



Parallel Breadth First Search on GPU Clusters

<http://mapgraph.io>

Zhisong Fu
SYSTAP, LLC
fuzhisong@systap.com

Harish Kumar Dasari
University of Utah
hdasari@sci.utah.edu

Bradley Bebee
SYSTAP, LLC
beeb@systap.com

Martin Berzins
University of Utah
mb@sci.utah.edu

Bryan Thompson
SYSTAP, LLC
bryan@systap.com
(Presenting)

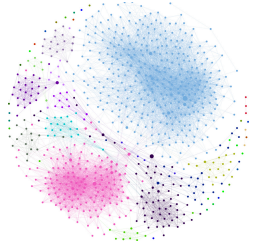
This work was (partially) funded by the DARPA XDATA program under AFRL Contract #FA8750-13-C-0002. This material is based upon work supported by the Defense Advanced Research Projects Agency (DARPA) under Contract No. D14PC00029. The authors would like to thank Dr. White, NVIDIA, and the MVAPICH group at Ohio State University for their support of this work.

Z. Fu, H.K. Dasari, B. Bebee, M. Berzins, B. Thompson. Parallel Breadth First Search on GPU Clusters. IEEE Big Data. Bethesda, MD. 2014.

Many-Core is Your Future

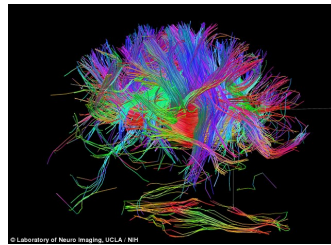


Graphs are everywhere and need for graph analytics is growing rapidly.



Communication and Social Networks

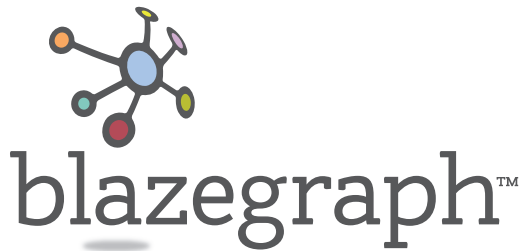
Human Brain and Biological Networks



E-Commerce and Online Service Analytics
Real-time Fraud Detection

- *Facebook has ~ 1 trillion edges in their graph.*
- *Over 30 minutes per iteration using 200 machines.*
- *All computations require multiple iterations (6-50).*

- *We could do it in seconds on a cluster of GPUs.*



SYSTAP, LLC



Small Business, Founded 2006

100% Employee Owned

Graph Database

- High performance, Scalable
 - 50B edges/node
 - High level query language
 - Efficient Graph Traversal
 - High 9s solution
- Open Source
 - Subscriptions

GPU Analytics

- Extreme Performance
 - 100s of times faster than CPUs
 - 10,000x faster than graph databases
 - 30,000,000,000 edges/sec on 64 GPUs
- DARPA funding
- Disruptive technology
 - Early adopters
 - Huge ROIs

Customers Powering Their Graphs with SYSTAP

Information
Management /
Retrieval



Genomics / Precision
Medicine



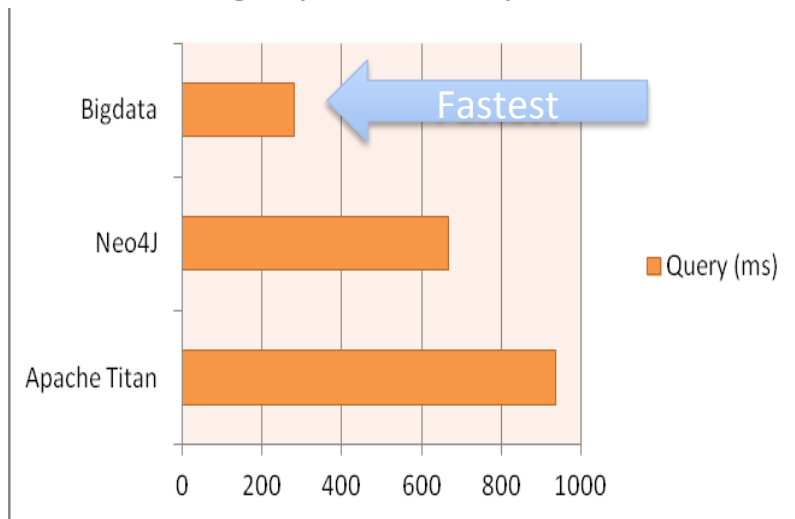
Defense, Intel, Cyber



SYSTAP is focused on building software that enable graphs at speed and scale.

- Blazegraph™ for Property and RDF Graphs
 - High Availability (HA) Architecture with Horizontal Scaling

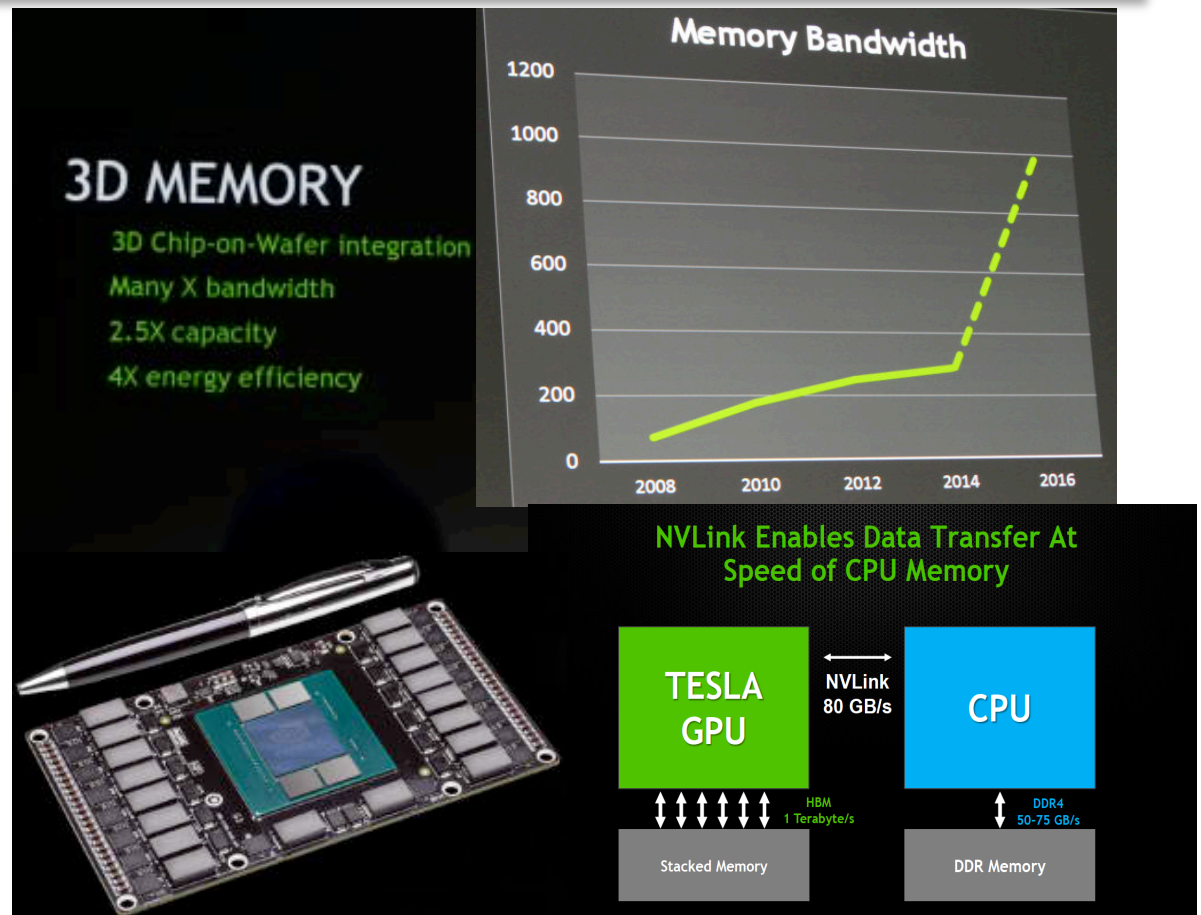
Blazegraph™ Comparison



- Mapgraph™ GPU-accelerated data parallel graph analytics
 - Vertex-Centric API.
 - Single or multi-GPU.
 - 10,000X faster than Accumulo, Neo4J, Titan, etc.
 - 3x faster than Cray XMT-2 at 1/3rd the price.




GPU Hardware Trends

- K40 GPU (today)
 - 12G RAM/GPU
 - 288 GB/s bandwidth
 - PCIe Gen 3
- Pascal GPU (Q1 2016)
 - 32G RAM/GPU
 - 1 TB/s bandwidth
 - Unified memory across CPU, GPUs
- NVLINK
 - High bandwidth access to host memory



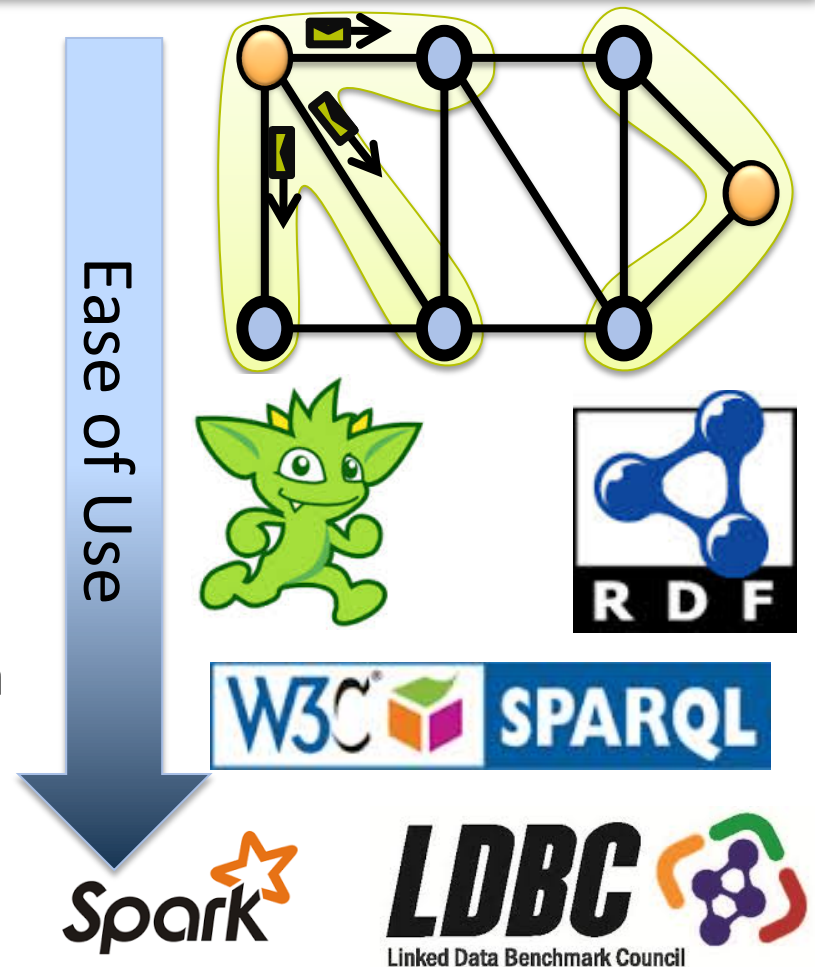
MapGraph: Extreme performance

- GTEPS is Billions (10^9) of Traversed Edges per Second.
 - This is the basic measure of performance for graph traversal.




