# Native-Task Performance Test Report

## Intel Software

Wang, Huafeng, Huafeng.wang@intel.com

Zhong, Xiang, xiang.zhong@intel.com

# 1. Background

# 2. Related Work

# 3. Preliminary Experiments

## 3.1 Experimental Environment

| Workbench | Peculiarity |
|---|---|
| **Wordcount** | CPU-intensive |
| **Sort** | IO-intensive |
| **DFSIO** | IO-intensive |
| **Pagerank** | Map :CPU-intensive<br>Reduce :IO-intensive |
| **Hivebench-Aggregation** | Map :CPU-intensive<br>Reduce :IO-intensive |
| **Hivebench-Join** | CPU-intensive |
| **Terasort** | Map :CPU-intensive<br>Reduce : IO-intensive |
| **K-Means** | Iteration stage: CPU-intensive<br>Classification stage: IO-intensive |
| **Nutchindexing** | CPU-intensive & IO-intensive |

**Cluster settings**

| Hadoop version | 1.0.3-Intel (patched with native task) |
|---|---|
| Cluster size | 4 |
| Disk per machine | 7 SATA Disk per node |
| Network | GbE network |
| CPU | E5-2680(32 core per node) |
| L3 Cache size | 20480 KB |
| Memory | 64GB per node |
| Map Slots | 3*32+1*26=122 |
| Reduce Slots | 3*16+1*13=61 |

**Job Configuration**

| io.sort.mb | 1GB |
|---|---|
| compression | Enabled |
| Compression algo | snappy |
| Dfs.block.size | 256MB |
| Io.sort.record.percent | 0.2 |
| Dfs replica | 3 |

## 3.2 Performance Metrics

| Data size before compression | Data size after compression | Native job run time(s | Original job run time(s) | Job performance improve- | Map stage performance nce |
|---|---|---|---|---|---|

| | | | ) | | ment | improvement |
|---|---|---|---|---|---|---|
| **Wordcount** | 1TB | 500GB | 1523.43 | 3957.11 | 159.8% | 159.8% |
| **Sort** | 500GB | 249GB | 2662.43 | 3066.97 | 15.2% | 45.4% |
| **DFSIO-Read** | 1TB | NA | 1249.68 | 1384.52 | 10.8% | 26% |
| **DFSIO-Write** | 1TB | NA | 6639.22 | 7165.97 | 7.9% | 7.9% |
| **Pagerank** | Pages:500M Total:481GB | 217GB | 6105.71 | 11644.63 | 90.7% | 133.8% |
| **Hive-Aggregation** | Uservisits:5 G Pages:600M Total:820GB | 345GB | 1113.82 | 1662.74 | 49.3% | 76.2% |
| **Hive-Join** | Uservisits:5 G Pages:600M Total:860GB | 382GB | 1107.55 | 1467.08 | 32.5% | 42.8% |
| **Terasort** | 1TB | NA | 4203.35 | 6360.49 | 51.3% | 109.1% |
| **K-Means** | Clusters:5 Samples:2G Inputfilesec ondample:4 00M Total:378GB | 350GB | 5734.11 | 7045.62 | 22.9% | 22.9% |
| **Nutchindexing** | Pages:40M 22G | NA | 4387.6 | 4600.56 | 4.9% | 13.2% |



**Hi-Bench Job Execution Time Analysis**

second

Legend:
- Original mapred job runtime
- Native-Task job runtime

## 3.3 Results

### 3.3.1 Wordcount

➢ Job Details:

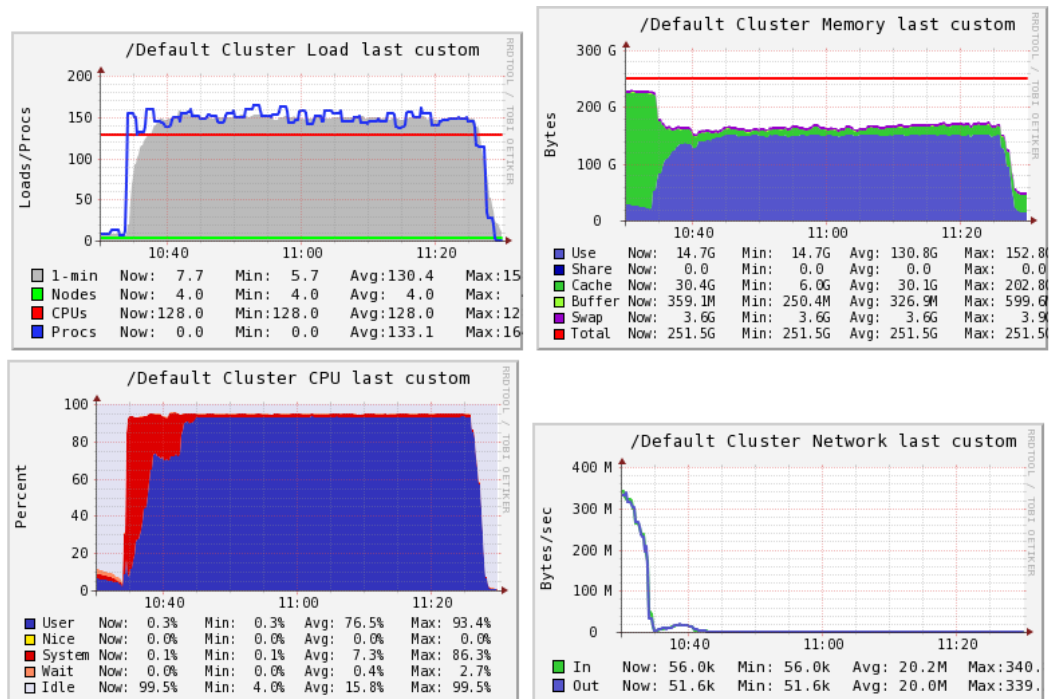| Name | Maps | Reducers |
|------|------|----------|
| wordcount | 1984 | 22 |





➢ Native-Task running state:

Start time: 9:14

Finish time: 9:37

➢ Original running state:



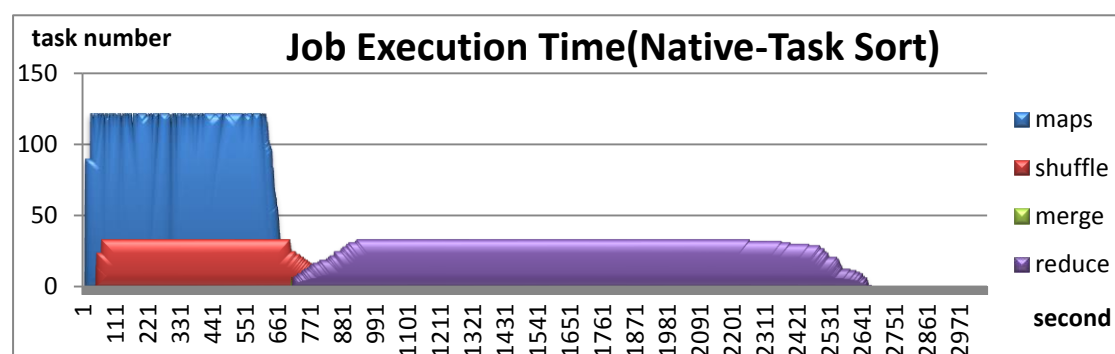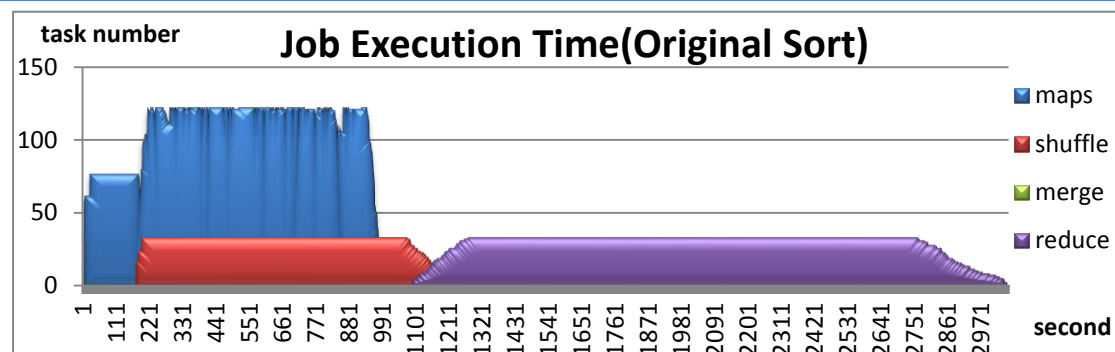Start time: 10:32

