



Seminar: Massive-Scale Graph Analysis
Summer Semester 2015

MapReduce for Graph Algorithms

Modeling & Approach

Ankur Sharma
ankur@stud.uni-saarland.de

May 8, 2015



Map-Reduce Framework

Big Data

MapReduce

Graph Algorithms

Graphs

Graph Algorithms

Design pattern

Summary & Conclusion

Big Data: Definition



- ▶ Sources? : Web pages, social network, tweets, messages and many more
- ▶ How Big? : Google's Index is larger than 100 Peta Bytes and consumed over 1M computing hours to build ^a
- ▶ Number of webpages indexed: 48B (Google) & 5B (Bing)^b
- ▶ Big Data Market will account \$50 Billion by 2017^c

^a<https://www.google.com/search/about/insidesearch/howsearchworks/crawling-indexing.html>

^bwww.worldwidewebsize.com

^cwww.wikibon.org



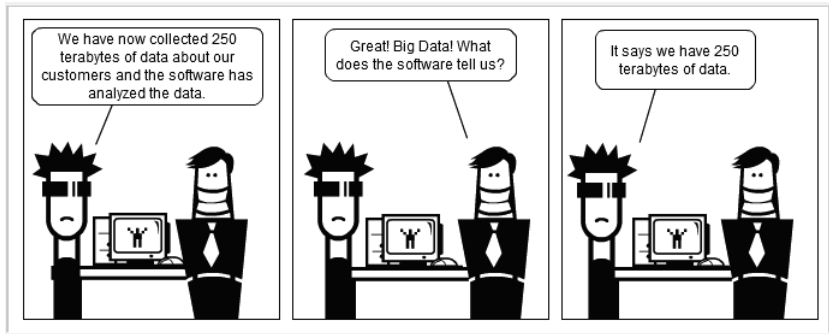


Figure: Comic from Socmedsean¹

¹Source: www.socmedsean.com



- ▶ Quick and efficient analysis of *Big Data* requires extreme parallelism.
- ▶ A lot of systems were developed in late 1990s for parallel data processing such as SRB^a
- ▶ These systems were efficient but lacked nice wrapper for average developer.
- ▶ In 2004, *Jeffrey Dean* and *Sanjay Ghemawat* from *Google* introduced *MapReduce*^b which is a part of company's proprietary infrastructure.
- ▶ Similar open source infrastructure *Hadoop*^c was developed in late 2005 by *Doug Cutting* & *Mike Cafarella*.

^aWan et. al. *The SDSC storage resource broker*, CASCON '98

^bGhemawat et. al. *MapReduce: Simplified data processing on large clusters*, OSDI'14

^cwww.hadoop.apache.org



Introduction

- ▶ MapReduce motivates from LISP's map and reduce operations
 - ▶ $(\text{map square } '(1\ 2\ 3\ 4\ 5)) \rightarrow (1\ 4\ 9\ 16\ 25)$
 - ▶ $(\text{reduce } +\ '(1\ 2\ 3\ 4\ 5)) \rightarrow (15)$
- ▶ Framework provides *automatic parallelism, fault-tolerance, I/O scheduling and monitoring* for data processing



Introduction

- ▶ MapReduce motivates from LISP's map and reduce operations
 - ▶ $(\text{map square } '(1\ 2\ 3\ 4\ 5)) \rightarrow (1\ 4\ 9\ 16\ 25)$
 - ▶ $(\text{reduce } + \ '(1\ 2\ 3\ 4\ 5)) \rightarrow (15)$
- ▶ Framework provides *automatic parallelism, fault-tolerance, I/O scheduling and monitoring* for data processing

MapReduce: Programming Model

- ▶ Input / Output are a set of key value pairs
- ▶ User only needs to provide simple functions `map()`, `reduce()` and `combine()`^a
 - ▶ $\text{map}(in_k, in_v) \rightarrow \text{list}(out_k, imt_v)$
 - ▶ $\text{combine}(in_k, \text{list}(in_v)) \rightarrow \text{list}(out_k, imt_v)$
 - ▶ $\text{reduce}(out_k, \text{list}(imt_v)) \rightarrow \text{list}(out_v)$

