

An Experimental Comparison of Pregel-like Graph Processing Systems

Minyang Han, Khuzaima Daudjee, Khaled Ammar,
M. Tamer Özsu, Xingfang Wang, Tianqi Jin

University of Waterloo

VLDB 2014

- 1 Motivation
- 2 Systems Tested
- 3 Methodology
- 4 Experimental Results

- 1 Motivation
- 2 Systems Tested
- 3 Methodology
- 4 Experimental Results

Many new graph processing systems...
...but existing studies lack scale and comprehensiveness.

Our study provides:

- Experiments with **scale**: up to 128 EC2 machines.
- **Comprehensive** combination of algorithms and datasets.
- Focus on **time**, **memory**, and **network**.
- Use of similar systems for an **apples-to-apples comparison**.

Many new graph processing systems...
...but existing studies lack scale and comprehensiveness.

Our study provides:

- Experiments with **scale**: up to 128 EC2 machines.
- **Comprehensive** combination of algorithms and datasets.
- Focus on **time**, **memory**, and **network**.
- Use of similar systems for an **apples-to-apples comparison**.

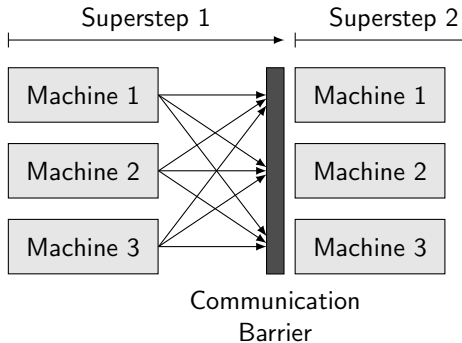
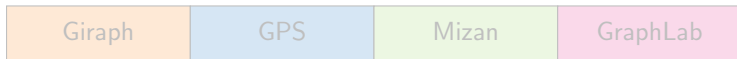
Many new graph processing systems...
...but existing studies lack scale and comprehensiveness.

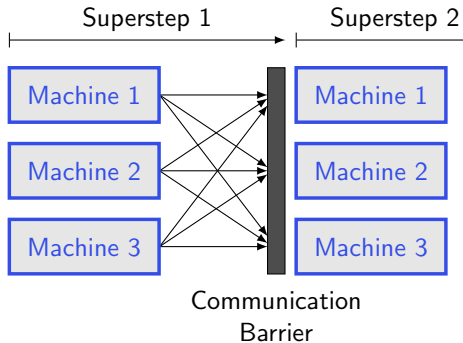
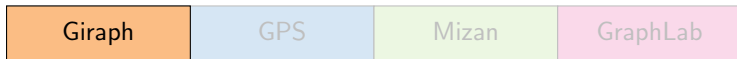
Our study provides:

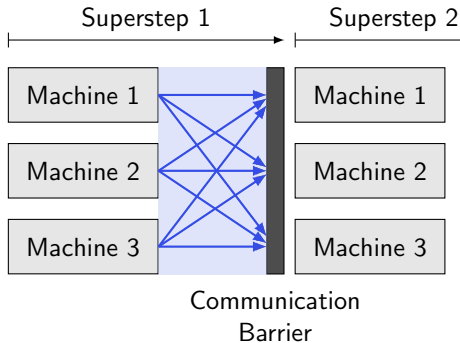
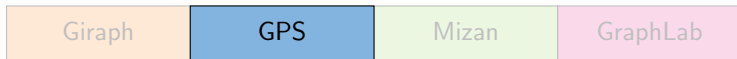
- Experiments with **scale**: up to 128 EC2 machines.
- **Comprehensive** combination of algorithms and datasets.
- Focus on **time**, **memory**, and **network**.
- Use of similar systems for an **apples-to-apples comparison**.

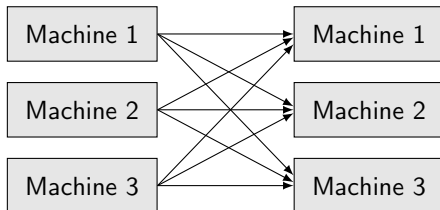
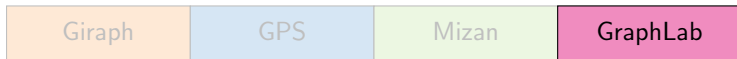
- 1 Motivation
- 2 Systems Tested**
- 3 Methodology
- 4 Experimental Results