Configuration	Cost	GTEPS	\$/GTEPS
4-Core CPU	\$4,000	0.2	\$5,333
 4-Core CPU + K20 GPU	\$7,000	3.0	\$2,333
XMT-2 (rumored price)	\$1,800,000	10.0	\$188,000
 64 GPUs (32 nodes with 2x K20 GPUs per node and InfiniBand DDRx4 – today)	\$500,000	30.0	\$16,666
 16 GPUs (2 nodes with 8x Pascal GPUs per node and InfiniBand DDRx4 – Q1, 2016)	\$125,000	>30.0	<\$4,166

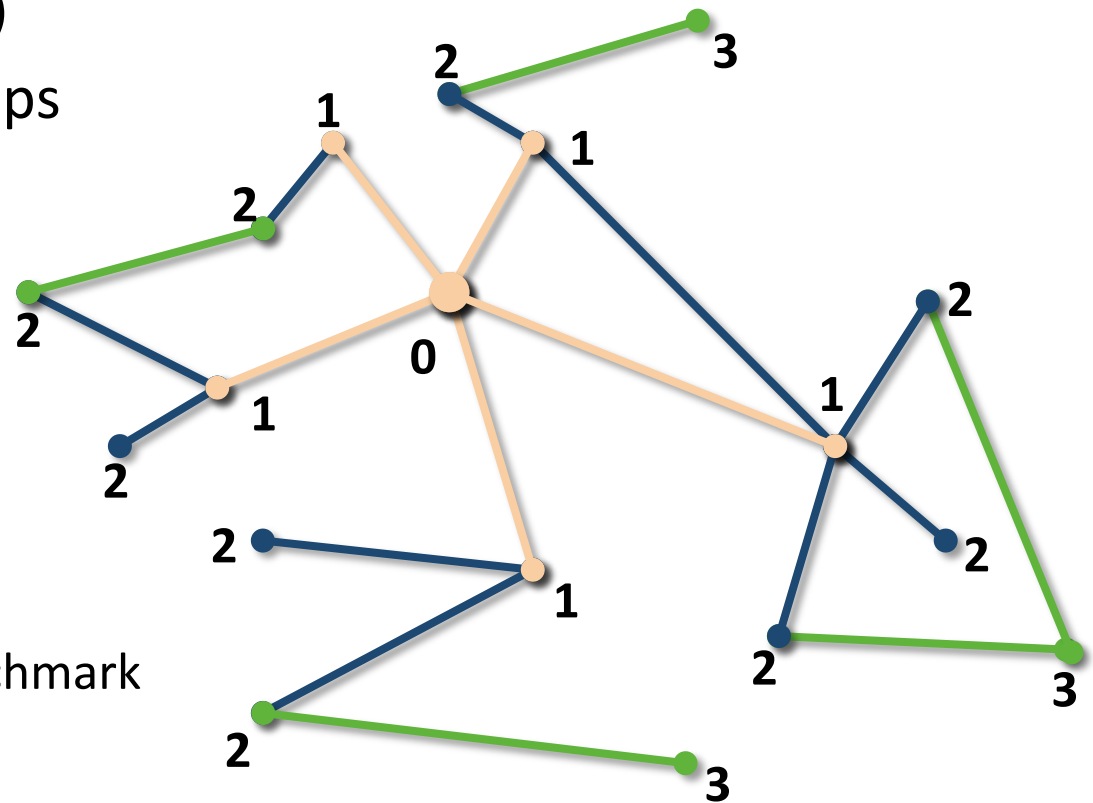
SYSTAP's MapGraph APIs Roadmap

- Vertex Centric API
 - Same performance as CUDA.
- Schema-flexible data-model
 - Property Graph / RDF
 - Reduce import hassles
 - Operations over merged graphs
- Graph pattern matching language
- DSL/Scala => CUDA code generation



Breadth First Search

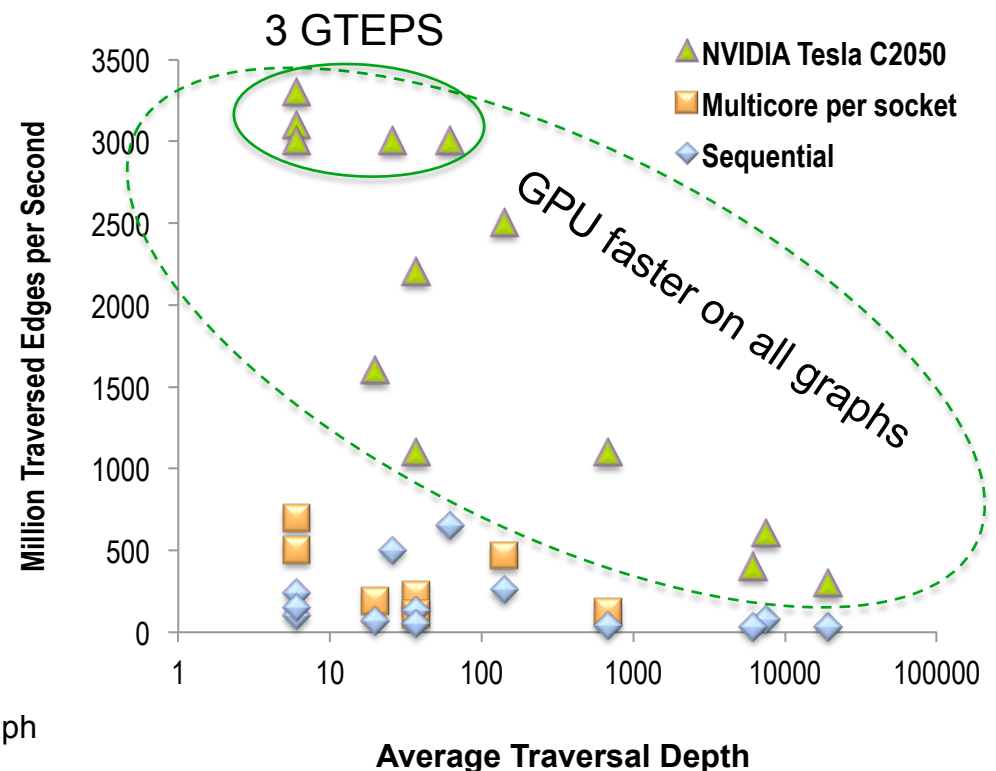
- Fundamental building block for graph algorithms
 - Including SPARQL (JOINS)
- In level synchronous steps
 - Label visited vertices
- For this graph
 - Iteration 0 
 - Iteration 1 
 - Iteration 2 
- Hard problem!
 - Basis for Graph 500 benchmark



GPUs – A Game Changer for Graph Analytics

- Graphs are a hard problem
 - Non-locality
 - Data dependent parallelism
 - Memory bus and communication bottlenecks
- GPUs deliver effective parallelism
 - 10x+ memory bandwidth
 - Dynamic parallelism