Finish time: 11:28

**Analysis**

Wordcount is a CPU-intensive workload and it's map stage run through the whole job. So the native-task has a huge performance improvement.
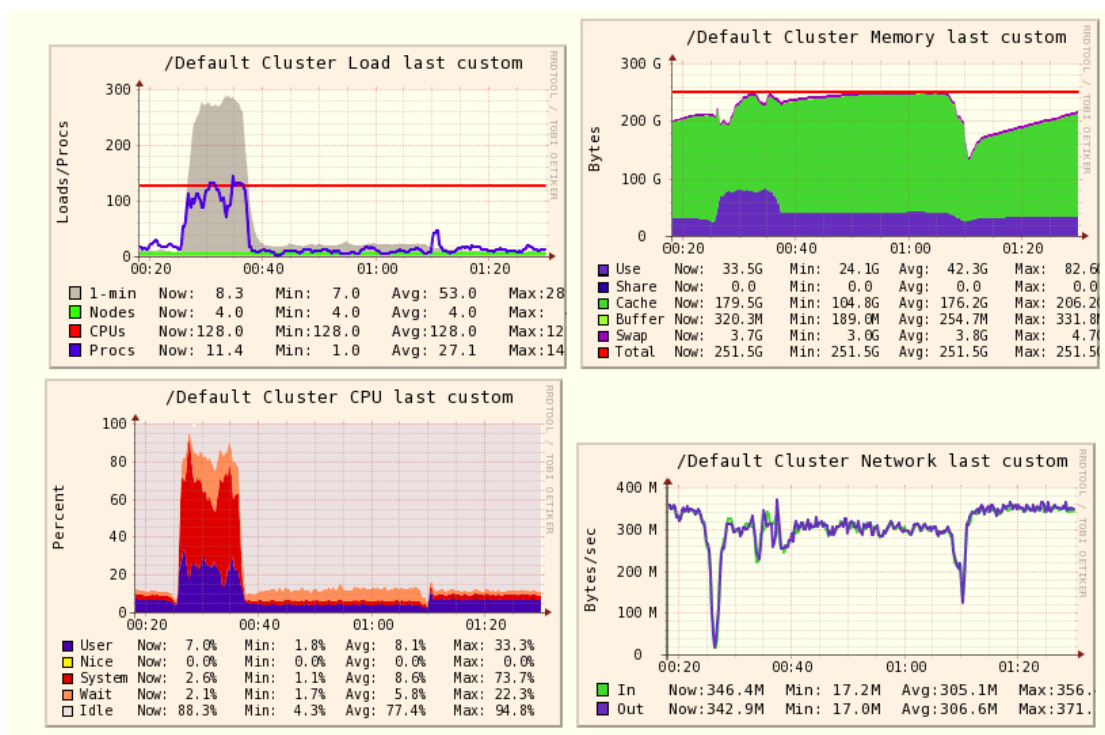
## 3.3.2 Sort

➢ Job Details:

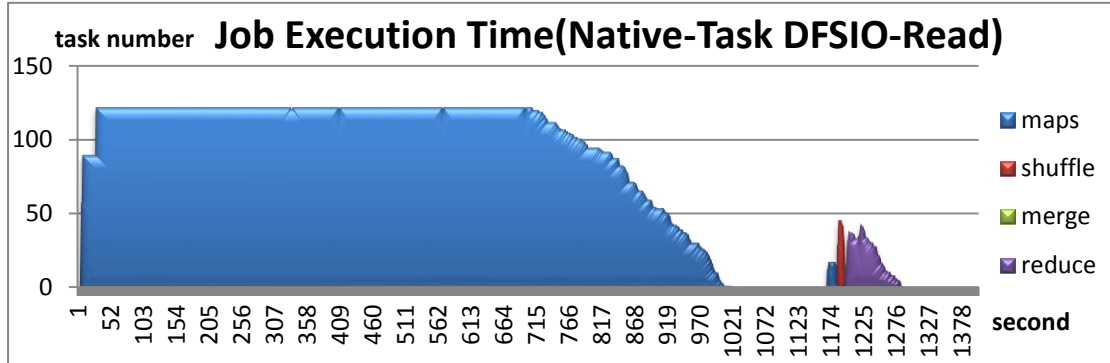| Name | Maps | Reducers |
| --- | --- | --- |
| **sorter** | 1024 | 33 |





➢ Native-Task running state:

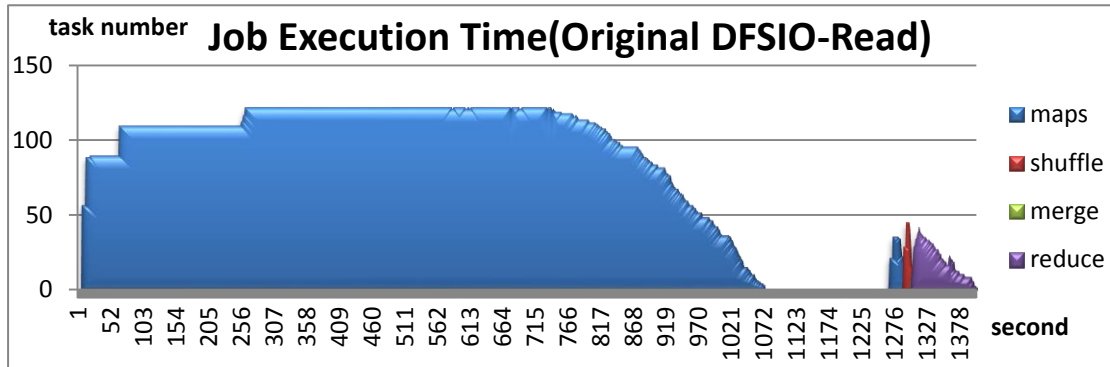Start time :00:25

Finish time :1:10

**Analysis**

Sort is IO-intensive at both map and reduce stage. We can see that it's reduce time occupy the most of whole job running time, because of that, the performance improvement is limited.
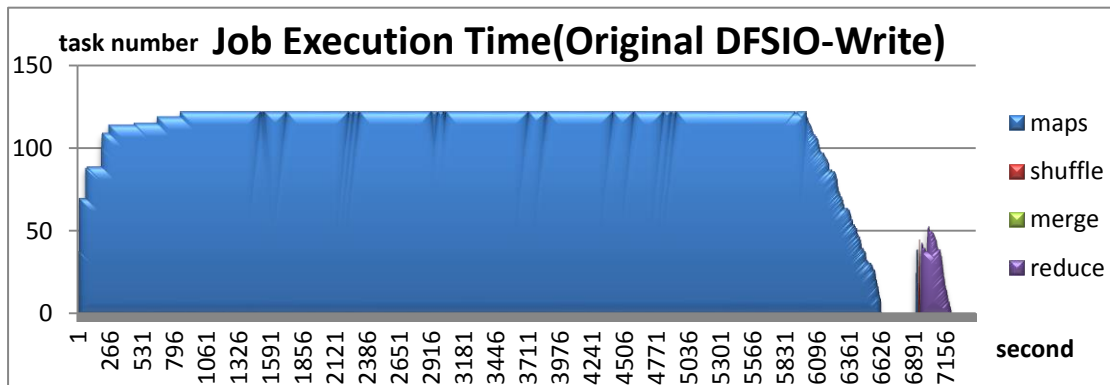
### 3.3.3 DFSIO-Read

➢ Job Details:

| Name | Maps | Reducers |
| --- | --- | --- |
| Datatools.jar | 256 | 1 |
| Result Analyzer | 50 | 63 |

Job Execution Time(Original DFSIO-Read)


Job Execution Time(Native-Task DFSIO-Read)

## 3.3.4 DFSIO-Write

➢ Job Details:

| Name | Maps | Reducers |
| --- | --- | --- |
| Datatools.jar | 512 | 1 |
| Result Analyzer | 50 | 63 |


Job Execution Time(Original DFSIO-Write)

Job Execution Time(Native-Task DFSIO-Write)

