

On the Design of Graph Analytical Software in Chapel

Oliver Alvarado Rodriguez, Zihui Du, David A. Bader

Department of Data Science

New Jersey Institute of Technology

Newark, NJ, USA



Supported by NSF grant: CCF-2109988



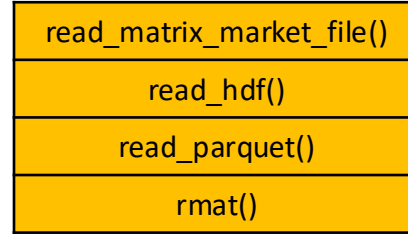
Introduction

- Graph analytical software consists of two main objectives: designing efficient **graph data structures** for fast data access and **algorithms** that exploit these efficient data accesses.
 - We have implemented an edge-based data structure based on a modified version of CSR we call the **Double-Index (DI)** data structure.
 - We have implemented algorithms for different graph analytical kernels such as **breadth-first search (BFS)**, triangle counting, connected components, etc.
 - All our functionality is bundled into the framework, Arachne, built on top of Arkouda.
- Firstly, this talk will present **DI** with a focus on new functionality to facilitate in-memory property graph analysis. Secondly, I will share our journey of optimizing **BFS** for distributed-memory execution in Chapel.

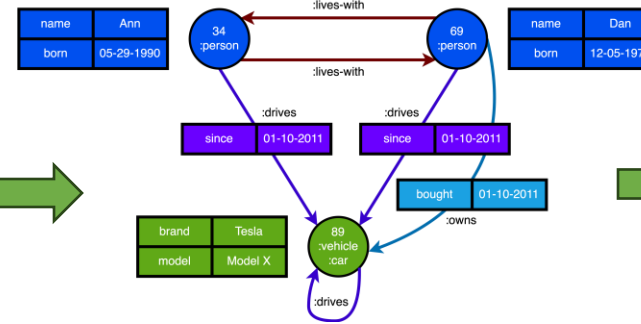
A Bird's Eye-View of Arachne+Arkouda

id	label	name	born	brand	model
34	person	Ann	1990	NULL	NULL
69					
	src id	dst id	relationship	since	bought
89					
	34	69	lives-with	NULL	NULL
	69	34	lives-with	NULL	NULL
	34	89	drives	2011	NULL
	69	89	drives	2011	NULL
	69	89	owns	NULL	2011
	89	89	drives	NULL	NULL

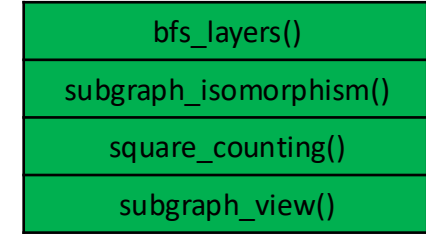
Load in CSVs, HDF5s, Parquets, etc.



Generate or load graphs in from various sources.

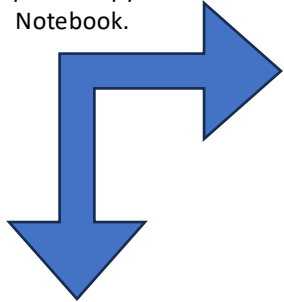


Convert tabular data to a property graph with all data closely maintained with vertex and edges.



Perform analysis or filter for NetworkX, iGraph, or graph-tool.

User edits a Python script or a Jupyter Notebook.

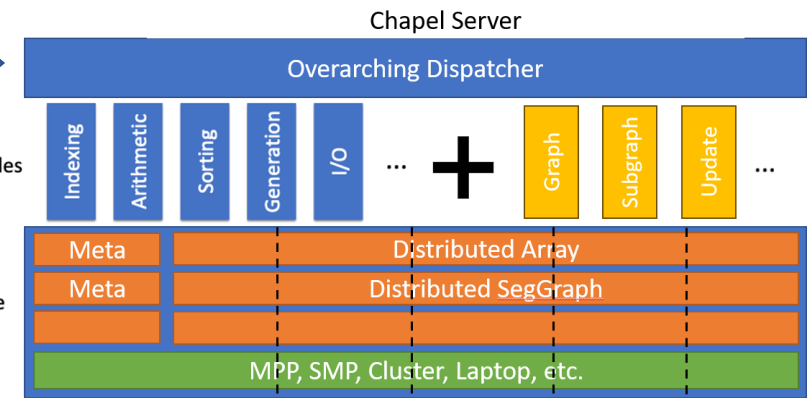


User

```

1. import arkouda as ak
2. import arachne as ar
3.
4. ## Get src and dst from input file.
5.
6. graph = ar.PropGraph()
7.
8. ## Generate label_df and relationships_df from input file.
9.
10. graph.load_edge_attributes(relationships_df)
11. graph.load_node_attributes(label_df)
12.
13. ## User generates labels_to_find and relationships_to_find.
14. returned_nodes = graph.node_attributes["column"] == 1
15. returned_edges = graph.edge_attributes["column"] == 2
16.
17. subgraph_src = ak.in1d(returned_edges[0], returned_nodes)
18. subgraph_dst = ak.in1d(returned_edges[1], returned_nodes)
19.
20. kept_edges = subgraph_src & subgraph_dst
21.
22. subgraph_src = subgraph_src[kept_edges]
23. subgraph_dst = subgraph_dst[kept_edges]
24.
25. subgraph = ar.Graph()
26. subgraph.add_edges_from(subgraph_src, subgraph_dst)
27. ## Run some other operations on subgraph! Reference our HPEC22 paper ©
    
```

Easily usable through NetworkX-like API.

