

Haloop

Efficient Iterative Data Processing on Large Clusters

University of Washington, 2010

Yingyi Bu / Bill Howe / Magdalena Balazinska / Michael D. Ernst



1. An opensource implemented base on Mapreduce
2. For performing large-scale data processing using “commodity computer” clusters
3. Can scale to thousands of nodes in a fault-tolerant manner

※ Because of version gap, There is some difference in Hadoop what you knew.



1. limit of Hadoop

- Not Support Iterative Data Processing,
(like PageRank, K-means)

2. Related work

Framework	Mahout	Twister	Pregel
Concept	Machine learning libraries on top of Hadoop.	Stream-based MapReduce framework	Processing graph data
Iteration	Not Support	Support	Not Support
Characteristic	Can use the iteration by injecting from user	Sensitive to failures	-



Haloop

For distribution processing big data **focused on cached data**



1. Architecture

2. Program Model

- Formula
- API

3. Loop Aware Task Scheduling

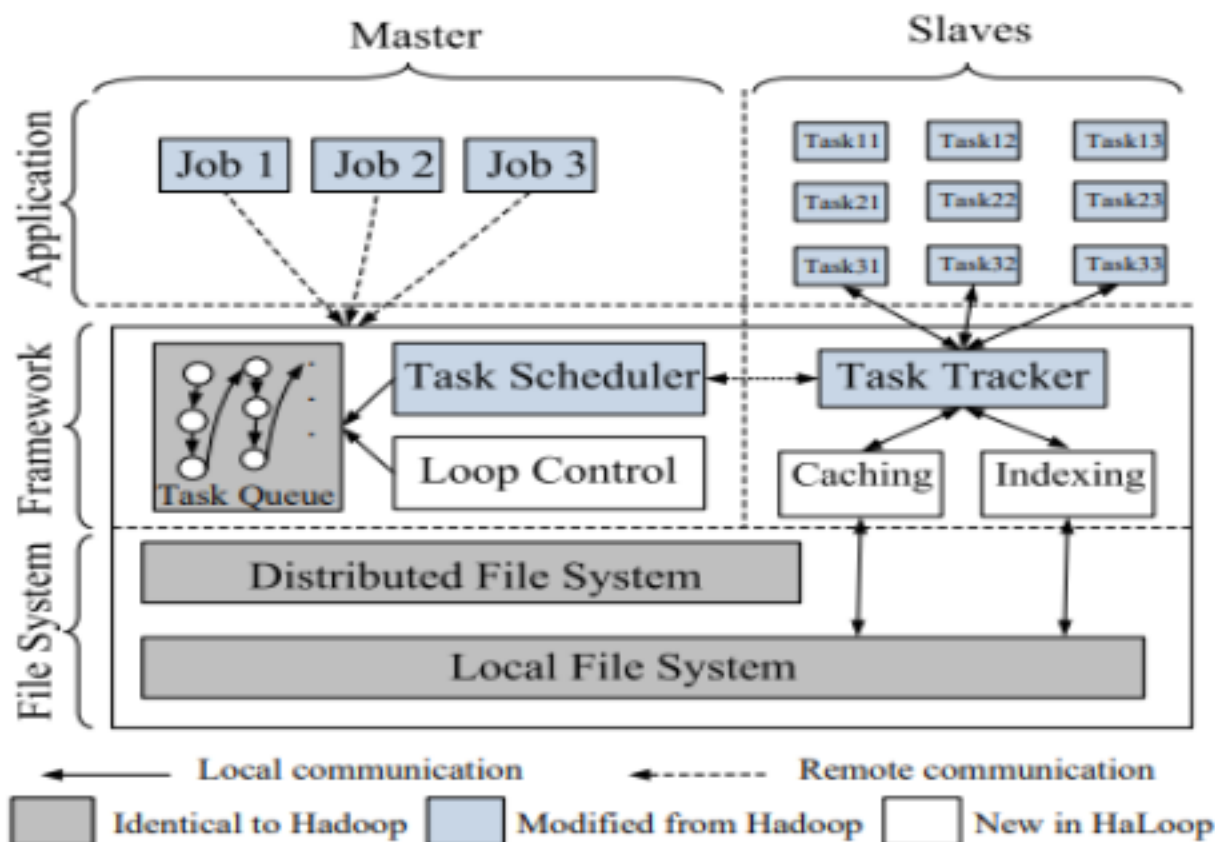
- Inter Iteration Locality
- Scheduling Algorithm

