# MR Job lifecycle on Yarn Cluster (cont'd)

| Client Word Count Job | | |
|---|---|---|
| | Node Manager Data Node | Blocks |
| | WordCount Map Task | |
| | Node Manager Data Node | Blocks |
| | Application Master | |
| | Node Manager Data Node | Blocks |
| | WordCount Map Task | |
| Resource Manager | Node Manager Data Node | Blocks |

Name Node

Job History Server

→ Application Master (AM):
  a. Is a temporary Service
  b. coordinator for jobs
  c. Can be on any DataNode containing requested resources
  d. Does not process data, but will coordinate the processing, while actual processing is performed by respective Nodemanager on each DataNode
  e. Data is processed locally by Nodemanager on each DataNode


→ Map-Reduce
  a. Mapper        : program to filter the data → key, value
  b. Reducer       : program to process/aggregate the mapper → key, value

```
Mapper : prog to filter the data -> key, val
Reducer : Prog to process/aggregate the mapper data -> k,v

Retail store example

store id, prod, qty_sold, price/unit,,,,,,

blk_1 [node 1]
1,201,5,80
1,202,10,20
2,201,20,80
2,202,30,20

blk_2 [node 3]
1,201,10,80
1,202,20,20
2,201,30,80
2,202,40,20

prob statement : find total qty sold for each prod

select prod, sum(qty_sold) from sales group by prod;

key = prod
value = qty
```

→ Map-Reduce: Disk Based processing (local disk storage - local node storage)
→ Spark: memory based processing (in-memory storage)

→ Three stages in Map-Reduce:
    a. Map
    b. Shuffle
    c. Reduce

→ Mappers are launched parallel for single job/query
→ Mapper:
    a. Data is filtered
    b. Input is text format, output in key-value format

→ Data shuffling:
    a. First resource container is created for Sort and shuffle stage
    b. Sort & shuffle is processed on separate node
    c. merging and then processing data on another node
    d. Format is key,[list of values from different nodes]
    e. Keys are sorted in ASC order, while values are in random order
    f. Output of shuffling in mapper stage is now input for reducer

```
2,202,30,20

blk_2 [node 3]
1,201,10,80
1,202,20,20
2,201,30,80
2,202,40,20

prob statement : find total qty sold for each prod

select prod, sum(qty_sold) from sales group by prod;

key = prod
value = qty

mapper1
201,5
202,10
201,20
202,30

mapper2
201,10
202,20
201,30
202,40
```
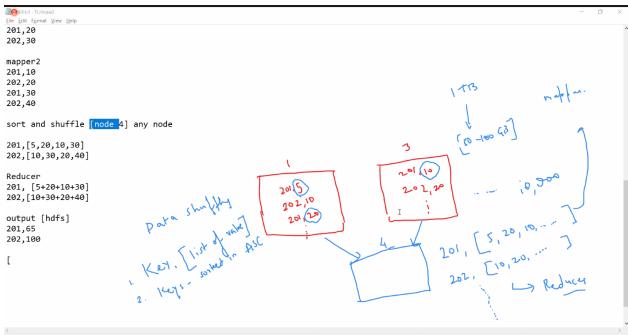
Windows (CRLF)    Ln 42, Col 1    100%

→ Reducer:
   a. First resource container is launched for reduce, Then reducer computes the final result
      And then it writes the final output to the HDFS
   b. Reducer is launched only once all the mapping-shuffling is done
   c.

```
201,20
202,30

mapper2
201,10
202,20
201,30
202,40

sort and shuffle [node 4] any node

201,[5,20,10,30]
202,[10,30,20,40]

Reducer
201, [5+20+10+30]
202,[10+30+20+40]

output [hdfs]
201,65
202,100

[
```