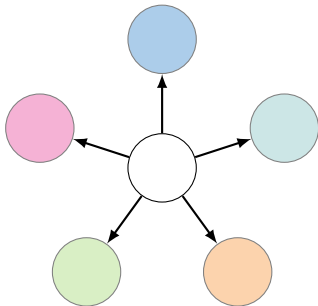


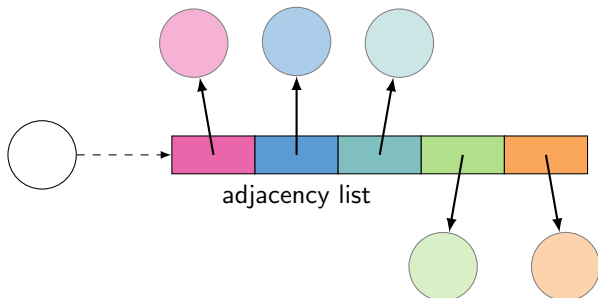


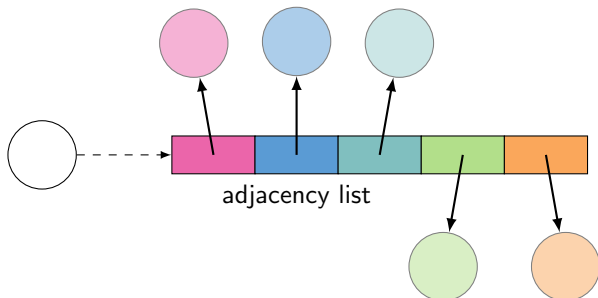






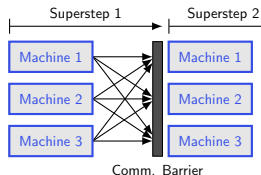






Can be implemented as:

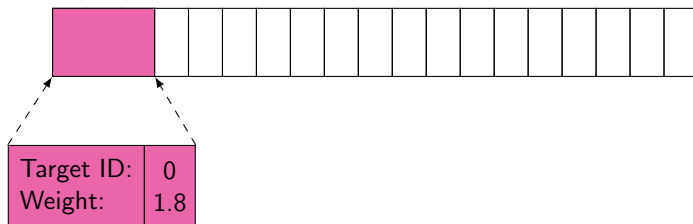
- Byte array (default).
- Hash map.



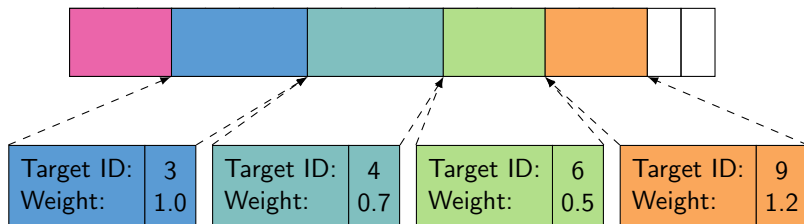
Byte array adjacency list:

--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

Byte array adjacency list:



Byte array adjacency list:



Byte array adjacency list:



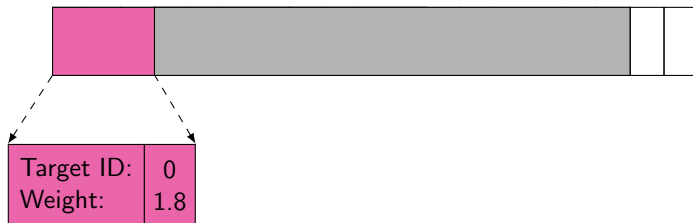
Byte array adjacency list:



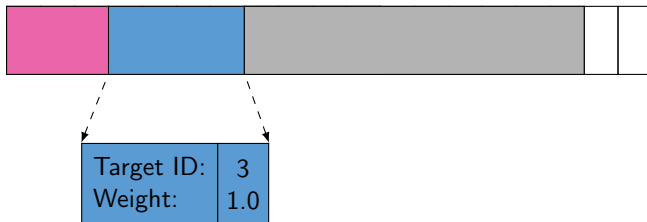
Byte array adjacency list:



Byte array adjacency list:



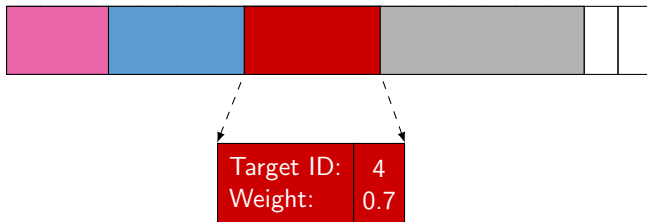
Byte array adjacency list:



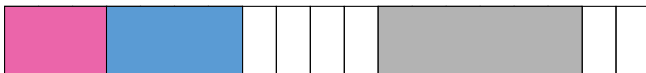
Byte array adjacency list:



Byte array adjacency list:



Byte array adjacency list:



Byte array adjacency list:



Hash map adjacency list:

Target ID

Weight

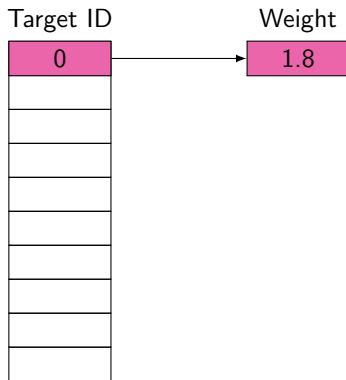
Hash map adjacency list:

Target ID:	0
Weight:	1.8

Target ID

Weight

Hash map adjacency list:



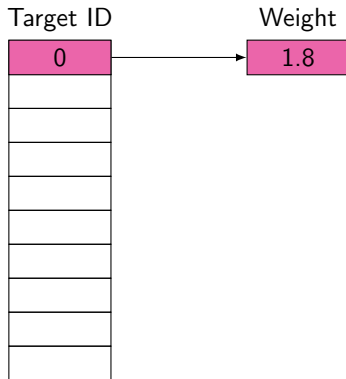
Hash map adjacency list:

Target ID:	3
Weight:	1.0

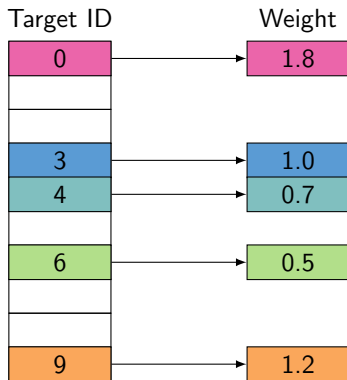
Target ID:	4
Weight:	0.7

Target ID:	6
Weight:	0.5

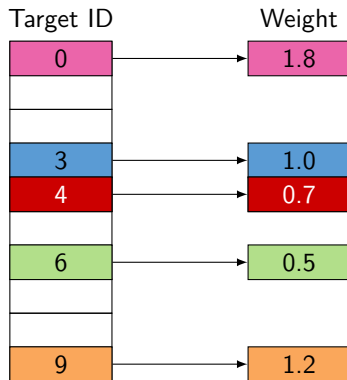
Target ID:	9
Weight:	1.2



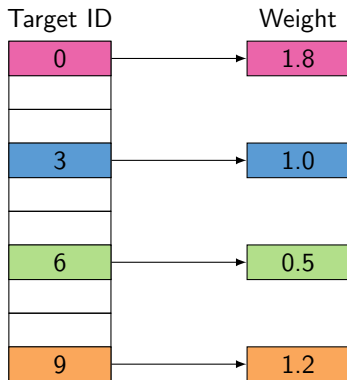
Hash map adjacency list:



Hash map adjacency list:



Hash map adjacency list:



Byte array:

- Space efficient. ✓
- Overheads for mutations. ✗

Hash map:

- Less space efficient. ✗
- Efficient for mutations. ✓

Byte array:

- Space efficient. ✓
- Overheads for mutations. ✗

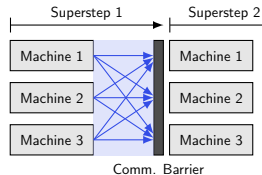
Hash map:

- Less space efficient. ✗
- Efficient for mutations. ✓

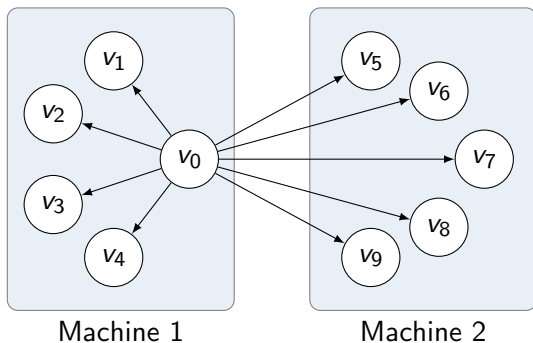
What's the tradeoff like?