4. Caching & Indexing

- Reducer Input Cache
- Reducer Output Cache
- Mapper Input Cache
- Cache Reloading

5. Experiment Evaluation

6. Conclusion



Based On : Hadoop Architecture

Modified : Task Scheduler / Task Tracker

New : Loop Control / Caching / Indexing



1. Suitable algorithms

Algorithms what using invariant data in each iterations

Because Hadoop get better performance when using caching and indexing.

Ex) Page Rank, K-means, Descendant query



1. Formula

$$R_{i+1} = R_0 \cup (R_i \bowtie L)$$

2. Terminate Conditions

- specified maximum number of iteration
- convergence condition
 - a) Default Fixpoint ($R_{i+1} = R_i$)
 - b) Approximate Fixpoint ($R_{i+1} - R_i \leq \text{Specified Threshold}$)



3. API

: Haloop provides an efficient foundation API for Iteration

1) Loop Body ; constructed a multistep MapReduce Job

a) AddMap

b) AddReduce

※ when constructing Loop Body, you must indicate that order of MapReduce

2) Termination Condition

a) Set Fixed Point Threshold

b) Result Distance

c) Set Max Num Of Iterations

3) Control Input

a) Set Iteration Input

b) Add Step Input

c) Add Invariant Table

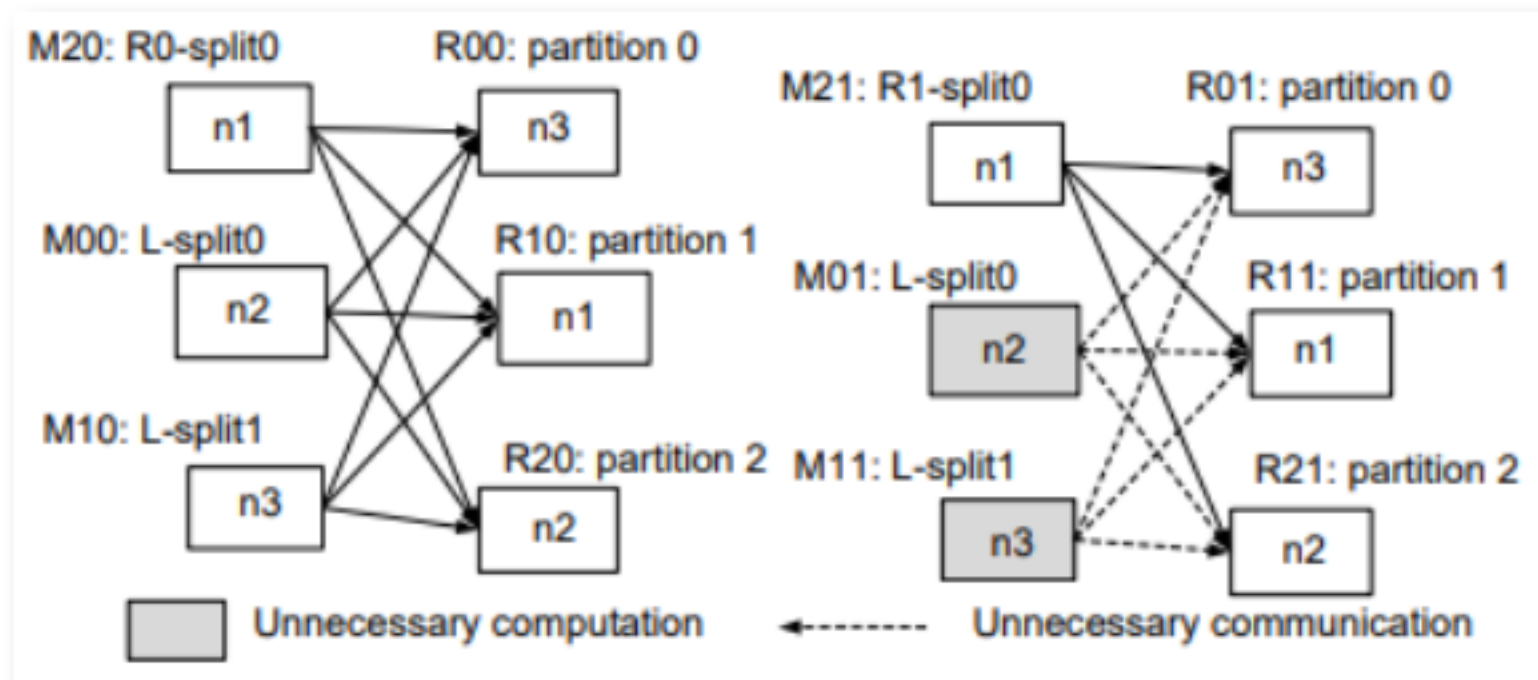


1. Goal of Haloop's Scheduler

To place on the **same physical machines**