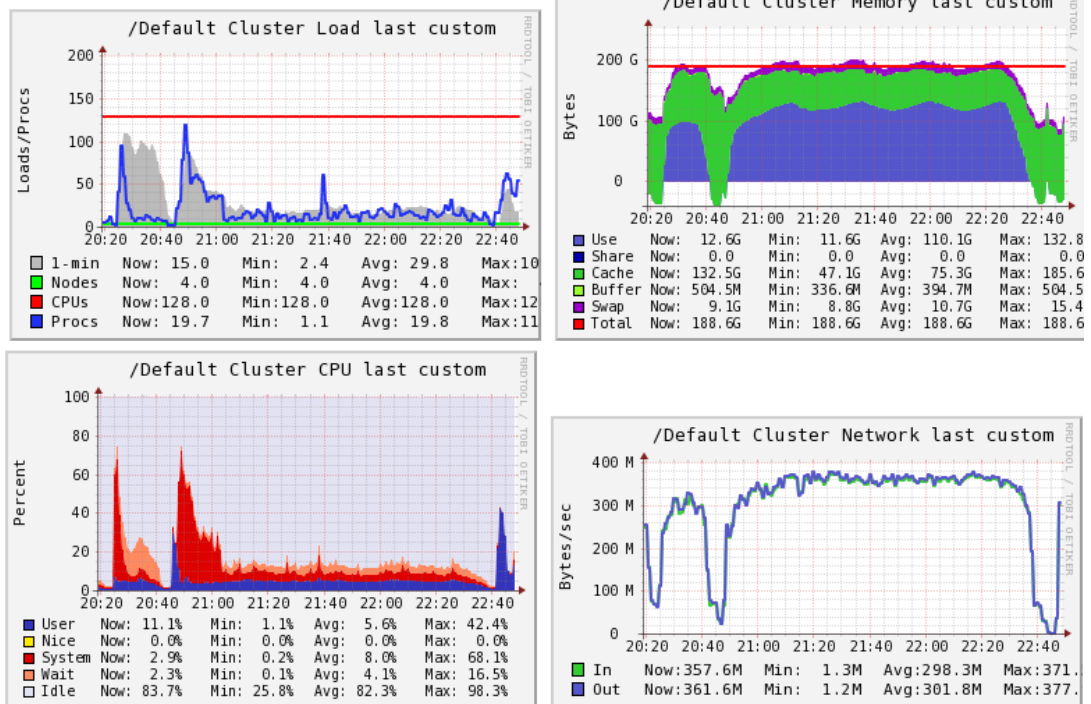> Native-Task running state:



Aggregation start time: 9:58

Aggregation finish time: 10:19

Join start time: 10:19

Join finish time: 12:10

> Original running state:

Aggregation start time: 20:22

Aggregation finish time: 20:46

Join start time: 20:46
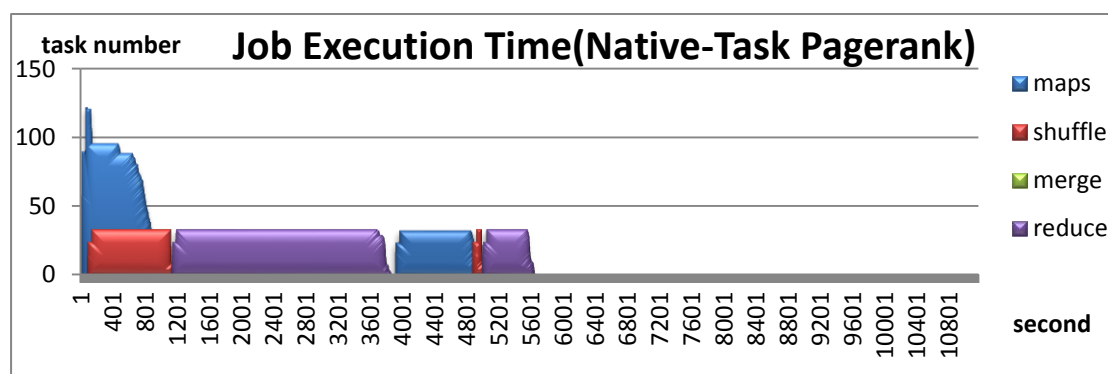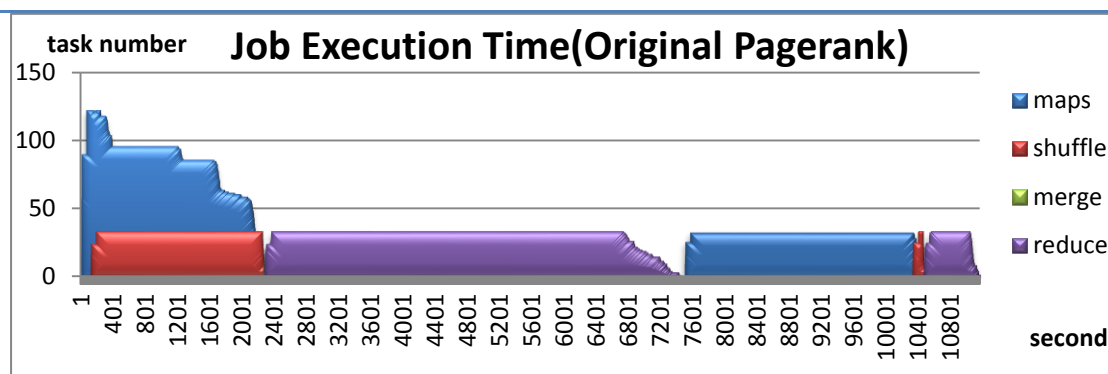
Join finish time: 22:45

**Analysis**

DFSIO is IO-intensive both at read and write stage. It's bottleneck is network bandwidth so the performance improvement is limited.
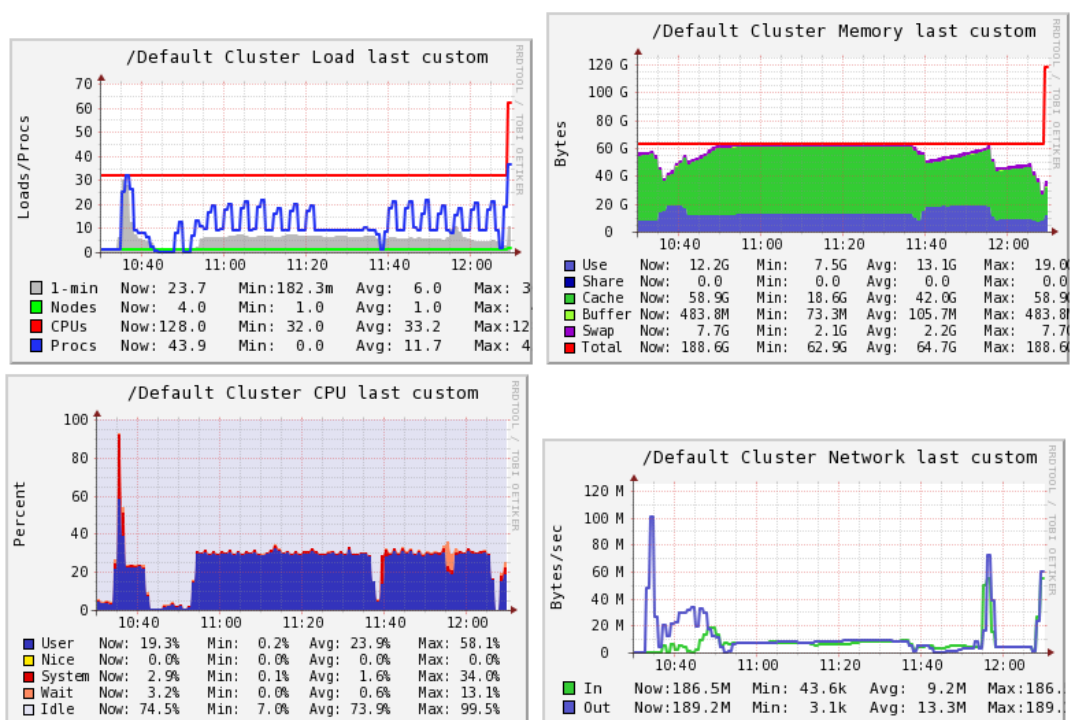
### 3.3.5 Pagerank

