

Scalable and Performant Graph Processing on GPU using Approximate Computing

Somesh Singh

CS14D406

Department of Computer Science and Engineering
Indian Institute of Technology Madras

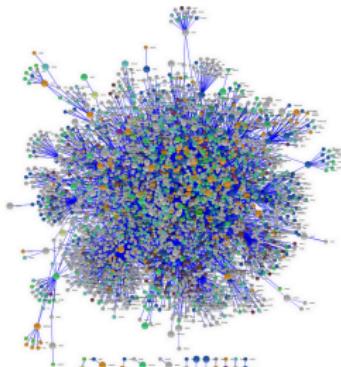
June 9, 2021



Timeline

- 
- July 14, 2014 Joined for PhD.
 - August 27, 2015 Comprehensive examination.
 - June 22, 2017 Seminar 1.
 - October 28, 2019 Seminar 2.
 - July 15, 2020 Synopsis.
 - June 9, 2021 PhD viva voce.

Graphs are Ubiquitous



Biological Network

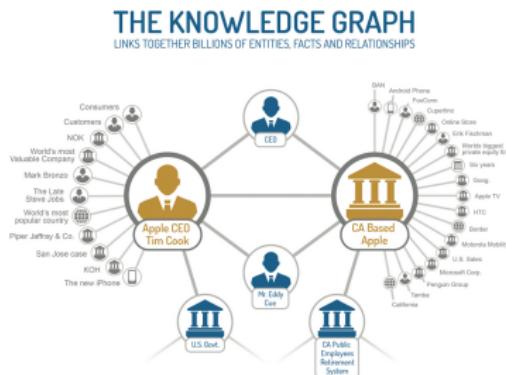


Social Network

Image Source: Google Images

Somesh Singh

Scalable and Performant Graph Processing on GPU using Approximate Computing

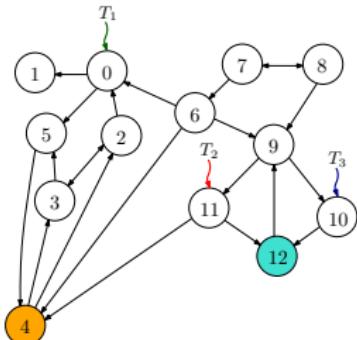


Knowledge Network



Road Network

Challenges in Parallel Graph Processing



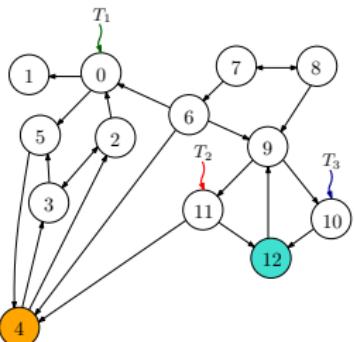
dest
1 5 0 3 2 5 2 3 4 0 4 9 6 8 7 9 10 11 12 13 14 15 16 17 18 19 20 21
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21

src
0 2 2 4 6 8 9 12 14 16 18 19 21 22
0 1 2 3 4 5 6 7 8 9 10 11 12 13

dist
0 1 2 3 4 5 6 7 8 9 10 11 12

CSR representation

Challenges in Parallel Graph Processing



dest	1	5	0	3	2	5	2	3	4	0	4	9	6	8	7	9	10	11	12	4	12	9
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	

src	0	2	2	4	6	8	9	12	14	16	18	19	21	22
0	1	2	3	4	5	6	7	8	9	10	11	12	13	

dist																				
0	1	2	3	4	5	6	7	8	9	10	11	12	13							

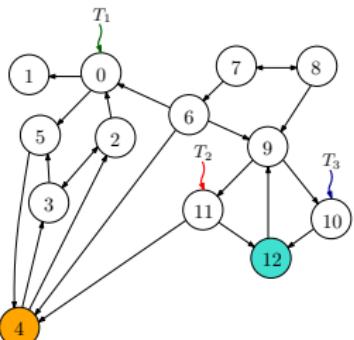
CSR representation

Assumptions

- vertex-centric model of parallelization
- propagation-based graph kernels

```
1 Graph G(V,E) = read_input(); // CPU
2 transfer_input(); // CPU → GPU
3 v.dist = ∞ ∀ v ∈ V;
4 source.dist = 0;
5 Worklist wl = {source};
6 do {
7     changed = false;
8     forall Node u : wl do {
9         for Node v : G.neighbors(u) do {
10            newVal = dist[u] + euv.wt();
11            if(newVal < dist[v]) {
12                oldVal = atomicMin(&dist[v], newVal);
13                if(newVal < oldVal) {
14                    wl.push(v);
15                    changed = true;
16                }
17            }
18        }
19    }
20 } while(changed);
21 transfer_output(); // GPU → CPU
```

