

Data-Intensive Distributed Computing

CS 451/651 (Fall 2018)

Part 8: Analyzing Graphs, Redux (1/2)

November 15, 2018

Jimmy Lin

David R. Cheriton School of Computer Science

University of Waterloo

These slides are available at <http://lintool.github.io/bigdata-2018f/>




This work is licensed under a Creative Commons Attribution-NonCommercial-Share Alike 3.0 United States
See <http://creativecommons.org/licenses/by-nc-sa/3.0/us/> for details

Graph Algorithms, again?

(srsly?)

What makes graphs hard?



Irregular structure

Fun with data structures!

Irregular data access patterns

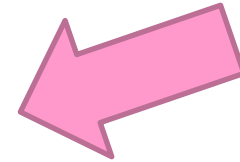
Fun with architectures!

Iterations

Fun with optimizations!

Characteristics of Graph Algorithms

Parallel graph traversals

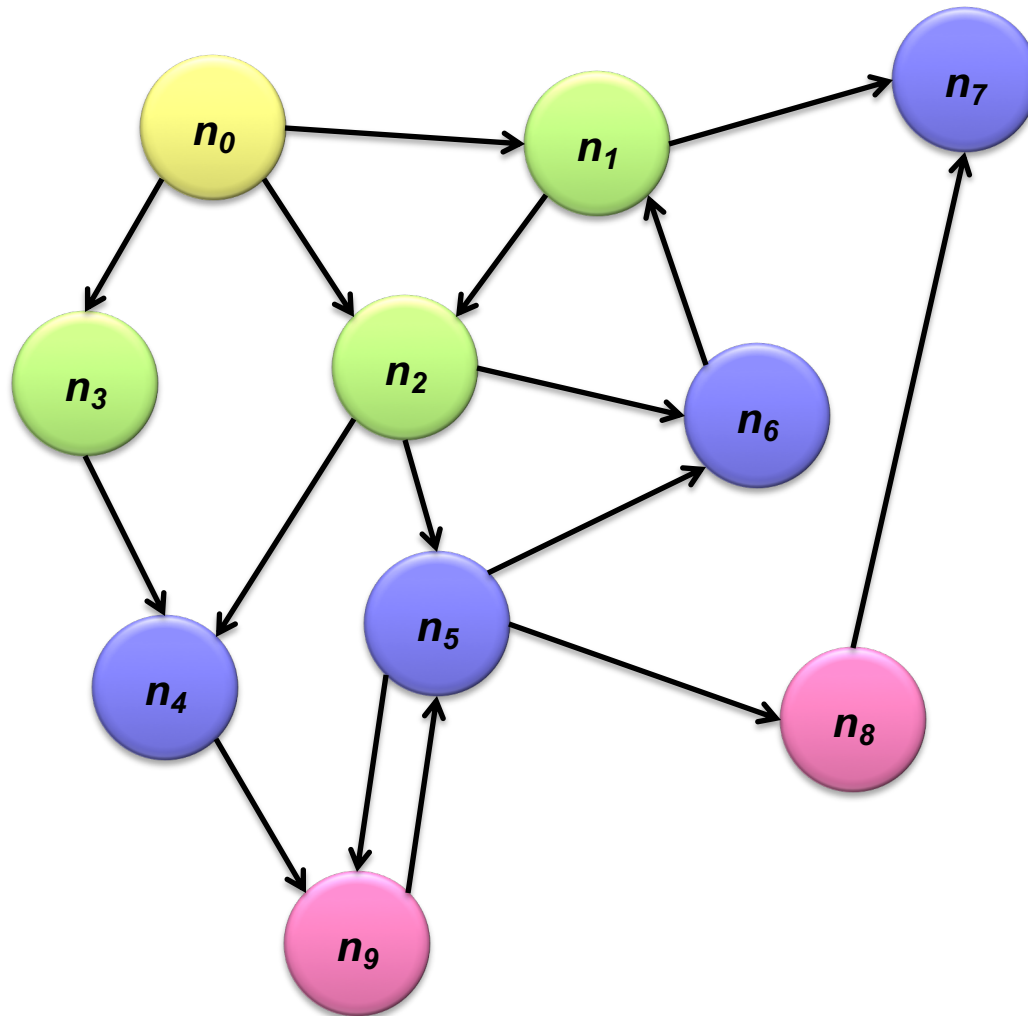


Local computations

Message passing along graph edges

Iterations

Visualizing Parallel BFS



PageRank: Defined

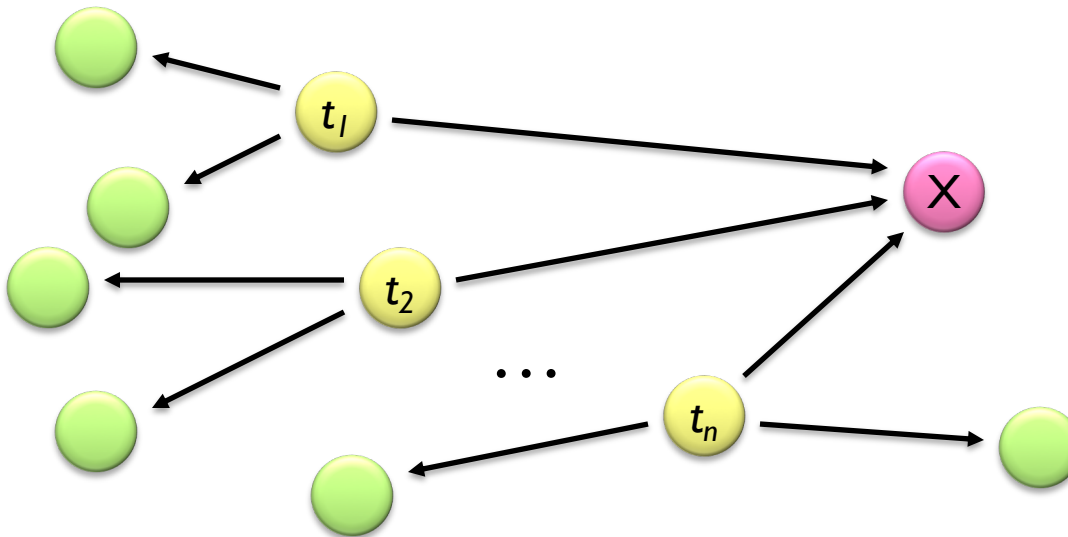
Given page x with inlinks $t_1 \dots t_n$, where

$C(t)$ is the out-degree of t

α is probability of random jump

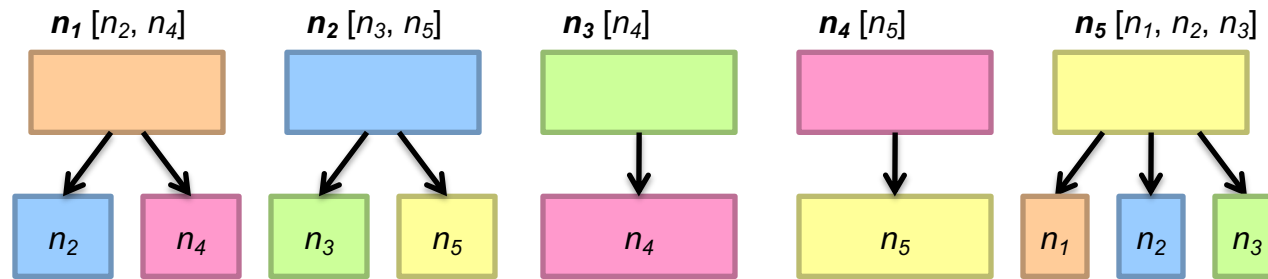
N is the total number of nodes in the graph

$$PR(x) = \alpha \left(\frac{1}{N} \right) + (1 - \alpha) \sum_{i=1}^n \frac{PR(t_i)}{C(t_i)}$$