Two optional optimizations:

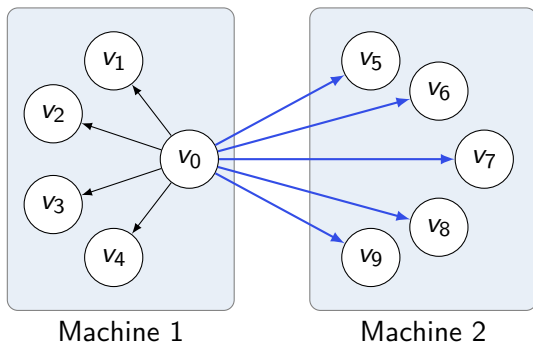
- LALP.
- Dynamic migration.



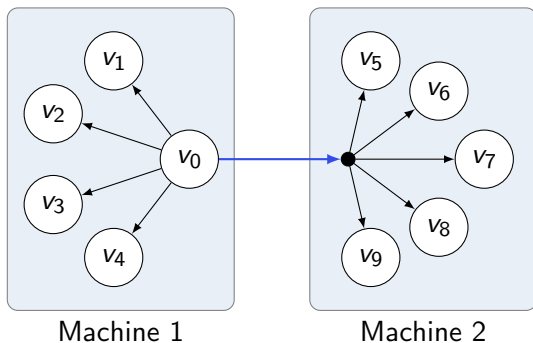
LALP (Large Adjacency List Partitioning):



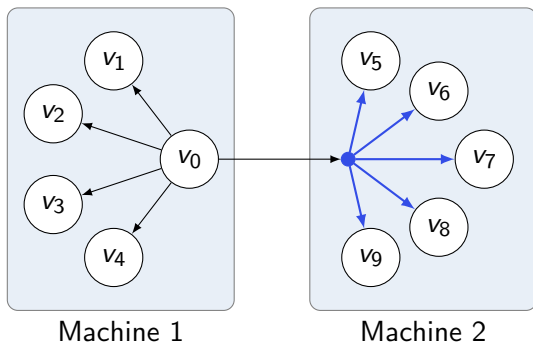
LALP (Large Adjacency List Partitioning):



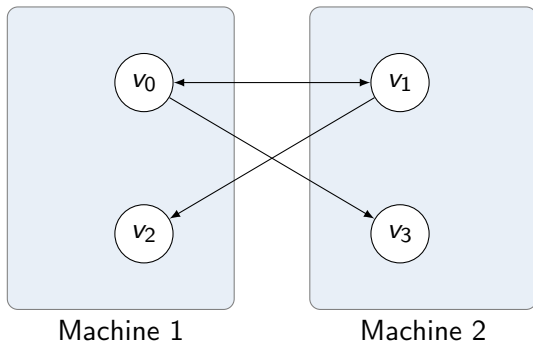
LALP (Large Adjacency List Partitioning):



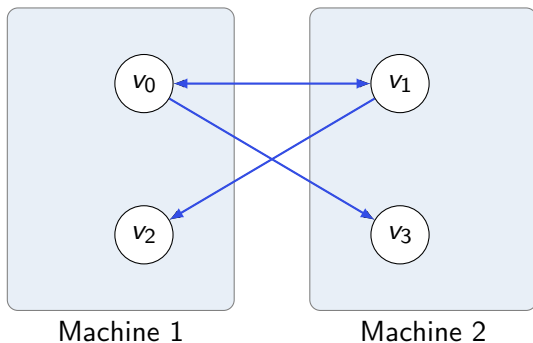
LALP (Large Adjacency List Partitioning):



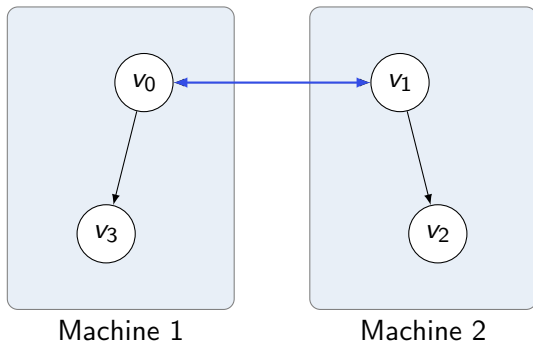
Dynamic migration:



Dynamic migration:



Dynamic migration:



LALP:

- Reduces network usage. ✓
- Sent messages must be same. ✗

Dynamic migration:

- Reduces network usage. ✓
- Incompatible with DMST. ✗

LALP:

- Reduces network usage. ✓
- Sent messages must be same. ✗

Dynamic migration:

- Reduces network usage. ✓
- Incompatible with DMST. ✗

Do these improve performance?