Challenges in Parallel Graph Processing



dest	1	5	0	3	2	5	2	3	4	0	4	9	6	8	7	9	10	11	12	4	12	9
src	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21

src	0	2	2	4	6	8	9	12	14	16	18	19	21	22
	0	1	2	3	4	5	6	7	8	9	10	11	12	13

dist																						
	0	1	2	3	4	5	6	7	8	9	10	11	12	13								

CSR representation

Assumptions

- vertex-centric model of parallelization
- propagation-based graph kernels

```
1 Graph G(V,E) = read_input(); // CPU
2 transfer_input(); // CPU → GPU
3 v.dist = ∞ ∀ v ∈ V;
4 source.dist = 0;
5 Worklist wl = {source};
6 do {
7     changed = false;
8     forall Node u : wl do {
9         for Node v : G.neighbors(u) do {
10            newVal = dist[u] + euv.wt();
11            if(newVal < dist[v]) {
12                oldVal = atomicMin(&dist[v], newVal);
13                if(newVal < oldVal) {
14                    wl.push(v);
15                    changed = true;
16                } } } }
17 } while(changed);
18 transfer_output(); // GPU → CPU
```

- Irregular accesses: The indirection “`dist[dest[id]]`”.
- Memory-latency bound.
- Load imbalance: Skew in vertex degrees.

Our Approach

- Combine *parallelization* with *approximate computing* to make graph processing more efficient at the expense of accuracy.
- Provide tunable knobs to control the performance-accuracy trade-off.

Our Approach

- Combine *parallelization* with *approximate computing* to make graph processing more efficient at the expense of accuracy.
- Provide tunable knobs to control the performance-accuracy trade-off.

Approximate computing techniques

- ① algorithm- and architecture-independent : *Graprox*
- ② algorithm-independent but architecture-specific : *Graffix*
- ③ algorithm-specific but architecture-independent : *ParTBC*

Graprox : Techniques for Approximate Parallel Graph Processing

1 Reduced Execution

- cut-short the number of outerloop iterations.

2 Partial Graph Processing

- process only a subset of the edges in each outerloop iteration.

3 Approximate Graph Representation

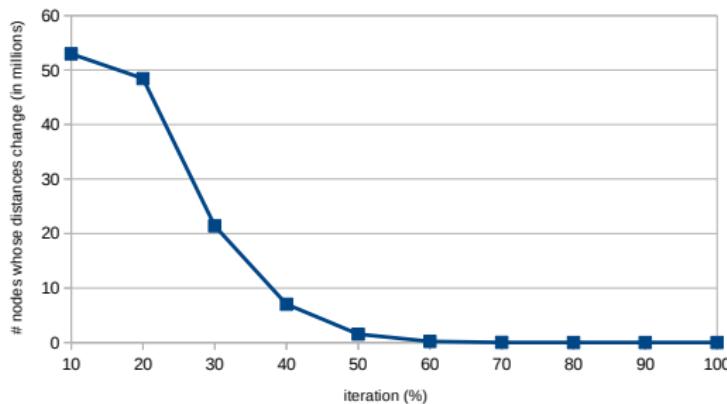
- merge nodes with overlapping neighbors, based on Jaccard's similarity.

4 Approximating Attribute Values

- round-off numeric attribute values to *discrete* values.

1 Reduced Execution

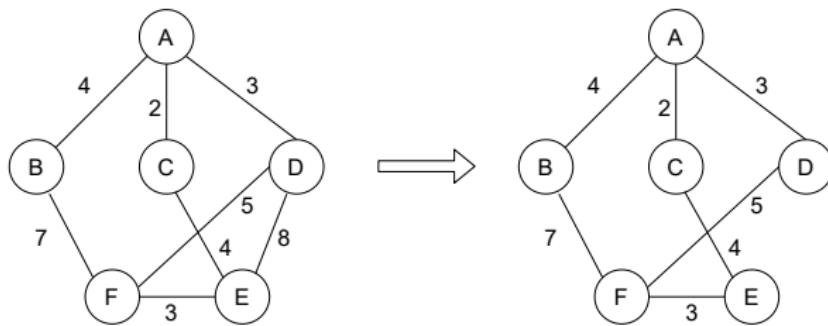
- Cut-short the number of outerloop iterations based on online stopping criteria.
- Helpful when majority of work gets done in the initial iterations.
- Results in fewer barriers in iterative parallel graph processing.



SSSP computation on rmat26 graph

2 Partial Graph Processing

- Process only a subset of the edges in each outerloop iteration.
- At each node, select the edges to be processed; ignore others.
- Helps improve performance since the work done per iteration (measured as number of edges traversed) is less.
- Processing fewer edges translates to lesser *synchronization* per iteration.



