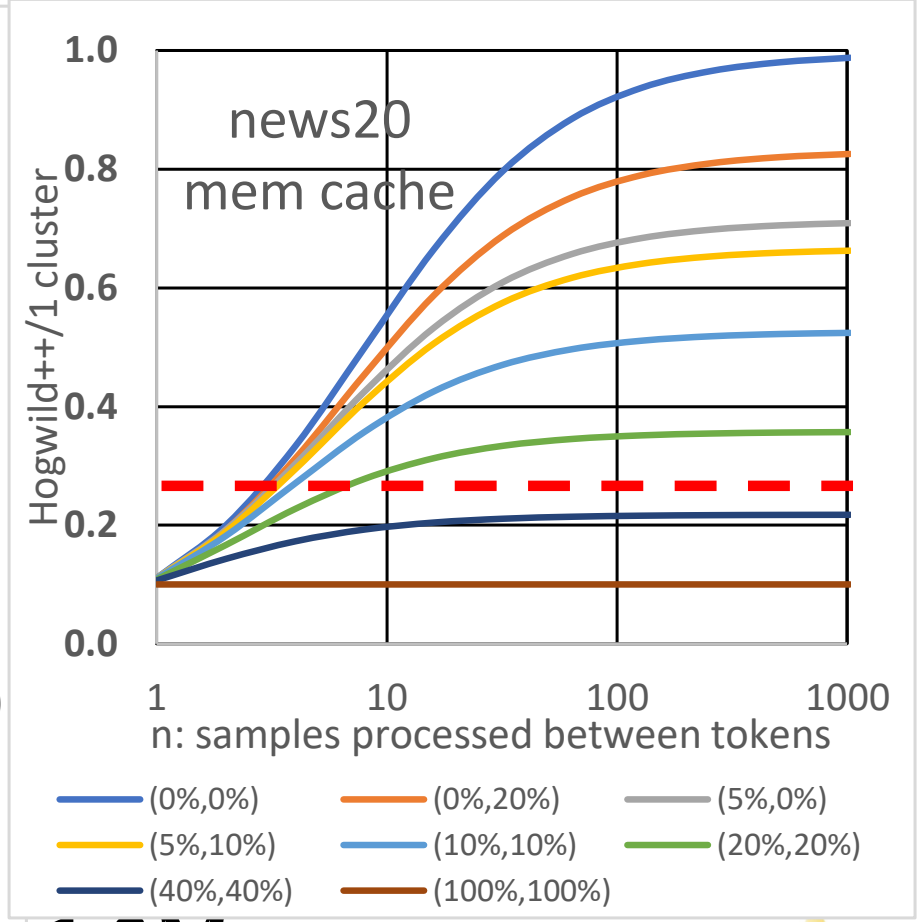
news20 Dataset



No Cache

$F = 1.3M$