Mizan:

- Lacks built-in system optimizations. ✗
- But competitive against older Giraph 0.1. ✓

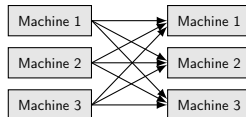
Mizan:

- Lacks built-in system optimizations. ✗
- But competitive against older Giraph 0.1. ✓

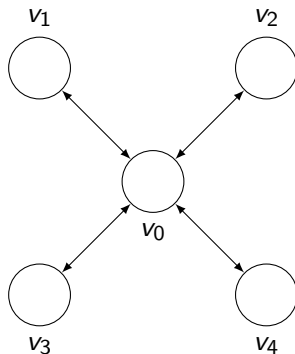
How does it compare now?

GraphLab features *asynchronous* execution:

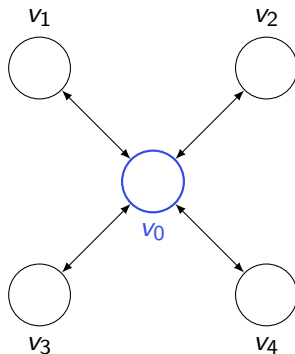
- No communication barriers. ✓
- Uses the *most recent* vertex values. ✓



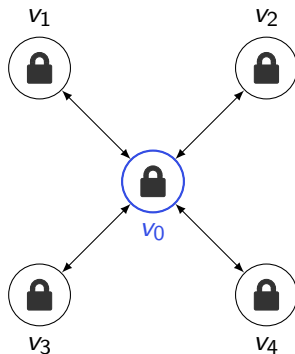
Implemented via distributed locking:



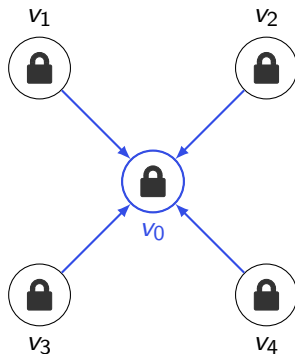
Implemented via distributed locking:



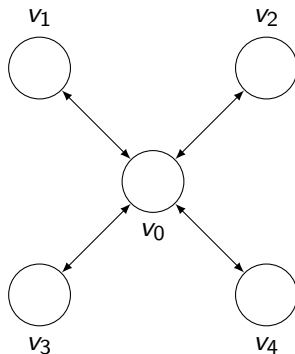
Implemented via distributed locking:



Implemented via distributed locking:



Implemented via distributed locking:



Asynchronous execution:

- No communication barriers. ✓
- Uses the *most recent* vertex values. ✓
- Potential overheads? ✗

Asynchronous execution:

- No communication barriers. ✓
- Uses the *most recent* vertex values. ✓
- Potential overheads? ✗

Performance of sync vs. async?

- 1 Motivation
- 2 Systems Tested
- 3 Methodology**
- 4 Experimental Results

Tackling Graph Processing Problems

Google recursion

Web Images Videos Maps Books More Search tools

About 2,040,000 results (0.28 seconds)

Did you mean: [recursion](#)

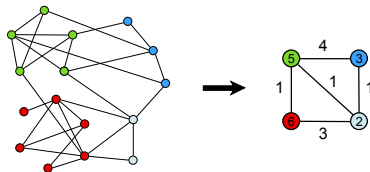
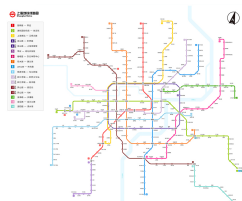
Recursion - Wikipedia, the free encyclopedia
[en.wikipedia.org/wiki/Recursion](#)

Recursion is the process of repeating items in a self-similar way. For instance, when the surfaces of two mirrors are exactly parallel with each other the nested ...

Recursion (computer science)
Recursion in computer science is a method where the solution to a ...

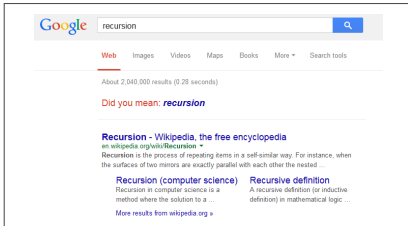
Recursive definition
A recursive definition (or inductive definition) in mathematical logic ...