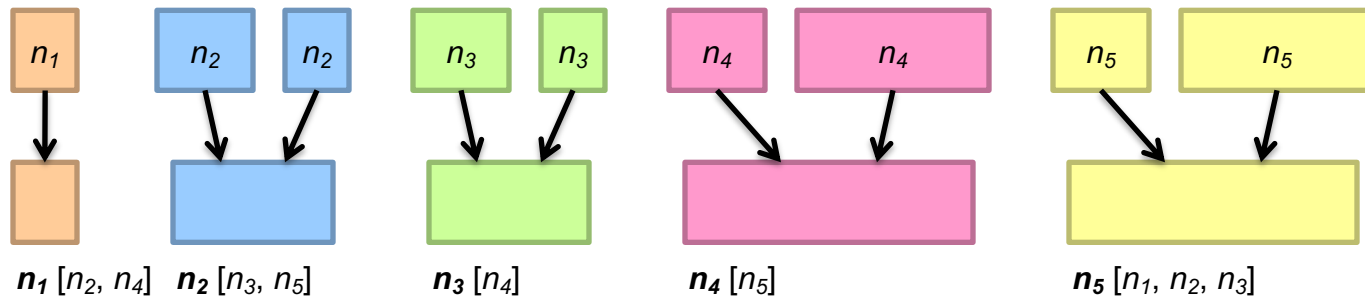


PageRank in MapReduce

Map



Reduce



PageRank vs. BFS

	PageRank	BFS
Map	PR/N	$d+1$
Reduce	sum	min

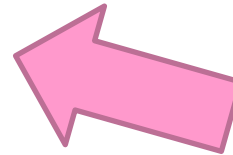
Characteristics of Graph Algorithms

Parallel graph traversals

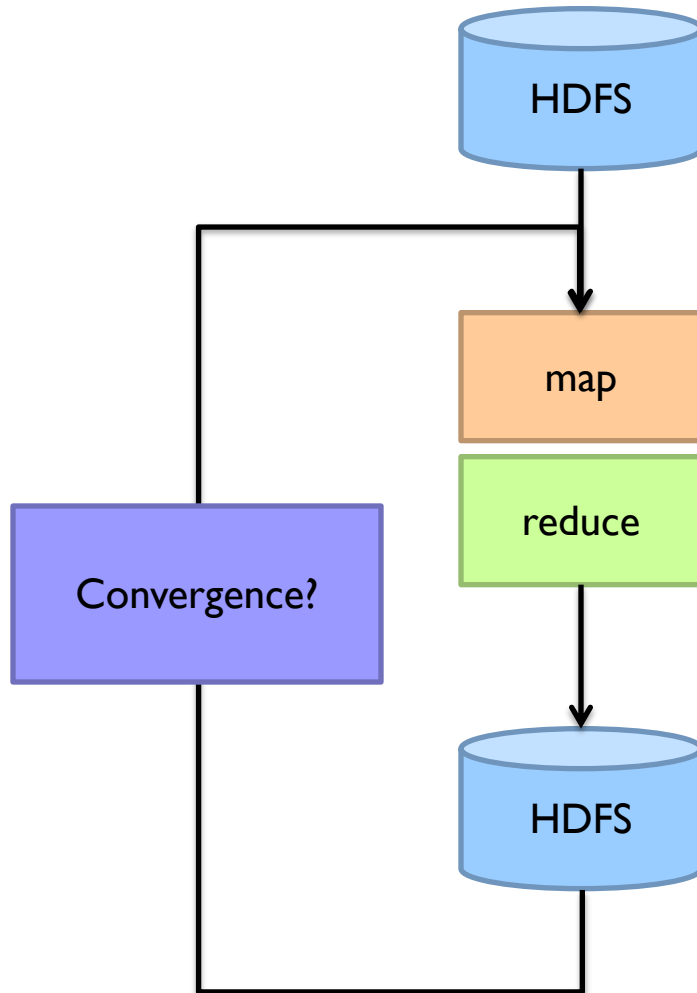
Local computations

Message passing along graph edges

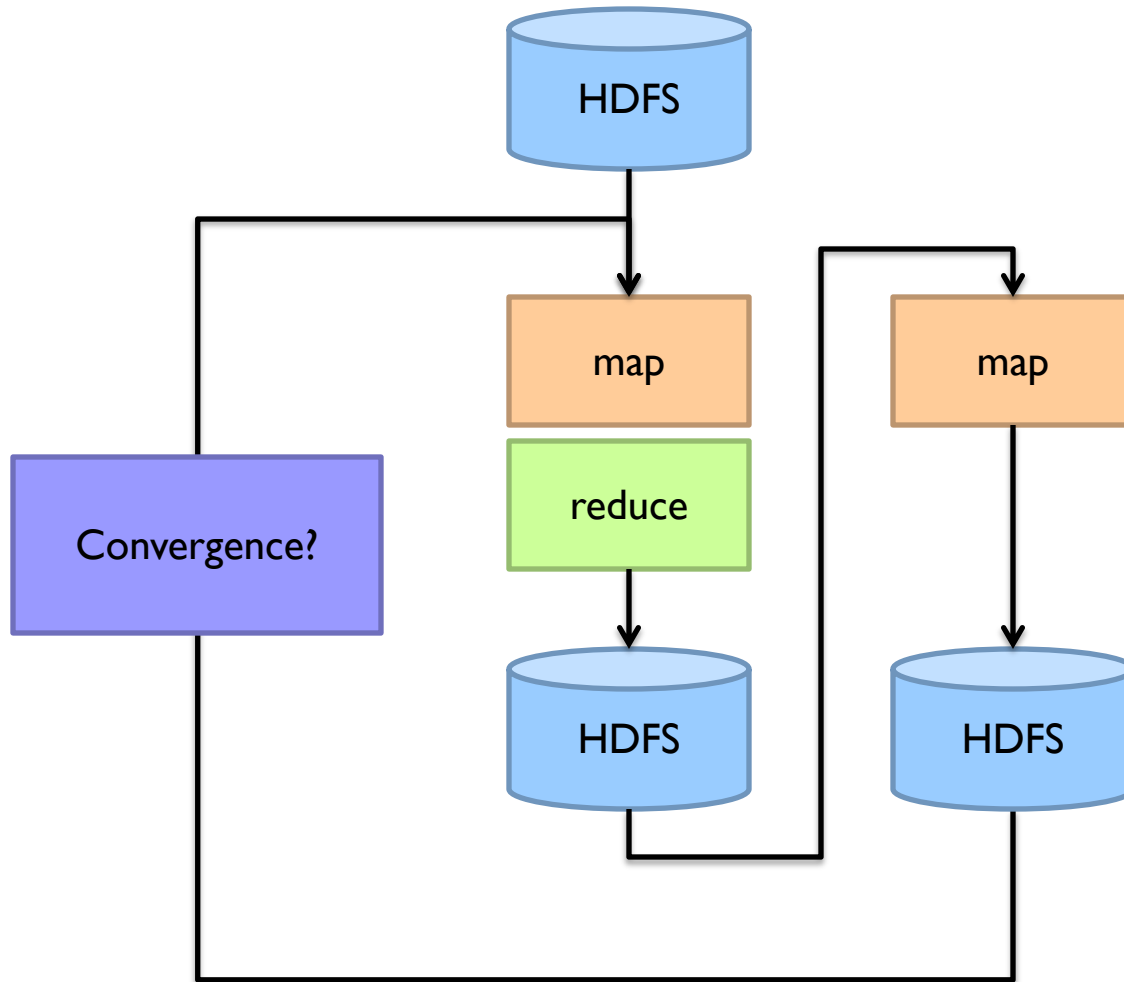
Iterations



BFS



PageRank



MapReduce Sucks

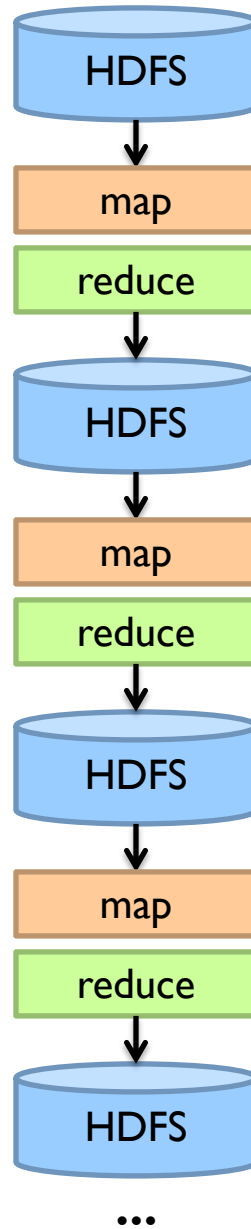
Hadoop task startup time

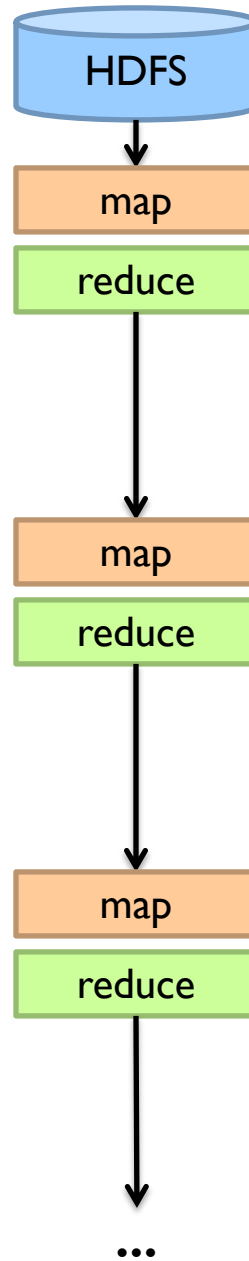
Stragglers

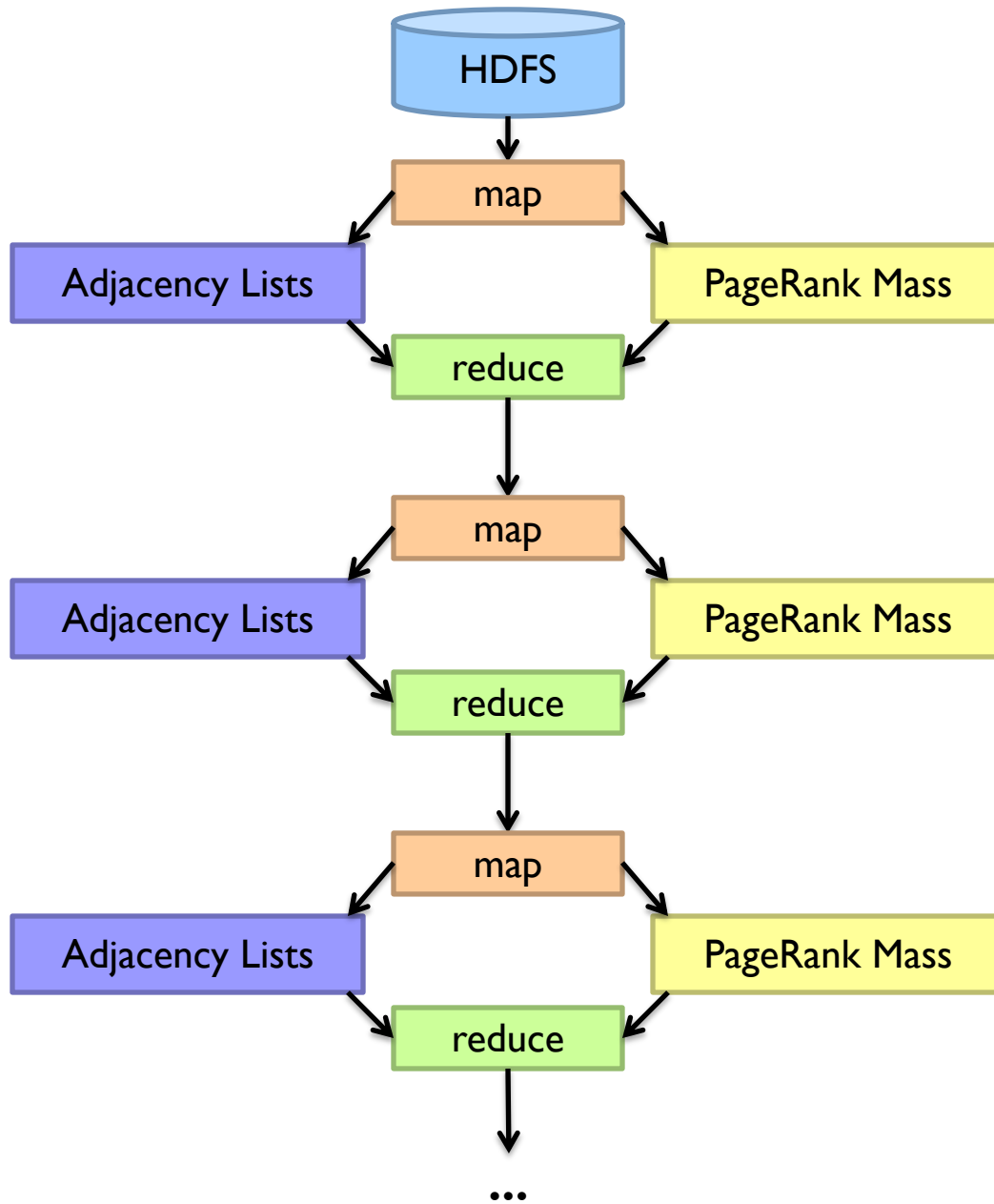
Needless graph shuffling

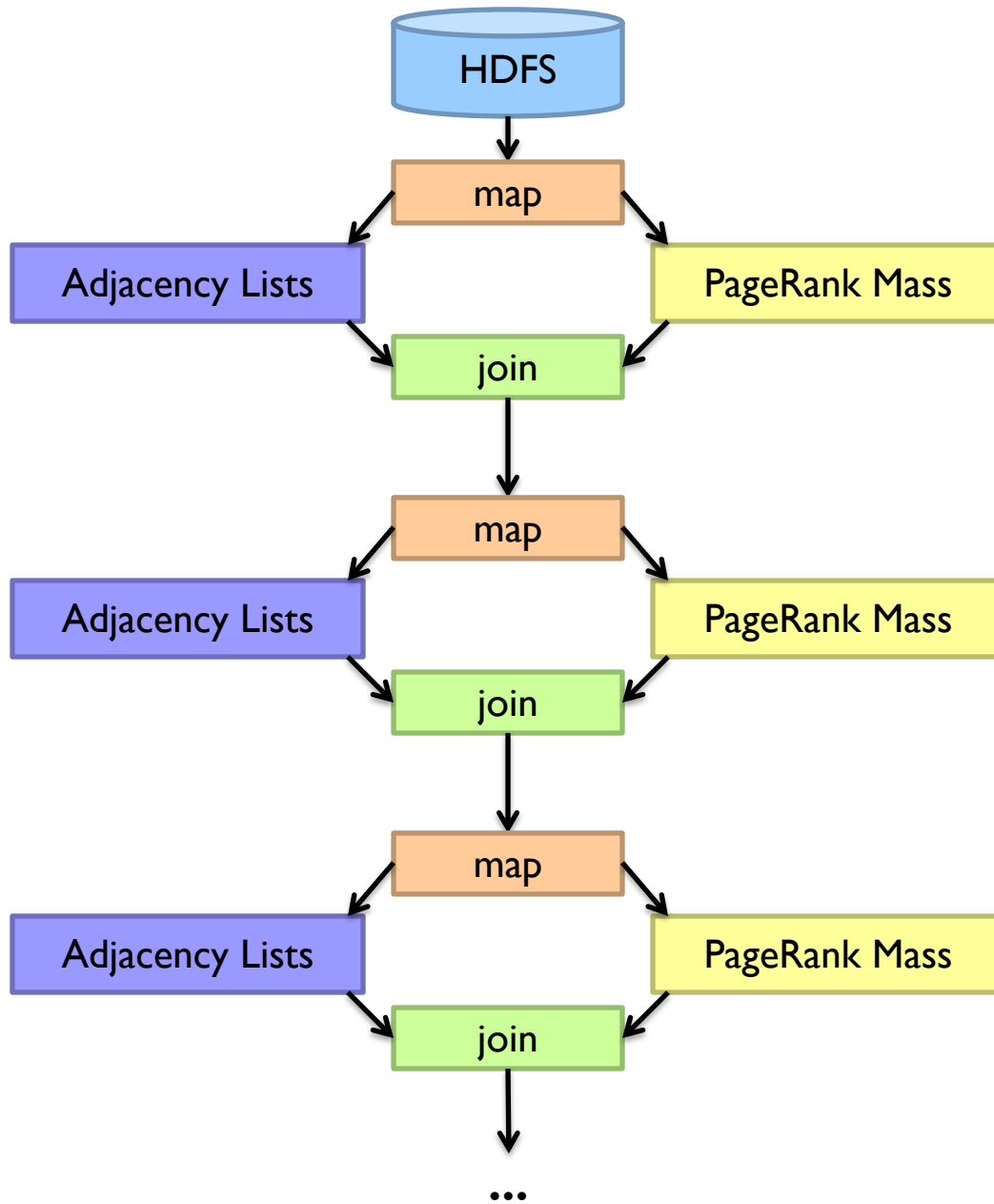
Checkpointing at each iteration

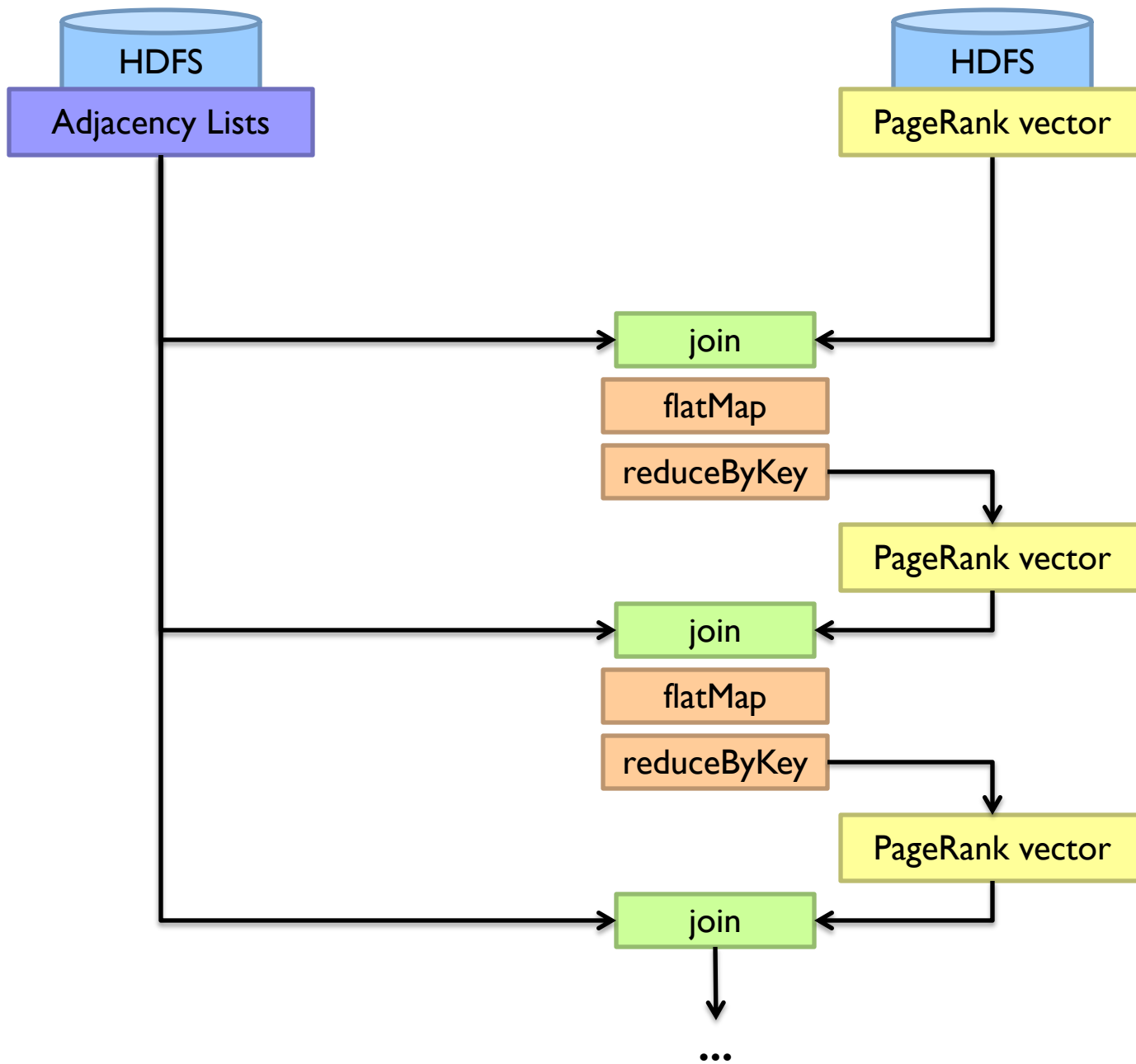
Let's Spark!