^aoptional

MapReduce: Execution Model

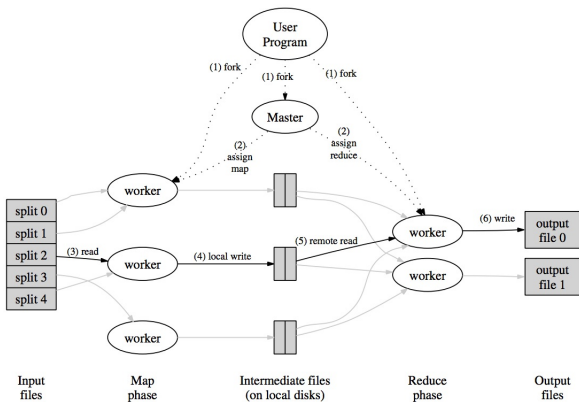


Figure: Overview: Map-Reduce execution model²

²Source: Ghemawat et. al. *MapReduce: Simplified data processing on large clusters*, OSDI'14

MapReduce: Example

Word Count



Map

```
function M(in_k, in_v)
  for each word in in_v:
    emit(word, 1)
end function
```

MapReduce: Example

Word Count



Map

```
function M(in_k, in_v)
  for each word in in_v:
    emit(word, 1)
end function
```

Combine/Reduce

```
function C_R(out_k, list{imt_v})
  emit(word,  $\sum imt_v$ )
end function
```

MapReduce: Example

Word Count

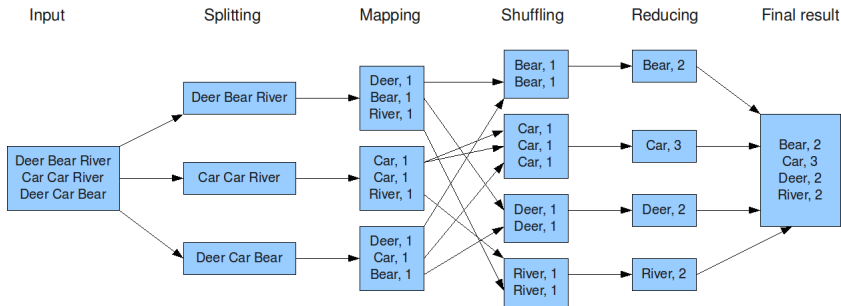


Map

```
function M(in_k, in_v)
  for each word in in_v:
    emit(word, 1)
end function
```

Combine/Reduce

```
function C_R(out_k, list{imt_v})
  emit(word,  $\sum imt_v$ )
end function
```





- ▶ Graph are everywhere: social graphs representing connections or citation graphs representing hierarchy in scientific research
- ▶ Due to massive scale, it is impractical to use conventional techniques for graph storage and in-memory analysis
- ▶ These constraints had driven the development of scalable systems such as distributed file systems like Google File System^a and Hadoop File System^b
- ▶ MapReduce provides a good way to partition and analyze graphs as suggested by Cohen et. al. ^c

^aGhemawat et. al. *The Google File System*, SOSP '03

^bwww.hadoop.apache.org

^cJonathan Cohen, Graph twiddling in a MapReduce world, Computing in Science & Engineering 2009

Augmenting edge with vertex degree



Input: All edges as *Key Value Pair* $\rightarrow \langle -, e_{ij} \rangle$: edge connecting v_i & v_j

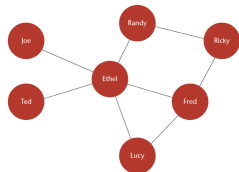
MapReduce Job I

- ▶ $\text{map}(\langle -, e_{ij} \rangle) \rightarrow \text{list}[\langle v_i, e_{ij} \rangle, \langle v_j, e_{ij} \rangle]$
- ▶ $\text{reduce}(v_k, \text{list}[e_{pk}, e_{qk} \dots]) \rightarrow [\langle e_{pk}, \text{deg}(v_k) \rangle, \langle e_{qk}, \text{deg}(v_k) \rangle, \dots]$
 - ▶ where $\text{deg}(v_k)$ = size of the list in the parameter