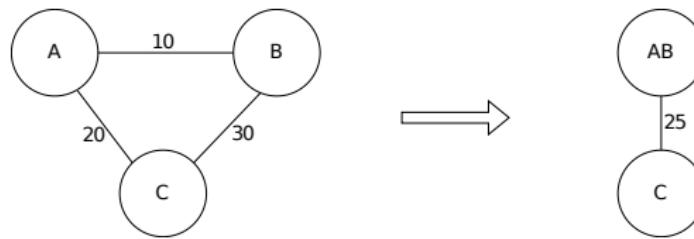
For SSSP computation

3 Approximate Graph Representation

- Lossy graph compression by merging nodes with overlapping neighbors.
- Jaccard's coefficient J_{ij} , for vertices v_i and v_j with sets of neighbors $N(v_i)$ and $N(v_j)$ respectively, is:

$$J_{ij} = \frac{|N(v_i) \cap N(v_j)|}{|N(v_i) \cup N(v_j)|}$$

- If there is a triangle a-b-c and a-b get merged:



- The compressed graph is fed as input to the exact parallel implementation of the algorithm.

4 Approximating Attribute Values

- Numeric attribute values (such as vertex attributes or edge weights) rounded-off to *discrete* values.
- Helps in reaching the termination criteria faster, in fewer rounds.
- For MST computation, rounding-off the edge-weights to the closest power of 2 helps reach the termination threshold in fewer iterations.
- SSSP transformed to BFS with careful discretization of edge weights for computing approximate shortest distance.
 - level-synchronous BFS can be implemented without explicit atomic instructions, reducing the synchronization.

1 Improving Memory Coalescing

- make the graph layout more *structured* to improve locality.
- *renumber* the graph vertices and *replicate* a select set of vertices.

2 Reducing Memory Latency

- process *well-connected* sub-graphs, iteratively, inside shared memory.

3 Reducing Thread Divergence

- normalize degrees across nodes assigned to a warp.

1 Improving Memory Coalescing

Vertex Renumbering

- Assign *nearby* ids to vertices to be accessed by warp-threads.

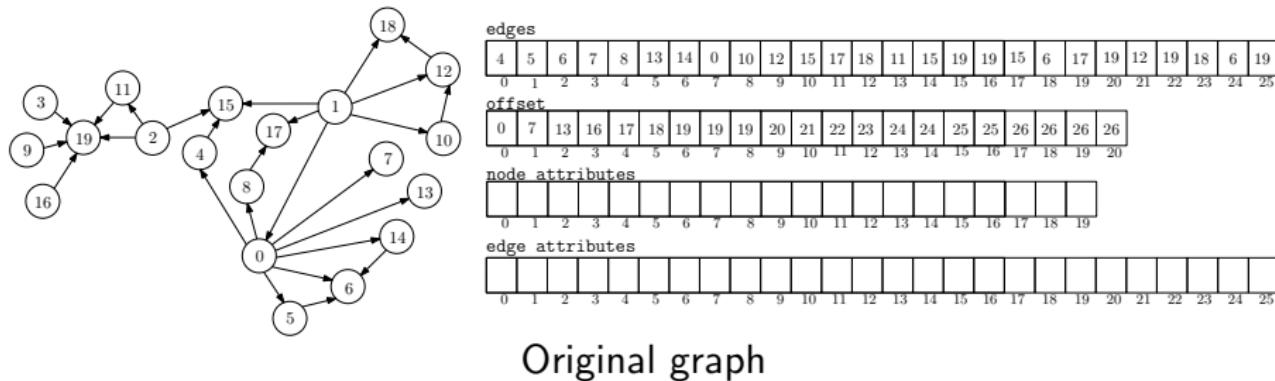
1 Improving Memory Coalescing

Vertex Renumbering

- Assign *nearby* ids to vertices to be accessed by warp-threads.

Approach

- Perform BFS from a highest outdegree node.
- Assign ids level-by-level; incrementally in a round-robin fashion at a level.

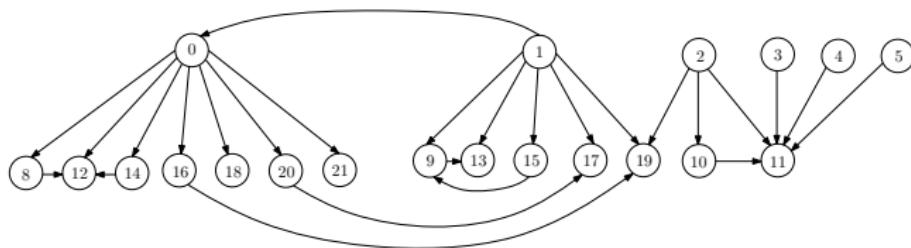


Improving Memory Coalescing

Approach

- Start a level at a multiple of k | $1 \leq k \leq \text{warp-size} \rightarrow \text{creates holes.}$
- Divide the node array (after renumbering) into chunks of size k .