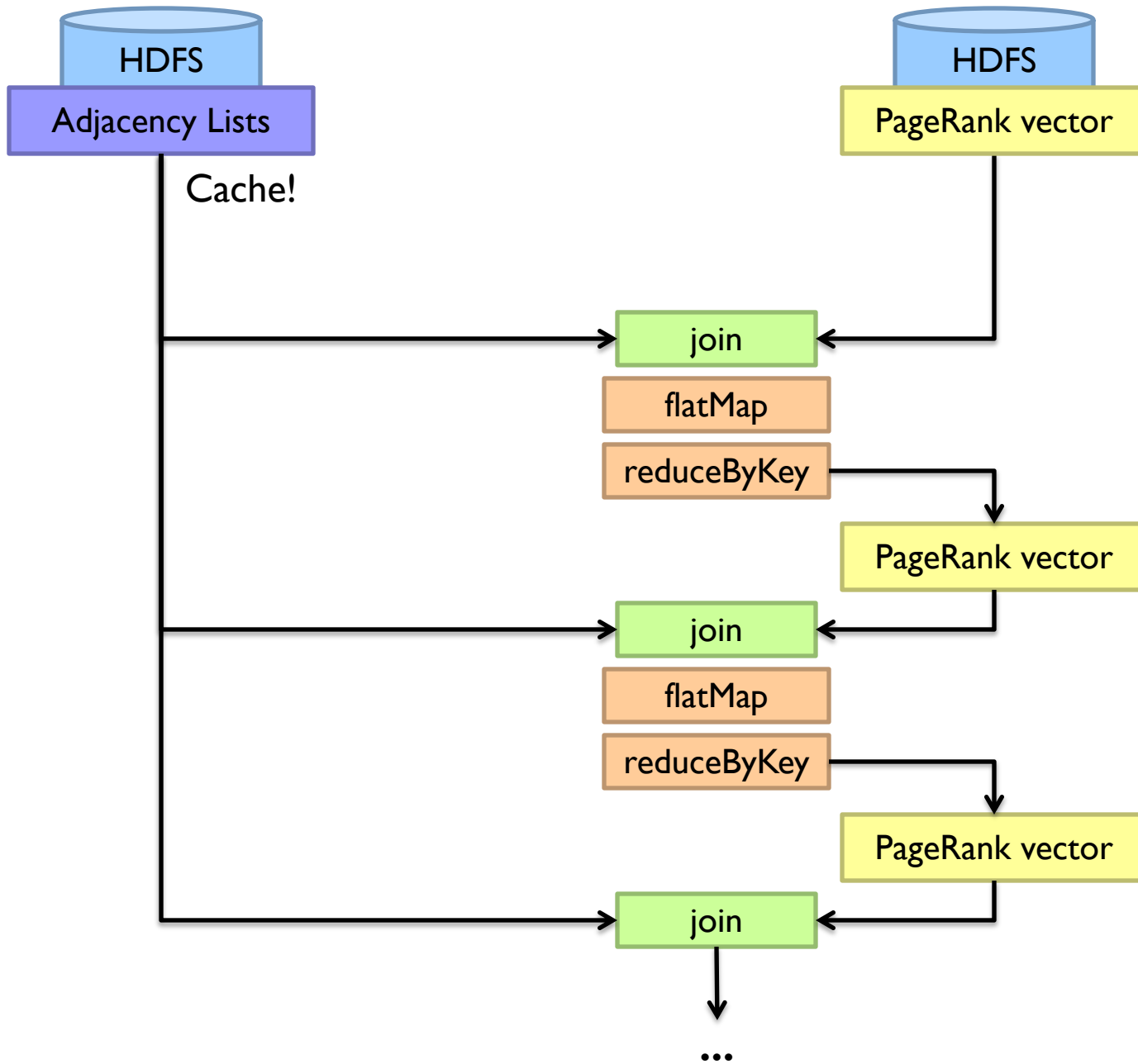




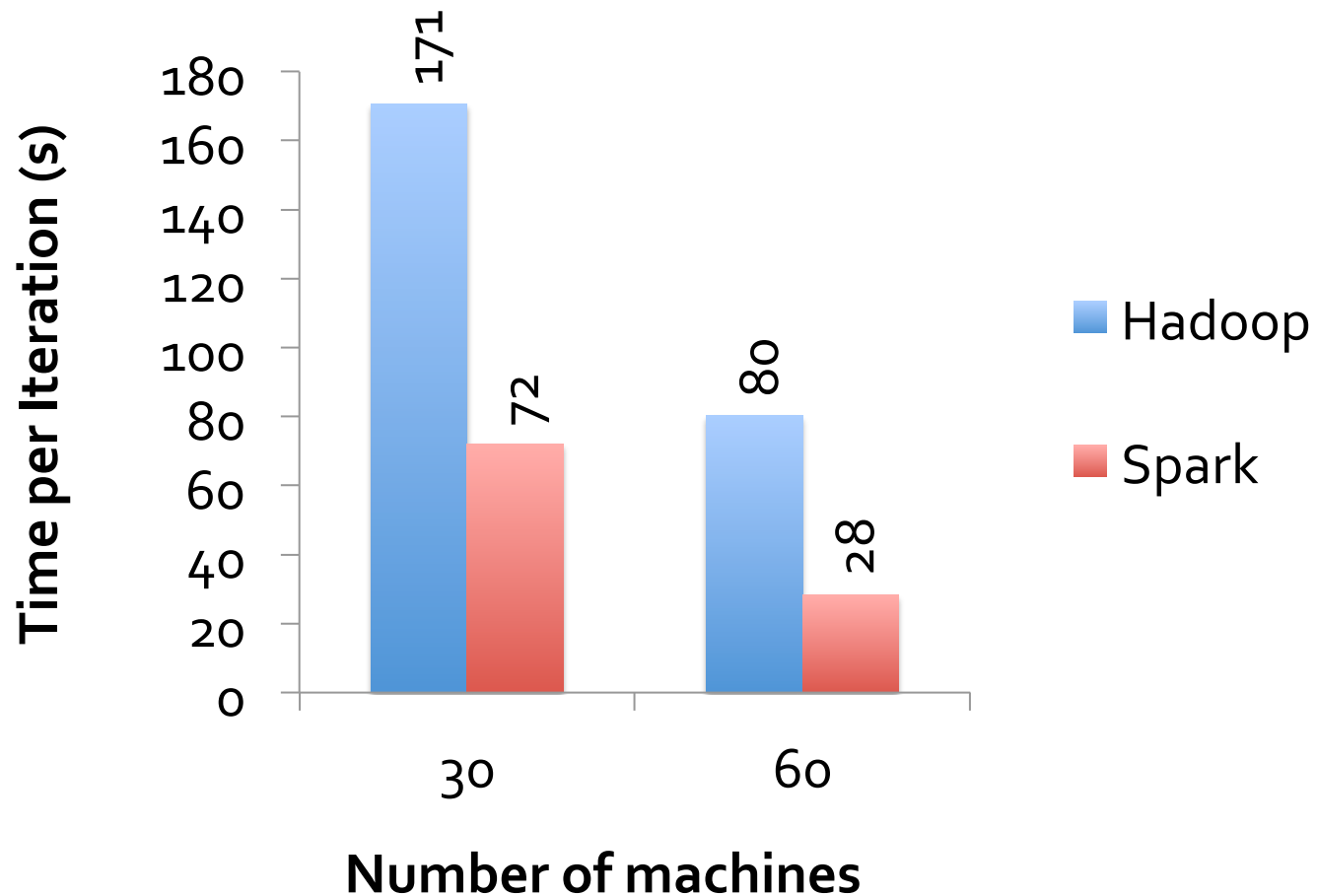








MapReduce vs. Spark



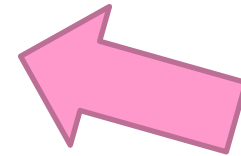
Characteristics of Graph Algorithms

Parallel graph traversals

Local computations

Message passing along graph edges

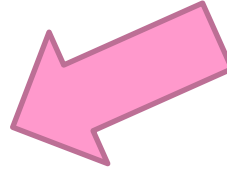
Iterations



Even faster?

Big Data Processing in a Nutshell

Let's be smarter about this!



Partition

Replicate

Reduce cross-partition communication

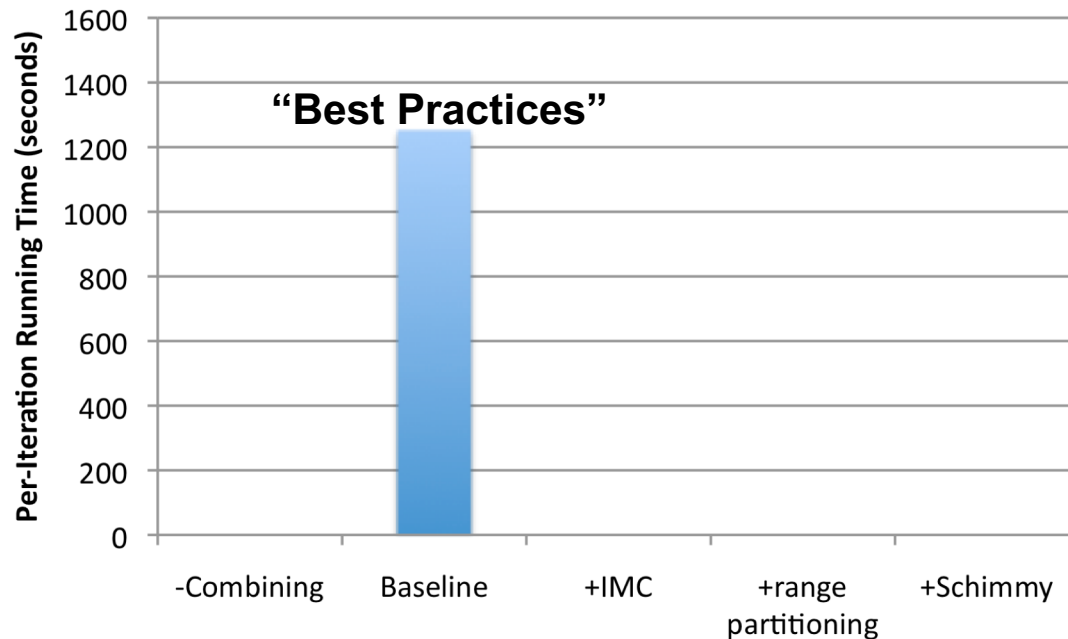
Simple Partitioning Techniques

Hash partitioning

Range partitioning on some underlying linearization

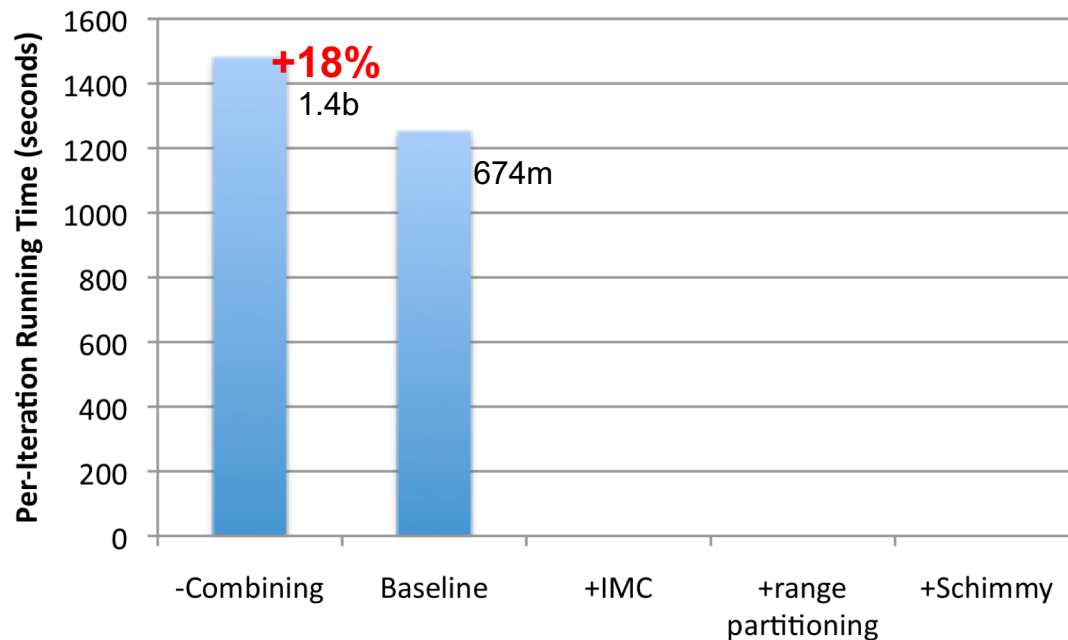
Web pages: lexicographic sort of domain-reversed URLs

How much difference does it make?



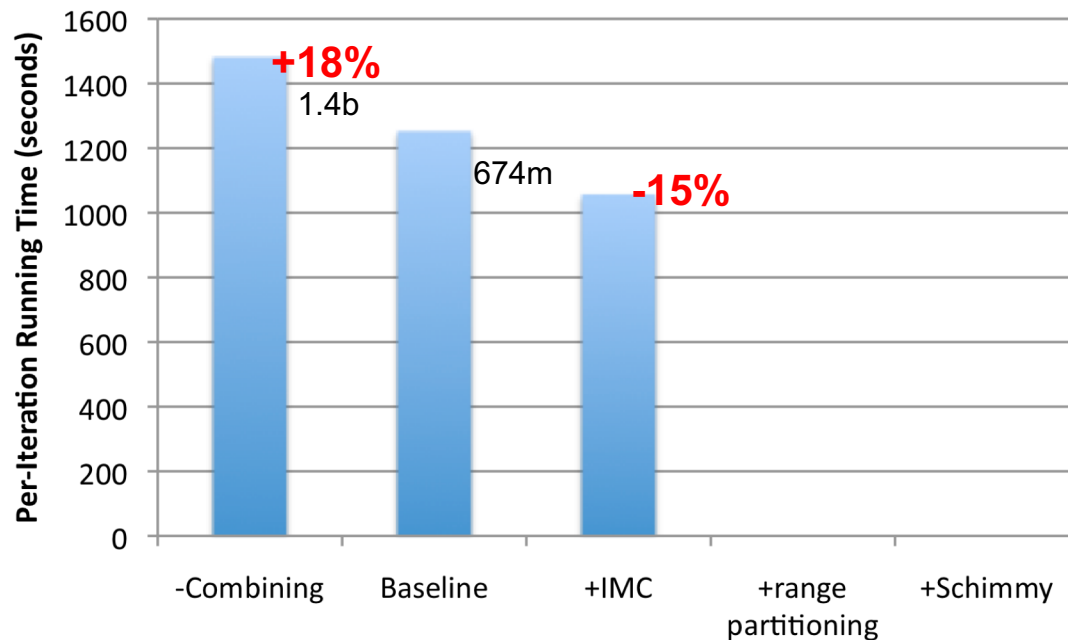
PageRank over webgraph
(40m vertices, 1.4b edges)

How much difference does it make?



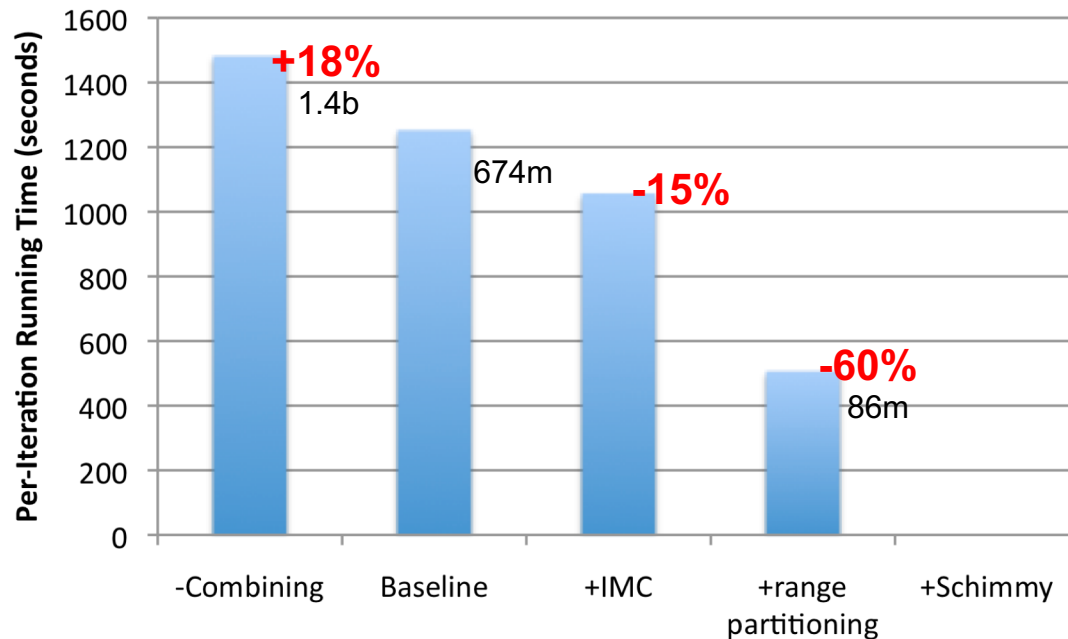
PageRank over webgraph
(40m vertices, 1.4b edges)

How much difference does it make?



PageRank over webgraph
(40m vertices, 1.4b edges)

How much difference does it make?



PageRank over webgraph
(40m vertices, 1.4b edges)

Schimmy Design Pattern

Basic implementation contains two dataflows:

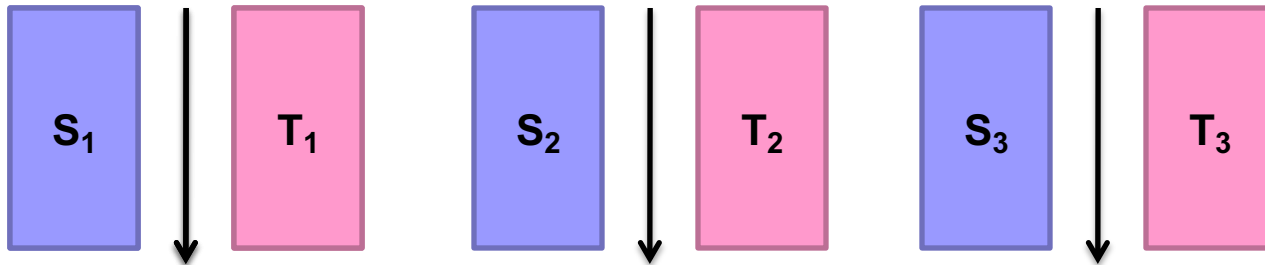
Messages (actual computations)