Input: Output of Job I

MapReduce Job II

- ▶ $\text{map}()$ is an *identity* function: does no modification to input
- ▶ $\text{reduce}(\langle e_{ij}, \text{list}[\text{deg}(v_i), \text{deg}(v_j)] \rangle) \rightarrow \langle e_{ij}, \text{deg}(v_i), \text{deg}(v_j) \rangle$

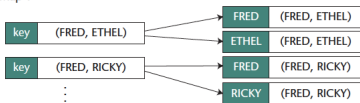


Augmenting edge ... contd

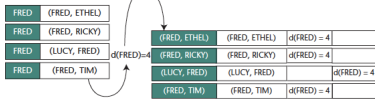
Example



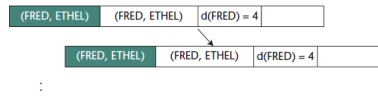
Map 1



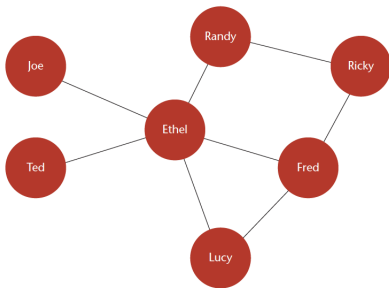
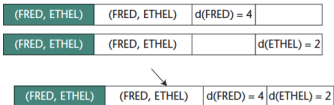
Reduce 1



Map 2 (Identity)



Reduce 2



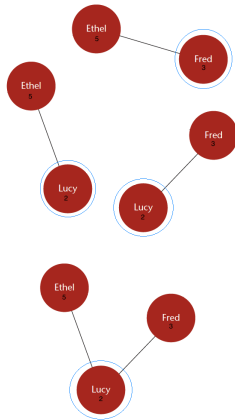
Enumerating Triangles



- ▶ The main idea behind enumerating the triangles is to find triad and an edge joining the open ends together
- ▶ **Input:** Augmented edges with vertex valency

MapReduce Job 1

- ▶ $\text{map}(\langle e_{ij}, \text{deg}(v_i), \text{deg}(v_j) \rangle) \rightarrow$
 - ▶ output: $\text{deg}(v_i) < \text{deg}(v_j) ? \langle v_i, e_{ij} \rangle : \langle v_j, e_{ij} \rangle$
 - ▶ If $\text{deg}(v_i) == \text{deg}(v_j)$, we can break the tie by any consistent method such as using vertex names
- ▶ $\text{reduce}(\langle v_i, \text{list}[\text{edges connecting } v_i] \rangle) \rightarrow$
 - ▶ For each pair of edge connected by v_i , emit $\langle e_1, e_2 \rangle$ with key as $\langle v_p, v_q \rangle$, where v_p and v_q represent the open end of triad
 - ▶ The order of v_p and v_q is decided using some consistent technique like vertex names



Enumerating Triangles ...contd

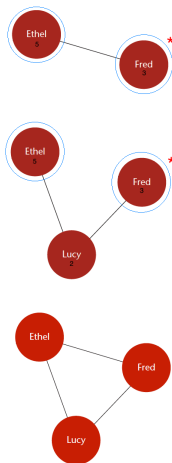


- ▶ **Input:** Augmented edge set and output of Job 1

MapReduce Job 2

- ▶ $\text{map}(\langle R \rangle) \rightarrow$
 - ▶ if R is a record from Augmented Edge Set, change key to $\langle v_i, v_j \rangle$ such that $\text{deg}(v_i) < \text{deg}(v_j)$
 - ▶ if R is a record from MR Job 1 and hence a **triad**, change key to $\langle v_i, v_j \rangle$ such that $\text{deg}(v_i) < \text{deg}(v_j)$, where v_i & v_j are vertices of open ends
- ▶ $\text{reduce}(\langle [v_i, v_j], \text{list}[\text{containing a triad, edge}]^a \rangle) \rightarrow$
 - ▶ If there is a record corresponding to an edge and a triad with same key, emit a triangle $\langle e_{ij}, e_{jk}, e_{ki} \rangle$