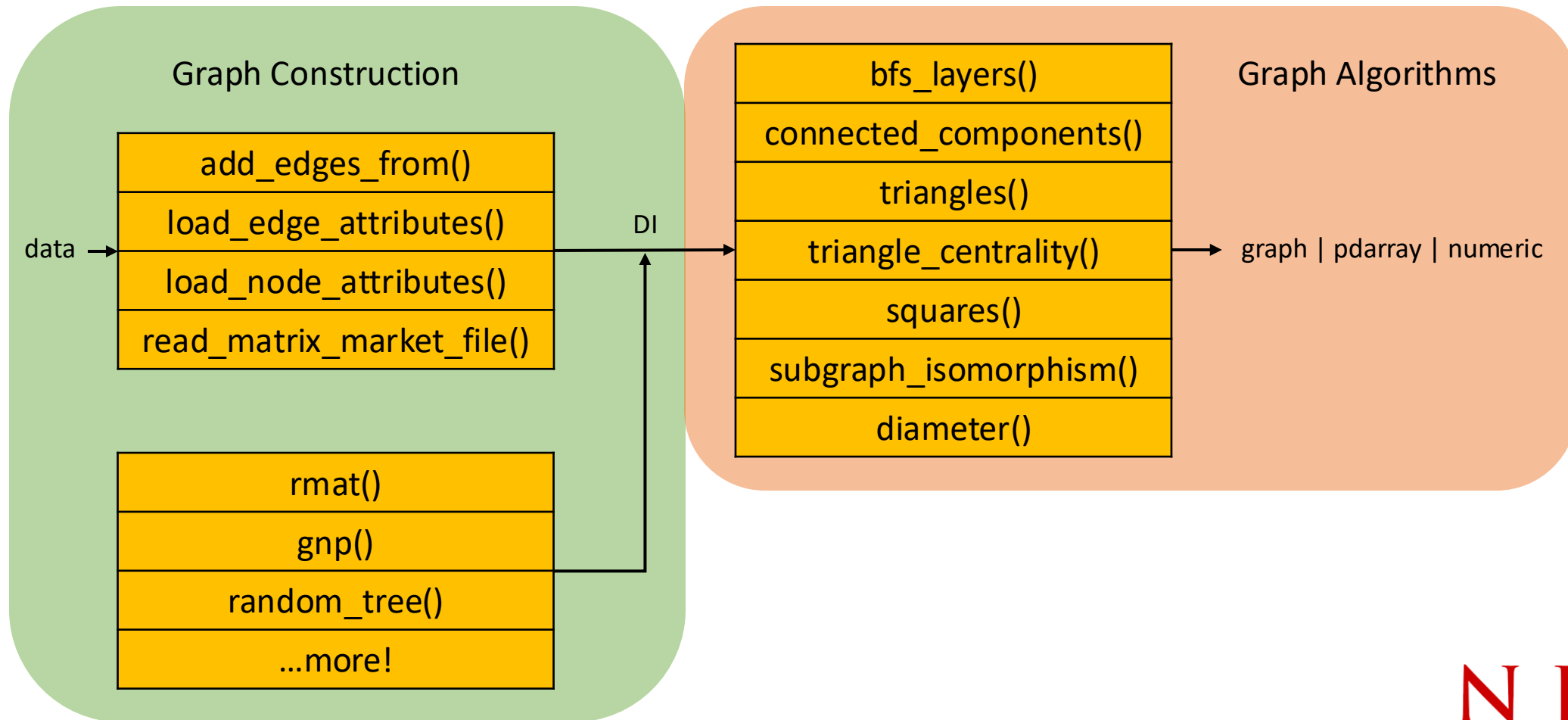


Original image source: <https://chapel-lang.org/CHIUIW/2020/Reus.pdf> was modified for this presentation

Runs on any hardware, data stays in the back-end, user calls API through Python: powerful and productive. Server can run on supercomputers; Python API usable locally.