Renumbered graph



Creation of *holes* after renumbering



$k = 8$

Improving Memory Coalescing

Vertex Replication

- A node occurs exactly once, so it cannot be nearby all its neighbors even after the renumbering.
- Replication brings such a node *close* to its otherwise *far* neighbors.

Improving Memory Coalescing

Vertex Replication

- A node occurs exactly once, so it cannot be nearby all its neighbors even after the renumbering.
- Replication brings such a node *close* to its otherwise *far* neighbors.

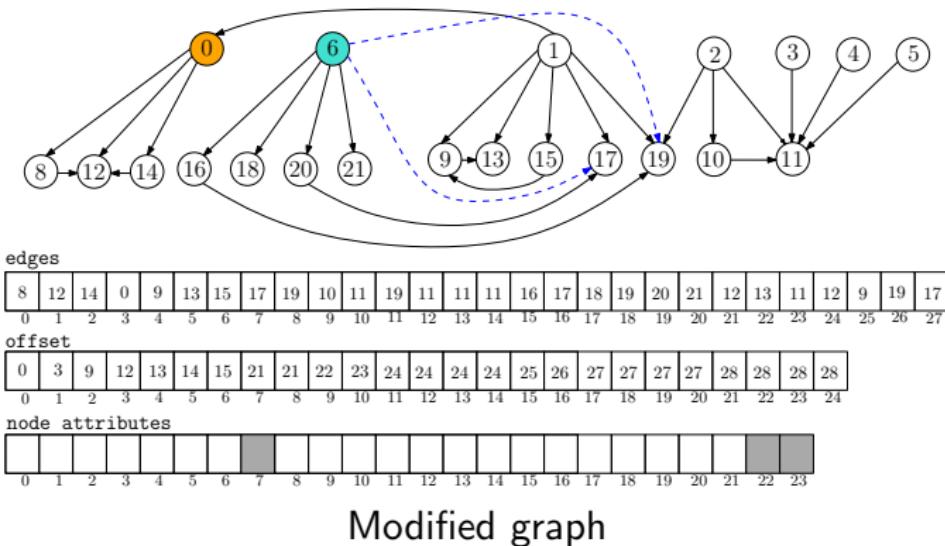
Approach

- If a node is well-connected to a chunk, replicate the node in a chunk in the previous BFS level.

$$\text{connectedness}_{\text{chunk}}^{\text{node}} \triangleq \left(\frac{\# \text{ edges to chunk from a node}}{\# \text{ non-hole nodes in chunk}} \right) \geq \text{threshold}$$
- Distribute the outgoing edges of a node among its copies.
- Add edges from node's replica to its 2-hop neighbors inside the chunk.
- Perform a merge operation on the values of the replicas after each iteration.

Improving Memory Coalescing

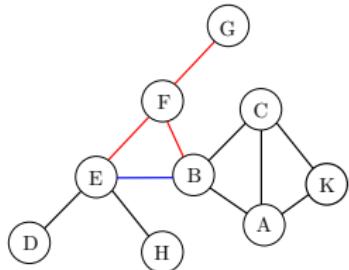
- Node 0 is well-connected to the chunk 16..23
- connectedness $_{16..23}^0 = \frac{4}{6} = 0.67 \geq 0.5$ (*threshold*)
- Node 0 is replicated in the chunk 0..7; its replica is assigned id 6.



2 Reducing Memory Latency using Shared Memory

- Clustering-coefficient measures the degree to which nodes in a graph tend to “cluster”.
- Local clustering-coefficient (LCC) of a node, X :

$$\text{LCC}_X = \frac{\# \text{ pairs of } X\text{'s neighbors that are neighbors}}{\# \text{ pairs of } X\text{'s neighbors}}$$



$$\begin{aligned}\# \text{ of pairs of } F\text{'s neighbors that are neighbors} &= 1 \\ \# \text{ of pairs of } F\text{'s neighbors} &= \binom{3}{2} = 3 \\ \text{LCC}_F &= \frac{1}{3}\end{aligned}$$

Reducing Memory Latency using Shared Memory

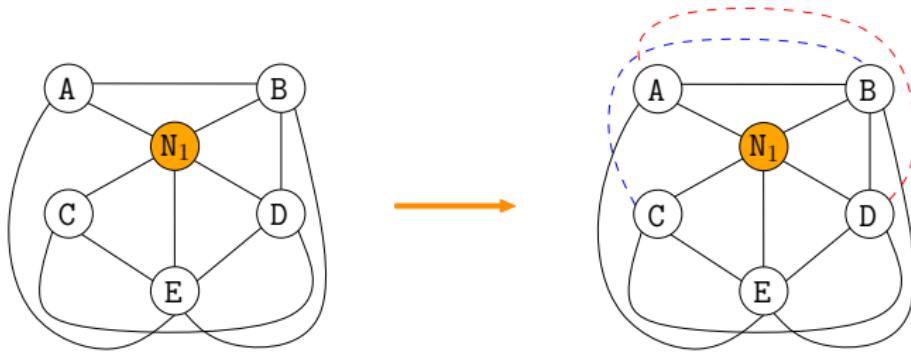
Vertices with *local clustering coefficient* (LCC) \geq threshold are more frequently accessed in iterative processing.