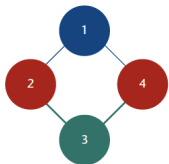
^a may or may not be there



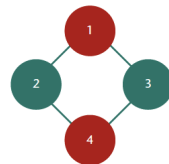
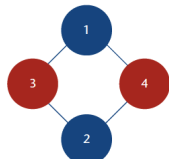
Enumerating Rectangles



Idea: Two triads joining with a common edge can be merged to form a rectangle



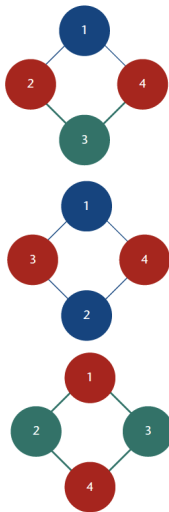
► $4! = 24$ different orderings. Since symmetry group(S_4) on 4 elements has $|S_4| = 8$, we have $\frac{24}{8} = 3$ distinct orderings



Enumerating Rectangles



Idea: Two triads joining with a common edge can be merged to form a rectangle



- ▶ $4! = 24$ different orderings. Since symmetry group(S_4) on 4 elements has $|S_4| = 8$, we have $\frac{24}{8} = 3$ distinct orderings

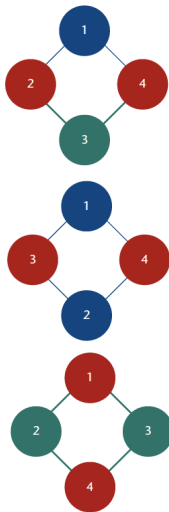
MapReduce Job 1

- ▶ $\text{map}(\langle e_{ij}, \text{deg}(v_i), \text{deg}(v_j) \rangle) \rightarrow$
 - ▶ $[\langle v_i, H, e_{ij} \rangle, \langle v_j, L, e_{ij} \rangle]$
- ▶ $\text{reduce}(\langle v_k, \text{list}[\text{edges incident on } v_k] \rangle) \rightarrow$
 - ▶ emit a triad for L-L & H-L record
- ▶ if we express $\text{deg}(v) = \text{deg}_L(v) + \text{deg}_H(v)$, each vertex appears in $O(\text{deg}_L^2(v))$ low triads and $O(\text{deg}_L(v)\text{deg}_H(v))$ mixed triads

Enumerating Rectangles



Idea: Two triads joining with a common edge can be merged to form a rectangle



- ▶ $4! = 24$ different orderings. Since symmetry group(S_4) on 4 elements has $|S_4| = 8$, we have $\frac{24}{8} = 3$ distinct orderings

MapReduce Job I

- ▶ $\text{map}(\langle e_{ij}, \text{deg}(v_i), \text{deg}(v_j) \rangle) \rightarrow$
 - ▶ $[\langle v_i, H, e_{ij} \rangle, \langle v_j, L, e_{ij} \rangle]$
- ▶ $\text{reduce}(\langle v_k, \text{list}[\text{edges incident on } v_k] \rangle) \rightarrow$
 - ▶ emit a triad for L-L & H-L record
- ▶ if we express $\text{deg}(v) = \text{deg}_L(v) + \text{deg}_H(v)$, each vertex appears in $O(\text{deg}_L^2(v))$ low triads and $O(\text{deg}_L(v)\text{deg}_H(v))$ mixed triads