OPEN SOURCE: <https://github.com/Bears-R-Us/arkouda-njit>
PUBLICATIONS & PRESENTATIONS AT: HPEC, HiPC, IPDPS, & PPoPP

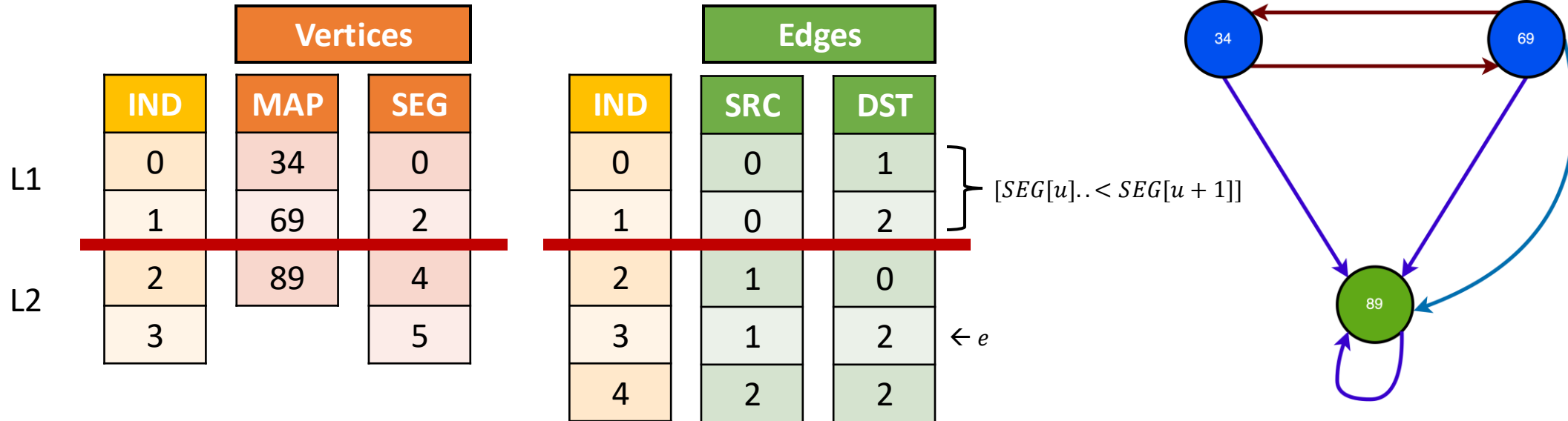
Modular View of Arachne Functionality



Double-Index (DI) Data Structure

Examples and Persistence

(Property) Graph Data Structure



- Allows for simple, compact, **distributable** storage of vertex and edge sets.
- Given an edge index, e , all vertices that make up that edge are found in **constant time**, avoiding a binary search into SRC (CSR offsets index equivalent).
- MAP allows explicitly storing original vertex labels, returning original graph involves index operations $SRC[MAP]$ and $DST[MAP]$.

(Property) Graph Data Structure

	Vertices						Edges						
	IND	MAP	SEG	INT	UINT	REAL	IND	SRC	DST	INT	UINT	BOOL	INT
L1	0	34	0	-1	1	1.1	0	0	1	-1	1	T	
	1	69	2	2	0	0.2	1	0	2	2	2	F	1
L2	2	89	4	5	4	4.1	2	1	0	0	0	T	
	3		5				3	1	2	-2	2	T	
							4	2	2	2	2	F	2

sparse
↑

- Same **distributable** storage of vertex and edge attributes as base DI.
- Given an edge or vertex index, all attribute data can be easily accessed.
- Same storage principles apply to strings, which are stored in an object containing a byte array for characters and segments for where each string starts in the byte array.
- Sparse attribute arrays maintaining locality can also be created to only store attribute values that belong to a subset of indices.

Persisting Graphs via Arkouda Symbol Table

- Graph is stored as a `GraphSymEntry` which is a wrapper to `SegGraph` that inherits from `CompositeSymEntry`.
- Sparse arrays are stored in a `SparseSymEntry` (shoutout to Vass from the Chapel team) that inherits from `GenSymEntry`.
- We have other special classes to persist data such as maps, replicated arrays, and associative arrays. Plans to store “sparse” Arkouda categoricals and strings.

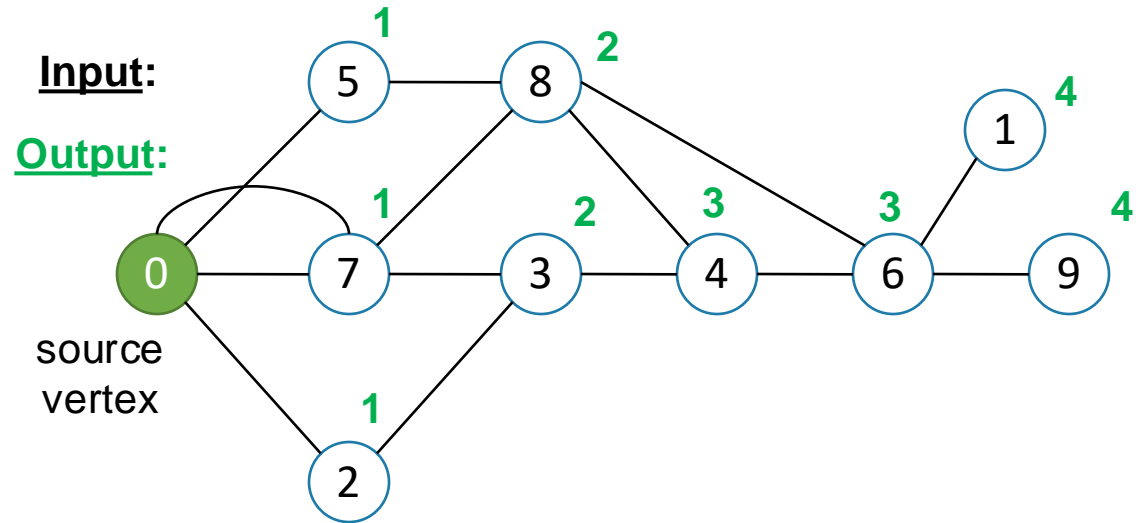
Breadth-First Search (BFS)