those map and reduce tasks that occur in **different iterations**
but **access the same data**



2. Inter-Iteration Locality



※ Example of Inter-Iteration Locality

Another Iteration / Same Data / Same Node



3. Scheduling Algorithm

Task Scheduling

Input: Node node

// The current iteration's schedule; initially empty

Global variable: Map(Node, List(Partition)) current

// The previous iteration's schedule

Global variable: Map(Node, List(Partition)) previous

```
1: if iteration == 0 then
2:   Partition part = hadoopSchedule(node);
3:   current.get(node).add(part);
4: else
5:   if node.hasFullLoad() then
6:     Node substitution = findNearestIdleNode(node);
7:     previous.get(substitution).addAll(previous.remove(node));
8:     return;
9:   end if
10:  if previous.get(node).size() > 0 then
11:    Partition part = previous.get(node).get(0);
12:    schedule(part, node);
13:    current.get(node).add(part);
14:    previous.remove(part);
15:  end if
16: end if
```



1. Overview

Thanks to the Inter-iteration locality, **Only one** physical node is needed for **processing Loop invariant data**

Types of caches

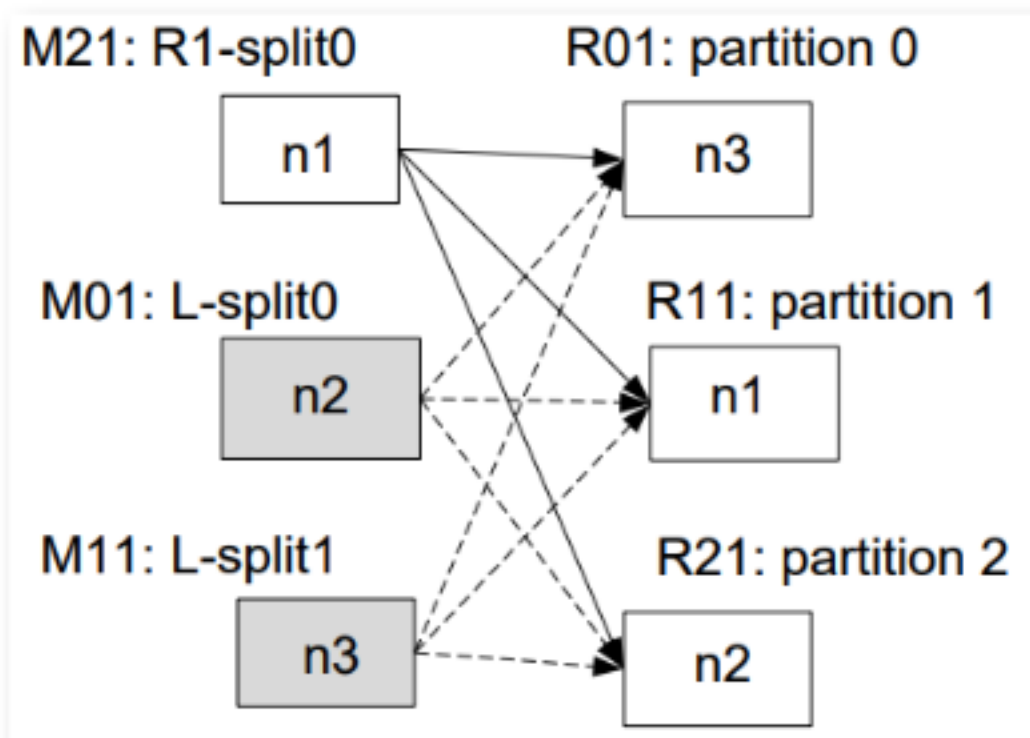
- Reducer Input Cache
- Reducer Output Cache
- Mapper Input Cache