➢ Job Details:

| Name | Maps | Reducers |
|---|---|---|
| **Pagerank_Stage1** | 163 | 33 |

### Job Execution Time(Original Pagerank)



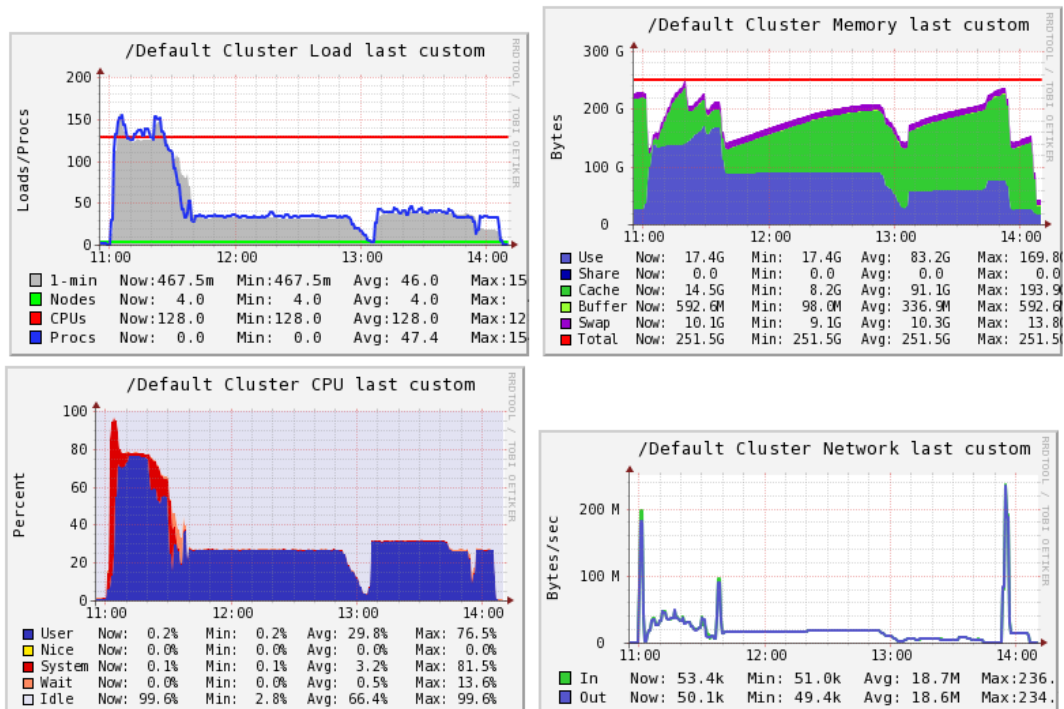### Job Execution Time(Native-Task Pagerank)



➢ Native-Task running state:



Start time: 10:33

Finish time: 12:08

> ➢ Original running state:
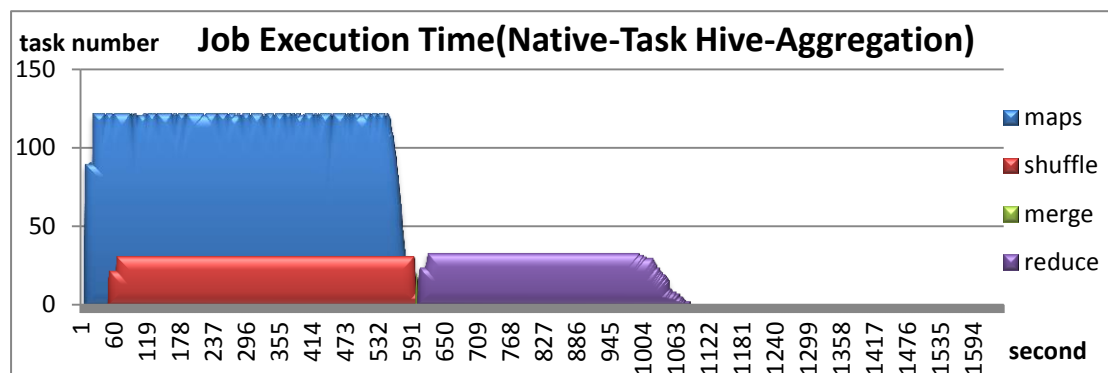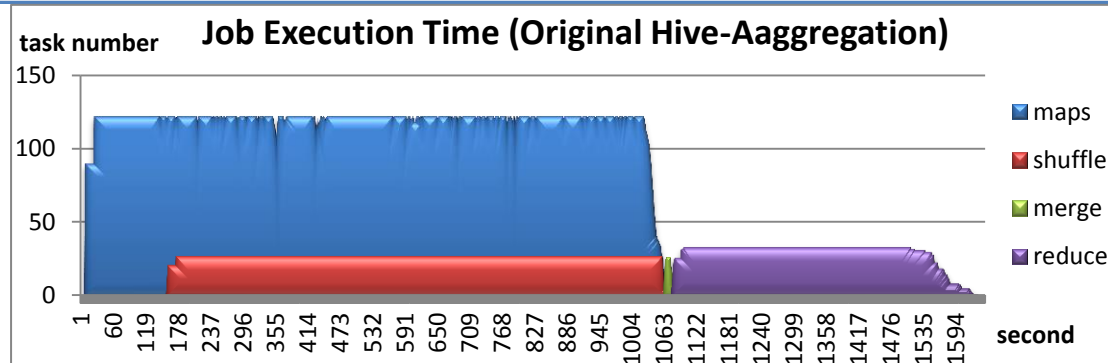


Start time: 10:59

Finish time: 14:06

**Analysis**

Pagerank is a CPU-intensive workload and it's map stage take about 50% of the whole job running time. So the performance improvement is obvious.
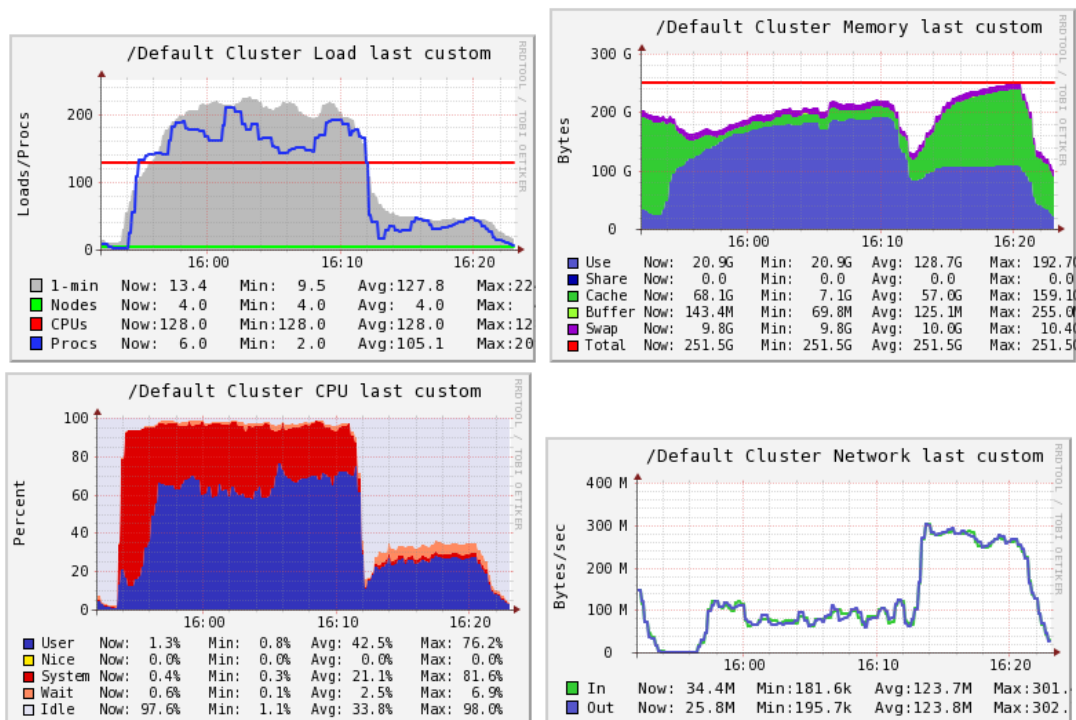
## 3.3.6 Hive-Aggregation

> ➢ Job Details:

| Name | Maps | Reducers |
|------|------|----------|
| **INSERT OVERWRITE TABLE** | 1386 | 33 |

## uservisits...sourceIP(Stage-1)



**Job Execution Time (Original Hive-Aaggregation)**

task number

maps
shuffle
merge
reduce

second



**Job Execution Time(Native-Task Hive-Aggregation)**

task number

maps
shuffle
merge
reduce

second

➢ Original running state:



/Default Cluster Load last custom

| | | | | |
|---|---|---|---|---|
| 1-min | Now: 13.4 | Min: 9.5 | Avg:127.8 | Max:22 |
| Nodes | Now: 4.0 | Min: 4.0 | Avg: 4.0 | Max: |
| CPUs | Now:128.0 | Min:128.0 | Avg:128.0 | Max:12 |
| Procs | Now: 6.0 | Min: 2.0 | Avg:105.1 | Max:20 |