A Journey of Optimizations

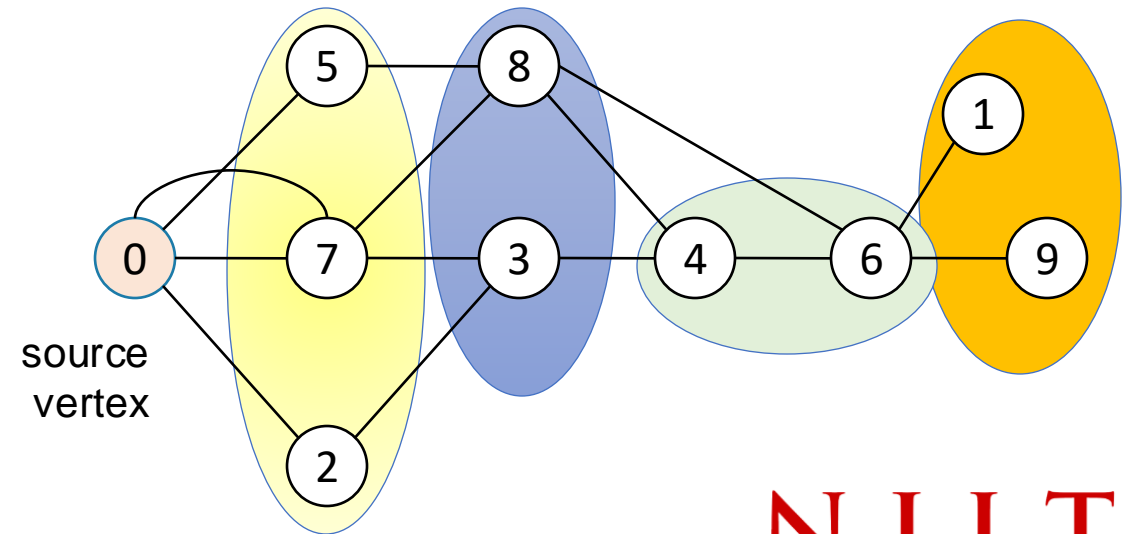
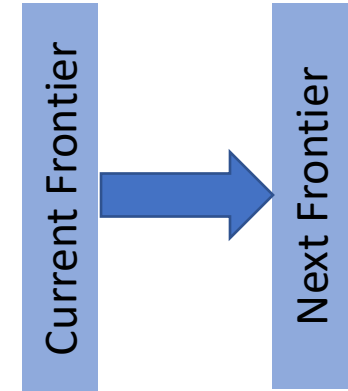
General Information

- Important algorithm for solving problems that requires a complete traversal of a graph: answer questions like “how far is every other vertex from our source?”
- One of the fundamental graph algorithms in computer science.
- Has a sequential complexity of $O(n + m)$ where n is the number of vertices and m is the number of edges.

Single Locale Parallel BFS (version 1.0)

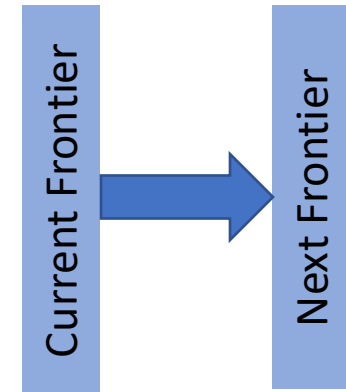
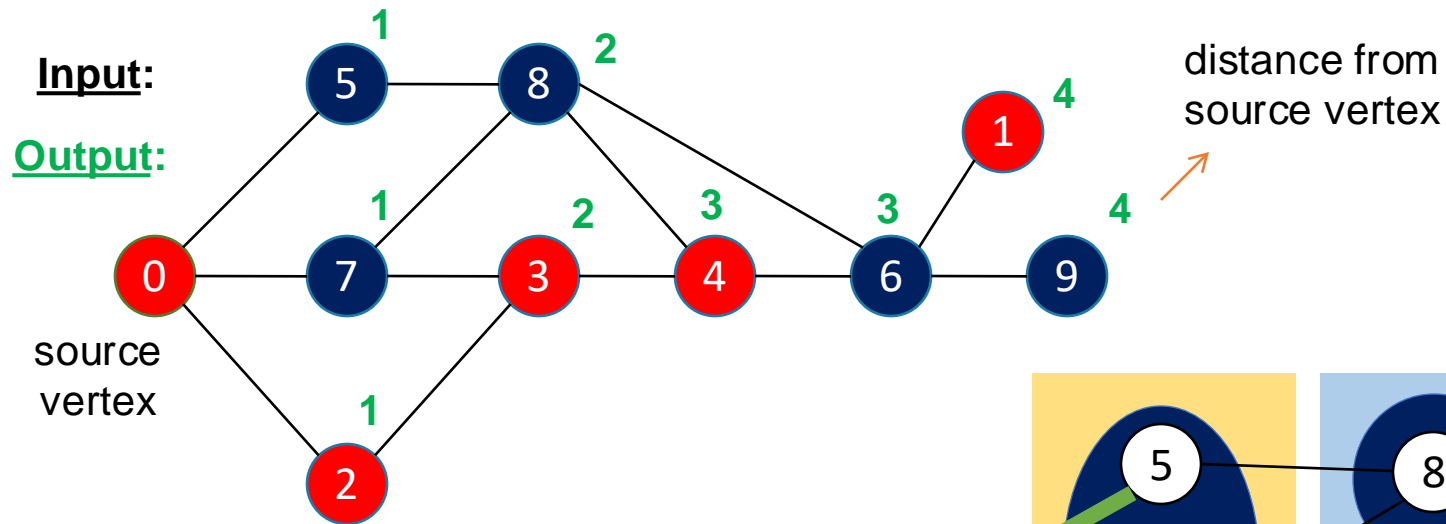