To reduce I/O cost, Haloop caches those data partitions

To accelerate processing, Haloop indexes cached data



2. Reducer Input Cache

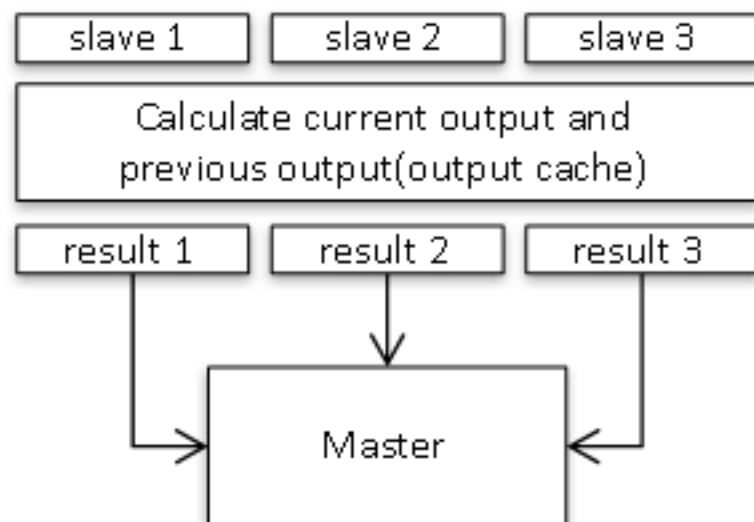


: There is no necessary to process M01 and M02 mapper, because reducers already have output of L-split0 and L-split1



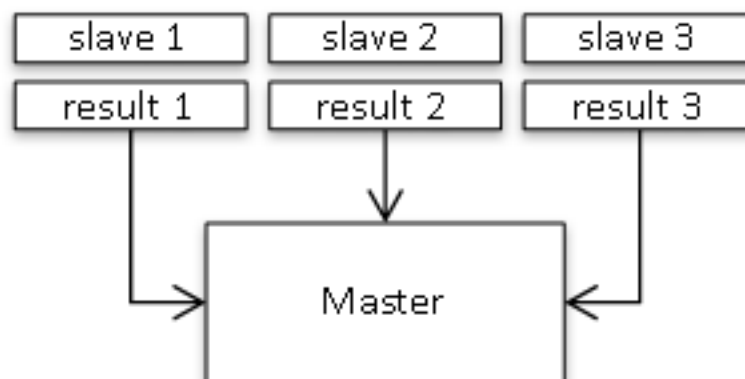
3. Reducer Output Cache

Using



1. compute those result
2. compare with computed result and threshold
3. decide to repeat next iteration

Non-Using



1. access previous result of iteration
2. group each result by key
3. Calculate current output and previous output
4. compute those result
5. compare with computed result and threshold
6. decide to repeat next iteration