MapReduce Job II

- ▶ $\text{map}(\langle v_k, \text{list}[\text{triads}] \rangle) \rightarrow$
 - ▶ $\text{list}[\langle \text{key: edge joining open ends, value: triad} \rangle]$
- ▶ $\text{reduce}(\langle e_{ij}, \text{list}[\text{triads open at same end}] \rangle) \rightarrow$
 - ▶ emit rectangle: $[e_{ij}, e_{jk}, e_{kl}, \& e_{li}]$

Finding k -truss



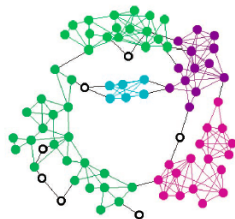
Definition: A k -truss is a sub-graph in which each edge is a part of atleast $k-2$ triangles

MapReduce Job

Input: Augmented edge set & enumerated triangles of the graph

- ▶ $\text{map}(\langle e_{ij}, \text{triangle}_i \rangle) \rightarrow$
 - ▶ if $(e_{ij} \in \text{triangle}_i) : \langle e_{ij}, 1 \rangle$
- ▶ $\text{reduce}(\langle e_{ij}, \text{list}[1, 1, \dots, 1] \rangle) \rightarrow$
 - ▶ if $(\sum \text{list}_i) \geq k : \text{emit } e_{ij}$

All edges emitted by map-reduce job belongs to the k -truss





Clusters are sub-graphs that partitions graph without losing the co-relevant information. **Barycentric Clustering** is a highly scalable approach of finding tightly connected sub-graphs.

- ▶ w_{ij} : weight of edge e_{ij} and $d_i = \sum w_{ij}$: weighted degree of vertex v_i
- ▶ If $x = (x_1, x_2, \dots, x_n)^T$ denotes the position of n vertices, then $x' = Mx$ modifies the vertex position by the average of its old position and position of its neighborhood.
- ▶ The process is repeated and literature³ suggests that 5 repetitions are adequate.
- ▶ If the length of edge > threshold, it can be cut

$$M = \begin{pmatrix} 1/(d_1 + 1) & w_{12}1/(d_1 + 1) & \dots & w_{1n}1/(d_1 + 1) \\ w_{21}1/(d_2 + 1) & 1/(d_2 + 1) & \dots & w_{2n}1/(d_2 + 1) \\ \vdots & \vdots & \ddots & \vdots \\ w_{n1}1/(d_n + 1) & w_{n2}1/(d_n + 1) & \dots & 1/(d_n + 1) \end{pmatrix}$$

³Cohen et. at. Barycentric Graph Clustering, 2008



MapReduce Job 1

- ▶ $\text{map}(\langle e_{ij} \rangle) \rightarrow [\langle v_i, e_{ij} \rangle, \langle v_j, e_{ij} \rangle]$
- ▶ $\text{reduce}(\langle v_i, \text{list}[e_{ij}, e_{ik}, \dots] \rangle) \rightarrow [\langle e_{ij}, P(v_i) \rangle, \langle e_{ik}, P(v_i) \rangle, \dots]$, where $P(v)$ denotes the position vector of vertex v .



MapReduce Job I

- ▶ $\text{map}(\langle e_{ij} \rangle) \rightarrow [\langle v_i, e_{ij} \rangle, \langle v_j, e_{ij} \rangle]$
- ▶ $\text{reduce}(\langle v_i, \text{list}[e_{ij}, e_{ik}, \dots] \rangle) \rightarrow [\langle e_{ij}, P(v_i) \rangle, \langle e_{ik}, P(v_i) \rangle, \dots]$, where $P(v)$ denotes the position vector of vertex v .

MapReduce Job II

- ▶ $\text{map} : (\text{identity})$
- ▶ $\text{reduce}(\langle e_{ij}, \text{list}[P(v_i), P(v_j)] \rangle) \rightarrow \langle e_{ij}, P(v_i), P(v_j) \rangle$



MapReduce Job I

- ▶ $\text{map}(\langle e_{ij} \rangle) \rightarrow [\langle v_i, e_{ij} \rangle, \langle v_j, e_{ij} \rangle]$
- ▶ $\text{reduce}(\langle v_i, \text{list}[e_{ij}, e_{ik}, \dots] \rangle) \rightarrow [\langle e_{ij}, P(v_i) \rangle, \langle e_{ik}, P(v_i) \rangle, \dots]$, where $P(v)$ denotes the position vector of vertex v .