→ find total volume for each stockID

```
Untitled - Notepad
File Edit Format View Help
202,100

Retail store example

store id, prod, qty_sold, price/unit,,,,,,

blk_1 [node 1]
1,201,5,80
1,202,10,20
2,201,20,80
2,202,30,20

blk_2 [node 3]
1,201,10,80
1,202,20,20
2,201,30,80
2,202,40,20

prob statement : find total sales for each store id

select store_id, sum(qty*price) from sales group by store_id

key = store_id
val = price*qty

m1
1
```

```
Untitled - Notepad
File Edit Format View Help
key = store_id
val = price*qty

m1
1,400
1,200
2,1600
2,600

m2
1,800
1,400
2,2400
2,800

1, [400,200,800,400]
2, [1600,600,2400,800]

Reducer
--------
1, [400+200+800+400]
2, [1600+600+2400+800]

output
1,1800
2,5400
```

```
# first upload myjar.jar file using ftp
[bigdatalab456422@ip-10-1-1-204 ~]$ ll
total 959100
-rw-rw-r-- 1 bigdatalab456422 bigdatalab456422 207106008 May 17 12:44
eclipse.gz
```

```
-rw-rw-r-- 1 bigdatalab456422 bigdatalab456422        50 May 16 12:19
file1.txt
-rw-rw-r-- 1 bigdatalab456422 bigdatalab456422        20 May 16 12:30
file2.txt
-rw-rw-r-- 1 bigdatalab456422 bigdatalab456422 209715200 May 17 09:16
myfile
-rw-rw-r-- 1 bigdatalab456422 bigdatalab456422 524288000 May 17 09:35
myfile2
-rw-rw-r-- 1 bigdatalab456422 bigdatalab456422      4088 May 18 11:41
```
**myjar.jar**
```
-rw-rw-r-- 1 bigdatalab456422 bigdatalab456422  40990862 May 17 09:21
NYSE.csv

[bigdatalab456422@ip-10-1-1-204 ~]$ jar tvf myjar.jar
   25 Thu May 18 17:09:20 UTC 2023 META-INF/MANIFEST.MF
  387 Thu May 18 15:53:20 UTC 2023 .project
 2408 Thu May 18 17:00:02 UTC 2023 StockVolume$MapClass.class
 2349 Thu May 18 17:00:02 UTC 2023 StockVolume$ReduceClass.class
 1697 Thu May 18 17:00:02 UTC 2023 StockVolume.class
  640 Thu May 18 17:00:00 UTC 2023 .classpath
```

# run map & reduce task
# generates part-r-0000 file, 'r' means reducer output

```
[bigdatalab456422@ip-10-1-1-204 ~]$ hadoop jar myjar.jar StockVolume
training/NYSE.csv training/out1
WARNING: Use "yarn jar" to launch YARN applications.
23/05/18 11:48:33 INFO client.RMProxy: Connecting to ResourceManager
at ip-10-1-1-204.ap-south-1.compute.internal/10.1.1.204:8032
23/05/18 11:48:33 WARN mapreduce.JobResourceUploader: Hadoop
command-line option parsing not performed. Implement the Tool
interface and execute your application with T
oolRunner to remedy this.
23/05/18 11:48:33 INFO mapreduce.JobResourceUploader: Disabling
Erasure Coding for path:
/user/bigdatalab456422/.staging/job_1684298513961_0041
23/05/18 11:48:34 INFO input.FileInputFormat: Total input files to
process : 1
23/05/18 11:48:34 INFO mapreduce.JobSubmitter: number of splits:1
23/05/18 11:48:34 INFO Configuration.deprecation:
yarn.resourcemanager.system-metrics-publisher.enabled is deprecated.
Instead, use yarn.system-metrics-publisher.enable
d
23/05/18 11:48:34 INFO mapreduce.JobSubmitter: Submitting tokens for
job: job_1684298513961_0041
23/05/18 11:48:34 INFO mapreduce.JobSubmitter: Executing with tokens:
[]
23/05/18 11:48:34 INFO conf.Configuration: resource-types.xml not
found
23/05/18 11:48:34 INFO resource.ResourceUtils: Unable to find
'resource-types.xml'.
23/05/18 11:48:34 INFO impl.YarnClientImpl: Submitted application
application_1684298513961_0041
23/05/18 11:48:34 INFO mapreduce.Job: The url to track the job:
http://ip-10-1-1-204.ap-south-1.compute.internal:6066/proxy/applicati
on_1684298513961_0041/
23/05/18 11:48:34 INFO mapreduce.Job: Running job:
job_1684298513961_0041
23/05/18 11:48:58 INFO mapreduce.Job: Job job_1684298513961_0041
running in uber mode : false
23/05/18 11:48:58 INFO mapreduce.Job:  map 0% reduce 0%
23/05/18 11:49:27 INFO mapreduce.Job:  map 67% reduce 0%
23/05/18 11:49:29 INFO mapreduce.Job:  map 100% reduce 0%
23/05/18 11:50:00 INFO mapreduce.Job:  map 100% reduce 100%
23/05/18 11:50:01 INFO mapreduce.Job: Job job_1684298513961_0041
completed successfully
23/05/18 11:50:01 INFO mapreduce.Job: Counters: 54
        File System Counters
```

```
        FILE: Number of bytes read=3584395
        FILE: Number of bytes written=7613961
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
        HDFS: Number of bytes read=40990986
        HDFS: Number of bytes written=2918
        HDFS: Number of read operations=8
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=2
        HDFS: Number of bytes read erasure-coded=0
    Job Counters
        Launched map tasks=1
        Launched reduce tasks=1
        Data-local map tasks=1
        Total time spent by all maps in occupied slots
(ms)=26310
        Total time spent by all reduces in occupied slots
(ms)=28739
        Total time spent by all map tasks (ms)=26310
        Total time spent by all reduce tasks (ms)=28739
        Total vcore-milliseconds taken by all map tasks=26310
        Total vcore-milliseconds taken by all reduce
tasks=28739
        Total megabyte-milliseconds taken by all map
tasks=26941440
        Total megabyte-milliseconds taken by all reduce
tasks=29428736
    Map-Reduce Framework
        Map input records=735026
        Map output records=735026
        Map output bytes=8781587
        Map output materialized bytes=3584391
        Input split bytes=124
        Combine input records=0
        Combine output records=0
        Reduce input groups=203
        Reduce shuffle bytes=3584391
        Reduce input records=735026
        Reduce output records=203
        Spilled Records=1470052
        Shuffled Maps =1
        Failed Shuffles=0
        Merged Map outputs=1
        GC time elapsed (ms)=575
```

```
            CPU time spent (ms)=8990
            Physical memory (bytes) snapshot=882450432
            Virtual memory (bytes) snapshot=5186535424
            Total committed heap usage (bytes)=1075314688
            Peak Map Physical memory (bytes)=621236224
            Peak Map Virtual memory (bytes)=2586968064
            Peak Reduce Physical memory (bytes)=261267456
            Peak Reduce Virtual memory (bytes)=2599567360
    Shuffle Errors
            BAD_ID=0
            CONNECTION=0
            IO_ERROR=0
            WRONG_LENGTH=0
            WRONG_MAP=0
            WRONG_REDUCE=0
    File Input Format Counters
            Bytes Read=40990862
    File Output Format Counters
            Bytes Written=2918
```

```
        HDFS: Number of bytes written=2918
        HDFS: Number of read operations=8
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=2
        HDFS: Number of bytes read erasure-coded=0
Job Counters
        Launched map tasks=1
        Launched reduce tasks=1
        Data-local map tasks=1
        Total time spent by all maps in occupied slots (ms)=20093
        Total time spent by all reduces in occupied slots (ms)=23248
        Total time spent by all map tasks (ms)=20093
        Total time spent by all reduce tasks (ms)=23248
        Total vcore-milliseconds taken by all map tasks=20093
        Total vcore-milliseconds taken by all reduce tasks=23248
        Total megabyte-milliseconds taken by all map tasks=20575232
        Total megabyte-milliseconds taken by all reduce tasks=23805952
Map-Reduce Framework
        Map input records=735026
        Map output records=735026
        Map output bytes=8781587
        Map output materialized bytes=3584391
        Input split bytes=123
        Combine input records=0
        Combine output records=0
        Reduce input groups=203
        Reduce shuffle bytes=3584391
        Reduce input records=735026
        Reduce output records=203
        Spilled Records=1470052
        Shuffled Maps =1
        Failed Shuffles=0
        Merged Map outputs=1
        GC time elapsed (ms)=538
        CPU time spent (ms)=8730
        Physical memory (bytes) snapshot=881442816
        Virtual memory (bytes) snapshot=5185540096
        Total committed heap usage (bytes)=1081081856
        Peak Map Physical memory (bytes)=619569152
        Peak Map Virtual memory (bytes)=2586906624
        Peak Reduce Physical memory (bytes)=261873664
        Peak Reduce Virtual memory (bytes)=2598633472
```