Output: $D = [0, 4, 1, 2, 3, 1, 3, 1, 2, 4]$



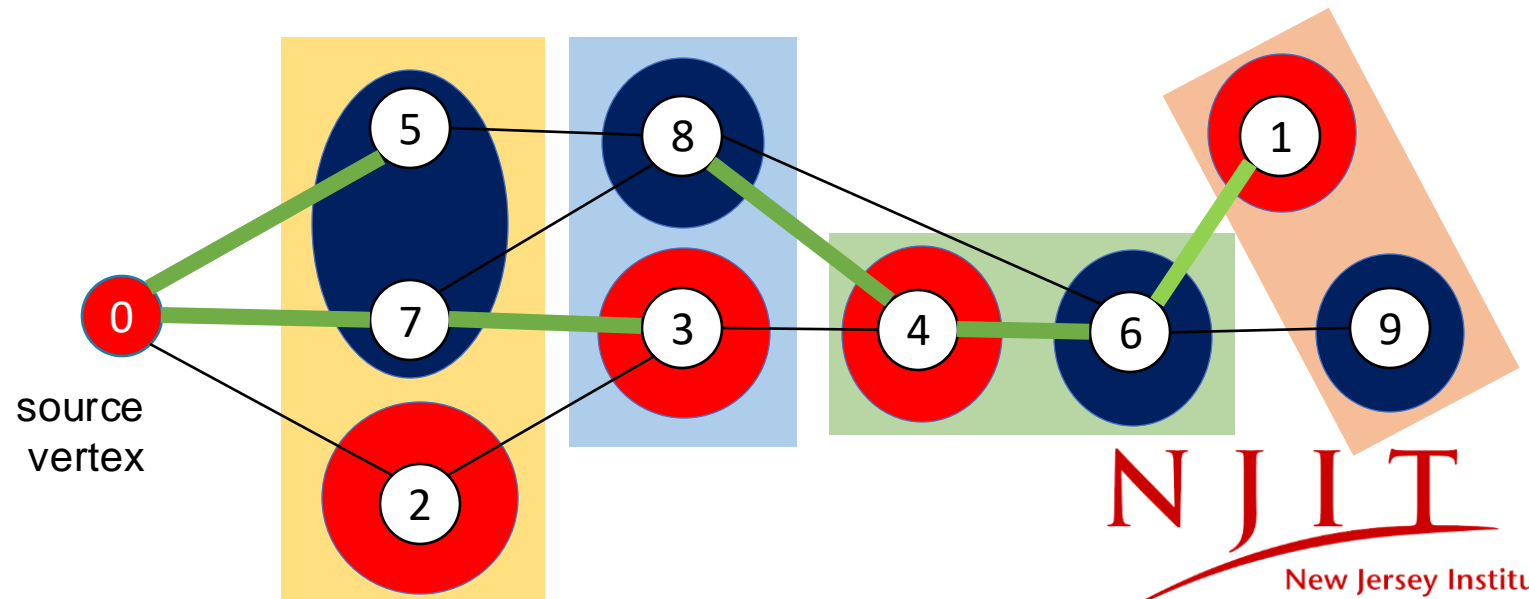
Multilocale Parallel BFS (version 1.5)

Assume our edge list is split down the middle, then the neighborhood of some vertices will live on one compute node while the rest live on another compute node.