Reducing Memory Latency using Shared Memory

Vertices with *local clustering coefficient* (LCC) $\geq \text{threshold}$ are more frequently accessed in iterative processing.

Approach

- Increase LCC of node if $\text{LCC} \leq \text{threshold}$ and $\text{LCC} \sim \text{threshold}$.
- Boost LCC of node if $\text{LCC} \geq \text{threshold}$.
- Cap on the total number of additional edges added in the graph.



3 Reducing Thread Divergence

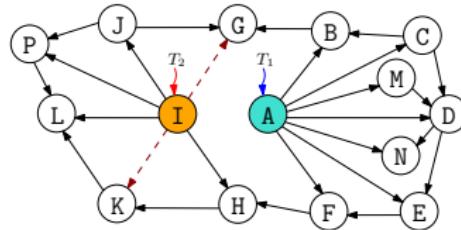
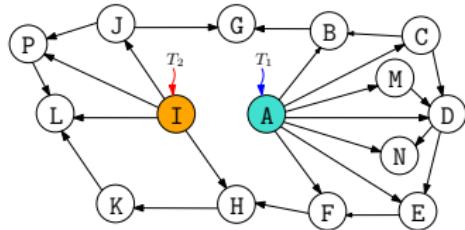
Make node degrees *nearly* uniform within each warp.

3 Reducing Thread Divergence

Make node degrees *nearly* uniform within each warp.

Approach

- Add edges to the nodes that are deficient in their connectivity.
- Add edges between 2-hop neighbors for faster convergence.
- Increase the degree of the candidate nodes to be close to $\alpha\%$ of max. degree (e.g., 85%); α is tunable.



ParTBC : Approximate top- k Betweenness Centrality Computation

Betweenness centrality (BC) is a metric that measures the significance of a vertex by the number of shortest paths leading through it.

Algorithm 1: Brandes' algorithm

Input: An unweighted graph $G(V, E)$

Output: Vertex betweenness centrality

```
1    $bc(v) = 0 \quad \forall v \in V$ 
2   for  $s \in V$  do
3       // Forward Pass: form BFS DAG  $D$ 
4       forall  $v : Node \in G$  do
5           compute  $\sigma_{sv}$ 
6           compute  $pred(s, v)$ 
7       // Backward Pass: backward traverse DAG  $D$ 
8       forall  $v : Node \in D$  do
9           compute  $\delta_s(v)$ 
10           $bc(v) += \delta_s(v)$ 
11      // Reset graph attributes
```

$$\text{Formally: } bc(v) = \sum_{s \neq t \neq v} \frac{\sigma_{st}(v)}{\sigma_{st}}$$

- σ_{sv} : # shortest paths from s to v .
 - $pred(s, w)$: immediate predecessors of w in $s \leadsto w$ path.
 - $\delta_s(v)$: dependency of v w.r.t. s .
- $$\delta_s(v) = \sum_{w | v \in pred(s, w)} \frac{\sigma_{sv}}{\sigma_{sw}} (1 + \delta_s(w))$$
- $bc(v) = \sum_{s \neq v \in V} \delta_s(v)$

ParTBC : Approximate top- k Betweenness Centrality Computation

Algorithm 2: ParTBC's approximate top- k computation

```
Input: An undirected, unweighted graph  $G(V, E)$ 
1  $bc(v) = 0 \quad \forall v \in V$                                 // initialization
2 // Phase-I
3  $G'(V, E) = \text{graphReordering}(G)$                       // vertex renumbering ( $G' \cong G$ )
4 // Phase-II
5 nextIter = true
6 while nextIter do
7     nextIter = false
8     s = getSource()                                         // pick source vertex
9     // Forward Pass: form BFS DAG D
10    forall v : Node  $\in G'$  do
11        compute  $\sigma_{sv}$ 
12        compute pred(s, v)
13
14    Let D be the DAG formed by the forward pass
15    // Backward Pass: backward traverse DAG D
16    forall v : Node  $\in D$  do
17        compute  $\delta_s(v)$ 
18         $bc(v) += \delta_s(v)$ 
19
20    // Reset graph attributes
21    if stopping criteria not met then
22        nextIter = true;
```