/Default Cluster Memory last custom

| | | | | |
|---|---|---|---|---|
| Use | Now: 20.9G | Min: 20.9G | Avg: 128.7G | Max: 192.7( |
| Share | Now: 0.0 | Min: 0.0 | Avg: 0.0 | Max: 0.0 |
| Cache | Now: 68.1G | Min: 7.1G | Avg: 57.0G | Max: 159.1( |
| Buffer | Now: 143.4M | Min: 69.8M | Avg: 125.1M | Max: 255.0( |
| Swap | Now: 9.8G | Min: 9.8G | Avg: 10.0G | Max: 10.4( |
| Total | Now: 251.5G | Min: 251.5G | Avg: 251.5G | Max: 251.5( |

/Default Cluster CPU last custom

| | | | | |
|---|---|---|---|---|
| User | Now: 1.3% | Min: 0.8% | Avg: 42.5% | Max: 76.2% |
| Nice | Now: 0.0% | Min: 0.0% | Avg: 0.0% | Max: 0.0% |
| System | Now: 0.4% | Min: 0.3% | Avg: 21.1% | Max: 81.6% |
| Wait | Now: 0.6% | Min: 0.1% | Avg: 2.5% | Max: 6.9% |
| Idle | Now: 97.6% | Min: 1.1% | Avg: 33.8% | Max: 98.0% |

/Default Cluster Network last custom

| | | | | |
|---|---|---|---|---|
| In | Now: 34.4M | Min:181.6k | Avg:123.7M | Max:301. |
| Out | Now: 25.8M | Min:195.7k | Avg:123.8M | Max:302. |

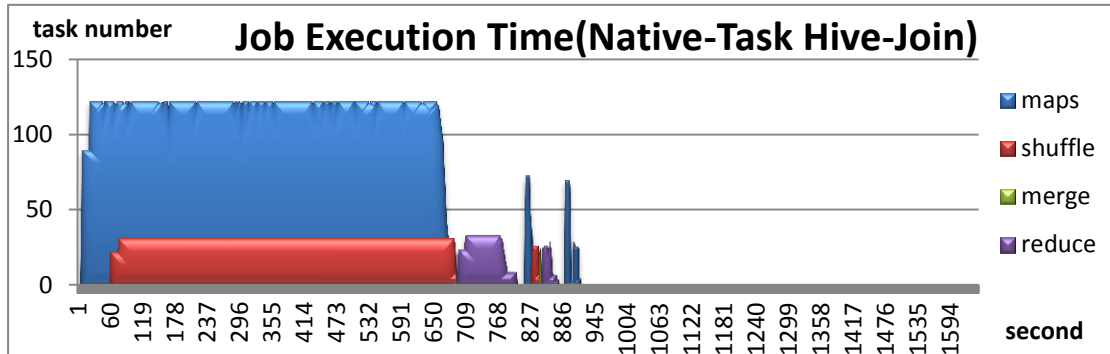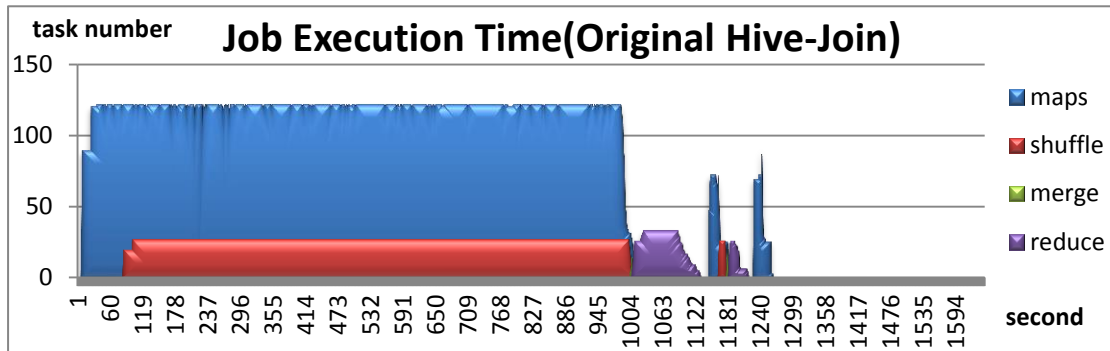Start time: 15:52

Finish time :16:22

**Analysis**

Hive-Aggregation is CPU-intensive at map stage and IO-intensive at reduce stage. It's map stage occupy the most of running time and when it comes to reduce stage, network bandwidth limits the performance. So the performance improvement at map stage is obvious.
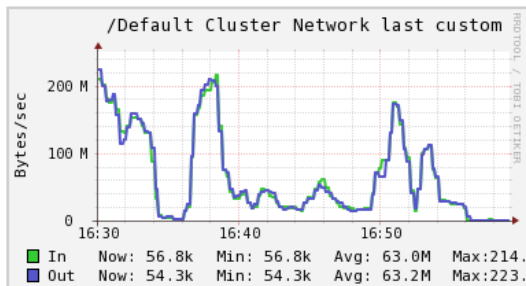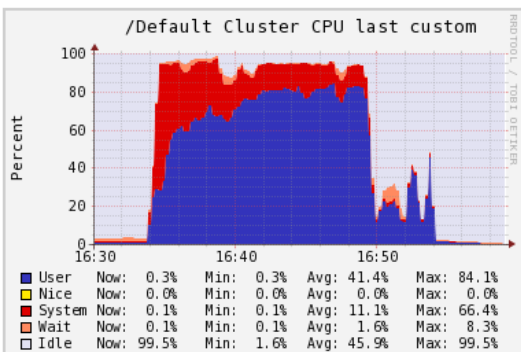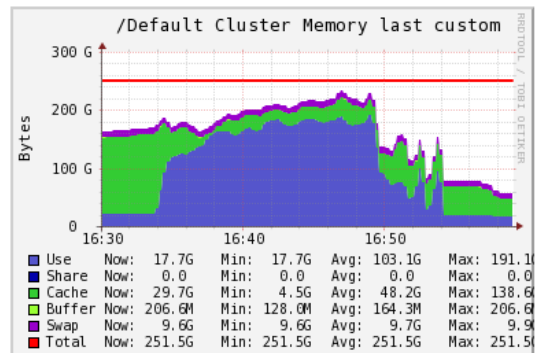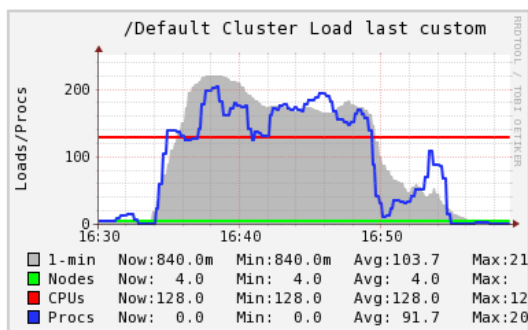
## 3.3.7 Hive-Join

➢ Job Details:

| Name | Maps | Reducers |
|---|---|---|
| **INSERT OVERWRITE TABLE rankings_uservisi...1(Stage-1)** | 1551 | 33 |
| **INSERT OVERWRITE TABLE rankings_uservisi...1(Stage-2)** | 99 | 33 |
| **INSERT OVERWRITE TABLE rankings_uservisi...1(Stage-3)** | 99 | 1 |
| **INSERT OVERWRITE TABLE rankings_uservisi...1(Stage-4)** | 1 | 1 |

Job Execution Time(Original Hive-Join)



Job Execution Time(Native-Task Hive-Join)

➢ Original running state:



Start time: 16:32

Finish time :16:58

**Analysis**

Hive-join is a CPU-intensive workload and it's map stage takes a high percent of whole running time. So we can see at map stage, the performance is improved by native-task.
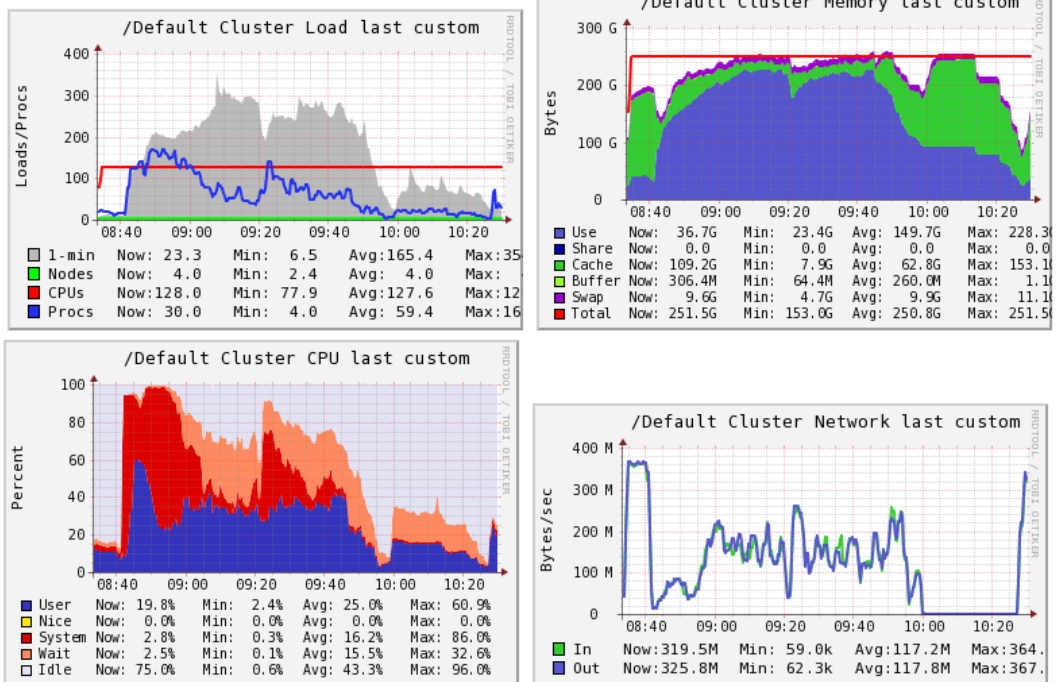
## 3.3.8 Terasort

➢ Job Details:

| Name | Maps | Reducers |
|------|------|----------|
| Terasort | 458 | 34 |



➢ Native-Task running state:

➢ Original running state:

Start time: 8:39

Finish time: 10:24

**Analysis**

Terasort is CPU-intensive at map stage and IO-intensive at reduce stage.It's map stage occupy the majority of the running time so there is a huge performance improvement at map stage.

## 3.3.9 K-Means

➢ Job Details:

| Name | Maps | Reducers |
|------|------|----------|
| **Cluster Iterator running iteration 1 …** | 1400 | 63 |
| **Cluster Iterator running** | 1400 | 63 |

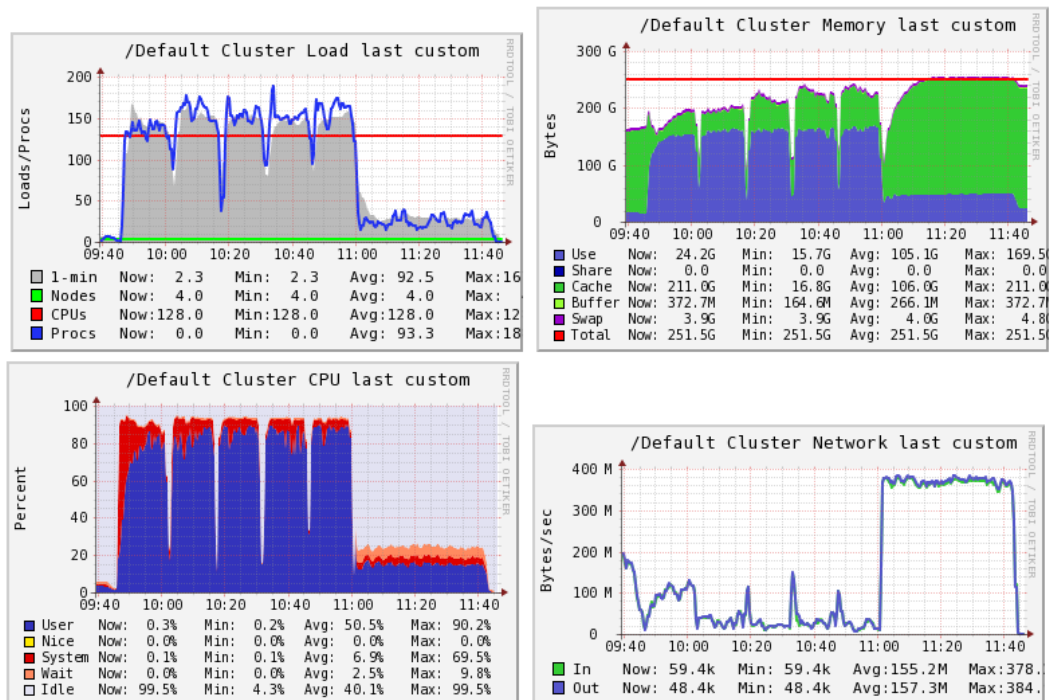| | | |
|---|---|---|
| **iteration 2 ...** | | |
| **Cluster Iterator running** | 1400 | 63 |
| **iteration 3 ...** | | |
| **Cluster Iterator running** | 1400 | 63 |
| **iteration 4 ...** | | |
| **Cluster Iterator running** | 1400 | 63 |
| **iteration 5 ...** | | |
| **Cluster     Classification** | 1400 | 0 |
| **Driver running** | | |





> ➢ Native-Task running state:

Start time: 20:38

Finish time: 22:23

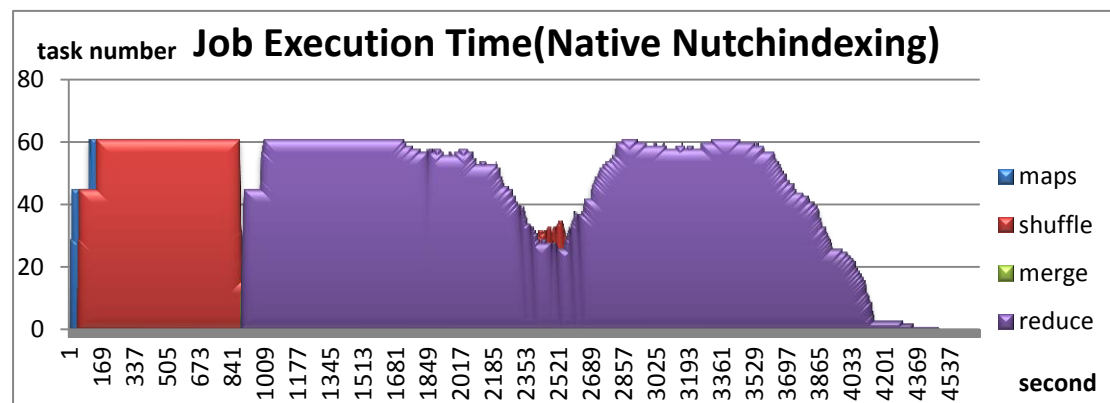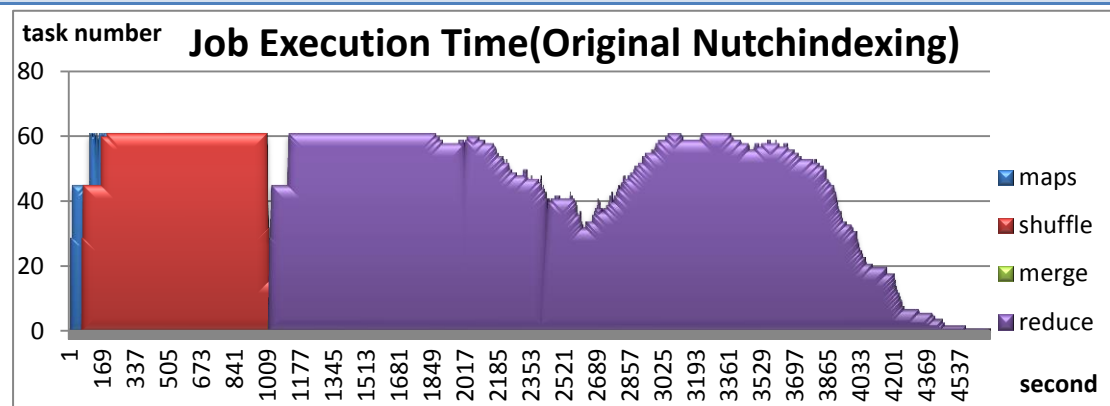> ➢ Original running state:



Start time: 9:43
Finish time: 11:41

**Analisys**

From the running state graph, we can see that the former 5 iteration is CPU-intensive and the last classification stage is IO-intensive. The two stages almost equally split the whole running time. So the performance improvement at map stage is evident.
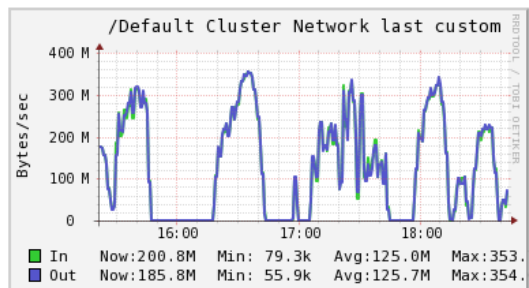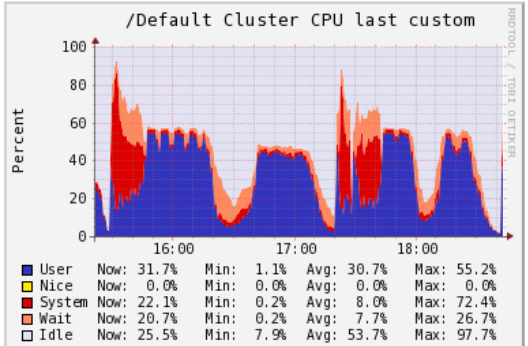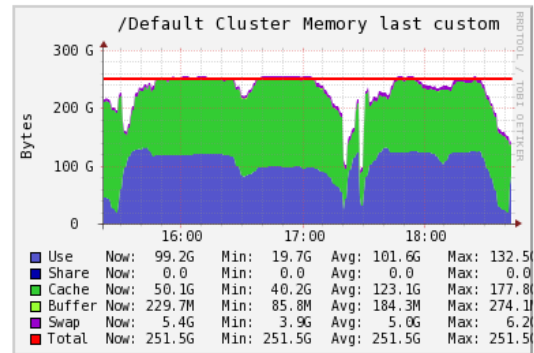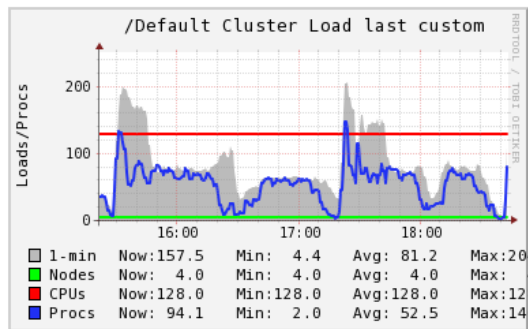
## 3.3.10 Nutchindexing

➢ Job Details:

| Name | Maps | Reducers |
|------|------|----------|
| index-lucene | 1171 | 133 |
| /HiBench/Nutch/Input/indexes | | |



Job Execution Time(Original Nutchindexing)
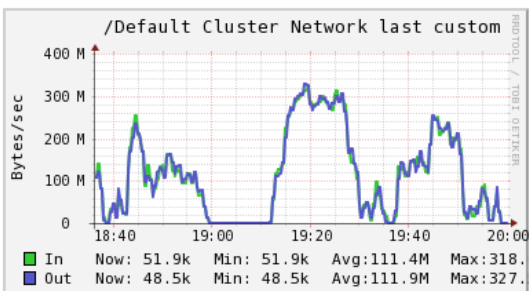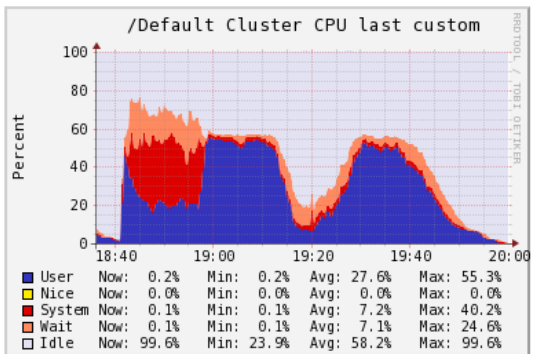


Job Execution Time(Native Nutchindexing)
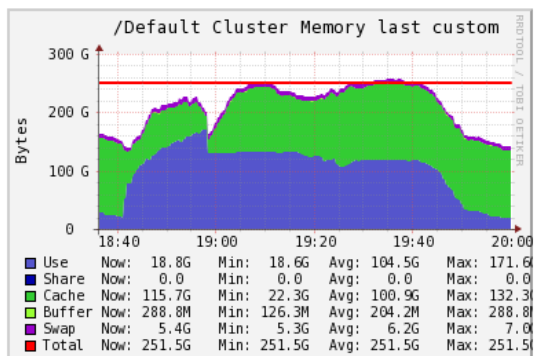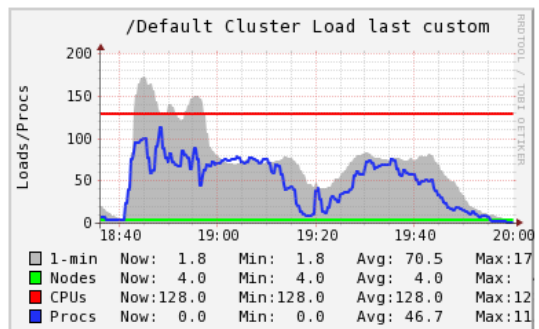
➢ Native-Task running state:

Start time: 17:26

Finish time: 18:40

> ➢ Original running state:



Start time: 18:40
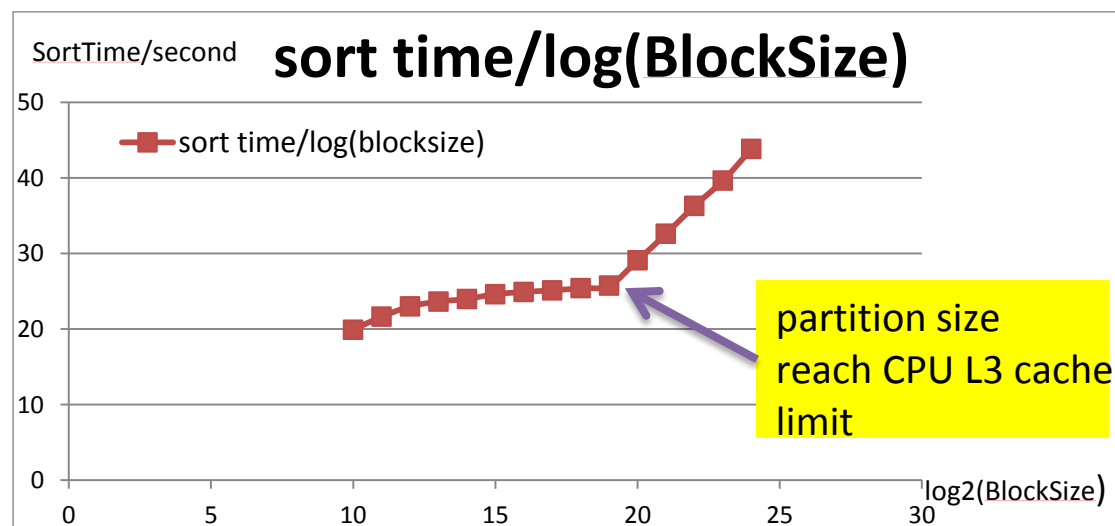
Finish time: 19:56

**Analysis**

Nutchindexing is CPU-intensive at map stage but the reduce stage take the majority of whole running time. So the performance improvement is not so huge.
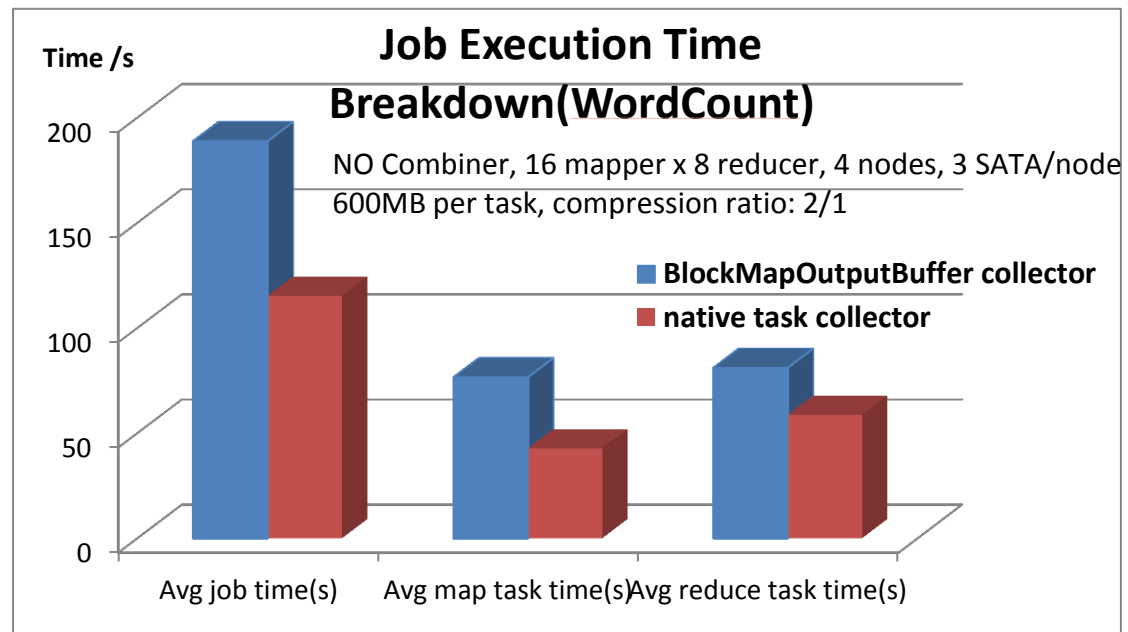
## 3.4 Other related results

### 3.4.1 Cache miss hurts Sorting performance

- Sorting time **increase rapidly** as cache miss rate increase

- We divide a large buffer into several memory unit.