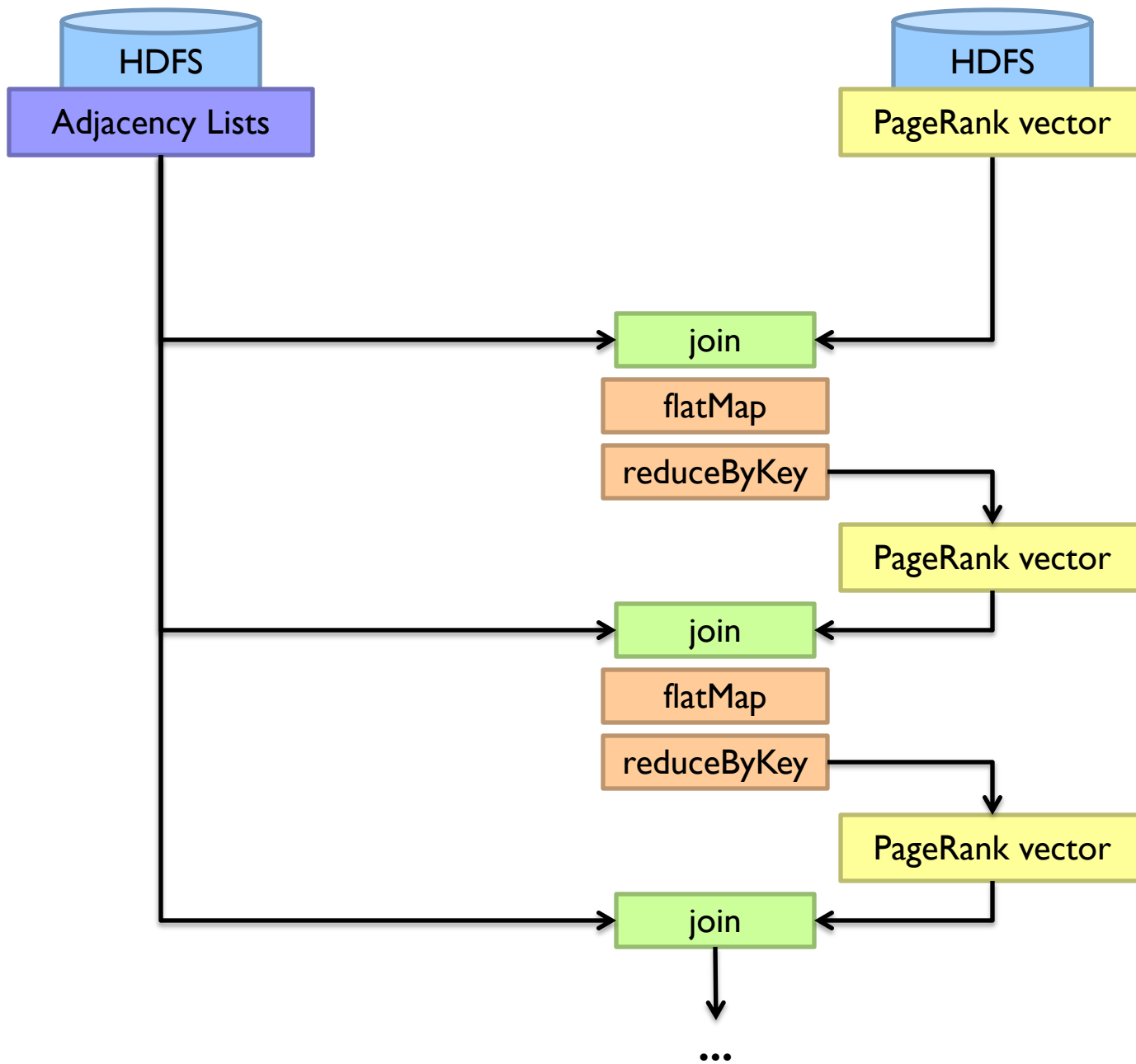
Graph structure (“bookkeeping”)

Schimmy: separate the two dataflows, shuffle only the messages

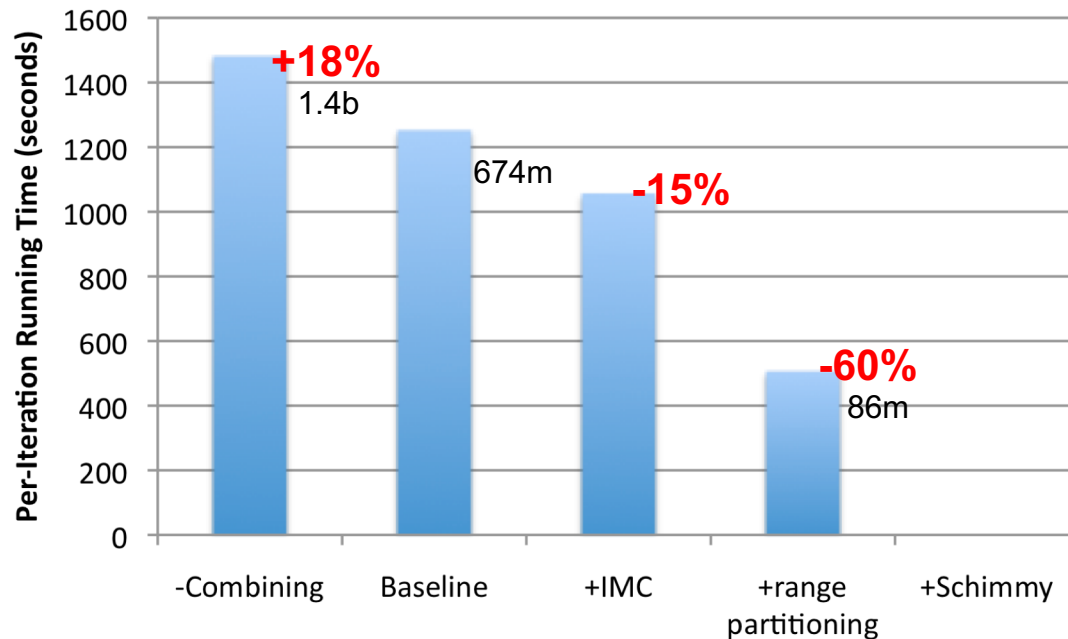
Basic idea: merge join between graph structure and messages

both relations are sorted by join key



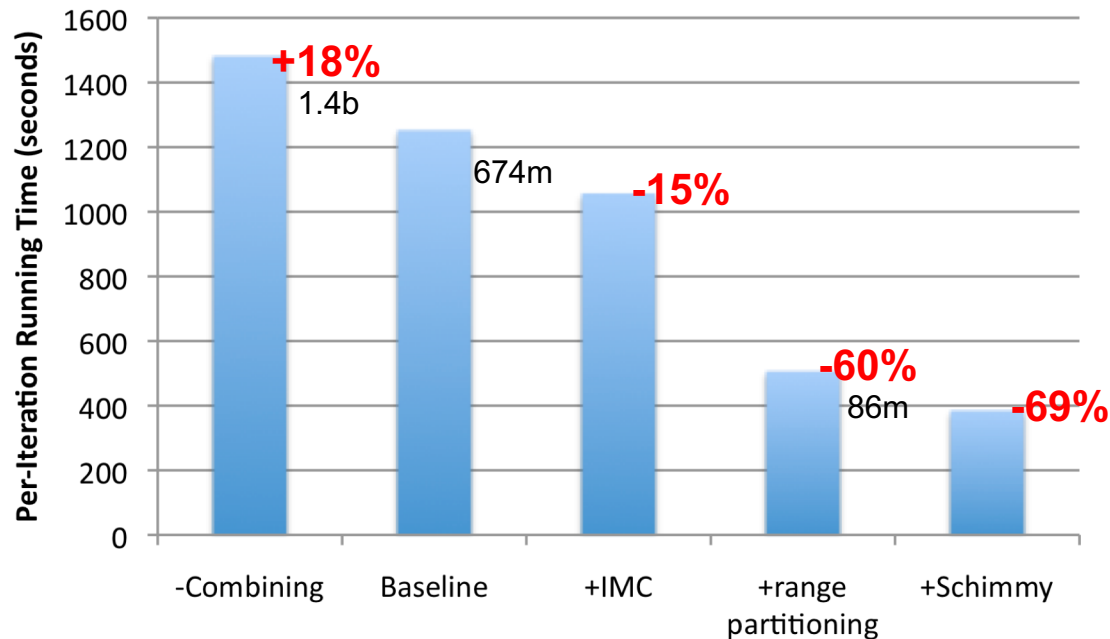


How much difference does it make?



PageRank over webgraph
(40m vertices, 1.4b edges)

How much difference does it make?



PageRank over webgraph
(40m vertices, 1.4b edges)

Simple Partitioning Techniques

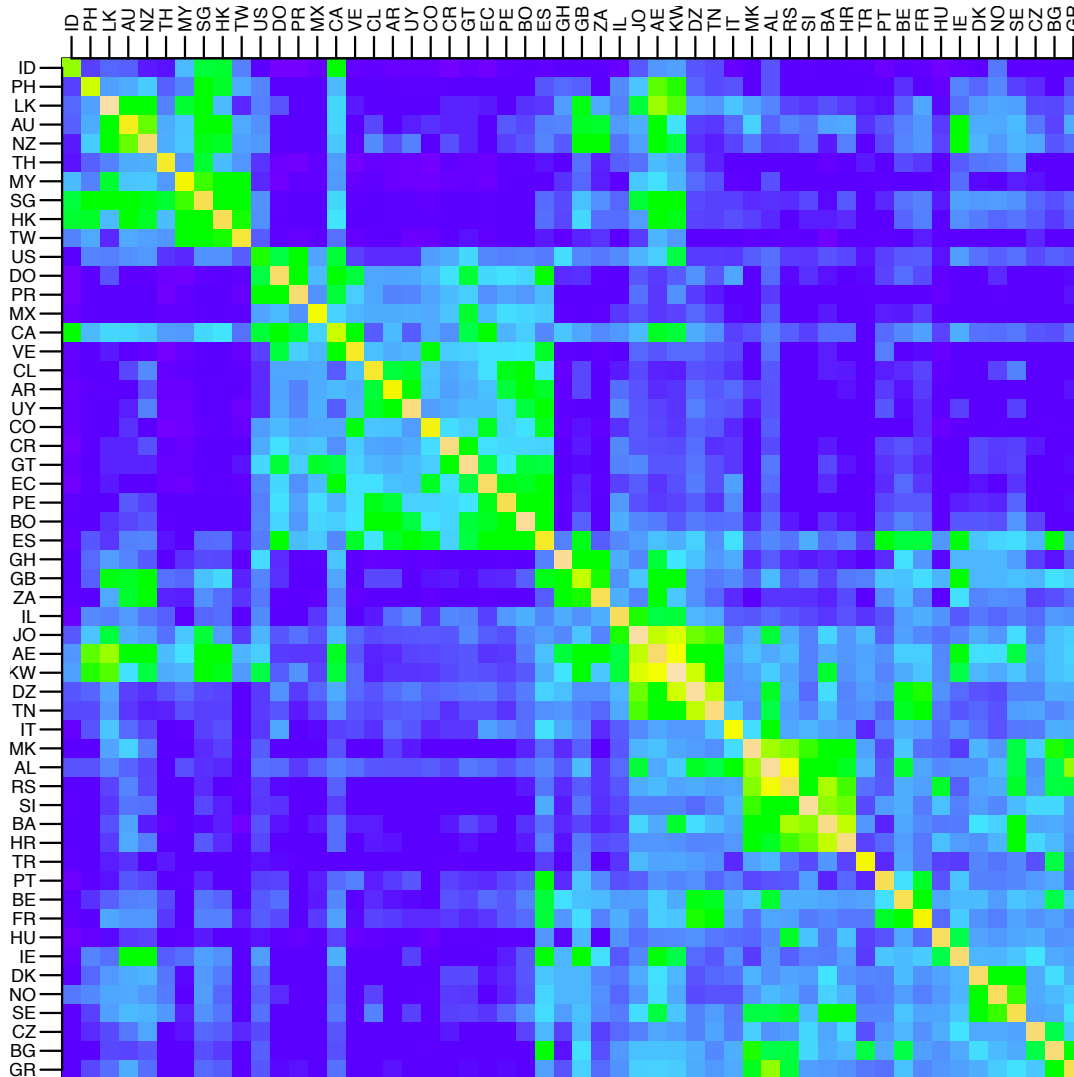
Hash partitioning

Range partitioning on some underlying linearization

Web pages: lexicographic sort of domain-reversed URLs

Social networks: sort by demographic characteristics

Country Structure in Facebook



Analysis of 721 million active users (May 2011)

54 countries w/ >1m active users, >50% penetration

Simple Partitioning Techniques

Hash partitioning

Range partitioning on some underlying linearization

Web pages: lexicographic sort of domain-reversed URLs

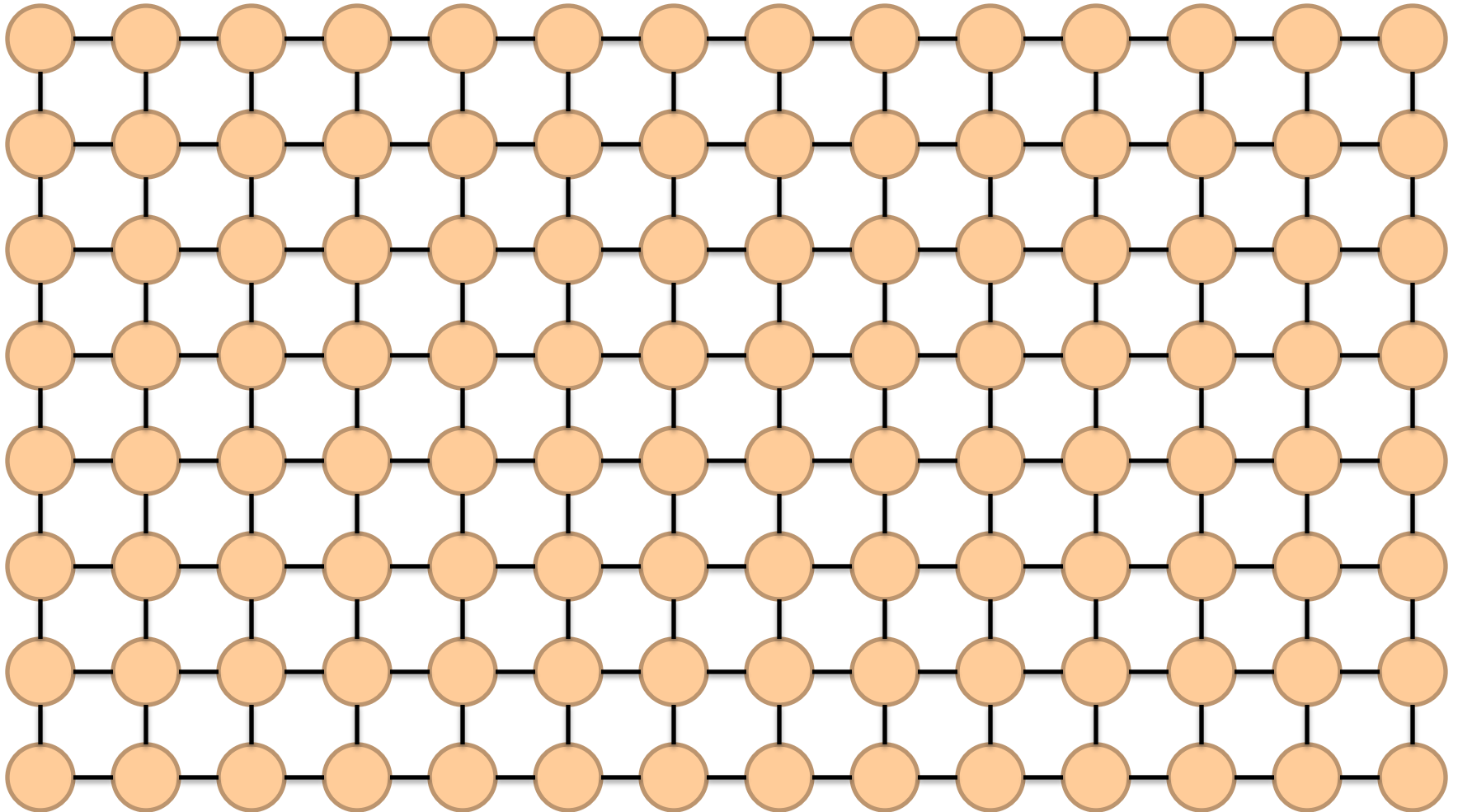
Social networks: sort by demographic characteristics

Geo data: space-filling curves

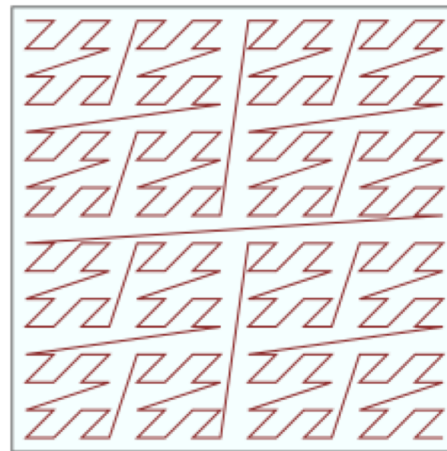
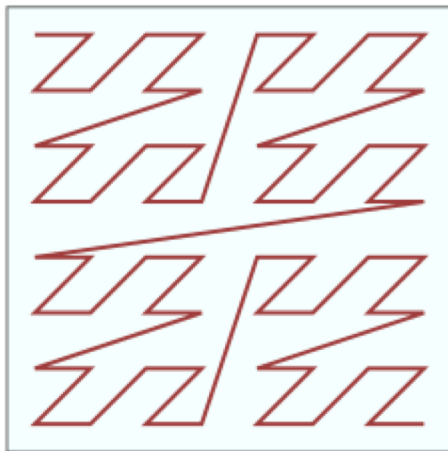
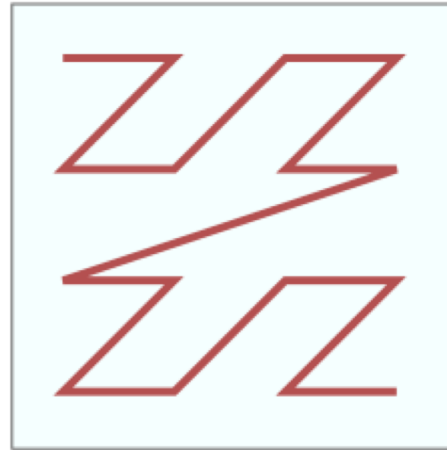
Aside: Partitioning Geo-data