ParTBC : Approximate top- k Betweenness Centrality Computation

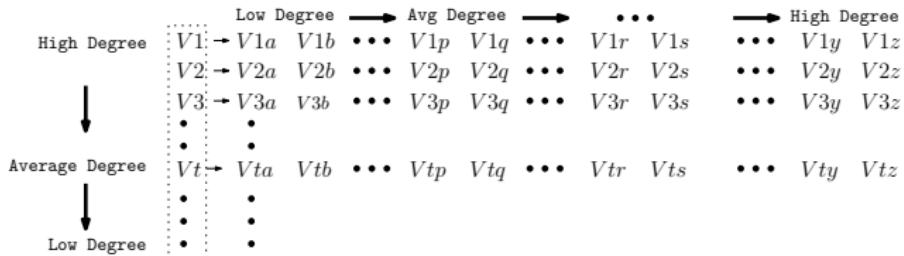
Observations

- Source does not contribute to its own BC source.
- High BC vertices are: i) high degree nodes, ii) nodes lying between large well-connected clusters.

Heuristics for selection of source vertices

- ① Random Selection
- ② Selection in Ascending and Descending Degree Order
- ③ Restricted Round Robin (RRR)
- ④ Dynamic (Dyn)
- ⑤ Dynamic Round Robin (DynRR)

Dynamic Round Robin (DynRR)



Arrangement of vertices for DynRR

- Select source nodes in a round-robin fashion for the initial 5% iterations.
- In subsequent iterations, pick the node with the least BC score as the next source.

Experimental Setup

CPU	Intel Xeon E5-2650 v2 (32 cores, 2.6 GHz, 96 GB RAM).
GPU	Nvidia (Kepler) Tesla K40C (15 SMXs, 2880 cores, 745 MHz, 12 GB global memory).
Software	CentOS 6.5, gcc 4.8.2, CUDA 8.0

Machine Configuration

Graph	$ V \times 10^6$	$ E \times 10^6$	Graph type
LiveJournal	4.8	68.9	Social network, small diameter
USA-road	23.9	57.7	Road network, large diameter
twitter	41.6	1468.3	Twitter graph 2010 snapshot
rmat26	67.1	1073.7	Synthetic scale-free graph
random26	67.1	1073.7	Synthetic random graph

Input Graphs for Graprox and Graffix

Experimental Setup

Graprox and Graffix

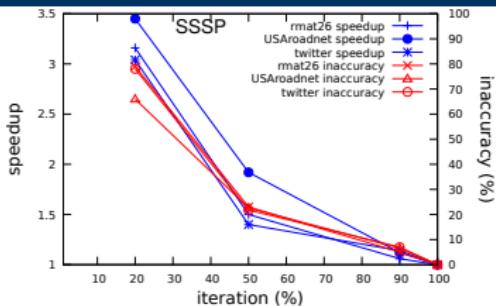
Graph Algorithms

- Single Source Shortest Path computation (SSSP)
- PageRank computation (PR)
- Strongly Connected Component computation (SCC)
- Minimum Spanning Tree computation (MST)
- Betweenness Centrality computation (BC)

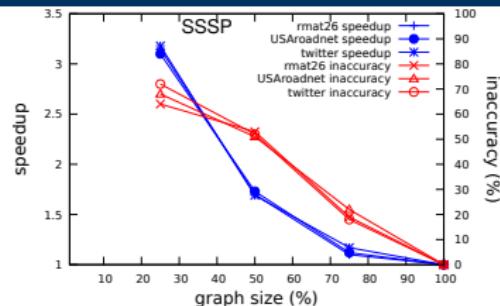
Baselines

- Baseline I:
 - SSSP, MST from [LonestarGPU](#);
 - SCC by [Devshatwar et al.](#);
 - BC by [McLaughlin and Bader](#);
 - [Our](#) exact parallel version of PR.
- Baseline II: SSSP, PR, BC from [Tigr](#).
- Baseline III: SSSP, PR, BC from [Gunrock](#).

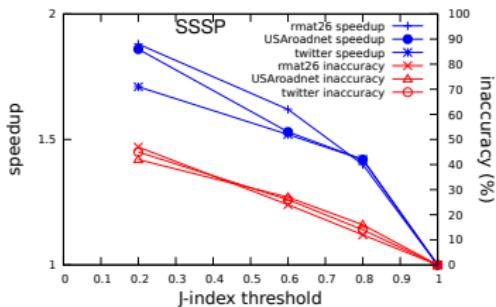
Graprox Results



Effect of reduced execution on performance and inaccuracy.



Effect of partial graph processing on performance and inaccuracy.