4. Mapper Input Cache

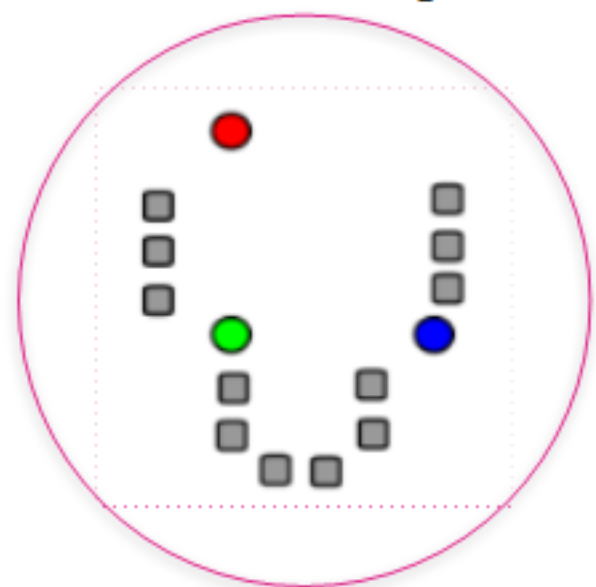
Because the **first iteration** get data from HDFS, local file system
and **cache** them

So there is **no necessary to reload** data next iterations

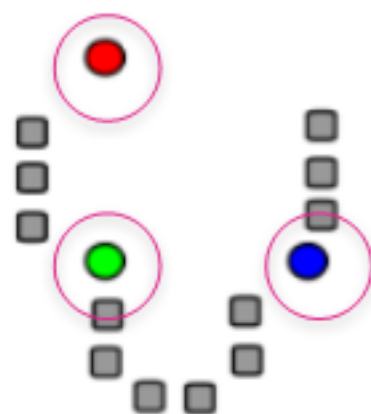


4. Mapper Input Cache

Non-Using



Using



- Datum point for clustering



- Data what will be clustered



- Data to read from HDFS or local file system



4. Cache Reloading

: if node is out of order or full, Cache reloading is necessary

1) Reducer Input Cache Reloading

: Duplicate Initial Mapper output result

2) Reducer output Cache & Mapper input Cache Reloading

: Reload from HDFS



- 1) Compared preformance of iterative data analysis on HaLoop and Hadoop
 - a) reducer input cache
 - b) reducer output cache
 - c) mapper input cache

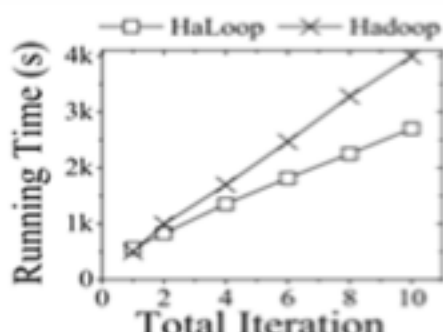
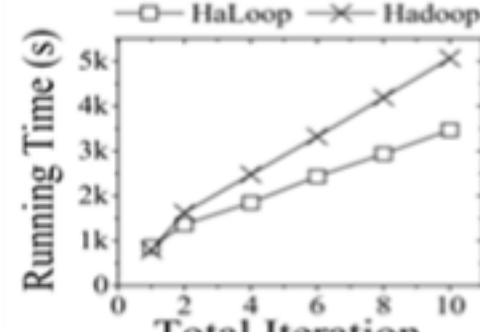
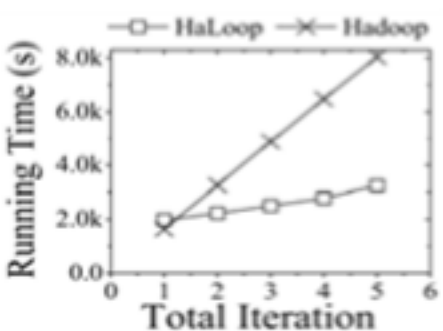
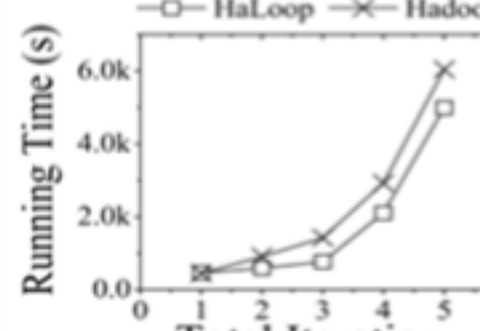
2) Using Data Set