Breadth-First Search on Graphs
10x Speedup on GPUs

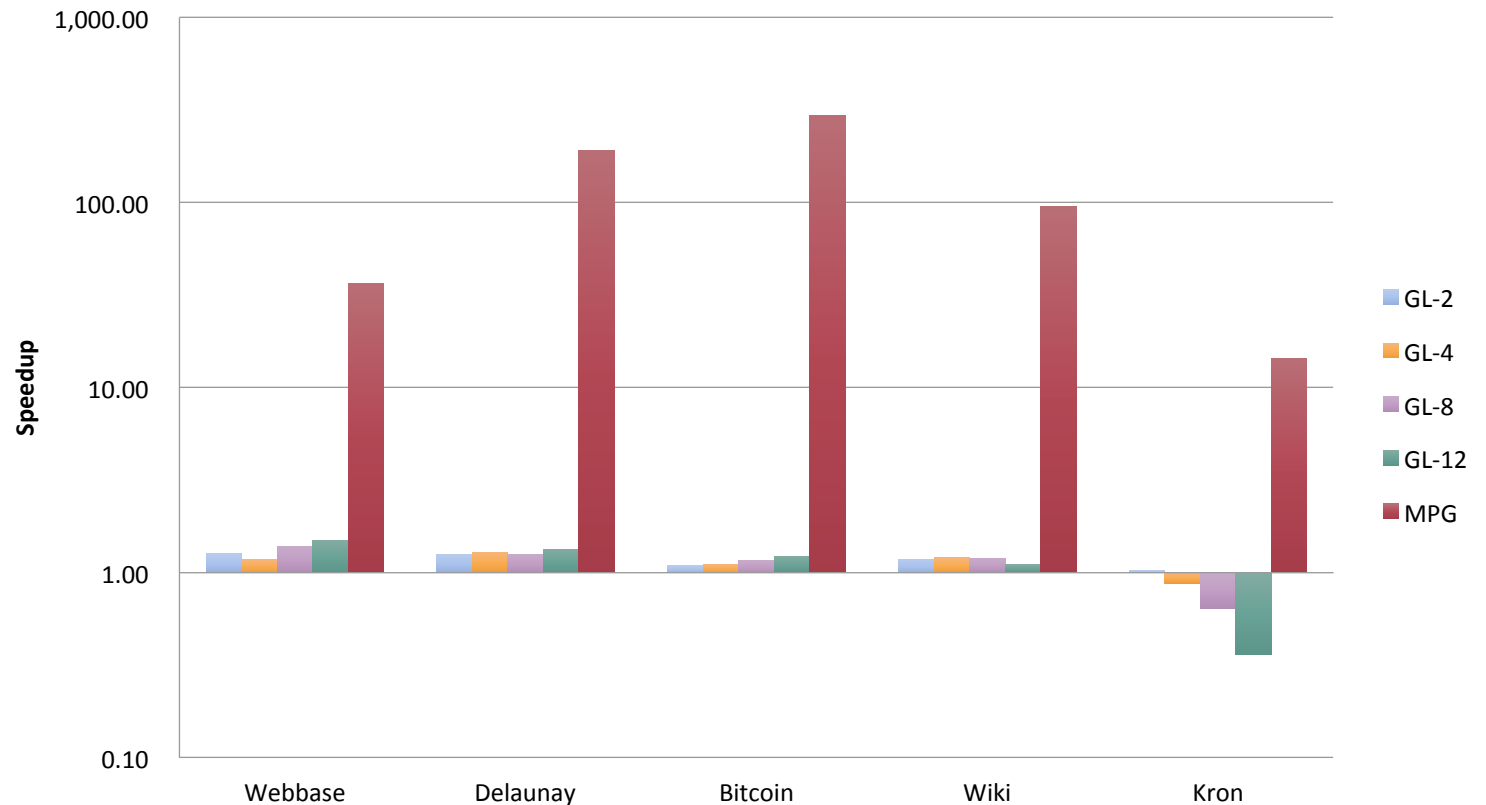


Graphic: Merrill, Garland, and Grimshaw, "GPU Sparse Graph Traversal", GPU Technology Conference, 2012

BFS Results : MapGraph vs GraphLab

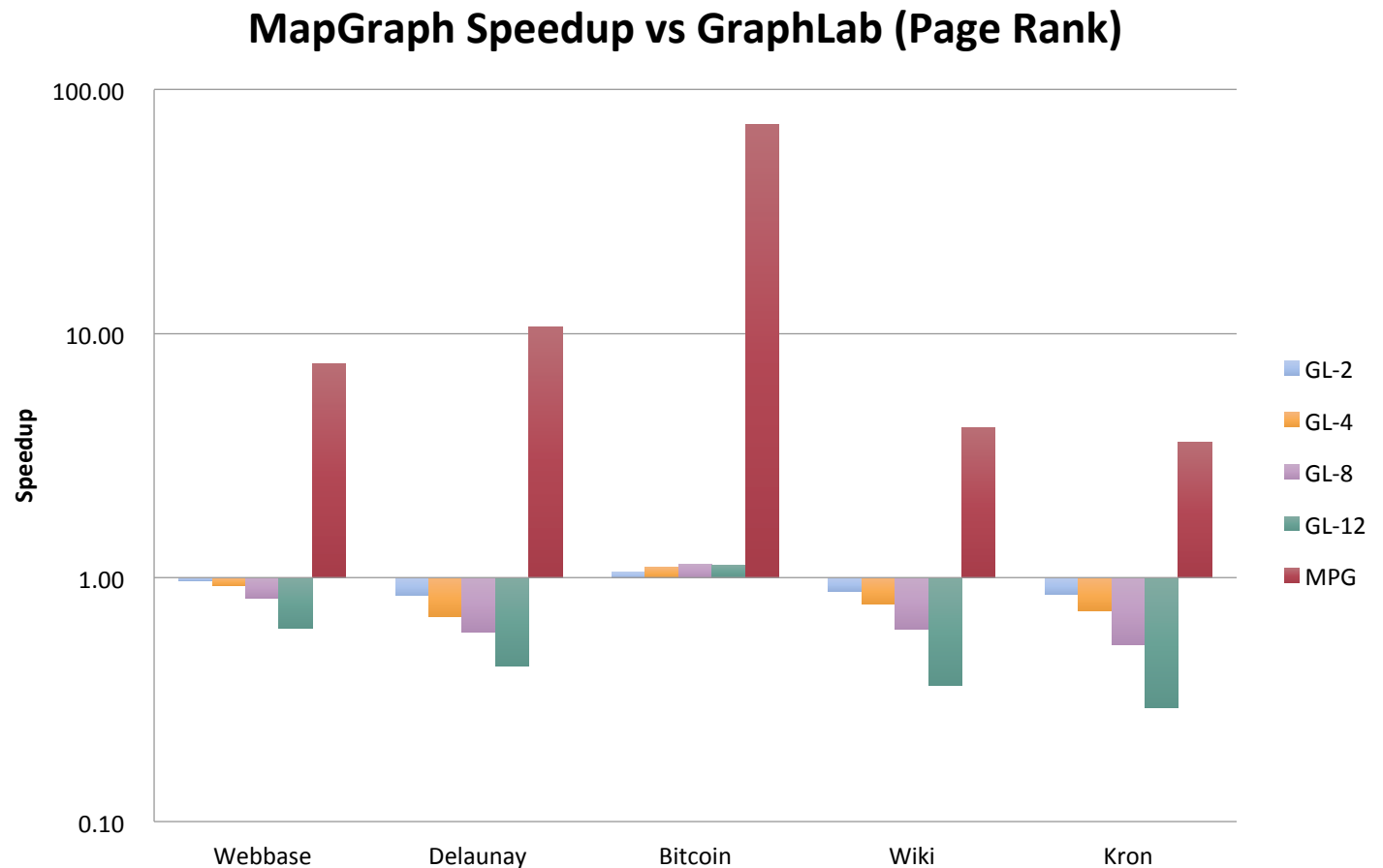
- CPU vs GPU
- GPU 15x-300x faster
- More CPU cores does not help

MapGraph Speedup vs GraphLab (BFS)



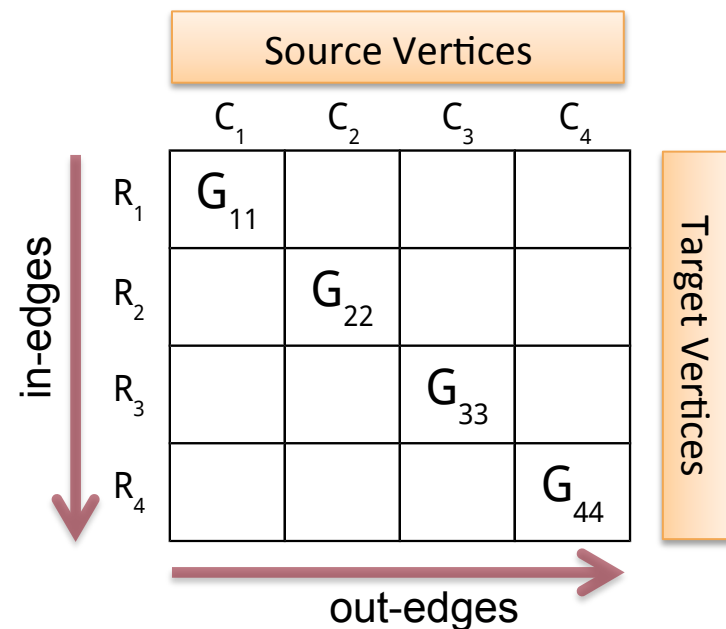
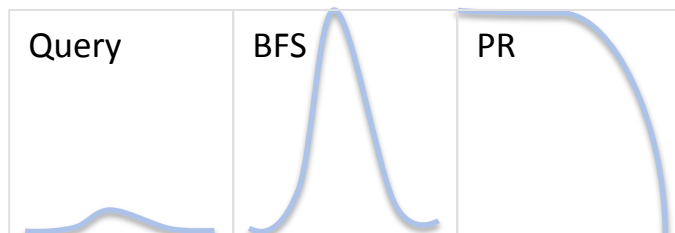
PageRank : MapGraph vs GraphLab

- CPU vs GPU
- GPU 5x-90x faster
- CPU slows down with more cores!



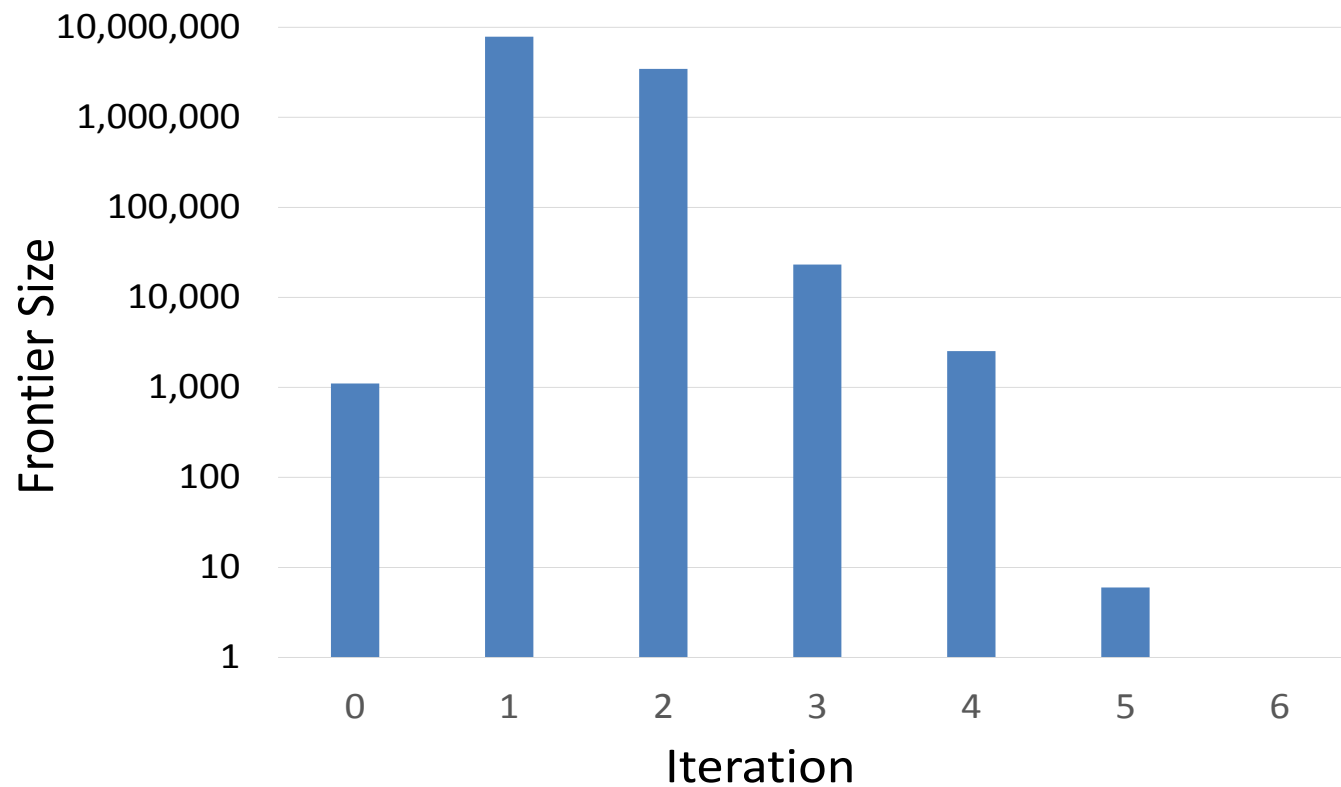
2D Partitioning (aka Vertex Cuts)

- $p \times p$ compute grid
 - Edges in rows/cols
 - Minimize messages
 - $\log(p)$ (versus p^2)
 - One partition per GPU
- Batch parallel operation
 - Grid row: out-edges
 - Grid column: in-edges
- Representative frontiers
- Parallelism – work must be distributed and balanced.
- Memory bandwidth – memory, not disk, is the bottleneck



Scale 25 Traversal

- Work spans multiple orders of magnitude.