Effect of varying J-index threshold on performance and inaccuracy.

Takeaway:

Approximate computation of graph algorithms is a robust way of dealing with irregularities.

Grprox Results

	Technique	Mean Speedup	Mean Inaccuracy
SSSP	Reduced execution	1.34 ×	6.07%
	Partial processing of graph	1.38 ×	16.19%
	Approx. graph representation	1.22 ×	13.87%
	Approx. attributes	1.92 ×	17.64%
MST	Reduced execution	1.18 ×	16.05%
	Partial processing of graph	1.65 ×	17.44%
	Approx. graph representation	1.44 ×	15.17%
	Approx. attributes	1.48 ×	19.07%
SCC	Reduced execution	1.25 ×	18.26%
	Partial processing of graph	1.32 ×	19.61%
	Approx. graph representation	1.45 ×	20.11%
	Approx. attributes	—	—
PR	Reduced execution	2.03 ×	2.75%
	Partial processing of graph	1.43 ×	15.74%
	Approx. graph representation	1.37 ×	13.70%
	Approx. attributes	—	—
BC	Reduced execution	1.74 ×	18.07%
	Partial processing of graph	1.42 ×	16.73%
	Approx. graph representation	1.33 ×	14.35%
	Approx. attributes	—	—

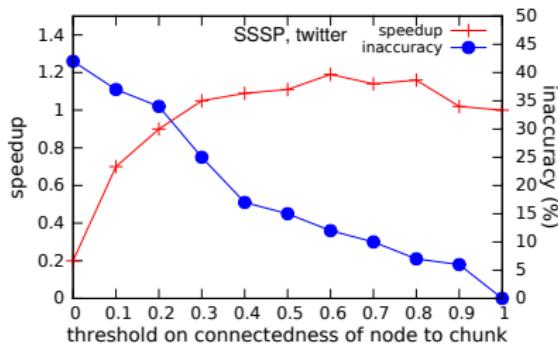
Takeaway:

Approximate computing techniques are consistently helpful in improving the execution performance of graph analytics in exchange for inaccuracy.

Graffix Results

Improving Memory Coalescing

	Baseline I	Baseline II	Baseline III
Mean Speedup	1.16×	1.10×	1.14×
Mean Inaccuracy	10%	9%	9%



Effect of varying the threshold for node replication on memory coalescing.
(Chunk size is set to 16.)

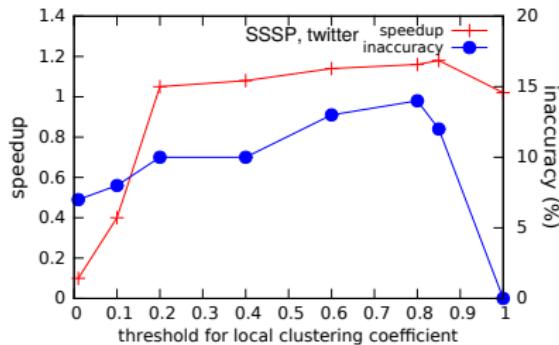
Takeaway:

Desired accuracy and performance for an algorithm – input graph pair can be achieved by tuning the chunk size and the threshold for node replication.

Graffix Results

Reducing Memory Latency

	Baseline I	Baseline II	Baseline III
Mean Speedup	1.20×	1.19×	1.19×
Mean Inaccuracy	13%	12%	12%



Effect of varying the LCC threshold on memory latency.

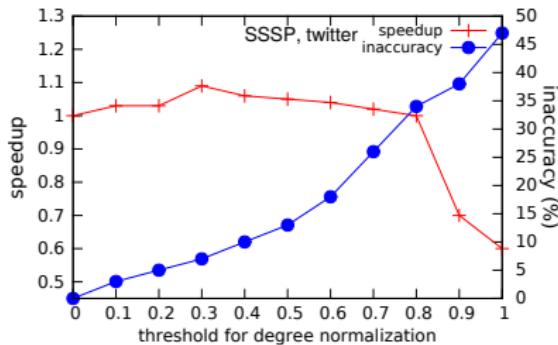
Takeaway:

Appreciable speedup, with low inaccuracy, can be achieved by processing well-connected subgraphs inside shared memory.

Graffix Results

Reducing Thread Divergence

	Baseline I	Baseline II	Baseline III
Mean Speedup	1.07×	1.03×	1.07×
Mean Inaccuracy	8%	8%	8%



Effect of varying the threshold for degree normalization.

Takeaway:

Small speedup with low inaccuracy can be achieved using a low threshold for degree normalization.