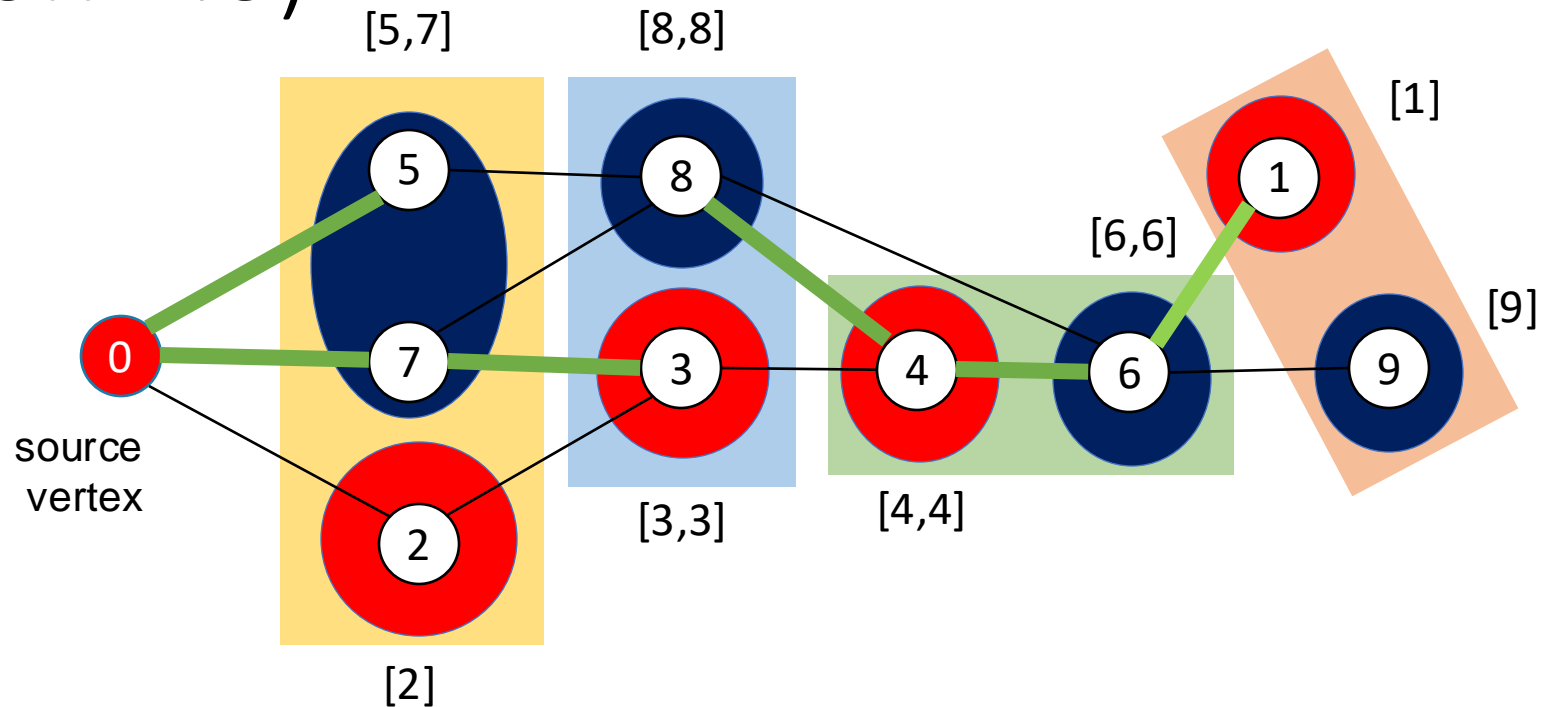


Output: $D = [0, 4, 1, 2, 3, 1, 3, 1, 2, 4]$

Any cross-color expansions are writes across the network; fine-grained writes hold up execution, large coarse-grained writes are better.



Multilocale Parallel BFS with Aggregators (version 2.0)



Each frontier is a list. Before we expand the frontiers in the following iteration, we aggregate them, and then write them to the appropriate frontier list.

Multilocale Parallel BFS Version 1.0

- Uses ideas of forward and reversed edges for undirected graphs. For example, $u-v$ is stored in SRC and DST and $v-u$ is stored in SRCr and DSTr.
- Use the “old” distributed bag to expand frontiers.

Expand frontier based of forward-edges

Expand frontier based of reversed-edges

07 June 2024

```
491 while (numCurF>0) {
492   coforall loc in Locales with (ref SetNextF,+ reduce topdown, + reduce bottomup, ref root, ref src, ref depth) {
493     on loc {
494       ref srcf=src;
495       ref df=dst;
496       ref nf=nei;
497       ref sf=start_i;
498
499       ref srcfR=srcR;
500       ref dfR=dstR;
501       ref nfr=neiR;
502       ref sfr=start_iR;
503
504       var edgeBegin=src.localSubdomain().lowBound;
505       var edgeEnd=src.localSubdomain().highBound;
506       var vertexBegin=src[edgeBegin];
507       var vertexEnd=src[edgeEnd];
508       var vertexBeginR=srcR[edgeBegin];
509       var vertexEndR=srcR[edgeEnd];
510
511       var switchratio=(numCurF:real)/nf.size:real;
512       if (switchratio<GivenRatio) { //top down
513         topdown:=1;
514         forall i in SetCurF with (ref SetNextF) {
515           if ((xlocal(i,vertexBegin,vertexEnd)) || (LF==0)) { // current edge has the vertex
516             var numNF=nf[i];
517             var edgeId=sf[i];
518             var nextStart=max(edgeId,edgeBegin);
519             var nextEnd=min(edgeEnd,edgeId+numNF-1);
520             ref NF=df[nextStart..nextEnd];
521             forall j in NF with (ref SetNextF){
522               if (depth[j]==-1) {
523                 depth[j]=cur_level+1;
524                 SetNextF.add(j);
525               }
526             }
527           }
528           if ((xlocal(i,vertexBeginR,vertexEndR)) || (LF==0)) {
529             var numNF=nfR[i];
530             var edgeId=sfR[i];
531             var nextStart=max(edgeId,edgeBegin);
532             var nextEnd=min(edgeEnd,edgeId+numNF-1);
533             ref NF=dfR[nextStart..nextEnd];
534             forall j in NF with (ref SetNextF) {
535               if (depth[j]==-1) {
536                 depth[j]=cur_level+1;
537                 SetNextF.add(j);
538               }
539             }
540           }
541         }
542       }
543     }
544   }
545 }
```