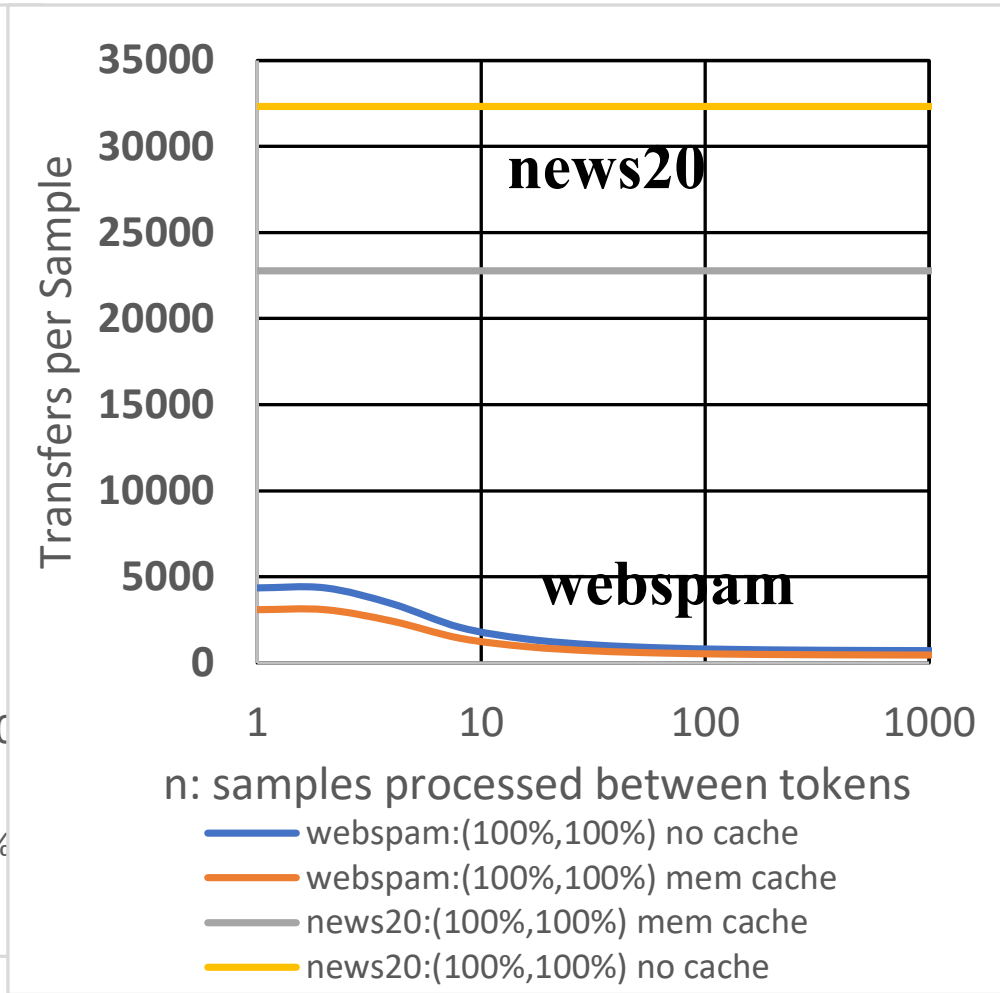
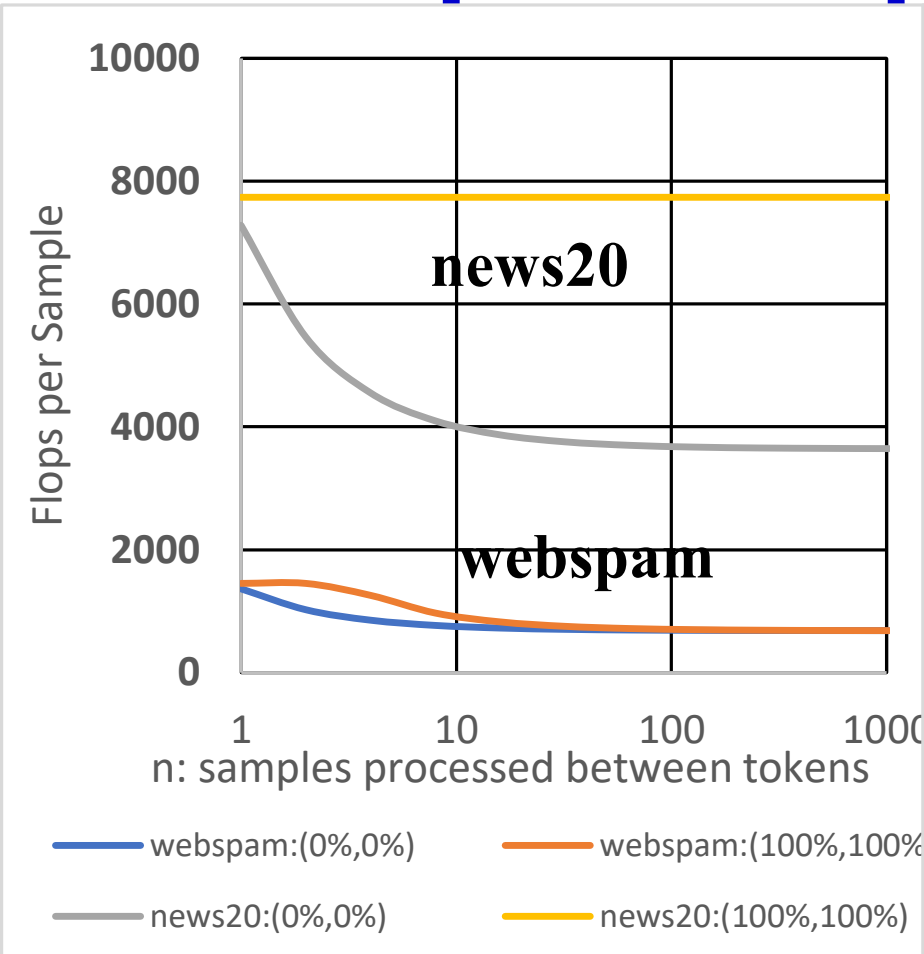
$F' = 455$