Distributed BFS Algorithm

```
1: procedure BFS(Root, Predecessor)
2:    $In_{ij}^0 \leftarrow \text{LocalVertex}(\text{Root})$  ← Starting vertex
3:   for  $t \leftarrow 0$  do
4:     Expand( $In_i^t$ ,  $Out_{ij}^t$ ) ← Data parallel 1-hop expand on all GPUs
5:      $LocalFrontier_t \leftarrow \text{Count}(Out_{ij}^t)$  ← Local frontier size
6:      $GlobalFrontier_t \leftarrow \text{Reduce}(LocalFrontier_t)$  ← Global frontier size
7:     if  $GlobalFrontier_t > 0$  then ← Termination check
8:       Contract( $Out_{ij}^t$ ,  $Out_j^t$ ,  $In_i^{t+1}$ ,  $Assign_{ij}$ ) ← Global Frontier contraction (“wave”)
9:       UpdateLevels( $Out_j^t$ ,  $t$ , level) ← Record level of vertices in the In /
10:      else ← Out frontier
11:        UpdatePreds( $Assign_{ij}$ ,  $Preds_{ij}$ , level) ← Compute predecessors from local
12:        break ← Done
13:      end if
14:       $t++$  ← Next iteration
15:    end for
16: end procedure
```


Distributed BFS Algorithm

```
1: procedure BFS(Root, Predecessor)
2:    $In_{ij}^0 \leftarrow \text{LocalVertex}(\text{Root})$   $\leftarrow$  Starting vertex
3:   for  $t \leftarrow 0$  do
4:     Expand( $In_i^t$ ,  $Out_{ij}^t$ )  $\leftarrow$  Data parallel 1-hop expand on all GPUs
5:      $LocalFrontier_t \leftarrow \text{Count}(Out_{ij}^t)$   $\leftarrow$  Local frontier size
6:      $GlobalFrontier_t \leftarrow \text{Reduce}(LocalFrontier_t)$   $\leftarrow$  Global frontier size
7:     if  $GlobalFrontier_t > 0$  then  $\leftarrow$  Termination check
8:       Contract( $Out_{ij}^t$ ,  $Out_j^t$ ,  $In_i^{t+1}$ ,  $Assign_{ij}$ )  $\leftarrow$  Global Frontier contraction (“wave”)
9:       UpdateLevels( $Out_j^t$ ,  $t$ , level)  $\leftarrow$  Record level of vertices in the In /
10:      else  $\leftarrow$  Out frontier
11:        UpdatePreds( $Assign_{ij}$ ,  $Preds_{ij}$ , level)  $\leftarrow$  Compute predecessors from local
12:        break  $\leftarrow$  Done  $\leftarrow$  In / Out levels (no communication)
13:      end if
14:       $t++$   $\leftarrow$  Next iteration
15:    end for
16: end procedure
```

- Key differences
 - $\log(p)$ parallel scan (vs sequential wave)
 - GPU-local computation of predecessors
 - 1 partition per GPU
 - GPUDirect (vs RDMA)

Expand

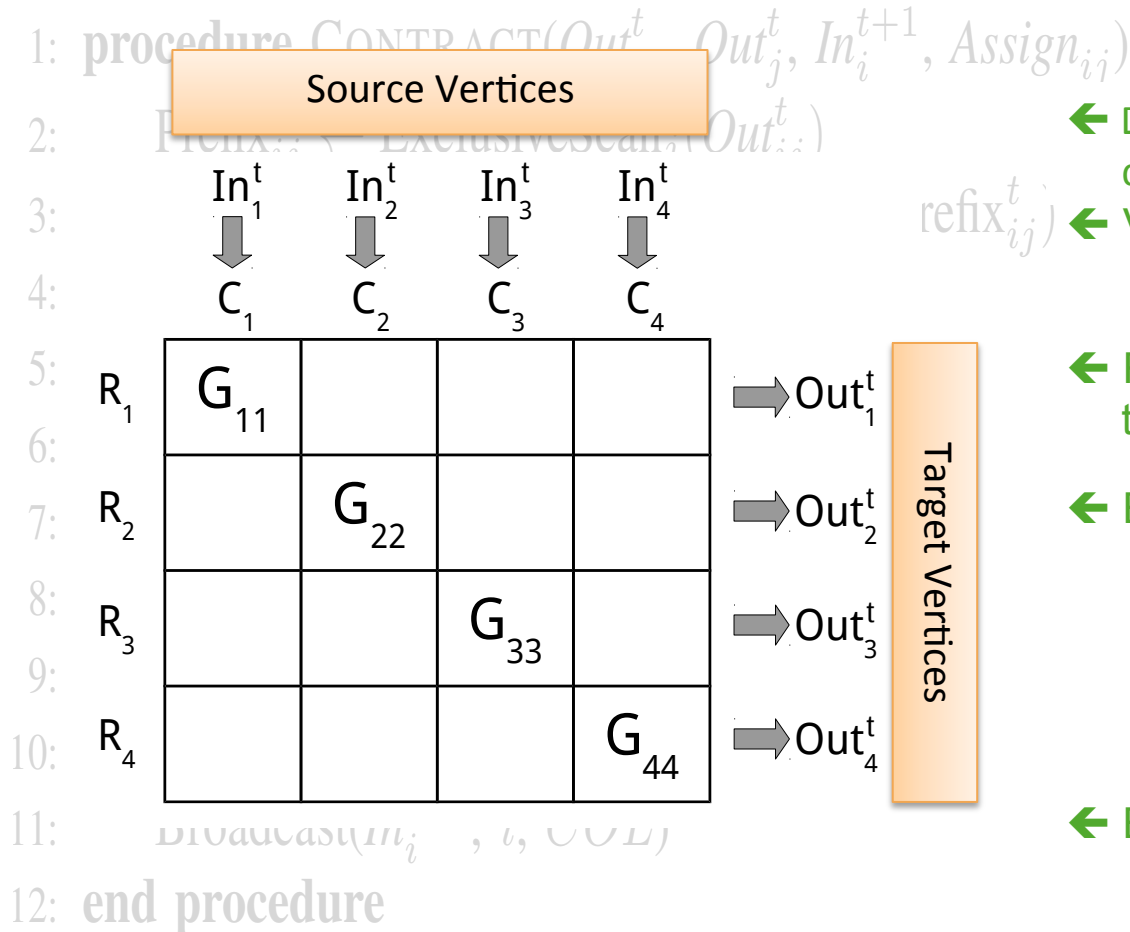
```
1: procedure EXPAND( $In_i^t$ ,  $Out_{ij}^t$ )
2:    $L_{in} \leftarrow \text{convert}(In_i^t)$ 
3:    $L_{out} \leftarrow \emptyset$ 
4:   for all  $v \in L_{in}$  in parallel do
5:     for  $i \leftarrow \text{RowOff}[v], \text{RowOff}[v + 1]$  do
6:        $c \leftarrow \text{ColIdx}[i]$ 
7:        $L_{out} \leftarrow c$ 
8:     end for
9:   end for
10:   $Out_{ij}^t \leftarrow \text{convert}(L_{out})$ 
11: end procedure
```

- The GPU implementation uses multiple strategies to handle data-dependent parallelism. See our SIGMOD 2014 paper for details.

Global Frontier Contraction and Communication

```
1: procedure CONTRACT( $Out_{ij}^t, Out_j^t, In_i^{t+1}, Assign_{ij}$ )
2:    $Prefix_{ij}^t \leftarrow ExclusiveScan_j(Out_{ij}^t)$  ← Distributed prefix sum. Prefix is vertices discovered by your left neighbors (log p)
3:    $Assigned_{ij} \leftarrow Assigned_{ij} \cup (Out_{ij}^t - Prefix_{ij}^t)$  ← Vertices first discovered by this GPU.
4:   if  $i = p$  then
5:      $Out_j^t \leftarrow Out_{ij}^t \cup Prefix_{ij}^t$  ← Right most column has global frontier for the row
6:   end if
7:   Broadcast( $Out_j^t, p, ROW$ ) ← Broadcast frontier over row.
8:   if  $i = j$  then
9:      $In_i^{t+1} \leftarrow Out_j^t$ 
10:  end if
11:  Broadcast( $In_i^{t+1}, i, COL$ ) ← Broadcast frontier over column.
12: end procedure
```

Global Frontier Contraction and Communication



← Distributed prefix sum. Prefix is vertices discovered by your left neighbors ($\log p$)