Experimental Setup

ParTBC

Graph	V	E	Graph type
fb-Friendships (FB)	63,731	817,035	Facebook friendship graph
usroad48 (RNUS)	102,615	147,656	Continental US road network
rmat17 (RMT)	130,977	2,091,451	Synthetic scale-free graph
random17 (RNM)	131,072	2,096,902	Synthetic random graph
roadnetSF (RNSF)	174,424	221,802	San Francisco road network
loc-Gowalla (LG)	196,591	950,327	Location-based social network
soc-Pokec (SP)	1,632,803	30,622,564	Online social network

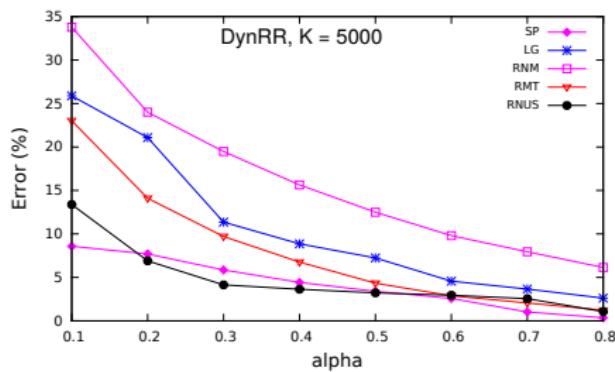
Input Graphs for ParTBC

Baseline: Exact GPU-parallel implementation of betweenness centrality.

ParTBC Results

Graph	Speedup w.r.t. exact parallel	Speedup breakdown (w.r.t. exact parallel)	
		VR	DynRR
fb-Friendships	2.80×	1.22×	2.29×
usroad48	2.68×	1.13×	2.37×
rmat17	2.65×	1.19×	2.22×
random17	1.67×	1.04×	1.61×
roadnetSF	2.71×	1.13×	2.40×
loc-Gowalla	2.48×	1.18×	2.10×
soc-Pokec	2.76×	1.20×	2.30×
geomean	2.5×	1.15×	2.17×

Performance of ParTBC w.r.t. exact parallel
Brandes' algorithm for error $\sim 6\%$
(VR: vertex renumbering)



Effect of the number of sources chosen
on error in top- k for DynRR

Takeaway:

The performance-accuracy trade-off can be controlled by carefully choosing the source nodes in top- k computation.

Conclusions

- ① Parallel graph processing is challenging due to *irregularity* in the data-access, control-flow, and communication patterns.
- ② We proposed
 - algorithm- and architecture-independent : Graprox
 - algorithm-independent but architecture-specific : Graffix
 - algorithm-specific but architecture-independent : ParTBCtechniques for approximate parallel graph processing.
- ③ The techniques provide *tunable knobs* to control the performance-accuracy trade-off in graph applications.
- ④ The techniques are generally applicable to a large class of parallel graph algorithms and input graphs of varying characteristics.
- ⑤ Approximate computing combined with parallelization promises to make heavy-weight graph computation practical, as well as scalable.

Publications

Conferences

- ① Somesh Singh and Rupesh Nasre. "Graffix: Efficient Graph Processing with a Tinge of GPU-Specific Approximations." In *Proceedings of the 49th International Conference on Parallel Processing (ICPP 2020)*. 23:1 – 23:11. <https://doi.org/10.1145/3404397.3404406>
- ② Somesh Singh and Rupesh Nasre. "Optimizing Graph Processing on GPUs Using Approximate Computing: Poster". In *Proceedings of the 24th Symposium on Principles and Practice of Parallel Programming (PPoPP 2019)*. 395 – 396. <https://doi.org/10.1145/3293883.3295736>

Journals

- ① Somesh Singh and Rupesh Nasre. "Scalable and Performant Graph Processing on GPUs Using Approximate Computing". *IEEE Transactions on Multi-Scale Computing Systems (TMSCS)* 4, 3 (2018), 190 – 203. <https://doi.org/10.1109/TMSCS.2018.2795543>

Publications

Conferences

- ① Somesh Singh and Rupesh Nasre. "Graffix: Efficient Graph Processing with a Tinge of GPU-Specific Approximations." In *Proceedings of the 49th International Conference on Parallel Processing (ICPP 2020)*. 23:1 – 23:11. <https://doi.org/10.1145/3404397.3404406>
- ② Somesh Singh and Rupesh Nasre. "Optimizing Graph Processing on GPUs Using Approximate Computing: Poster". In *Proceedings of the 24th Symposium on Principles and Practice of Parallel Programming (PPoPP 2019)*. 395 – 396. <https://doi.org/10.1145/3293883.3295736>

Journals

- ① Somesh Singh and Rupesh Nasre. "Scalable and Performant Graph Processing on GPUs Using Approximate Computing". *IEEE Transactions on Multi-Scale Computing Systems (TMSCS)* 4, 3 (2018), 190 – 203. <https://doi.org/10.1109/TMSCS.2018.2795543>

Thank You