[More results from wikipedia.org »](#)



Adapted from: Fortunato, arXiv:0906.0612 (2010)

Random Walk



Google recursion

Web Images Videos Maps Books More Search tools

About 2,040,000 results (0.28 seconds)

Did you mean: [recursion](#)

Recursion - Wikipedia, the free encyclopedia
[en.wikipedia.org/wiki/Recursion](#)

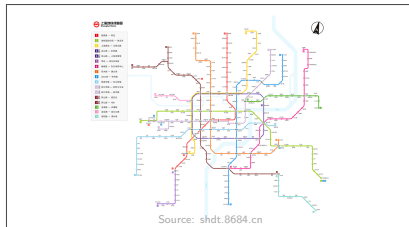
Recursion is the process of repeating items in a self-similar way. For instance, when the surfaces of two mirrors are exactly parallel with each other the nested ...

Recursion (computer science)
Recursion in computer science is a method where the solution to a ...

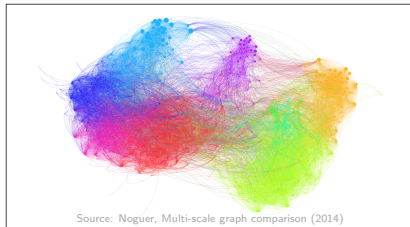
Recursive definition
A recursive definition (or inductive definition) in mathematical logic ...

[More results from wikipedia.org »](#)

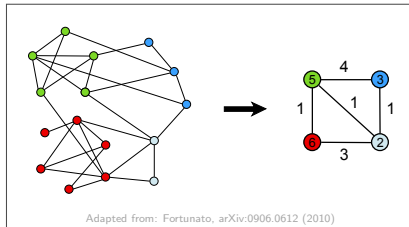
Sequential Traversal



Parallel Traversal



Graph Mutation



Category	Algorithm
Random walk	PageRank
Sequential traversal	SSSP (Single Source Shortest Path)
Parallel traversal	WCC (Weakly Connected Components)
Graph mutation	DMST (Distributed Minimum Spanning Tree)



LJ

orkut

OR

العربية

AR



TW



UK



LJ

orkut

OR

العربية

AR



TW



UK

$|E| = 1.46\text{B}$

$|E| = 3.73\text{B}$



LJ

orkut

OR

العربية

AR



TW



UK

$|E| = 1.46\text{B}$

$|E| = 3.73\text{B}$

Datasets are:

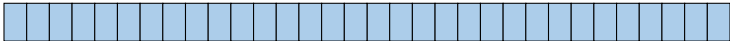
- Stored on HDFS as text files.
- Loaded via random hash partitioning.

- **Time:** *total* time = *setup* time + *computation* time
- **Memory:** *maximum* memory usage (across all machines)
- **Network:** *total* network usage (summed over all machines)