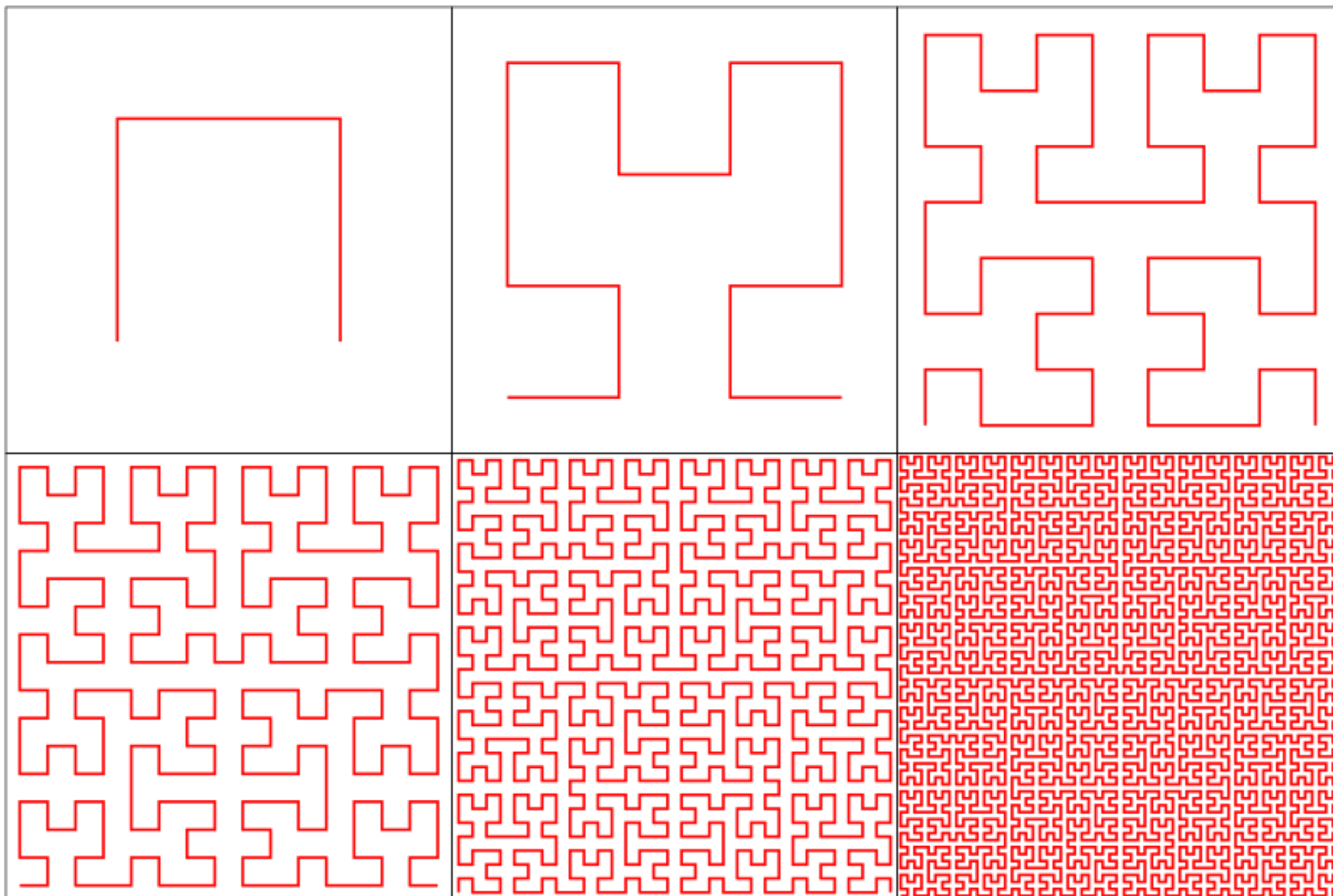
Geo-data = regular graph



Space-filling curves: Z-Order Curves



Space-filling curves: Hilbert Curves



Simple Partitioning Techniques

Hash partitioning

Range partitioning on some underlying linearization

Web pages: lexicographic sort of domain-reversed URLs

Social networks: sort by demographic characteristics

Geo data: space-filling curves

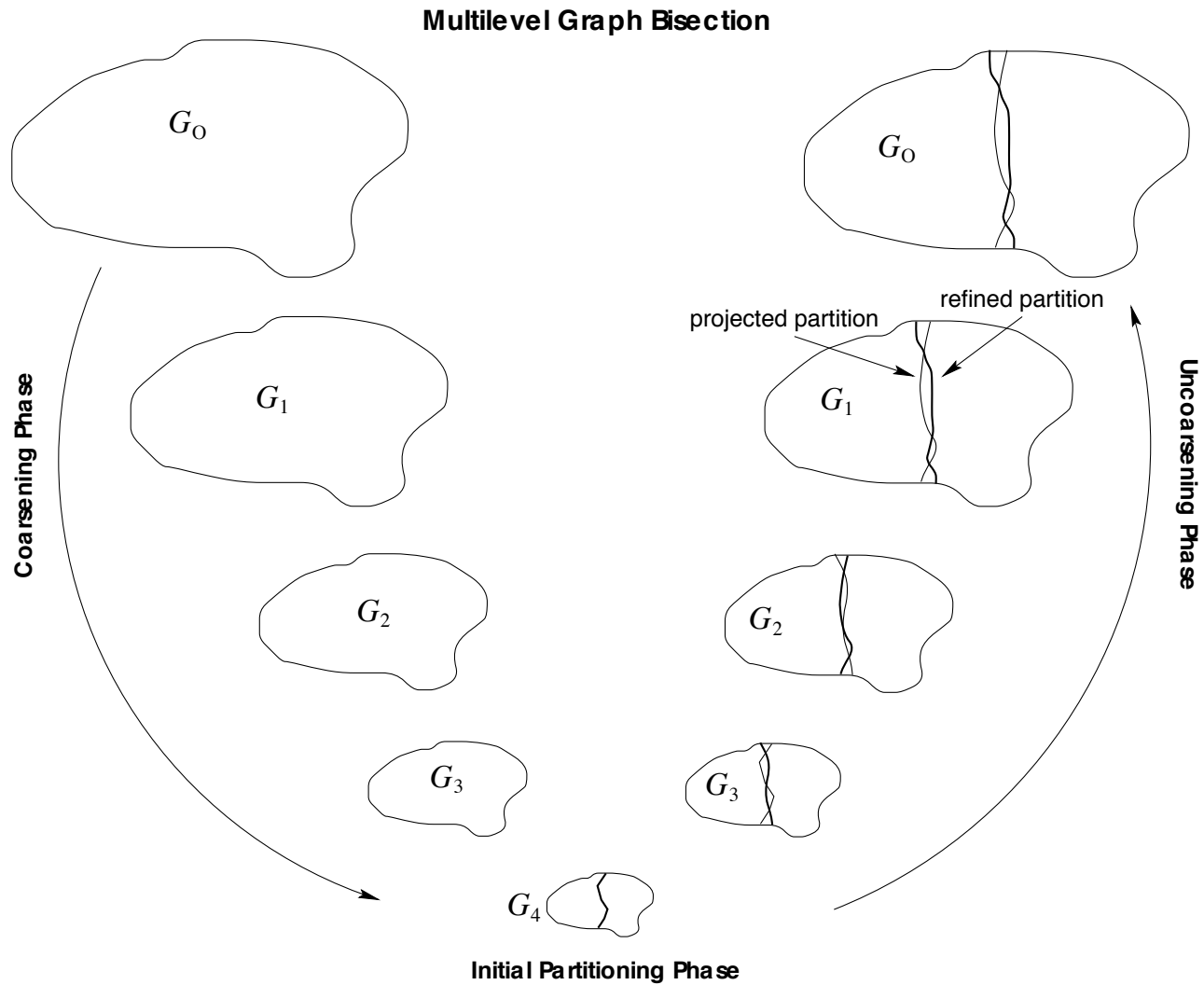
But what about graphs in general?



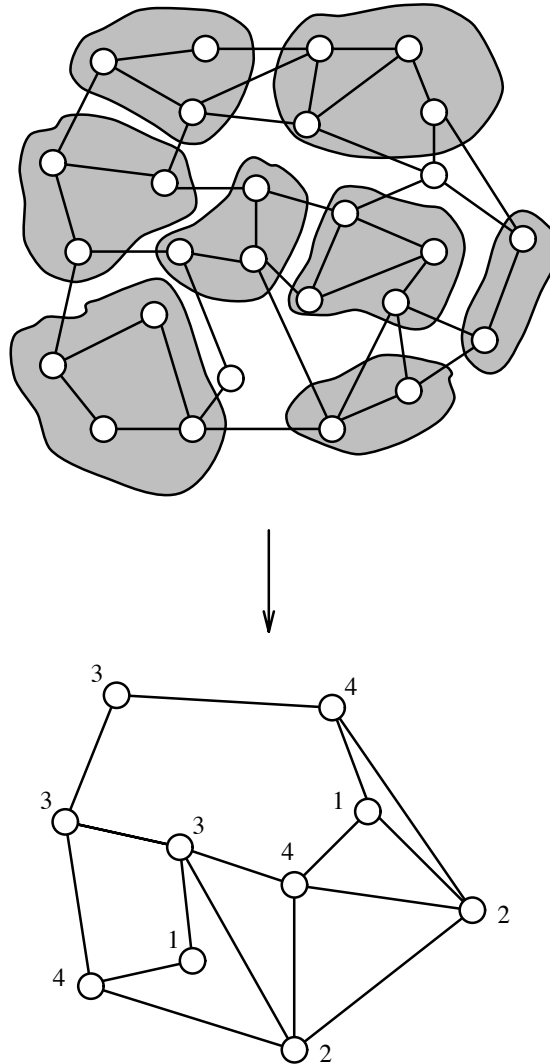
General-Purpose Graph Partitioning

Graph coarsening
Recursive bisection

General-Purpose Graph Partitioning



Graph Coarsening



Chicken-and-Egg

To coarsen the graph you need to identify dense local regions
To identify dense local regions quickly you need to traverse local edges
But to traverse local edges efficiently you need the local structure!

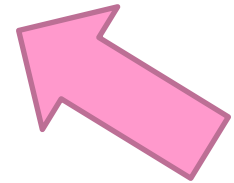
To efficiently partition the graph, you need to already know what the partitions are!
Industry solution?

Big Data Processing in a Nutshell

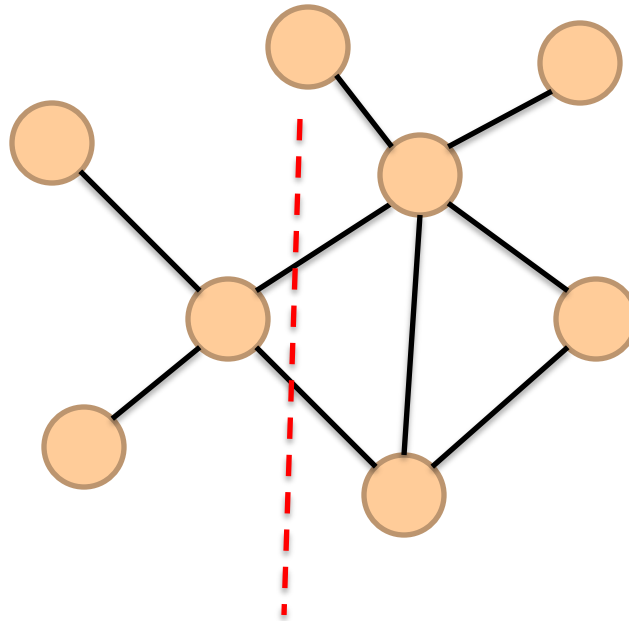
Partition

Replicate

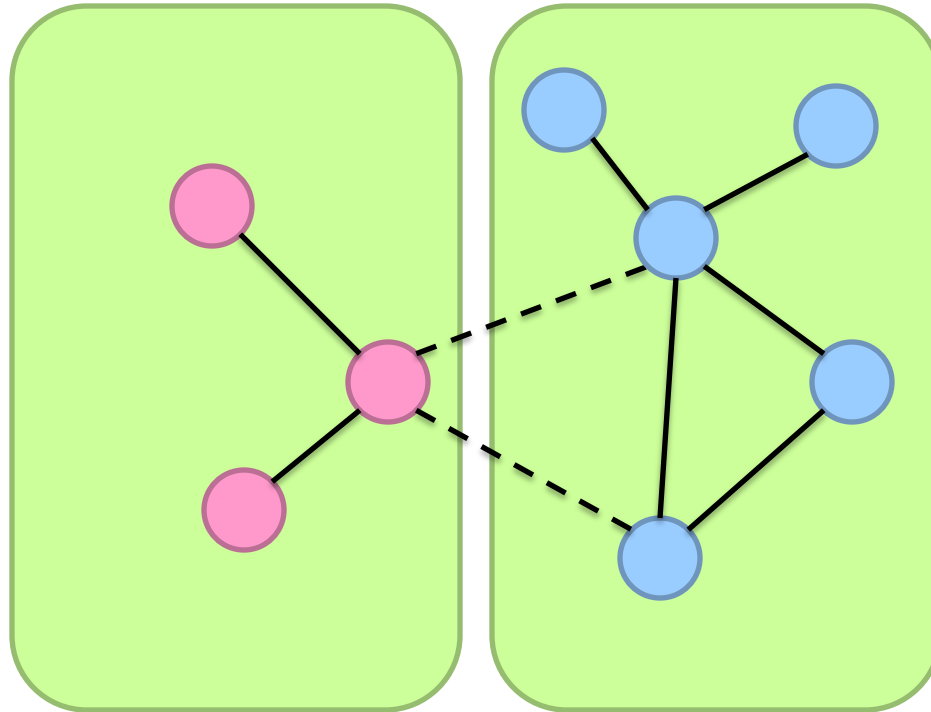
Reduce cross-partition communication



Partition



Partition



What's the fundamental issue?