Multilocale Parallel BFS Version 1.5

- Combines the forward and reversed arrays to ensure every vertex has full access to its neighbors instead of a split view.

Expand frontier based of symmetrized edges ←

```
491 while (numCurF>0) {
492   coforall loc in Locales with (ref SetNextF,+ reduce topdown,+ reduce bottomup, ref root, ref src, ref depth) {
493     on loc {
494       ref srcf=src;
495       ref df=dst;
496       ref nf=nei;
497       ref sf=start_i;
498
499       ref srcfR=srcR;
500       ref dfR=dstR;
501       ref nfr=neiR;
502       ref sfr=start_iR;
503
504       var edgeBegin=src.localSubdomain().lowBound;
505       var edgeEnd=src.localSubdomain().highBound;
506       var vertexBegin=src[edgeBegin];
507       var vertexEnd=src[edgeEnd];
508       var vertexBeginR=srcR[edgeBegin];
509       var vertexEndR=srcR[edgeEnd];
510
511       var switchratio=(numCurF:real)/nf.size:real;
512       if (switchratio<GivenRatio) { //top down
513         topdown=-1;
514         forall i in SetCurF with (ref SetNextF) {
515           if ((xlocal(i,vertexBegin,vertexEnd)) || (LF==0)) { // current edge has the vertex
516             var numNF=nf[i];
517             var edgeId=sf[i];
518             var nextStart=max(edgeId,edgeBegin);
519             var nextEnd=min(edgeEnd,edgeId+numNF-1);
520             ref NF=df[nextStart..nextEnd];
521             forall j in NF with (ref SetNextF){
522               if (depth[j]==-1) {
523                 depth[j]=cur_level+1;
524                 SetNextF.add(j);
525               }
526             }
527           }
528         }
529       }
530     }
531   }
532 }
```

Multilocale Parallel BFS Version 2.0

```
105     while true {
106         var pending_work:bool;
107         coforall loc in Locales with(|| reduce pending_work, ref depth, ref frontier_sets) {
108             on loc {
109                 var src_low = src.localSubdomain().low;
110                 var src_high = src.localSubdomain().high;
111                 forall u in frontier_sets[frontier_sets_idx] with (|| reduce pending_work,
112                     var adj_list_start = seg[u];
113                     var num_neighbors = seg[u+1] - adj_list_start;
114                     if (num_neighbors != 0) {
115                         var adj_list_end = adj_list_start + num_neighbors - 1;
116
117                         // Only pull the part of the adjacency list that is local.
118                         var actual_start = max(adj_list_start, src_low);
119                         var actual_end = min(src_high, adj_list_end);
120
121                         ref neighborhood = dst.localSlice(actual_start..actual_end);
122                         for v in neighborhood {
123                             if (depth[v] == -1) {
124                                 pending_work = true;
125                                 depth[v] = cur_level + 1;
126                                 var locs = find_locs(v, ranges);
127                                 for lc in locs do frontier_agg.copy(lc.id, v);
128                             }
129                         }
130                     }
131                 } //end forall
132                 frontier_sets[frontier_sets_idx].clear();
133             } // end on loc
134         } // end coforall loc
135         if !pending_work {
136             break;
137         }
138         cur_level += 1;
139         frontier_sets_idx = (frontier_sets_idx + 1) % 2;
140     } // end while
141     return "success";
```

Locale parallelism

Parallelism per-locale

Maintaining locality

Neighborhood expansion with aggregation

Aggregator per-task

Multilocale BFS Communication Volume Heatmap

delaunayn20	get			put		
locale	1.0	1.5	2.0	1.0	1.5	2.0
0	15672640	7873842	639827	5629422	2749193	138070
1	15834332	7939017	687156	1952226	1016946	127936
2	15715554	7722659	226754	1942839	962031	45217
3	15817879	7723971	226880	1951313	962201	45060
4	15964559	7724880	226691	1961552	962199	51217
5	15739226	7726504	230024	1940688	962439	52714
6	15569450	7727678	229096	1925536	962680	51977
7	15341933	7736094	225083	1904757	963418	48413