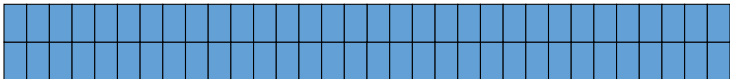
EC2 Instances



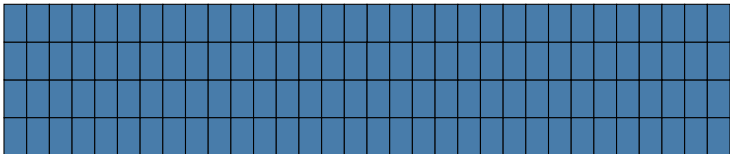
EC2 Instances



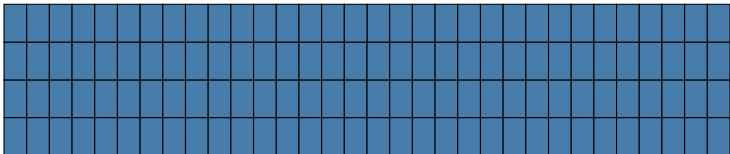
EC2 Instances



EC2 Instances

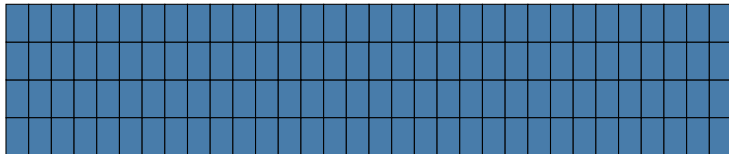


EC2 Instances



- m1.xlarge: 4 vCPUs, 15GB memory, Ubuntu 12.04

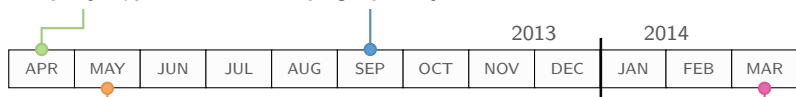
EC2 Instances



- m1.xlarge: 4 vCPUs, 15GB memory, Ubuntu 12.04

Mizan 0.1bu1

GPS rev 110



Giraph 1.0.0

GraphLab 2.2

- 1 Motivation
- 2 Systems Tested
- 3 Methodology
- 4 Experimental Results**

For memory, we define per-machine memory usage as the difference between the maximum and minimum memory used by a single machine during an experiment. This excludes the background usage of Hadoop and OS processes, which is typically between 200 to 300MB. We focus on *maximum memory usage*, the maximum per-machine usage across all worker machines. This gives the minimum memory resources all machines need to run an experiment without failure.

Similarly, per-machine network usage is the total number of bytes sent or received by a single machine for an experiment. We focus on *total network usage*, the sum of per-machine usage across all worker machines, as it best enables a high-level comparison of each system's network I/O. Additionally, we distinguish between total outgoing (sent) and total incoming (received) network usage.

Lastly, using different datasets and number of machines enables us to investigate the *scalability* of each system: how performance scales with the same graph on more machines or with larger graphs on the same number of machines.

To track these metrics, we use the total and setup times reported by all systems. For network usage, some systems have built-in message counters but they all differ. Message count is also a poor measure of network traffic, as messages can be of different sizes (such as in DMST) and header sizes are ignored. Further, the notion of messages is ill-defined for GraphLab's asynchronous mode. For these reasons, we use `/proc/net/dev` to track the total bytes sent and received at each machine, before and after every experiment.

For more fine-grained statistics, and to compute memory usage, we rely on one-second interval reports from `sar` and `free`. These are started locally on both the master and worker machines to record CPU, network, and memory utilization in a decentralized manner. They are started before and killed after each experiment to ensure a small but constant overhead across all experiments. Finally, Giraph 1.0.0 has a bug that prevents proper clean up of Java processes after a successful run. We solve this by killing the offending processes after each run.

For each experiment, we perform five runs and report both the mean and 95% confidence intervals for computation time, setup time, maximum memory usage, and total incoming network usage. We show only incoming network I/O as it is similar to total outgoing network I/O and exhibits identical patterns.

6. EXPERIMENTAL RESULTS

In this section, we present and analyze our experimental results. Due to space constraints and in the interests of readability, we present only a subset of all our experimental data. In particular, we do not show data for the master, since its memory and network usage are smaller and less revealing of system performance. All data and plots, along with our code and EC2 images, are available online.⁴

Table 4: Performance for AR, TV, UK.

	Computation	Setup	Total Time	Memory	Network
No Mutations					
Gi-BA	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆
Gi-HM	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆
GPS	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆
Mizan	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆
GL-S	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆
GL-A	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆
With Mutations (DMST)					
Gi-BA	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆
Gi-HM	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆
GPS	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆	☆☆☆☆

6.1 Summary of Results

A relative ranking of system performance, for computation, setup, and total time, maximum memory usage, and total network I/O, are presented for non-mutation and mutation (DMST) algorithms in Table 4. System names are abbreviated as Gi-BA and Gi-HM for Giraph byte array and hash map respectively; GPS for GPS with no optional optimizations; and GL-S and GL-A for GraphLab's synchronous and asynchronous modes. We exclude GPS's other modes as they provide little performance benefits. We focus on the larger graphs (AR, TV, and UK) and summarize our findings for LJ and UR in Section 6.6.

For each performance attribute, we provide a 4 star relative ranking, with 1 for poor performance, 2 for suboptimal, 3 for good, and 4 for excellent. The ranking of a system is roughly an average of its performance across algorithms.

Overall, for the non-mutation algorithms tested, Giraph and GraphLab are close: we recommend Giraph byte array over GraphLab for very large graphs on a limited number of machines, but GraphLab's synchronous mode over Giraph otherwise. In each case, the recommended system has the best all-around performance. When mutations are required, we recommend Giraph hash map over byte array if memory is not a constraint. If memory is the primary constraint, then GPS is the best option for both mutation and non-mutation algorithms.

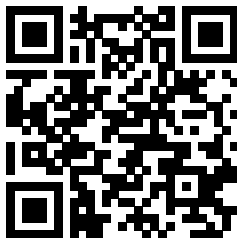
Next, we discuss results for each system, focusing on the larger graphs, before highlighting our findings for the LJ and UR graphs.