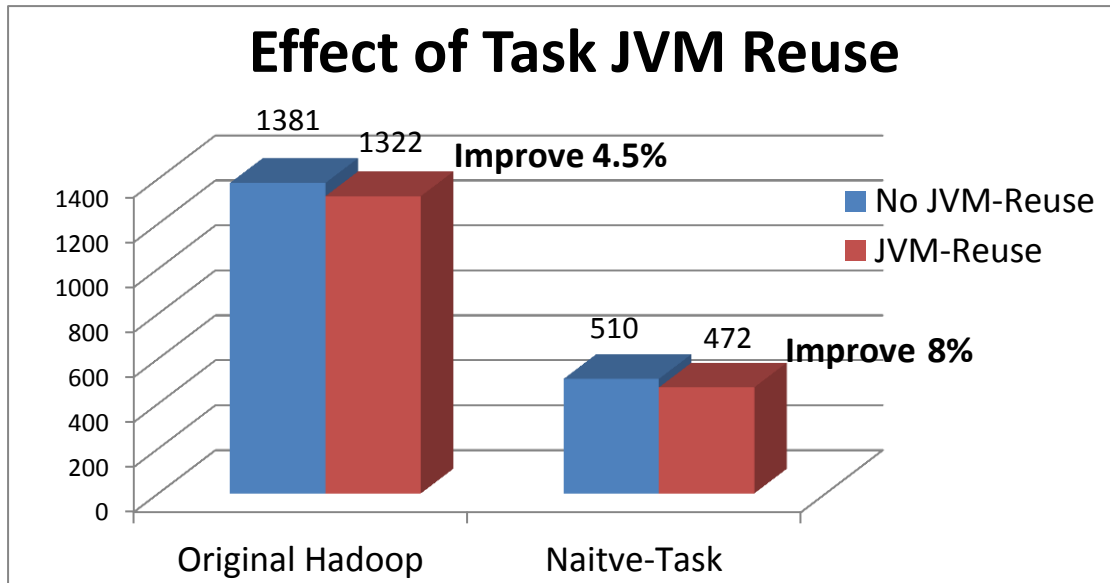- BlockSize is the size of memory unit we doing the sorting.

## 3.4.2 Compare with BlockMapOutputBuffer



- **70% faster** than BlockMapOutputBuffer collector.

- BlockMapOutputBuffer **supports ONLY BytesWritable**
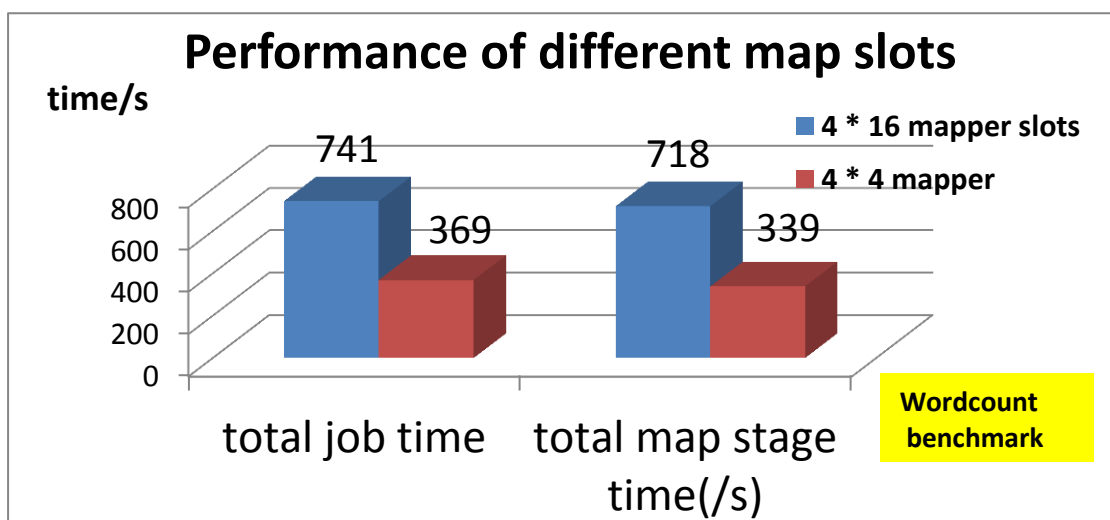
## 3.4.3 Effect of JVM reuse

- 4.5% improvement for Original Hadoop, 8% improvement for Native-Task

**Effect of Task JVM Reuse**

1381
1322 **Improve 4.5%**

■ No JVM-Reuse
■ JVM-Reuse

1400
1200
1000
800   510   472   **Improve 8%**
600
400
200
0

Original Hadoop          Naitve-Task
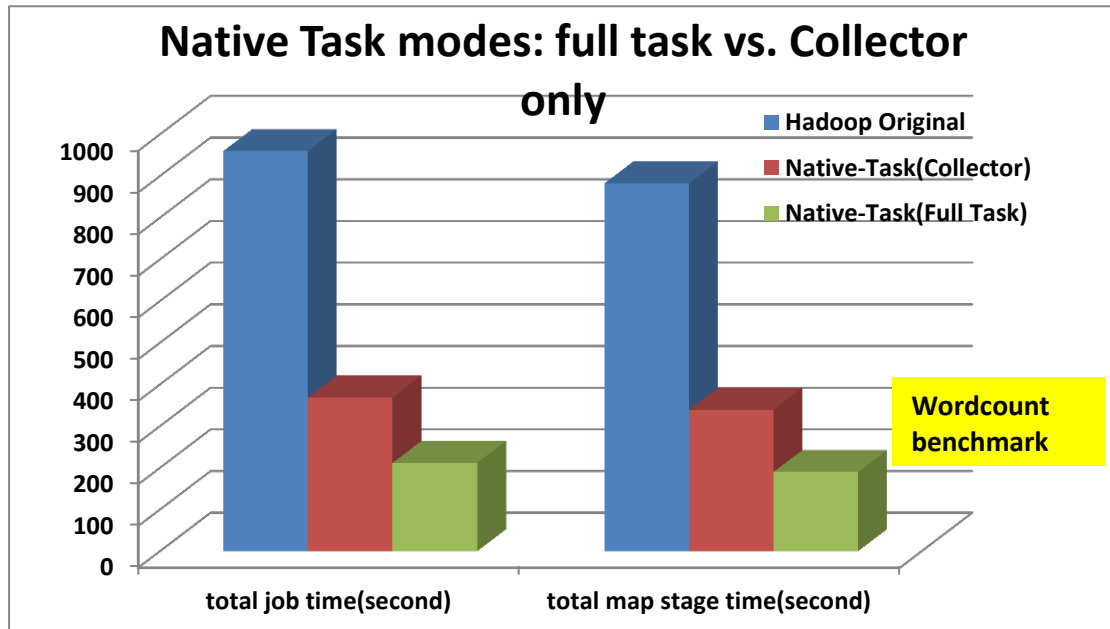
4 nodes, 4 map slots per node.

## 3.4.4 Hadoop don't scale well when slots number increase

- 4 nodes(32 core per node), 16 map slots max, **CPU, memory, disk are NOT fully used.**

- **Performance drops unexpectedly when slots# increase**

**Performance of different map slots**

time/s

■ 4 * 16 mapper slots
■ 4 * 4 mapper

741          718
369          339

800
600
400
200
0

total job time   total map stage   **Wordcount**
time(/s)          **benchmark**

### 3.4.5 Native-Task mode: full task optimization

- **2x faster further for Native-Task full time optimization**, compared with native collector



# 4. Conclusions