Name	Data capacity	Detail
Livejournal	18GB	social network data
Triples	120GB	semantic web data
Freebase	12GB	concept linkage graph

- 3) Environment setting
 - virtual machine nodes: 50
 - slave nodes(EC2) : 90
 - master node : 1

Evaluation of Reducer Input Cache



Over all	
Page Rank	
 <p>Running Time (s)</p> <p>Total Iteration</p> <p>—□— HaLoop —×— Hadoop</p> <p>Livejournal Dataset, 50 nodes</p>	 <p>Running Time (s)</p> <p>Total Iteration</p> <p>—□— HaLoop —×— Hadoop</p> <p>Freebase Dataset, 90 nodes</p>
Descendant query	
 <p>Running Time (s)</p> <p>Total Iteration</p> <p>—□— HaLoop —×— Hadoop</p> <p>Triples Dataset, 90 nodes</p>	 <p>Running Time (s)</p> <p>Total Iteration</p> <p>—□— HaLoop —×— Hadoop</p> <p>Livejournal Dataset, 50 nodes</p>

Evaluation of Reducer Input Cache

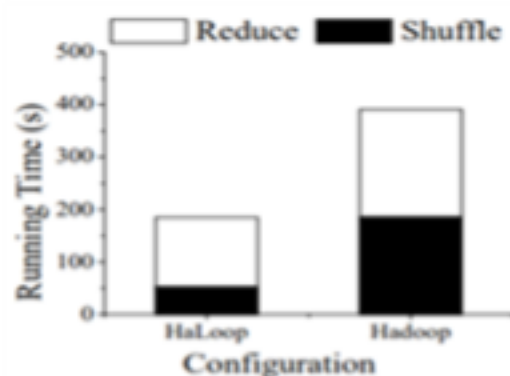


Join Step	
Page Rank	
<p>Running Time (s)</p> <p>Iteration</p> <p>Legend: HaLoop (square), Hadoop (cross)</p>	<p>Running Time (s)</p> <p>Iteration</p> <p>Legend: HaLoop (square), Hadoop (cross)</p>
Livejournal Dataset, 50 nodes	Freebase Dataset, 90 nodes
Descendant query	
<p>Running Time (s)</p> <p>Iteration</p> <p>Legend: HaLoop (square), Hadoop (cross)</p>	<p>Running Time (s)</p> <p>Iteration</p> <p>Legend: HaLoop (square), Hadoop (cross)</p>
Triples Dataset, 90 nodes	Livejournal Dataset, 50 nodes