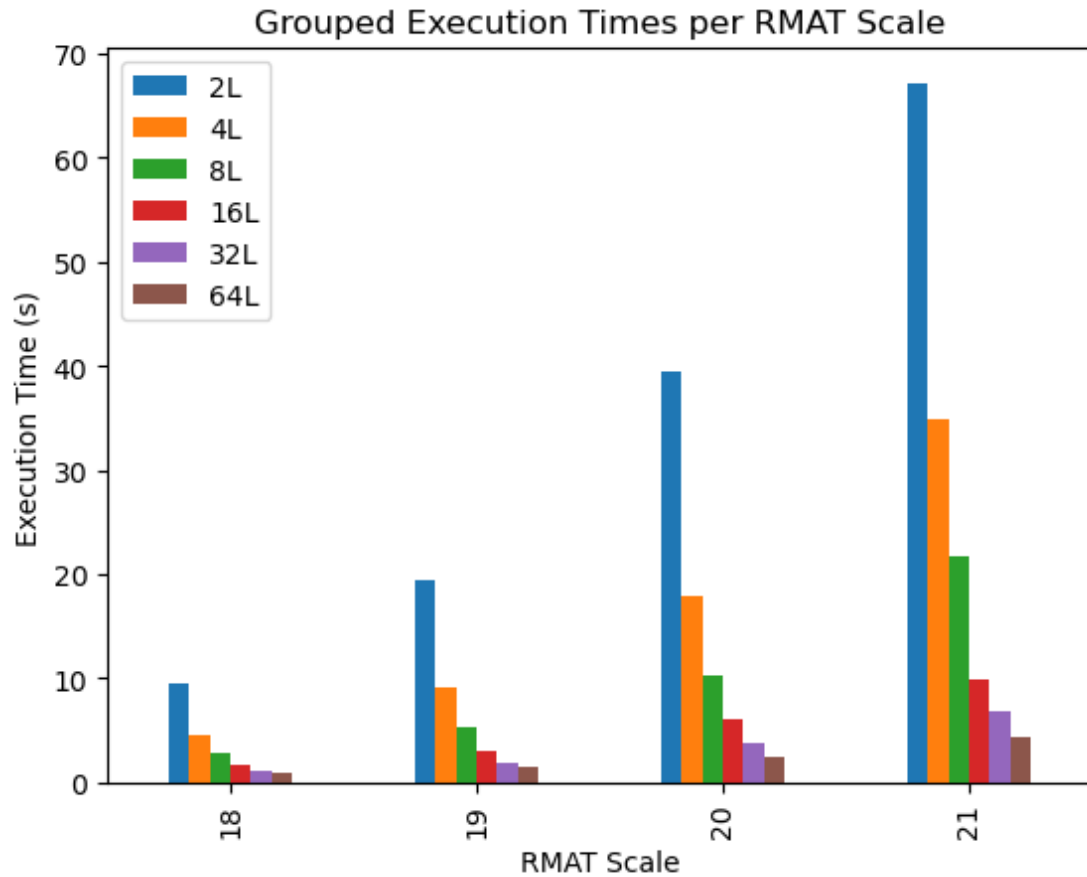
1.0: 84 seconds (HPEC 21')

2.0: 3.36 seconds

delaunayn20 is a graph with 3 million edges and a large diameter

Takeaway: Aggregating writes drastically reduces communication volumes, improving performance, all done easily through Chapel by adapting aggregators for different uses.

2.0 BFS Scalability



Speed-Up Over 2 Locales

	4L	8L	16L	32L	64L
18	2.11	3.43	5.87	8.10	9.66
19	2.14	3.69	6.35	10.28	13.04
20	2.20	3.84	6.41	10.60	15.90
21	1.93	3.09	6.84	9.86	15.56

Takeaway: As the number of locales increased, we see a good speed-up for distributed-memory breadth-first search.

Lessons Learned

- Using Chapel (or any PGAS-based languages and frameworks) don't magically get rid of the complications of parallelizing and distributing graph operations.
- Adapting communication-aware optimizations, such as being aware of how neighborhoods are split across locales, can help improve graph-based performances.

Conclusion

- Using a programming language like Chapel allows us to quickly implement both shared-memory and distributed-memory algorithms to enable highly productive large-scale graph analysis.
- Using an existing framework like Arkouda allows us to focus more on graph algorithms while offloading tasks such as object persistence and array sorting.

Future Work

- Not everything needs to be distributed – large queries can be done in a distributed manner and smaller graphs analyzed on one compute node; can we **hybridize** our graph tools?
- Performance, performance, performance. Array-based operations are wonderful in Chapel, but do we need to build harnesses in Arachne to **call out to external programs** written in MPI, YGM, or other massively distributed tools?
- How can we **dynamically optimize** during runtime? For example, code regions that perform a lot of reads or writes on GASNet+Infiniband suffer when multiple parallel threads are writing since those values are transmitted sequentially. Chapel currently doesn't allow for forall loops to dynamically use a runtime-given thread count.
- **There isn't one data structure to rule them all.** Add capabilities in Arachne to build at runtime the data structure that is best for a given problem.

Thank You 😊
Questions?