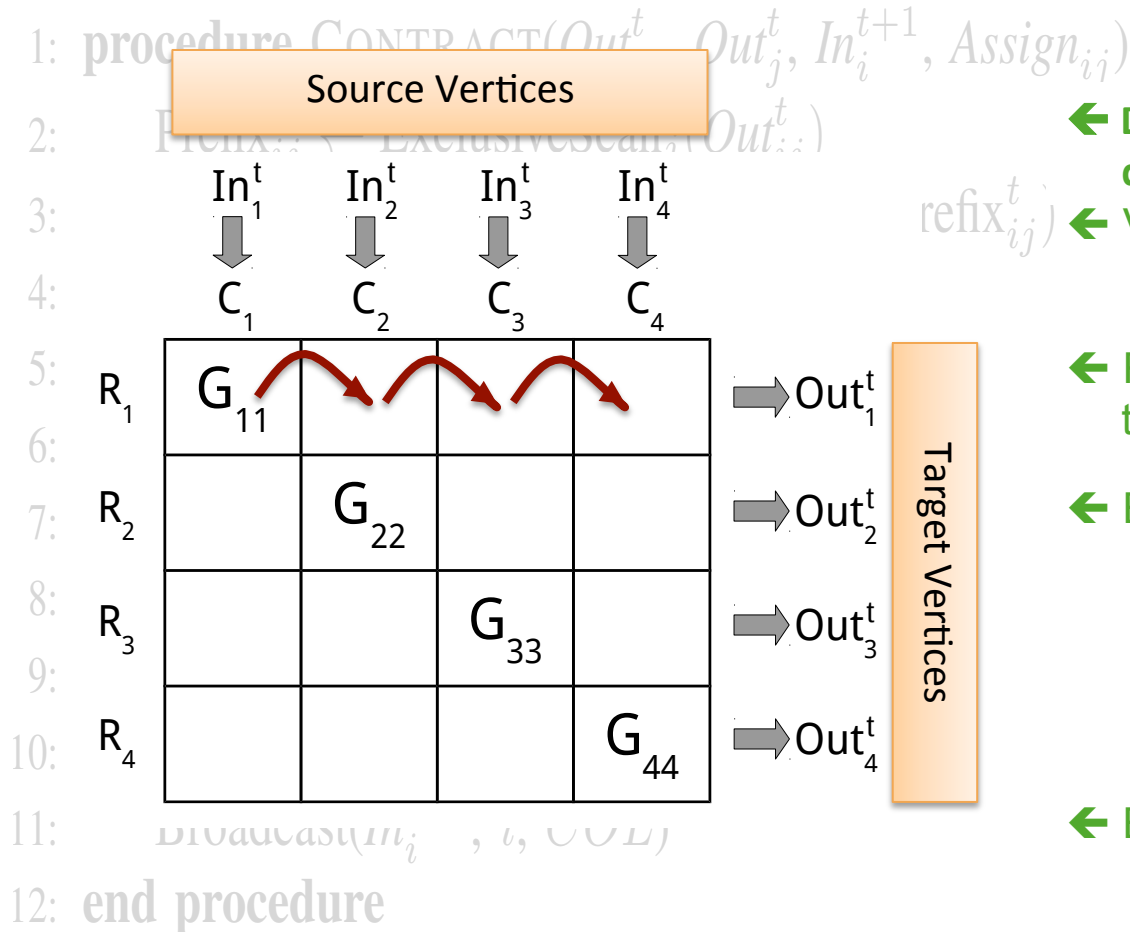
← Vertices first discovered by this GPU.

← Right most column has global frontier for the row

← Broadcast frontier over row.

← Broadcast frontier over column.

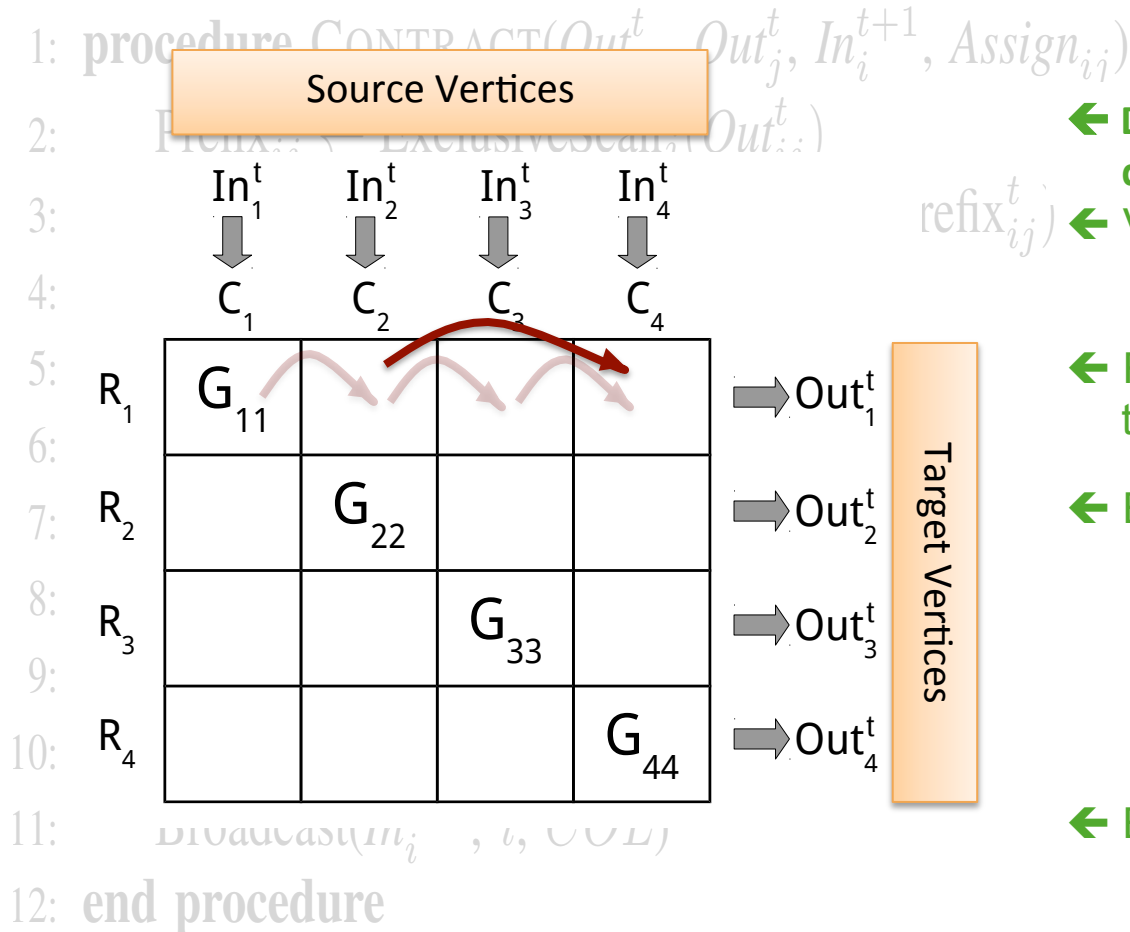
Global Frontier Contraction and Communication



- ← Distributed prefix sum. Prefix is vertices discovered by your left neighbors ($\log p$)
- ← Vertices first discovered by this GPU.
- ← Right most column has global frontier for the row
- ← Broadcast frontier over row.
- ← Broadcast frontier over column.

- We use a parallel scan that minimizes communication steps (vs work).
- This improves the overall scaling efficiency by 30%.

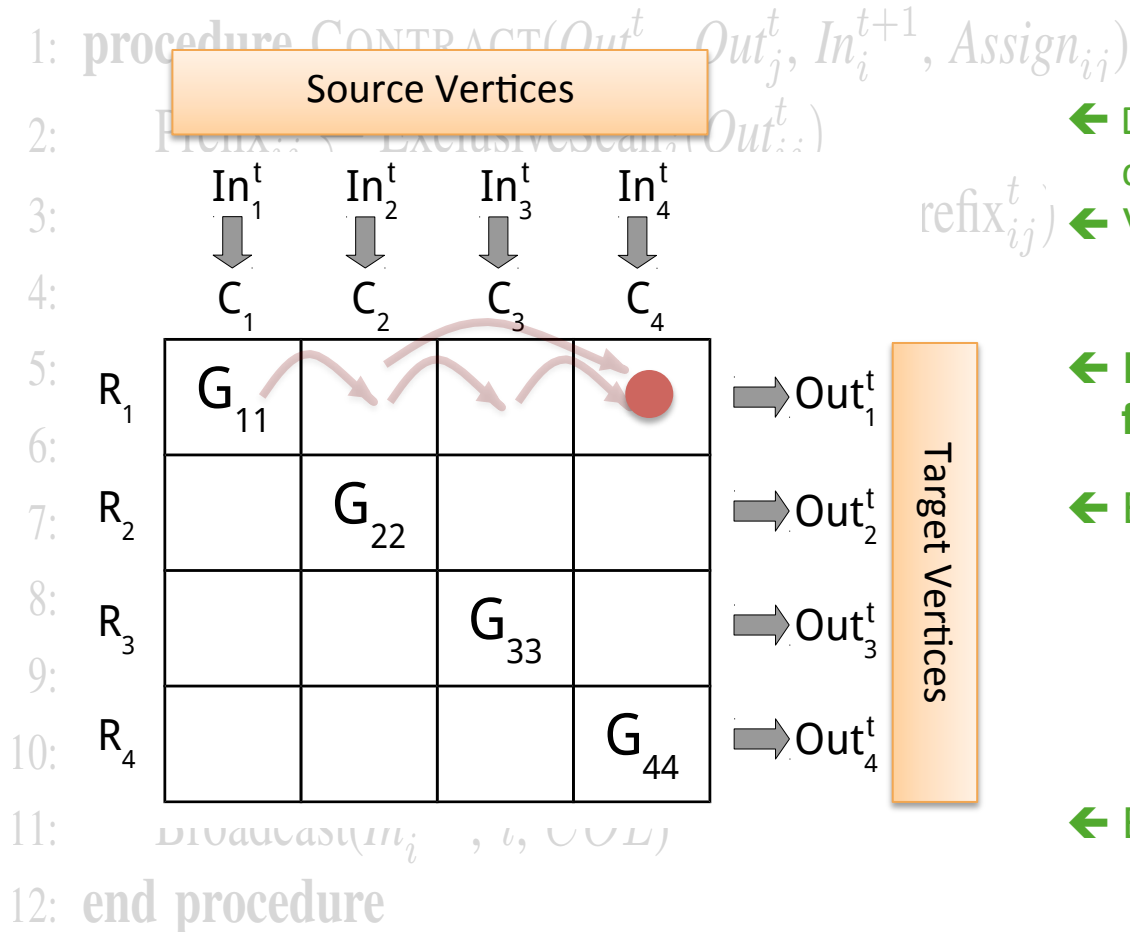
Global Frontier Contraction and Communication



- ← Distributed prefix sum. Prefix is vertices discovered by your left neighbors ($\log p$)
- ← Vertices first discovered by this GPU.
- ← Right most column has global frontier for the row
- ← Broadcast frontier over row.
- ← Broadcast frontier over column.