6.2 Giraph

Compared to GPS, Giraph's computation times are longer for PageRank and comparable or shorter for SSSP and WCC (Figure 2). Since Giraph's setup times are much shorter than GPS's, especially on 64 and 128 machines, Giraph is faster overall: up to 3× shorter total time on PageR-

Code, results, wiki, and EC2 images:

`http://xvz.github.io/graph-processing`



1. Graph storage should be memory and mutation efficient.

No Mutations		
	Time	Memory
Byte array	✓	✓
Hash map	✗	✗

With Mutations (DMST)		
	Time	Memory
Byte array	✗✗	✓
Hash map	✓	✗

2. LALP and dynamic migration provide little gains.

- **Time:** negative to no improvements. ✗
- **Memory:** large overheads on large graphs. ✗
- **Network:** little to no improvements. ✓

3. Message *processing* optimizations are very important.

	Computation	Setup	Memory	Network
No Mutations				
Giraph (byte array)	★ ★ ★ ☆	★ ★ ★ ★	★ ★ ★ ☆	★ ★ ★ ★
Giraph (hash map)	★ ★ ★ ☆	★ ★ ★ ☆	★ ★ ☆ ☆	★ ★ ★ ★
GPS (none)	★ ★ ★ ☆	★ ★ ☆ ☆	★ ★ ★ ★	★ ★ ☆ ☆
Mizan	★ ☆ ☆ ☆	★ ☆ ☆ ☆	★ ☆ ☆ ☆	★ ★ ★ ☆
GraphLab (sync)	★ ★ ★ ★	★ ★ ★ ☆	★ ★ ★ ☆	★ ★ ★ ☆
GraphLab (async)	★ ☆ ☆ ☆	★ ★ ★ ☆	★ ★ ☆ ☆	★ ☆ ☆ ☆

4. Distributed locking for asynchronous execution is not scalable.
- Performance degrades as more machines are used. **X**
 - Due to: lock contention, termination scheme, lack of message batching.

5. Giraph scales better across graphs;
GraphLab scales better across more machines.

5. Giraph scales better across graphs;
GraphLab scales better across more machines.

64 machines	TW	UK
Giraph (byte array)	5.8GB	7.0GB
GraphLab (sync)	4.5GB	14GB

5. Giraph scales better across graphs;
GraphLab scales better across more machines.

64 machines	TW	UK
Giraph (byte array)	5.8GB	7.0GB
GraphLab (sync)	4.5GB	14GB

TW	16 machines	128 machines
Giraph (byte array)	8.5GB	5.8GB
GraphLab (sync)	11GB	3.3GB

Comparison Summary

	Computation	Setup	Memory	Network
No Mutations				
Giraph (byte array)	★ ★ ★ ☆	★ ★ ★ ★	★ ★ ★ ☆	★ ★ ★ ★
Giraph (hash map)	★ ★ ★ ☆	★ ★ ★ ☆	★ ★ ☆ ☆	★ ★ ★ ★
GPS (none)	★ ★ ★ ☆	★ ★ ☆ ☆	★ ★ ★ ★	★ ★ ☆ ☆
Mizan	★ ☆ ☆ ☆	★ ☆ ☆ ☆	★ ☆ ☆ ☆	★ ★ ★ ☆
GraphLab (sync)	★ ★ ★ ★	★ ★ ★ ☆	★ ★ ★ ☆	★ ★ ★ ☆
GraphLab (async)	★ ☆ ☆ ☆	★ ★ ★ ☆	★ ★ ☆ ☆	★ ☆ ☆ ☆
With Mutations (DMST)				
Giraph (byte array)	★ ☆ ☆ ☆	★ ★ ★ ★	★ ★ ☆ ☆	★ ★ ★ ★
Giraph (hash map)	★ ★ ★ ★	★ ★ ★ ★	★ ★ ☆ ☆	★ ★ ★ ★
GPS (none)	★ ★ ★ ★	★ ☆ ☆ ☆	★ ★ ★ ★	★ ☆ ☆ ☆

- Comprehensive experimental comparison of four graph processing systems, with scale.
- System optimizations are critical, but some optimizations can degrade performance.
- Performance and scalability depend heavily on metric used: there is no clear winner.

- Comprehensive experimental comparison of four graph processing systems, with scale.
- System optimizations are critical, but some optimizations can degrade performance.
- Performance and scalability depend heavily on metric used: there is no clear winner.

- Comprehensive experimental comparison of four graph processing systems, with scale.
- System optimizations are critical, but some optimizations can degrade performance.
- Performance and scalability depend heavily on metric used: there is no clear winner.