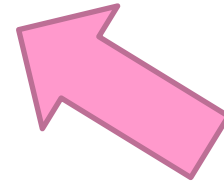
Characteristics of Graph Algorithms

Parallel graph traversals

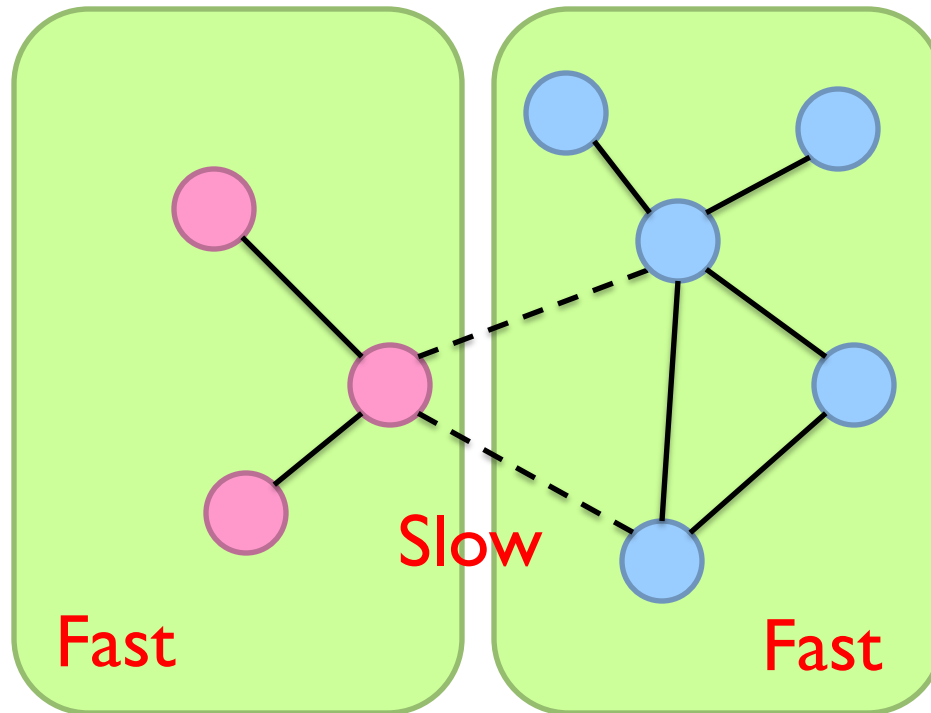
Local computations

Message passing along graph edges

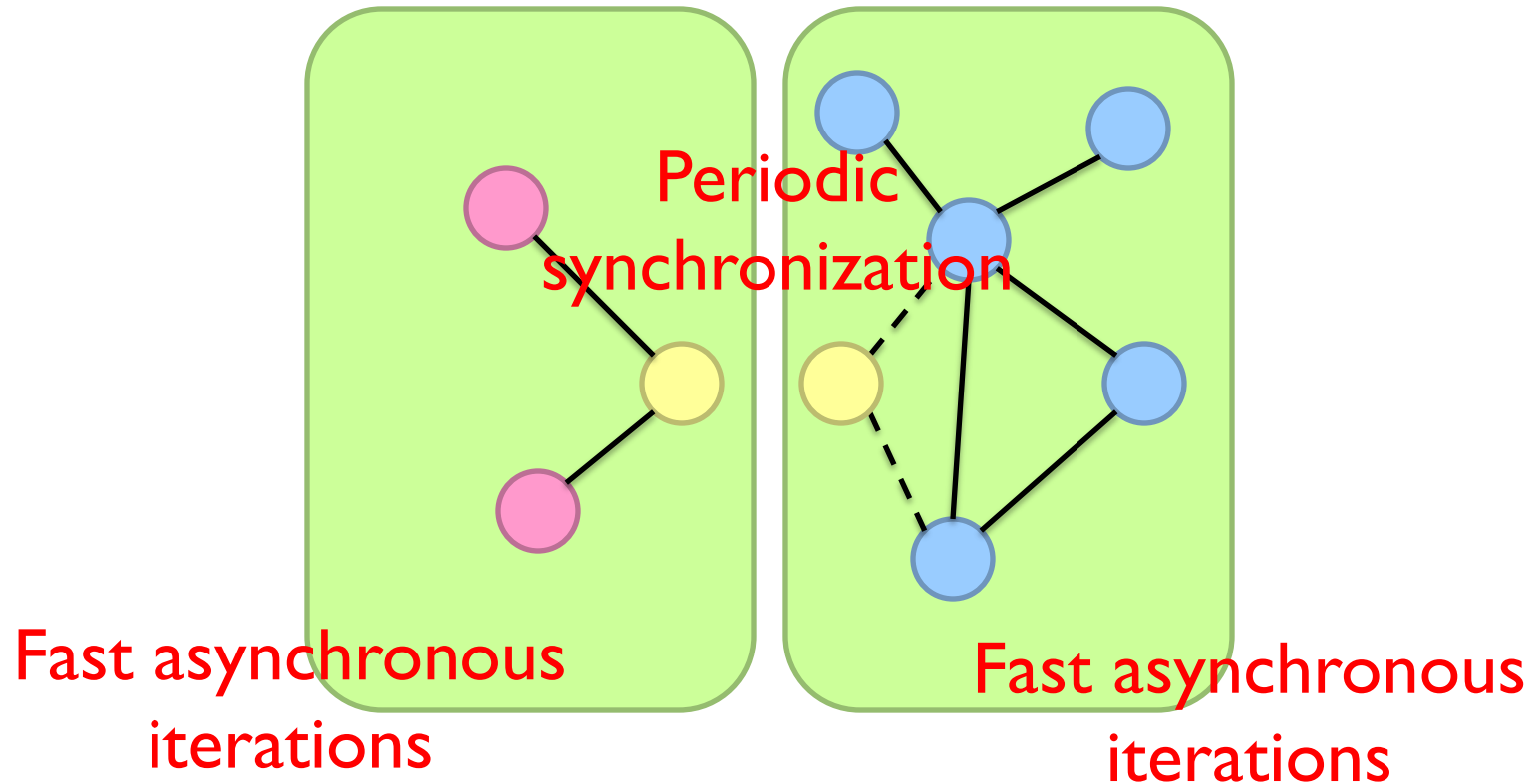
Iterations



Partition

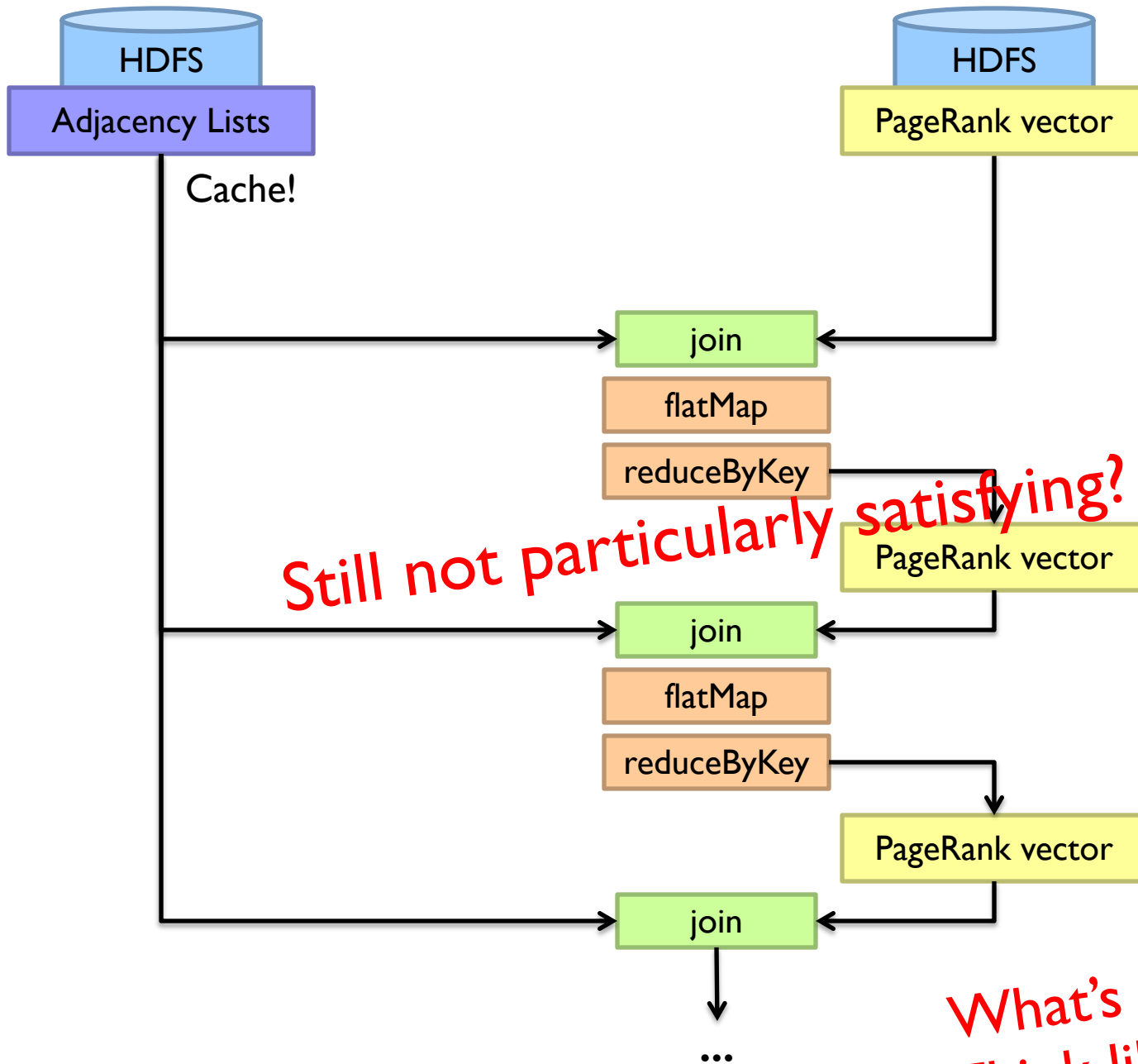


State-of-the-Art Distributed Graph Algorithms



Graph Processing Frameworks





Still not particularly satisfying?

*What's the issue?
Think like a vertex!*

Pregel: Computational Model

Based on Bulk Synchronous Parallel (BSP)

Computational units encoded in a directed graph

Computation proceeds in a series of supersteps

Message passing architecture

Each vertex, at each superstep:

Receives messages directed at it from previous superstep

Executes a user-defined function (modifying state)

Emits messages to other vertices (for the next superstep)

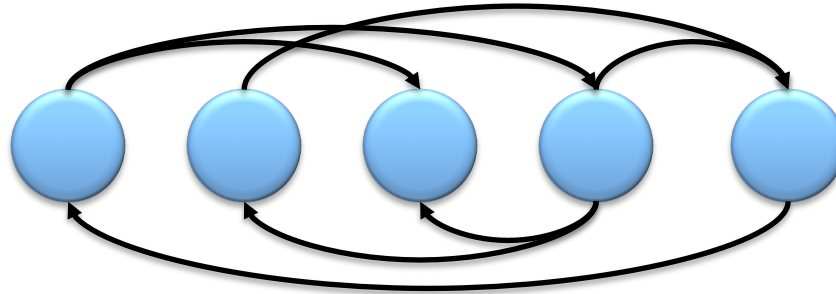
Termination:

A vertex can choose to deactivate itself

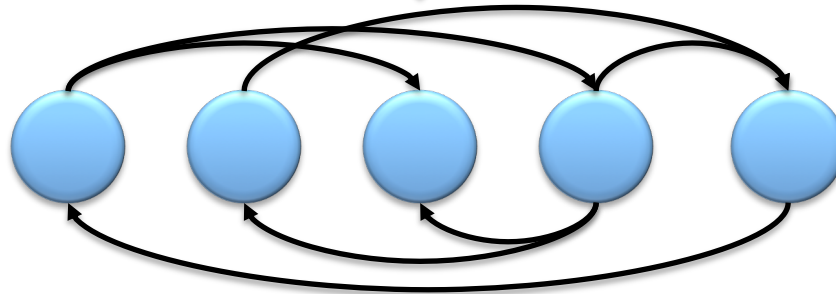
Is “woken up” if new messages received

Computation halts when all vertices are inactive

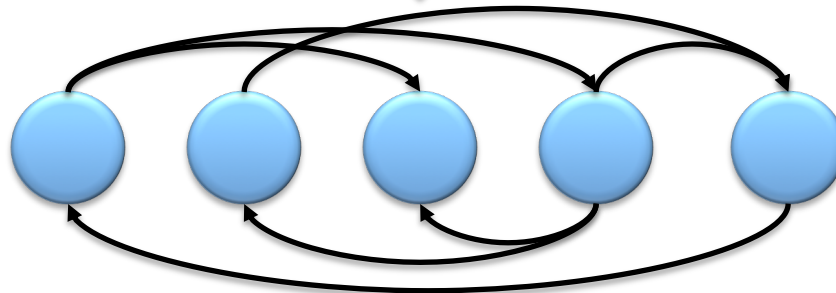
superstep t



superstep $t+1$



superstep $t+2$



Pregel: Implementation

Master-Worker architecture

Vertices are hash partitioned (by default) and assigned to workers
Everything happens in memory

Processing cycle:

Master tells all workers to advance a single superstep
Worker delivers messages from previous superstep, executing vertex computation
Messages sent asynchronously (in batches)
Worker notifies master of number of active vertices

Fault tolerance

Checkpointing
Heartbeat/revert

Pregel: SSSP

```
class ShortestPathVertex : public Vertex<int, int, int> {
    void Compute(MessageIterator* msgs) {
        int mindist = IsSource(vertex_id()) ? 0 : INF;
        for (; !msgs->Done(); msgs->Next())
            mindist = min(mindist, msgs->Value());
        if (mindist < GetValue()) {
            *MutableValue() = mindist;
            OutEdgeIterator iter = GetOutEdgeIterator();
            for (; !iter.Done(); iter.Next())
                SendMessageTo(iter.Target(),
                             mindist + iter.GetValue());
        }
        VoteToHalt();
    }
};
```

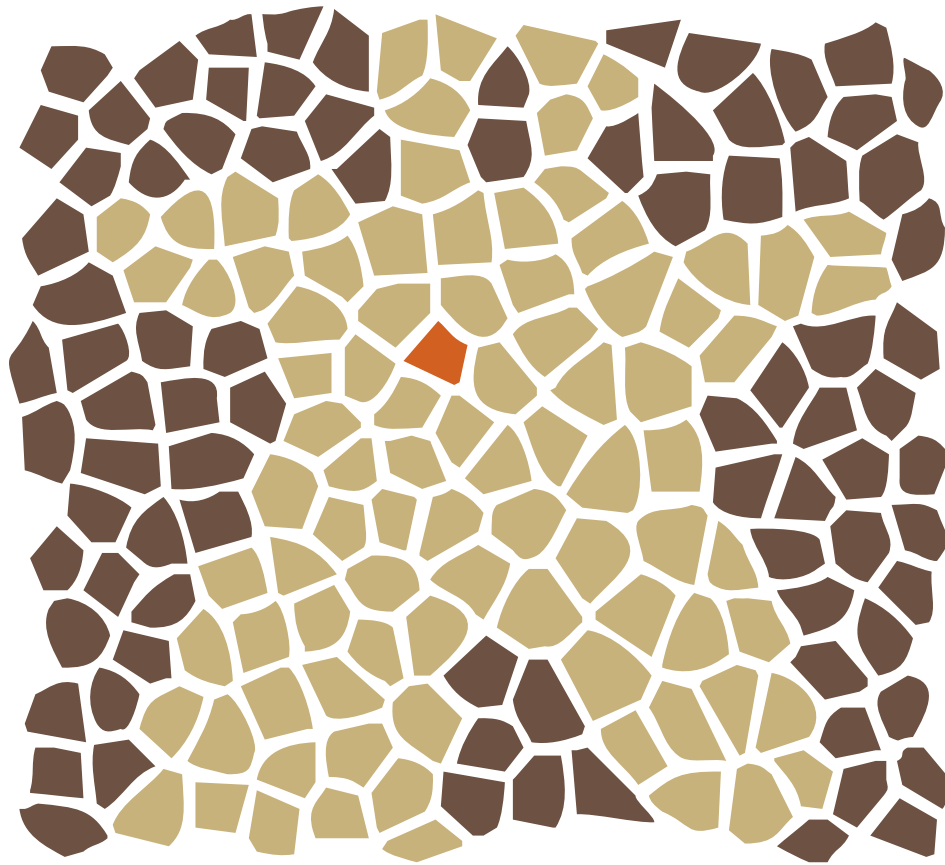
Pregel: PageRank

```
class PageRankVertex : public Vertex<double, void, double> {
public:
    virtual void Compute(MessageIterator* msgs) {
        if (superstep() >= 1) {
            double sum = 0;
            for (; !msgs->Done(); msgs->Next())
                sum += msgs->Value();
            *MutableValue() = 0.15 / NumVertices() + 0.85 * sum;
        }

        if (superstep() < 30) {
            const int64 n = GetOutEdgeIterator().size();
            SendMessageToAllNeighbors(GetValue() / n);
        } else {
            VoteToHalt();
        }
    }
};
```

Pregel: Combiners

```
class MinIntCombiner : public Combiner<int> {  
    virtual void Combine(MessageIterator* msgs) {  
  
        int mindist = INF;  
        for (; !msgs->Done(); msgs->Next())  
            mindist = min(mindist, msgs->Value());  
        Output("combined_source", mindist);  
    }  
  
};
```