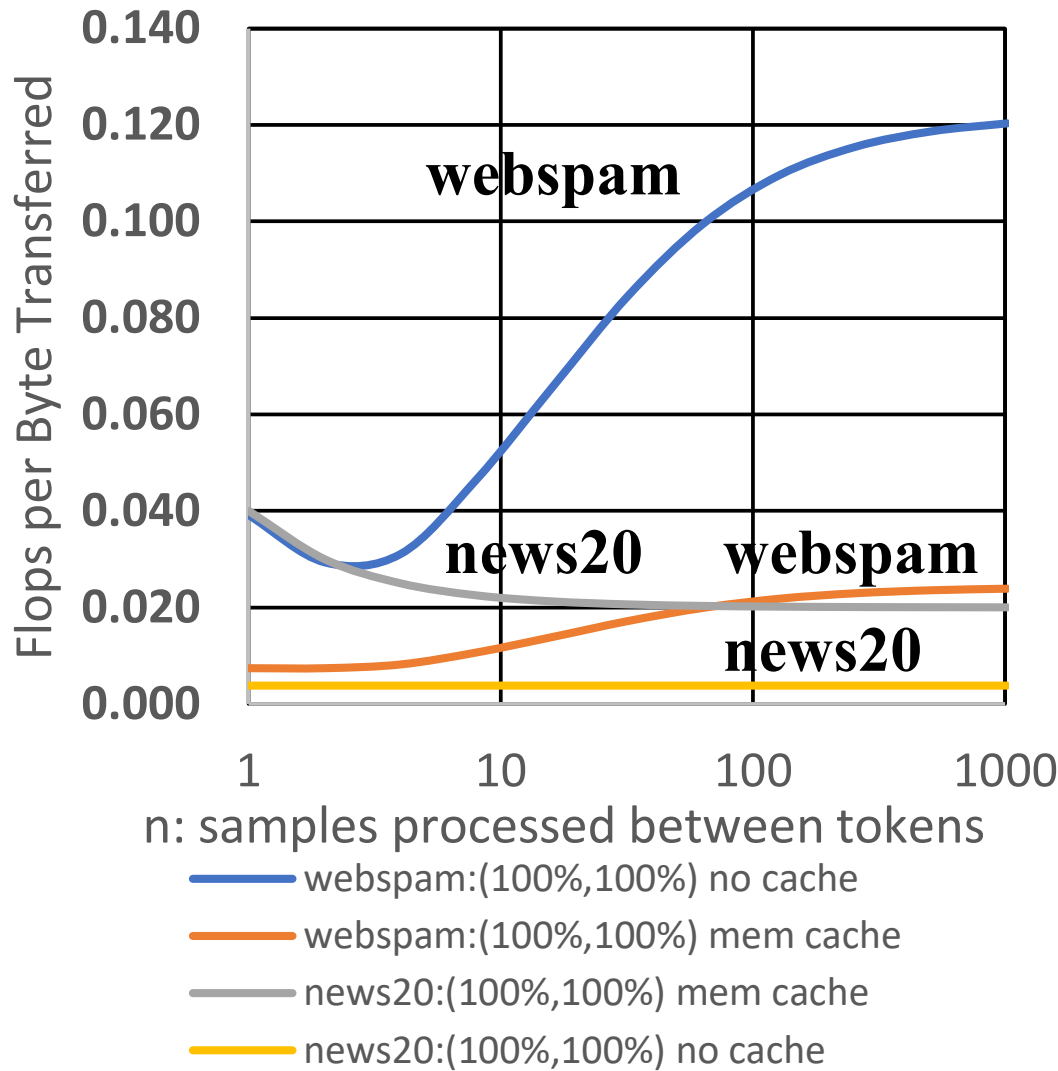
Mem-Cache



"per Sample" Metrics



Intensity



**Remember:
Today's Systems
Are 4-8**



Take-aways

- SVM very memory-bound
 - Very low flops/byte of memory bandwidth
 - No reason to have numeric-intensive accelerators
- Sparsity present in many data sets
 - And if recognized can have huge effect on time
- Multi-threading in cluster training a good match
 - Add enough threads to saturate memory
- Code complexity is in cross cluster updates
 - Ideal match to handle via migrating threads
- Efficient AMOs essential – esp. floating point
 - 40% of transfers are involved in AMOs
 - Would be much more for conventional architectures
- **All told: migrating thread good match**



Next Steps

- Consider integrating in code to track changes to objective function
- Consider effects of single thread token processing in more detail
- Develop demo code & look at news20 scaling
- Instrument to understand real-world F'' and F'''
- Look at other SGD applications



References

F. Niu, B. Recht, C. Re, and S. J. Wright. Hogwild!: A lock-free approach to parallelizing stochastic gradient descent. In *Proceedings of the 24th International Conference on Neural Information Processing Systems*, NIPS'11, pages 693–701, USA, 2011. Curran Associates Inc.

C. D. Sa, C. Zhang, K. Olukotun, and C. Ré. Taming the wild: A unified analysis of hog wild! -style algorithms. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'15, page 2674–2682, Cambridge, MA, USA, 2015. MIT Press.

H. Zhang, C. J. Hsieh, and V. Akella. Hogwild++: A new mechanism for decentralized asynchronous stochastic gradient descent. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*, pages 629–638, Dec 2016.



Projections