MapReduce Job II

- ▶ $\text{map} : (\text{identity})$
- ▶ $\text{reduce}(\langle e_{ij}, \text{list}[P(v_i), P(v_j)] \rangle) \rightarrow \langle e_{ij}, P(v_i), P(v_j) \rangle$

Now we execute 5 iterations of next two map reduce jobs

MapReduce Job III

- ▶ $\text{map}(\langle e_{ij}, P(v_i), P(v_j) \rangle) \rightarrow [\langle v_i, e_{ij}, P(v_i), P(v_j) \rangle, \langle v_j, e_{ij}, P(v_i), P(v_j) \rangle]$
- ▶ $\text{reduce}(\langle v_i, \text{list}[\langle e_{ij}, P(v_i), P(v_j) \rangle, \langle e_{ik}, P(v_i), P(v_k) \rangle, \dots] \rangle) \rightarrow$ calculates modified $P'(v_i) = M \cdot P(v_i)$ and emits $[\langle e_{ij}, P'(v_i), P(v_j) \rangle, \langle e_{ik}, P'(v_i), P(v_k) \rangle, \dots]$



MapReduce Job IV

- ▶ map: (identity)
- ▶ reduce($\langle e_{ij}, \text{list}[\langle e_{ij}, P'(v_i), P'(v_j) \rangle, \langle e_{ij}, P(v_i), P'(v_j) \rangle] \rangle$) $\rightarrow \langle e_{ij}, P'(v_i), P'(v_j) \rangle$



MapReduce Job IV

- ▶ map: (identity)
- ▶ reduce($\langle e_{ij}, \text{list}[\langle e_{ij}, P'(v_i), P(v_j) \rangle, \langle e_{ij}, P(v_i), P'(v_j) \rangle] \rangle) \rightarrow \langle e_{ij}, P'(v_i), P'(v_j) \rangle$

MapReduce Job V

- ▶ map($e_{ij}, \text{list}[P_1, P_2, \dots, P_r]$) $\rightarrow [\langle v_i, \text{AVG}(P_i), \text{AVG}(P_j) \rangle, \langle v_j, \text{AVG}(P_j) \rangle]$
- ▶ reduce($\langle v_i, \text{list}[\langle e_{ij}, P_j \rangle, \langle e_{ik}, P_k \rangle \dots] \rangle) \rightarrow [\langle e_{ij}, \dots, \langle e_{ik}, \sum P_j, P_k, \dots \rangle \dots], \sum P_j, P_k, \dots, a_{ij} = \sum P_j, P_k, \dots$ represents the weight of edge e_{ij} .



MapReduce Job IV

- ▶ map: (identity)
- ▶ reduce($\langle e_{ij}, \text{list}[\langle e_{ij}, P'(v_i), P(v_j) \rangle, \langle e_{ij}, P(v_i), P'(v_j) \rangle] \rangle \rightarrow \langle e_{ij}, P'(v_i), P'(v_j) \rangle$)

MapReduce Job V

- ▶ map($e_{ij}, \text{list}[P_1, P_2, \dots, P_r]$) $\rightarrow \langle v_i, \text{AVG}(P_i), \text{AVG}(P_j) \rangle, \langle v_j, \text{AVG}(P_j) \rangle$
- ▶ reduce($\langle v_i, \text{list}[\langle e_{ij}, P_j \rangle, \langle e_{ik}, P_k \rangle \dots] \rangle \rightarrow \langle e_{ij}, \dots, \langle e_{ik}, \sum P_j, P_k, \dots \rangle \dots \rangle, \sum P_j, P_k, \dots, a_{ij} = \sum P_j, P_k, \dots$ represents the weight of edge e_{ij} .