A P A C H E
G I R A P H

Giraph Architecture

Master – Application coordinator

Synchronizes supersteps

Assigns partitions to workers before superstep begins

Workers – Computation & messaging

Handle I/O – reading and writing the graph

Computation/messaging of assigned partitions

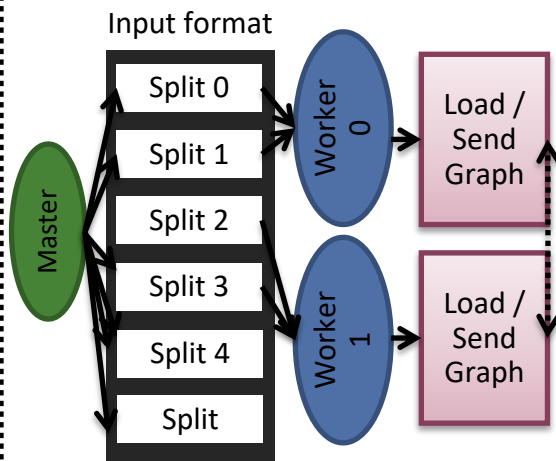
ZooKeeper

Maintains global application state

Giraph Dataflow

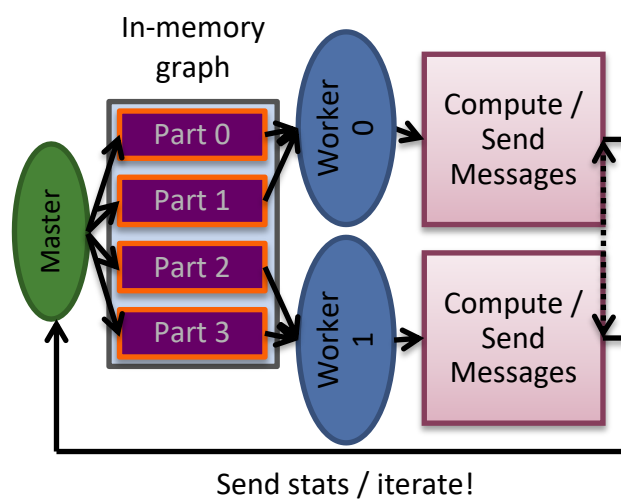
1

Loading the graph



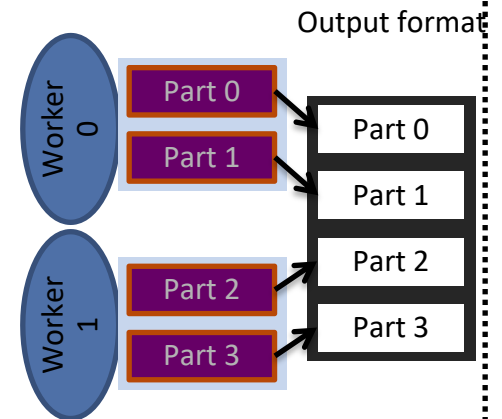
2

Compute/Iterate



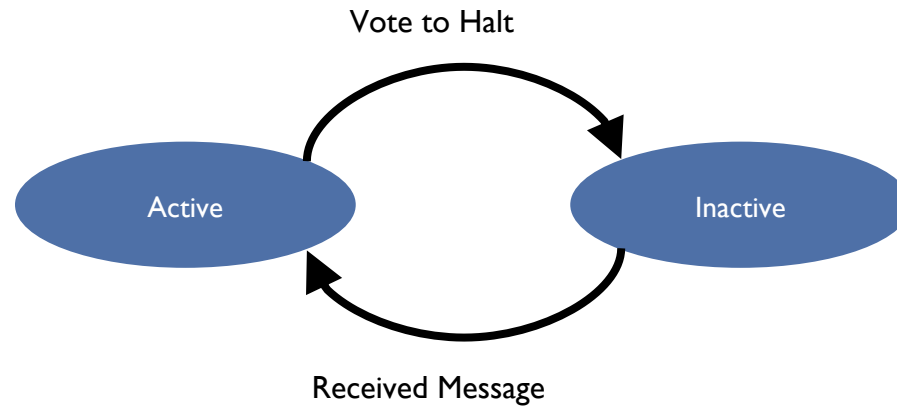
3

Storing the graph

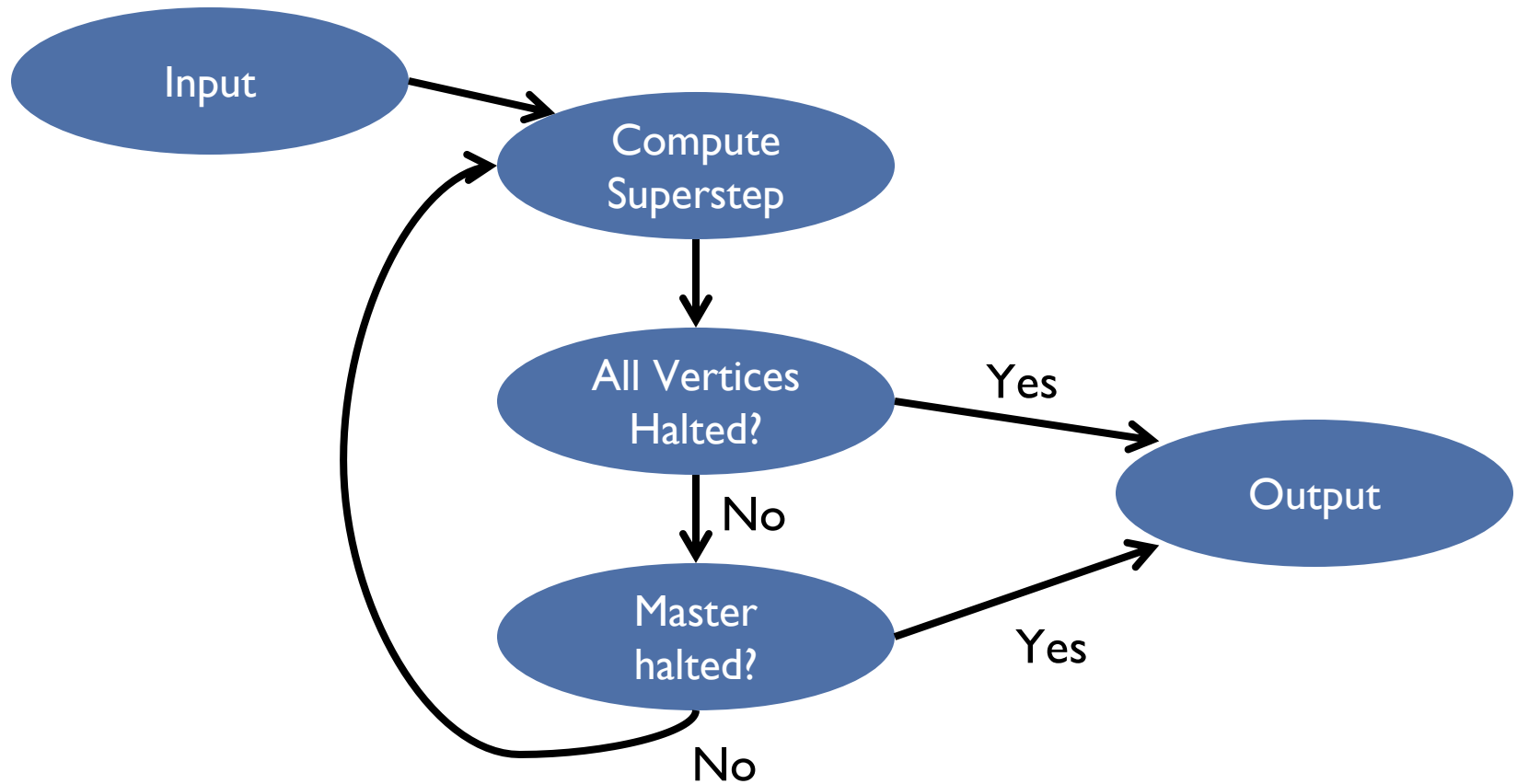


Giraph Lifecycle

Vertex Lifecycle



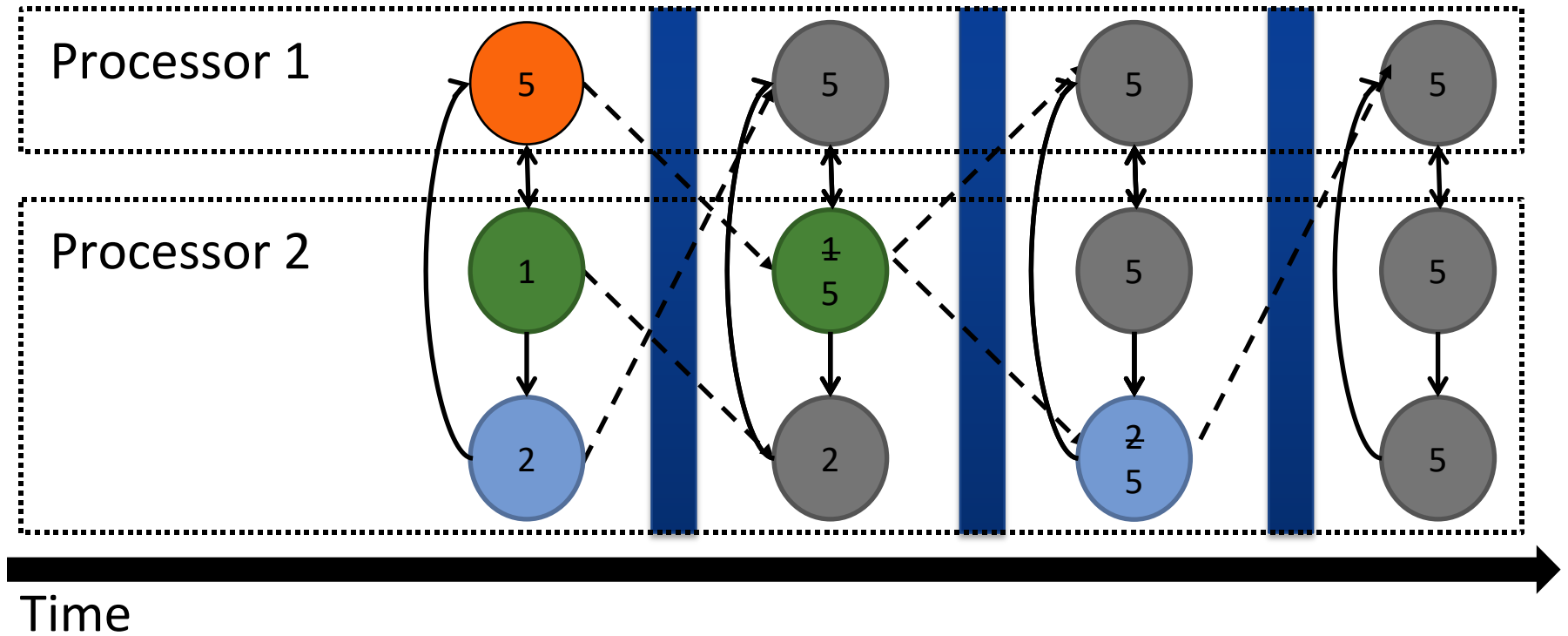
Giraph Lifecycle

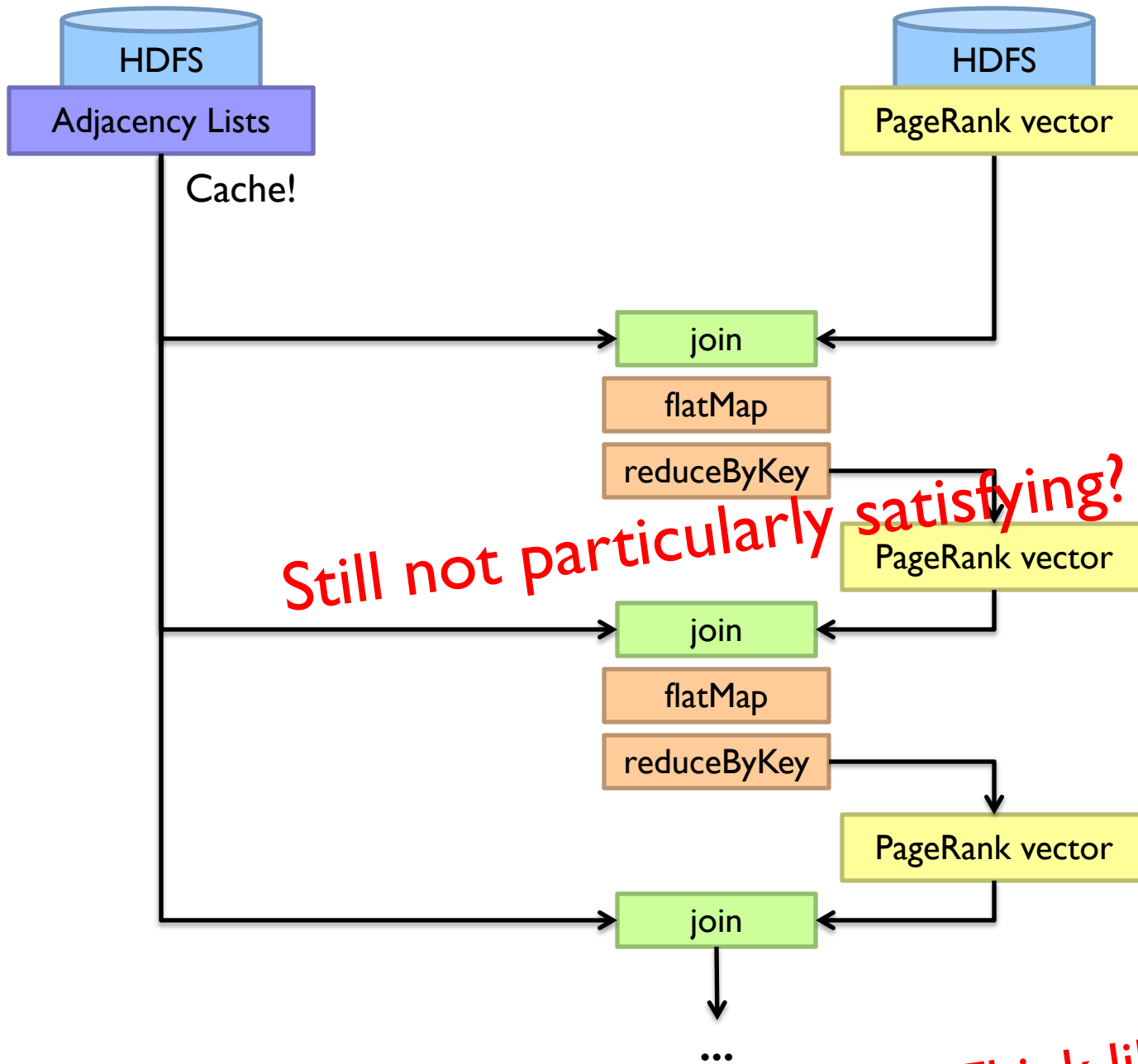


Giraph Example

```
public class MaxComputation extends BasicComputation<IntWritable, IntWritable,
    NullWritable, IntWritable> {
    @Override
    public void compute(Vertex<IntWritable, IntWritable, NullWritable> vertex,
        Iterable<IntWritable> messages) throws IOException
    {
        boolean changed = false;
        for (IntWritable message : messages) {
            if (vertex.getValue().get() < message.get()) {
                vertex.setValue(message);
                changed = true;
            }
        }
        if (getSuperstep() == 0 || changed) {
            sendMessageToAllEdges(vertex, vertex.getValue());
        }
        vertex.voteToHalt();
    }
}
```

Execution Trace

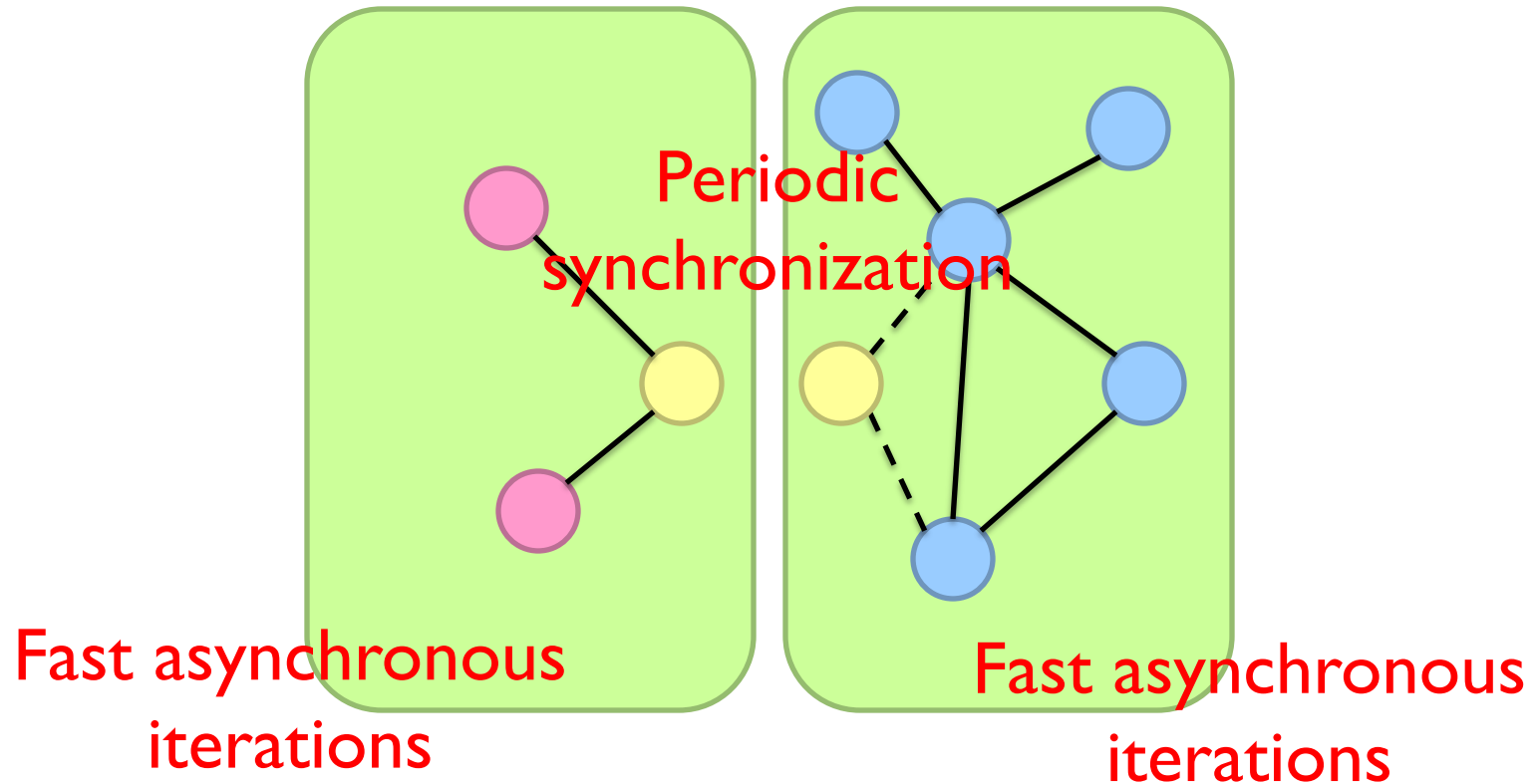




Still not particularly satisfying?