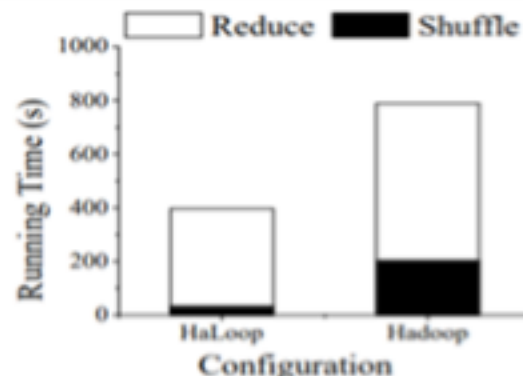


Cost Distribution for Join Step : iteration 3

Page Rank

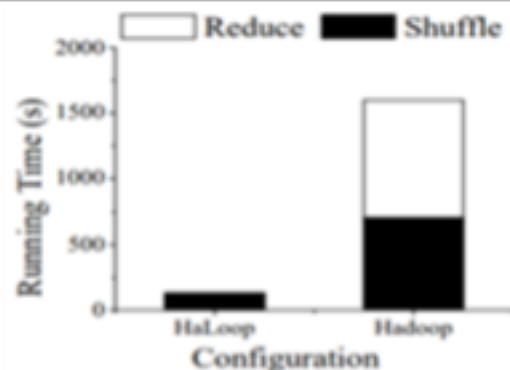


Livejournal Dataset, 50 nodes

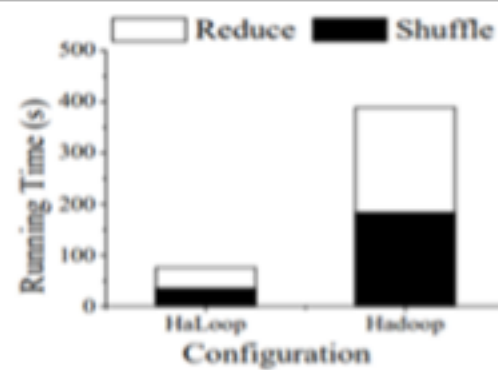


Freebase Dataset, 90 nodes

Descendant query



Triples Dataset, 90 nodes



Livejournal Dataset, 50 nodes

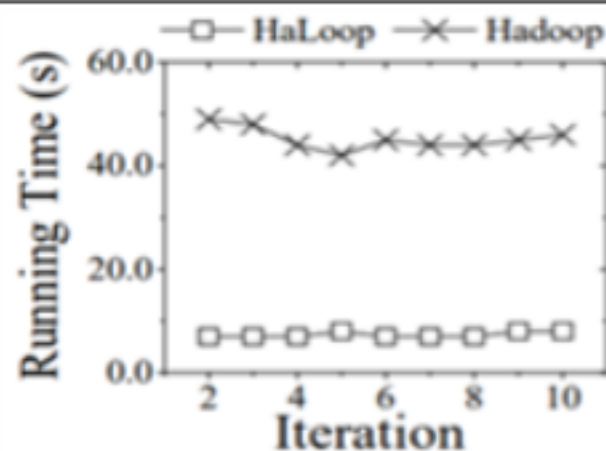


I/O in Shuffle Phase of Join Step	
Page Rank	
<p>Shuffled Data (Bytes)</p> <p>Iteration</p> <p>—○— HaLoop —×— Hadoop</p>	<p>Shuffled Data (Bytes)</p> <p>Iteration</p> <p>—○— HaLoop —×— Hadoop</p>
Livejournal Dataset, 50 nodes	Freebase Dataset, 90 nodes
Descendant query	
<p>Shuffled Data (Bytes)</p> <p>Iteration</p> <p>—○— HaLoop —×— Hadoop</p>	<p>Shuffled Data (Bytes)</p> <p>Iteration</p> <p>—○— HaLoop —×— Hadoop</p>
Triples Dataset, 90 nodes	Livejournal Dataset, 50 nodes

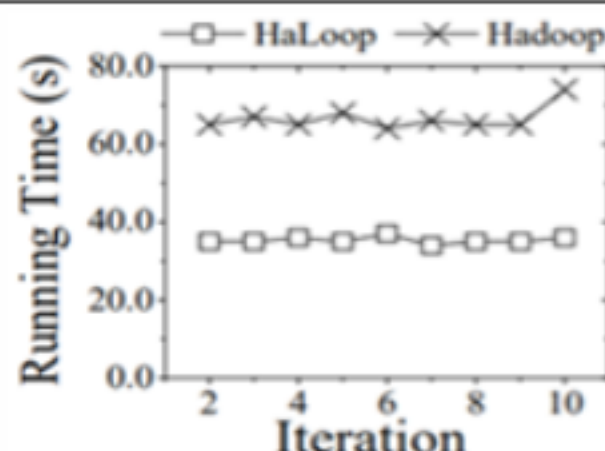


Compared the cost of fixpoint evaluation

Page Rank



Livejournal Dataset, 50 nodes



Freebase Dataset, 90 nodes

Analysis

- HaLoop reduces the cost of this step to 40%
- Hadoop need an extra MapReduce job