- We use a parallel scan that minimizes communication steps (vs work).
- This improves the overall scaling efficiency by 30%.

Global Frontier Contraction and Communication



← Distributed prefix sum. Prefix is vertices discovered by your left neighbors ($\log p$)

← Vertices first discovered by this GPU.

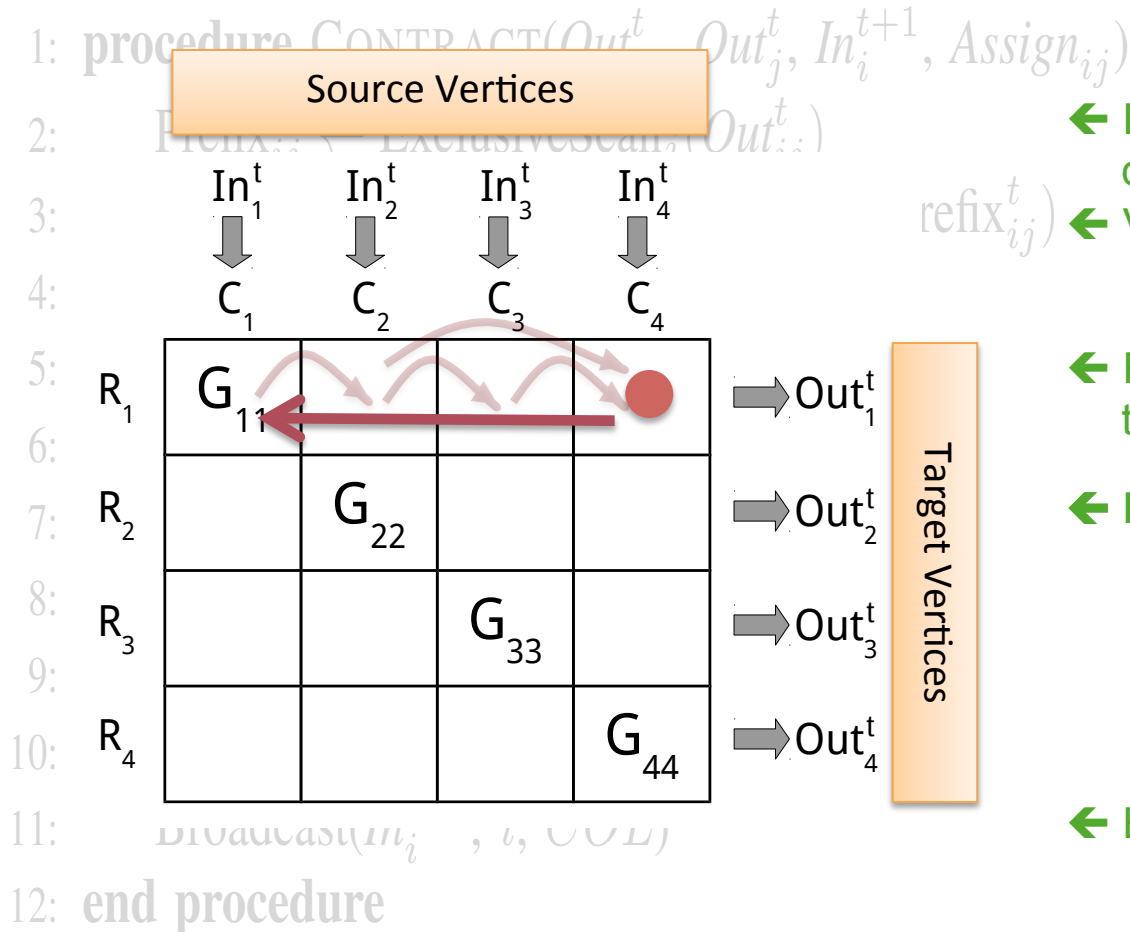
← Right most column has global frontier for the row

← Broadcast frontier over row.

← Broadcast frontier over column.

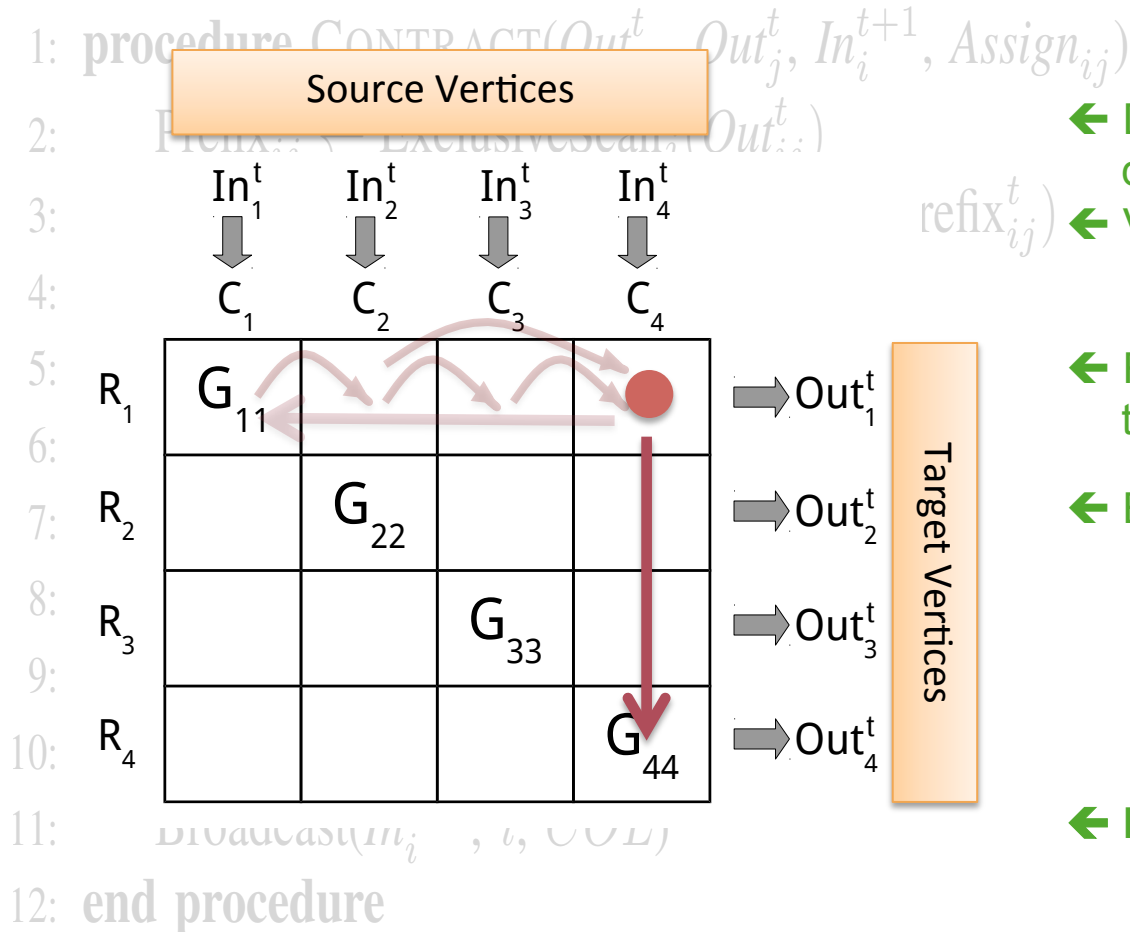
- We use a parallel scan that minimizes communication steps (vs work).
- This improves the overall scaling efficiency by 30%.

Global Frontier Contraction and Communication



- ← Distributed prefix sum. Prefix is vertices discovered by your left neighbors ($\log p$)
- ← Vertices first discovered by this GPU.
- ← Right most column has global frontier for the row
- ← **Broadcast frontier over row.**
- ← Broadcast frontier over column.

Global Frontier Contraction and Communication



- ← Distributed prefix sum. Prefix is vertices discovered by your left neighbors ($\log p$)
- ← Vertices first discovered by this GPU.
- ← Right most column has global frontier for the row
- ← Broadcast frontier over row.
- ← Broadcast frontier over column.

• All GPUs now have the new frontier.

Update Levels

```
1: procedure UPDATELEVELS( $Out_j^t$ ,  $t$ ,  $level$ )
2:   for all  $v \in Out_j^t$  in parallel do
3:      $level[v] \leftarrow t$ 
4:   end for
5: end procedure
```

- We store both the In and Out levels.
- This allows us to compute the predecessors in a GPU local manner.

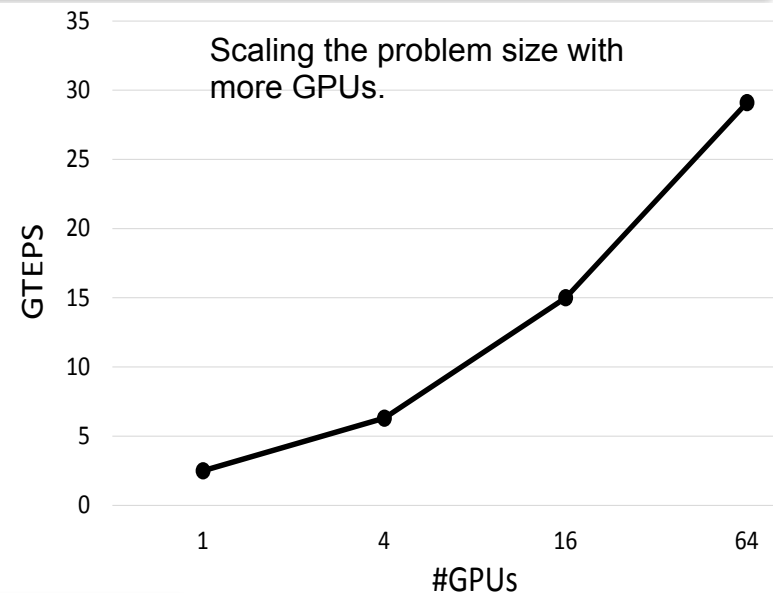
Predecessor Computation

```
1: procedure UPDATEPREDS(Assignedij, Predsij, level)
2:   for all  $v \in \text{Assigned}_{ij}$  in parallel do
3:      $\text{Pred}[v] \leftarrow -1$ 
4:     for  $i \leftarrow \text{ColOff}[v], \text{ColOff}[v + 1]$  do
5:       if  $\text{level}[v] == \text{level}[\text{RowIndex}[i]] + 1$  then
6:          $\text{Pred}[v] \leftarrow \text{RowIndex}[i]$ 
7:       end if
8:     end for
9:   end for
10: end procedure
```

- Predecessors are computed after the traversal is complete using node-local In/Out levels.
- No inter-node communication is required.