Think like a vertex!

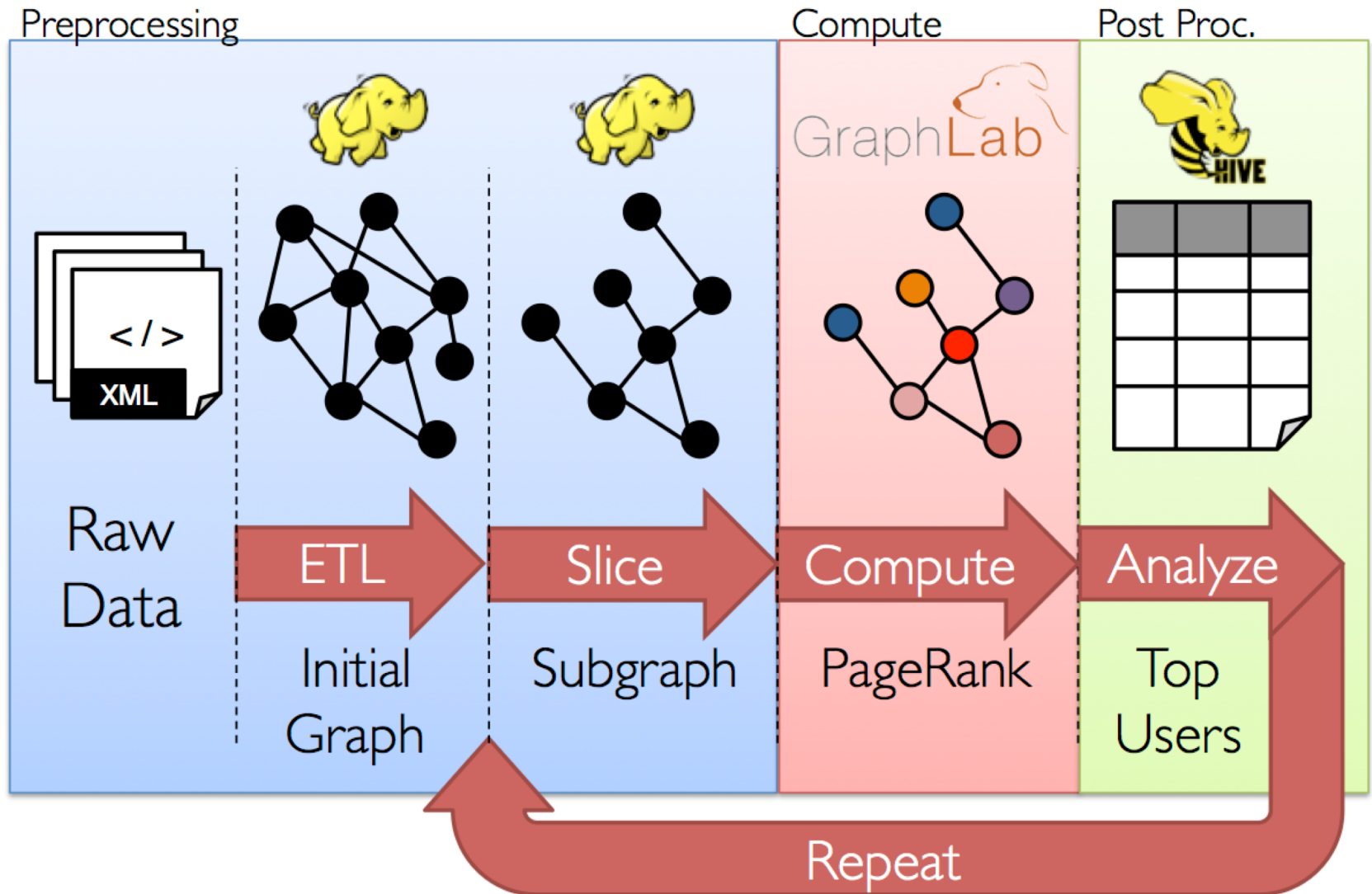
State-of-the-Art Distributed Graph Algorithms



Graph Processing Frameworks



GraphX: Motivation



GraphX = Spark for Graphs

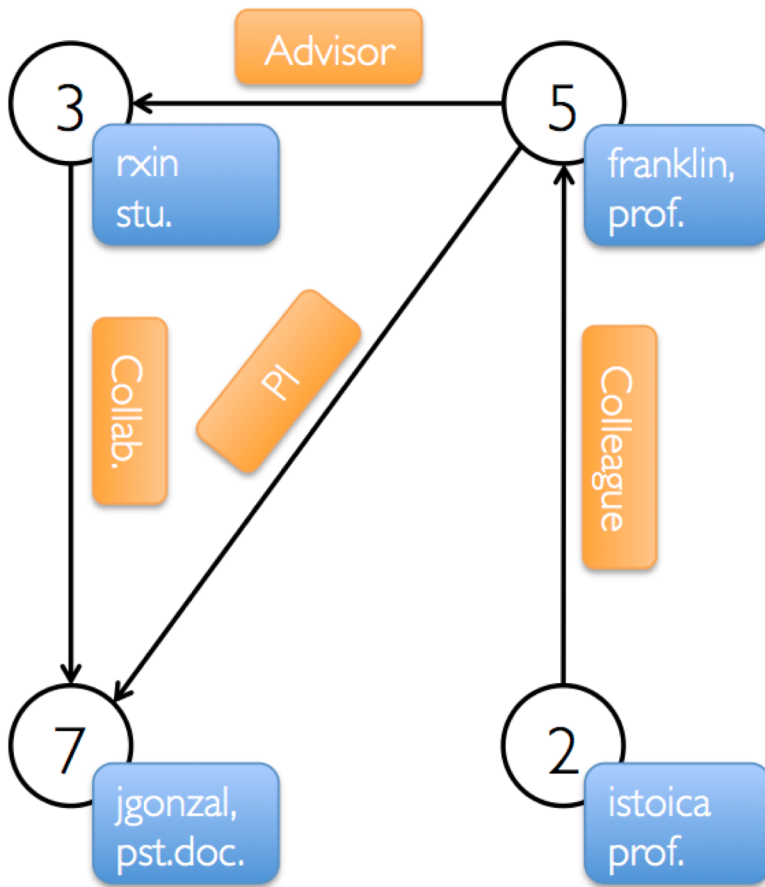
Integration of record-oriented and graph-oriented processing

Extends RDDs to Resilient Distributed Property Graphs

```
class Graph[VD, ED] {  
  val vertices: VertexRDD[VD]  
  val edges: EdgeRDD[ED]  
}
```

Property Graph: Example

Property Graph



Vertex Table

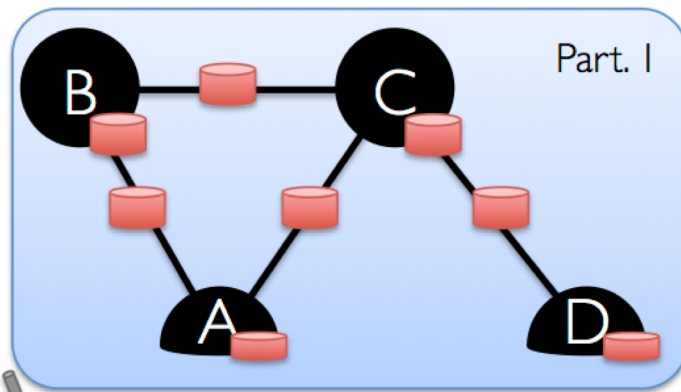
Id	Property (V)
3	(rxin, student)
7	(jgonzal, postdoc)
5	(franklin, professor)
2	(istoica, professor)

Edge Table

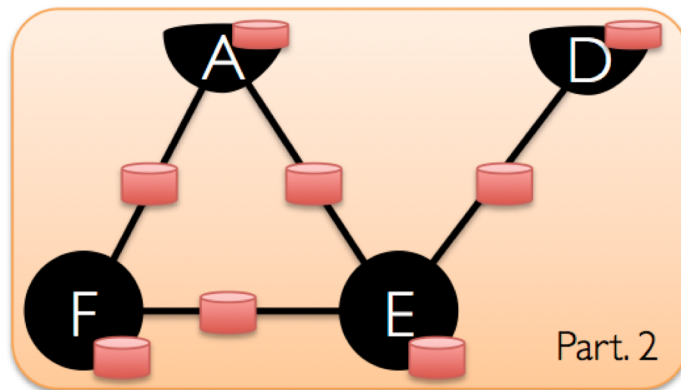
SrcId	DstId	Property (E)
3	7	Collaborator
5	3	Advisor
2	5	Colleague
5	7	PI

Underneath the Covers

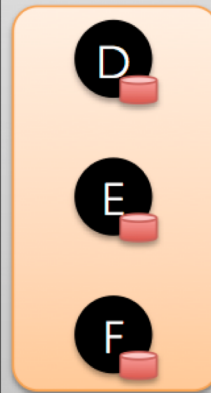
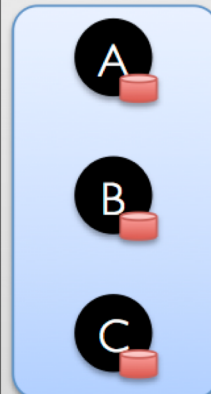
Property Graph



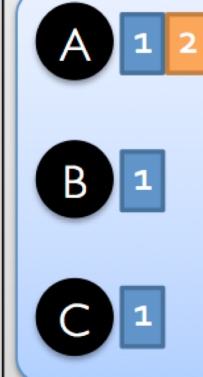
2D Vertex Cut Heuristic



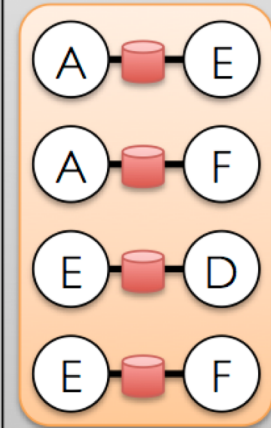
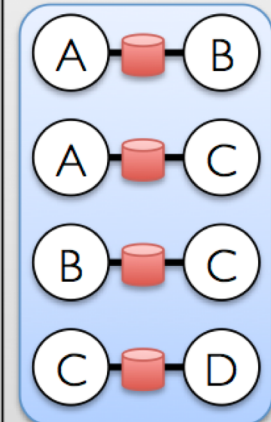
Vertex Table
(RDD)



Routing
Table
(RDD)



Edge Table
(RDD)



GraphX Operators

“collection” view

```
val vertices: VertexRDD[VD]  
val edges: EdgeRDD[ED]  
val triplets: RDD[EdgeTriplet[VD, ED]]
```

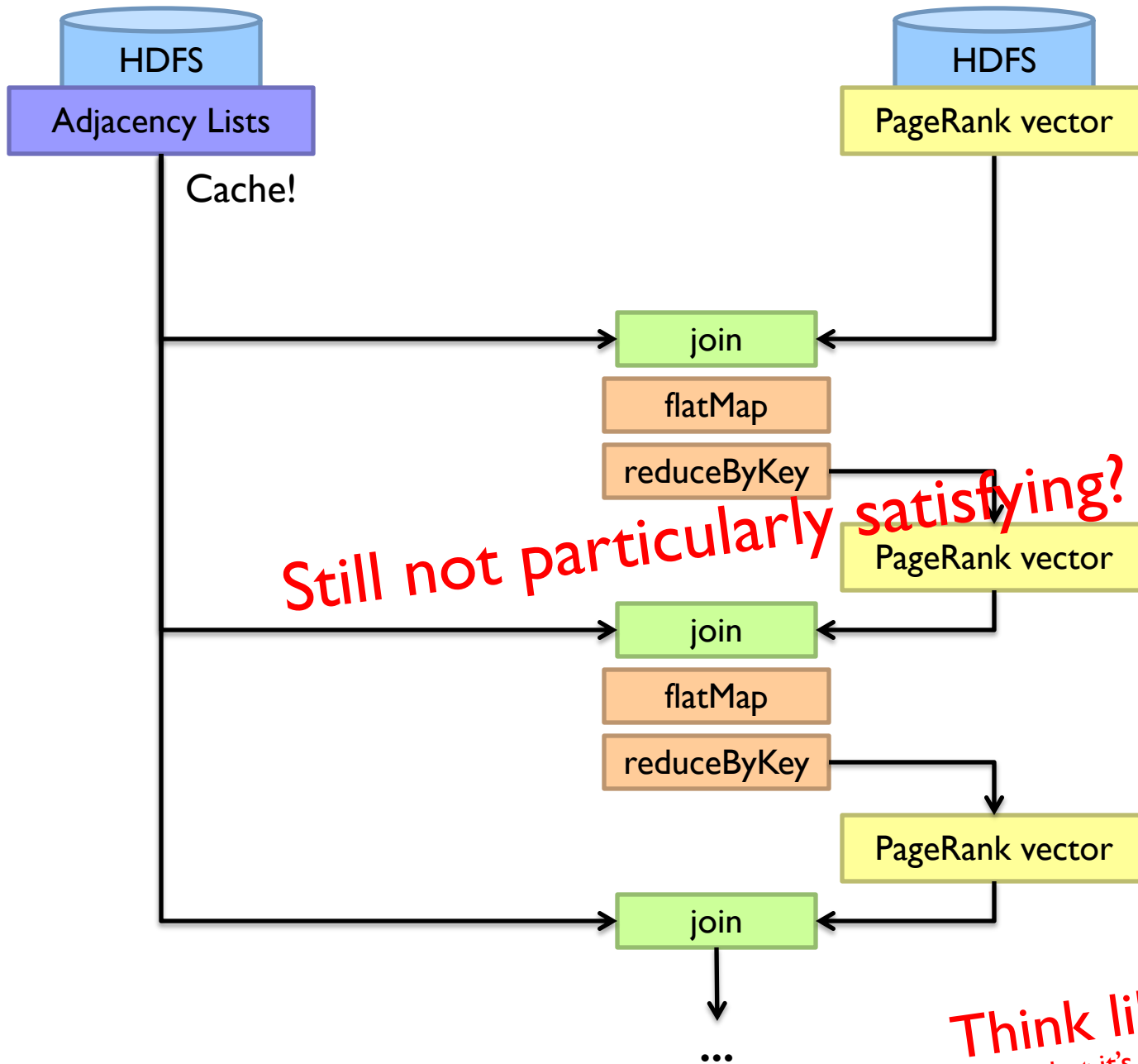
Transform vertices and edges

```
mapVertices  
mapEdges  
mapTriplets
```

Join vertices with external table

Aggregate messages within local neighborhood

Pregel programs



Still not particularly satisfying?

Think like a vertex!
(Yeah, but it's still really all just RDDs)