1. Using Application

k-means clustering algorithm

2. Environment Setting

8-node physical machine cluster

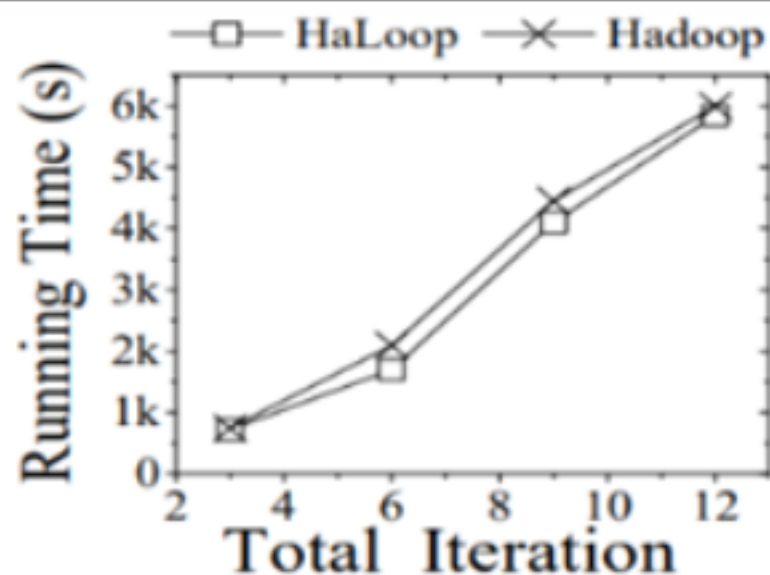
3. Data Set

Name	Data capacity
cosmo-dark	46GB
cosmo-gas	54GB

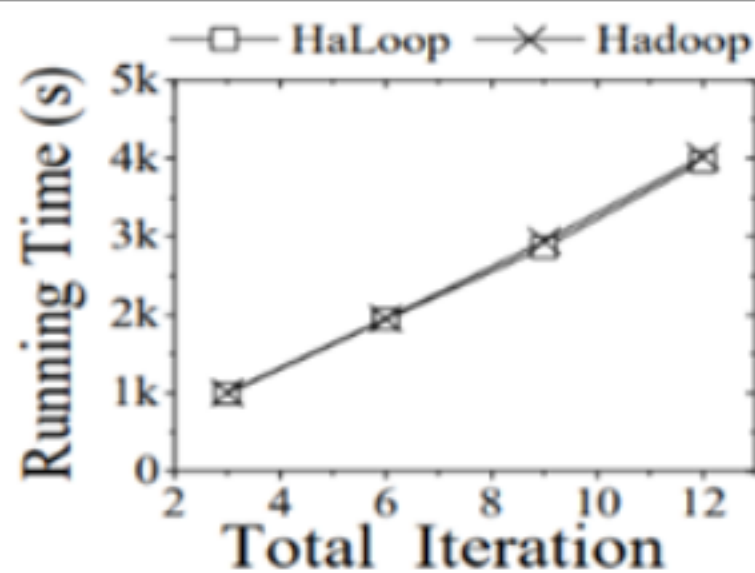


Running Time

k-means



Cosmo-dark



Cosmo-gas



Haloop is

New Programming Model and Architecture for Iterative Programs

have advantages

- a) Data can reuse across iterations by Loop-Aware Task Scheduler
- b) Reduce the I/O cost in next iterations by Caching loop-invariant data
- c) Reduce useless Mapreduce step for checking fixpoint or convergence by Caching Reducer's output

outperforms in Iterative Programs

- a) Reduce query runtimes by 1.85
- b) Shuffles only 4% of the data between mappers and reducers



Q&A