Weak Scaling

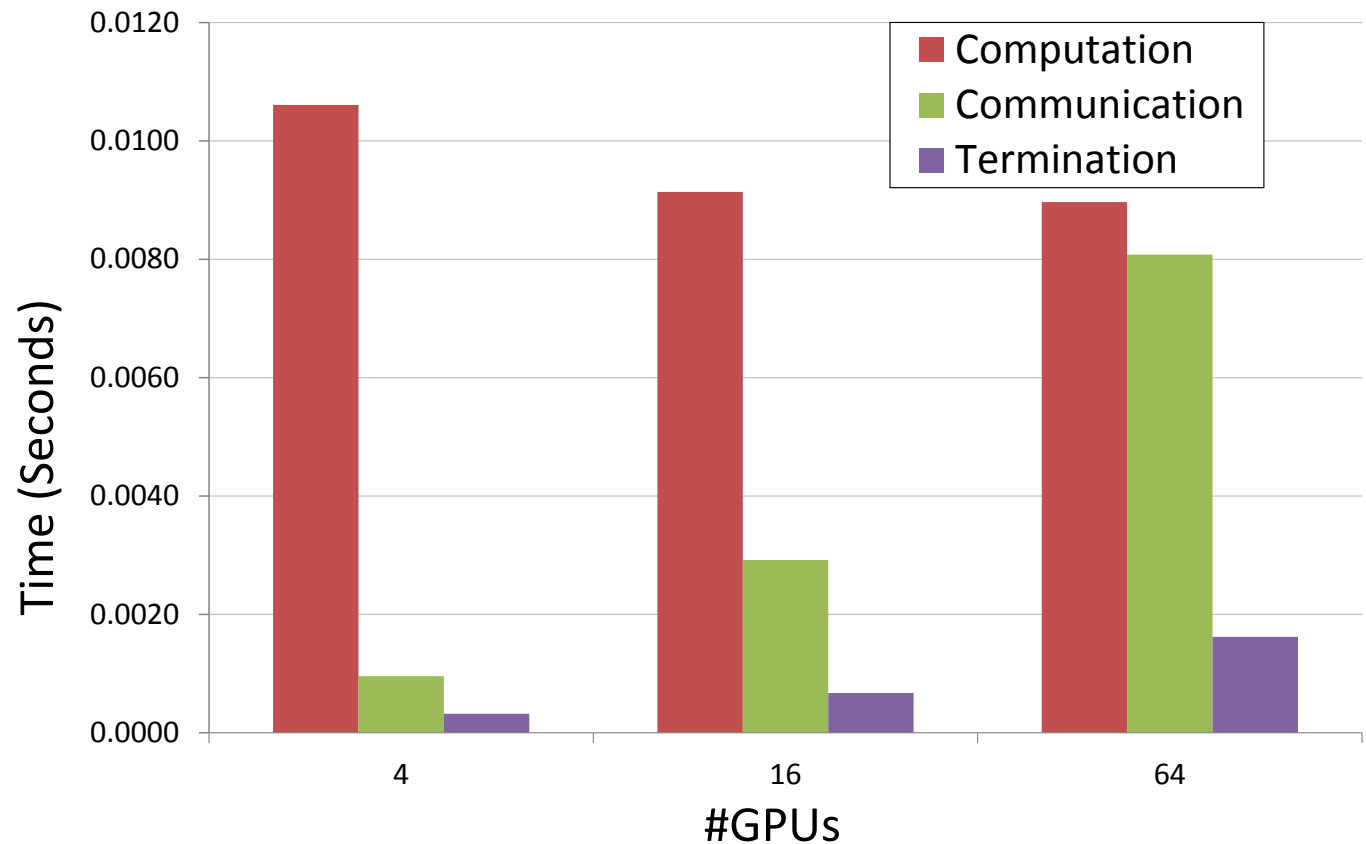
- Scaling the problem size with more GPUs
 - Fixed problem size per GPU.
- Maximum scale 27 (4.3B edges)
 - 64 K20 GPUs => .147s => 29 GTEPS
 - 64 K40 GPUs => .135s => 32 GTEPS
 - K40 has faster memory bus.



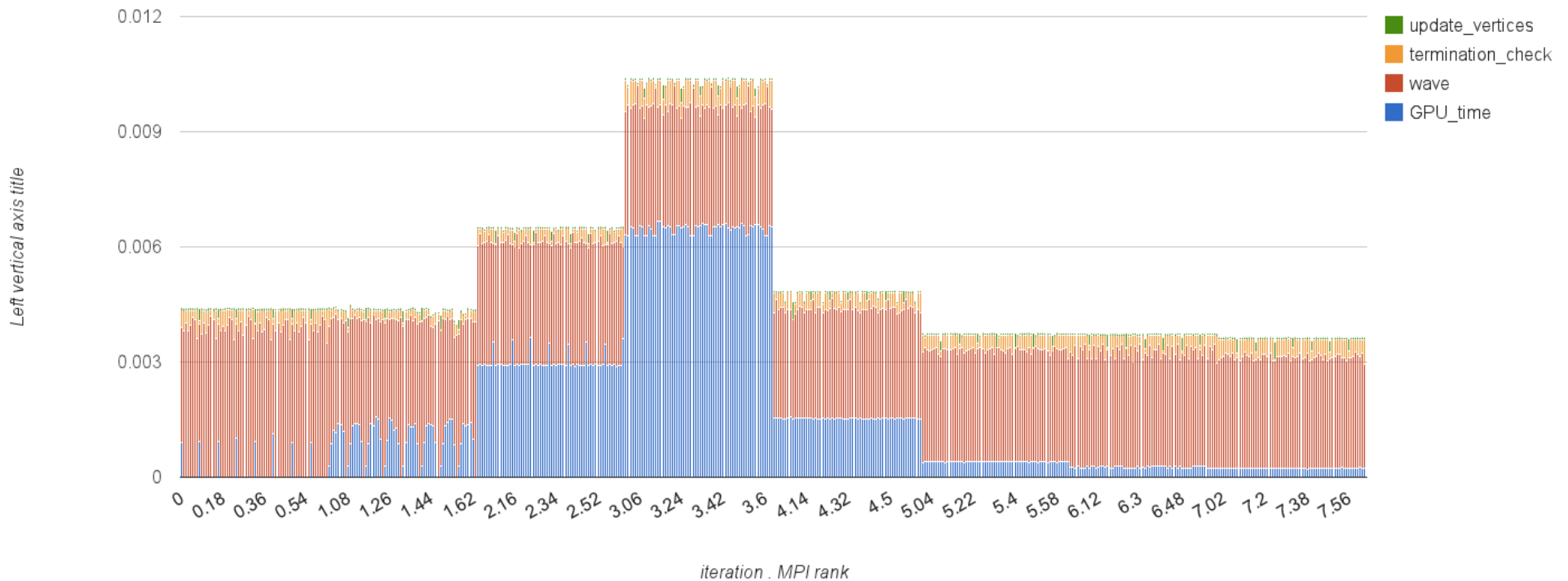
GPUs	Scale	Vertices	Edges	Time (s)	GTEPS
1	21	2,097,152	67,108,864	0.0254	2.5
4	23	8,388,608	268,435,456	0.0429	6.3
16	25	33,554,432	1,073,741,824	0.0715	15.0
64	27	134,217,728	4,294,967,296	0.1478	29.1

Central Iteration Costs (Weak Scaling)

- Communication costs are not constant.
- 2D design implies cost grows as $2 \log(2p)/\log(p)$
- How to scale?
 - Overlapping
 - Compression
 - Graph
 - Message
 - Hybrid partitioning
 - Heterogeneous computing



Costs During BFS Traversal

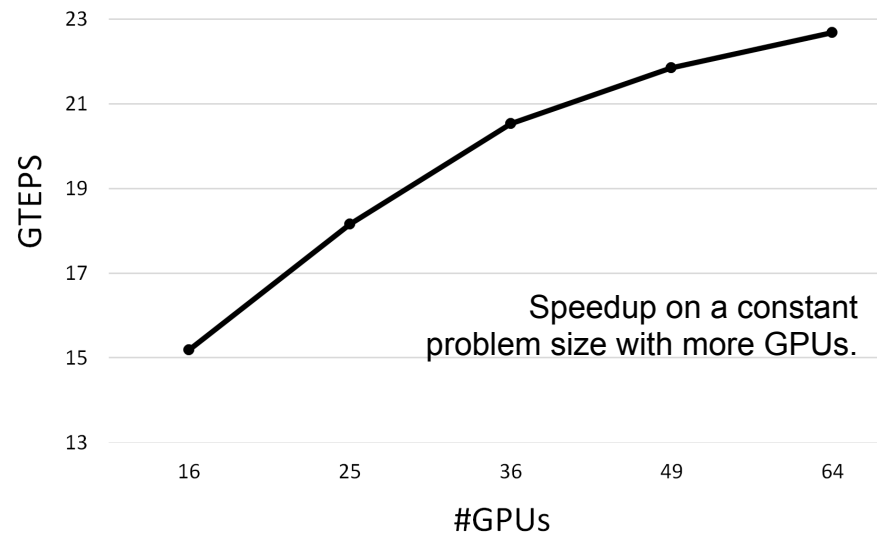


- Chart shows the different costs for each GPU in each iteration (64 GPUs).
- Wave time is essentially constant, as expected.
- Compute time peaks during the central iterations.
- Costs are reasonably well balanced across all GPUs after the 2nd iteration.