MapReduce Job VI

- ▶ map : (identity)
- ▶ reduce($\langle e_{ij}, \text{list}[\text{weights emitted by previous map reduce job}] \rangle$) : computes \bar{a}_{ij} , if $\bar{a}_{ij} > a_{ij}$, emit the edge

$$\bar{a}_{ij} = \frac{\sum_{k \in N_i} a_{jk} + \sum_{k \in N_j} a_{ik} - a_{ij}}{|N_i| + |N_j| - 1}$$

- ▶ where N_k represents the set of vertices in the neighborhood of v_k

Algorithm using MapReduce

- ▶ **Step 1** → consists of two MapReduce Jobs focussing on translating the vertex-vertex information into zone zone information
 - ▶ **MapReduce 1** assigns edges with zones, producing 1 record for each vertex in edge and then translating this information to assign 1 zone to both the vertices
 - ▶ **MapReduce 2** merges the result of Job 1. If distinct zones are connected by an edge, then the zone with lower zone_id captures the other
- ▶ **Step 2** → consists of one MapReduce Job
 - ▶ **MapReduce 1** It record any updates to zone information for each vertex
- ▶ Step 1 and 2 alternate until step 1 produces no new record, which marks the completion of zone assignment

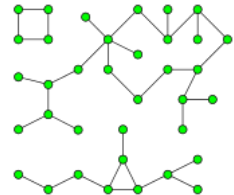


Figure: Graph with three components



Lin et. al^a suggested following design patterns for efficiently implementing graph algorithms using MapReduce.

- ▶ **Local aggregation** done using combiners reduce the amount intermediate data. It is effective only in case of high duplication factor of intermediate key-value pair.
- ▶ **In-Mapper Combining** performs better than than regular combiner. Combiner often have to create and destroy objects and perform object serialization and de-serialization which is often costly.
- ▶ **Schimmy Pattern** divides the graph into n partitions each sorted by *vertex id*. Number of reducers are then set to n and it is guaranteed that reducer R_i processes same *vertex ids* as in partition P_i in sorted manner. Hence reducing the cost of shuffle and sort.
- ▶ **Range Partitioning** proposes to maximize intra-block connectivity and minimize inter-block connectivity. Dividing a graph into blocks with such a property helps reducing the amount of communication over network.

^aJimmy Lin, Michael Schatz, Design patterns for efficient graph algorithms in MapReduce, MLG'10

Design Patterns ... contd

Performance Comparison

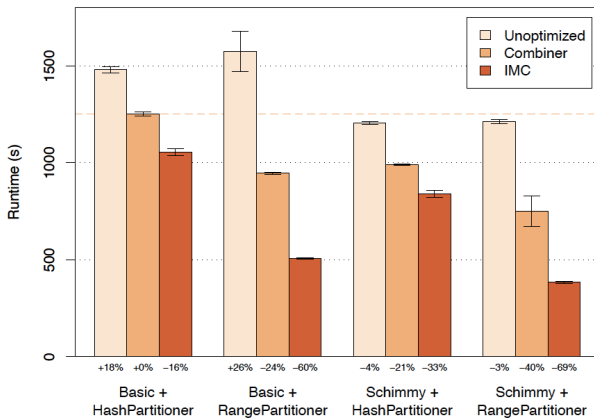


Figure: Performance comparison for PageRank using suggested strategies.⁴

⁴Jimmy Lin, Michael Schatz, Design patterns for efficient graph algorithms in MapReduce, MLG'10



- ▶ In general, it is a good contribution to attract MapReduce community to develop distributed systems for graph processing.
- ▶ Author focused on breaking down smaller problems in map-reduce jobs and merging them together to solve bigger problem.
- ▶ The contribution provides no information about the performance of the implementation. No comparison to standard approach.
- ▶ Most of the algorithms need complete re-transformation for implementation using MR framework. Some may become easy and some like finding component becomes complex.
- ▶ Lin et. al. introduced some really good strategies for efficient MapReduce implementation which can be used even for non-graph algorithms.

Thank You
Questions?