Strong Scaling

- Speedup on a constant problem size with more GPUs
- Problem scale 25
 - 2^{25} vertices (33,554,432)
 - 2^{26} directed edges (1,073,741,824)
- Strong scaling efficiency of 48%
 - Versus 44% for BG/Q

GPUs	GTEPS	Time (s)
16	15.2	0.071
25	18.2	0.059
36	20.5	0.053
49	21.8	0.049
64	22.7	0.047



Directions to Improve Scaling

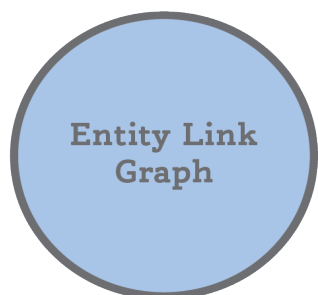
- Overlap computation with communication
 - Multiple partitions per GPU
 - Frontier compression
- Hybrid partitioning
 - Degree aware data layout + bottom up search optimization
 - This also requires asynchronous communications and per-target node buffers.
 - Graph aware partitioning plus 2D data layout
- Uintah style data warehouse
 - Hand off tasks to workers (Uintah)
 - Hybrid CPU/GPU computation strategies (TOTEM)

Concept: Accelerate Key Value Stores

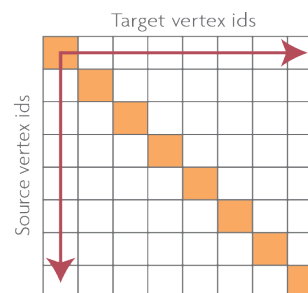
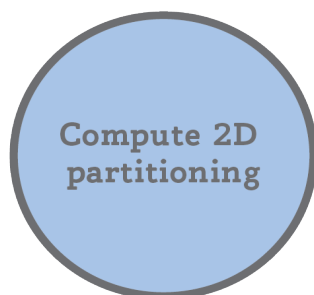
ACCUMULO CLOUD

MAP/REDUCE CLOUD

GPU CLUSTER



Parallel Scan



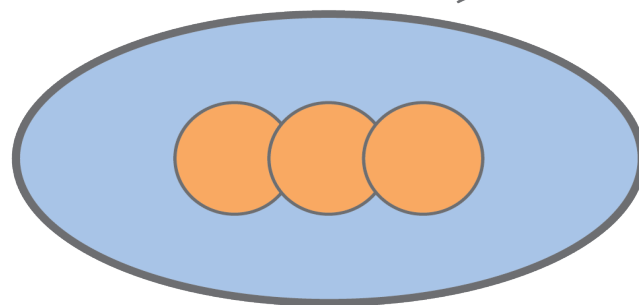
1. Graph data in Accumulo



1 BFS step @ 270 MTEPS => 2 Hours

- Accumulo + MR
- 1200 nodes
- 1T edges

2. Parallel reduce writes 2D edge partitioning



Parallel File System (HDFS)

3. Parallel read of 2D data GPU DRAM

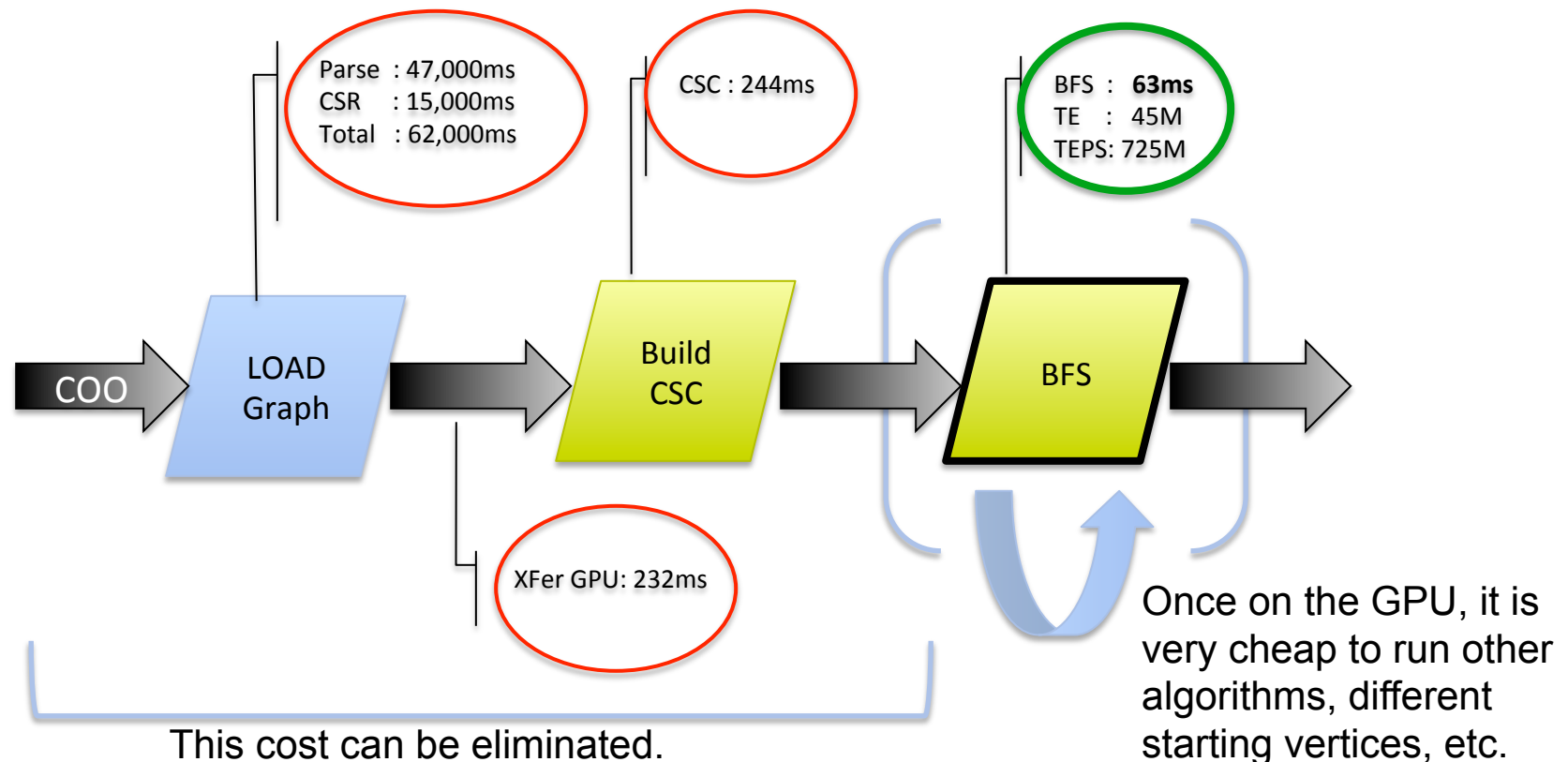
5. Parallel reduction writes output to HDFS or other target.

Assuming 10000x Speedup => .72s/step

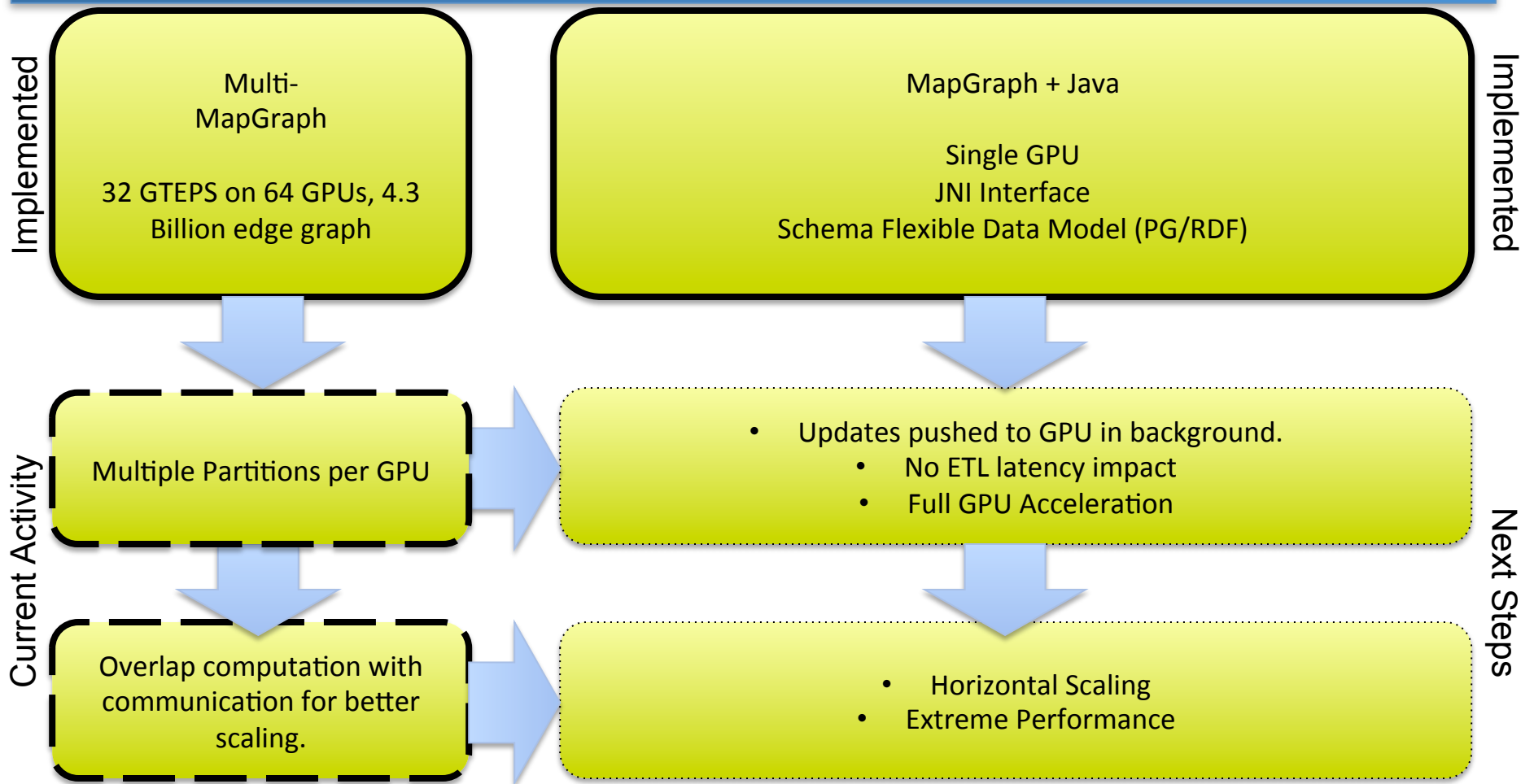
4. Extreme speed for graph analytics.

MapGraph Timings (single GPU)

- Orkut social network 2.3M vertices, 92M edges. Most time is *load* on CPU.
- Next step eliminates overhead: 62500ms => 63ms (1000x faster)



Current and Future Code Streams



MapGraph Beta Customer?

Contact
Bradley Bebee
beeb@systap.com