*(handwritten annotations: "unique keys" pointing to Reduce input groups=203; "output" pointing to Reduce output records=203)*

**Hue** | Query | Search saved documents... | Jobs | bigdatalab45644

**File Browser**

Search for file name | Actions | Move to trash | Upload | New

Home / user / bigdatalab45644 / training / out1 | Trash

| Name | Size | User | Group | Permissions | Date |
|---|---|---|---|---|---|
| ↰ | | bigdatalab45644 | bigdatalab45644 | drwxr-xr-x | May 18, 2023 04:48 AM |
| . | | bigdatalab45644 | bigdatalab45644 | drwxr-xr-x | May 18, 2023 04:49 AM |
| _SUCCESS | 0 bytes | bigdatalab45644 | bigdatalab45644 | -rw-r--r-- | May 18, 2023 04:49 AM |
| part-r-00000 | 2.8 KB | bigdatalab45644 | bigdatalab45644 | -rw-r--r-- | May 18, 2023 04:49 AM |

Show 45 of 2 items | Page 1 of 1

Tables (5000)

1415_m6_dataset_credit_data_rx3_u95wx2jg
1a_bank3
1myemids
1stjuly
2008_data
5000sales
555555555555555555d
5july
5thjulyy
6thjuly
7jul2022
7th_july
7thjulaman
7thjuly
7thjuly007
7thjuly1
7thjuly11
7thjuly111
7thjuly2022
7thjuly2022apaxis
7thjuly2022db
7thjuly_fbn
7thjuly_silpa

---

**Hue** | Query | Search saved documents... | Jobs | bigdatalab45644

**File Browser**

Back | Home | Page 1 to 1 of 1

Edit file | / user / bigdatalab45644 / training / out1 / part-r-00000

Refresh
View as binary
Download

Last modified
05/18/2023 5:19 PM +05:30

User
bigdatalab45644

Group
bigdatalab45644

Size
2.85 KB

Mode
100644

```
AA      42061448400  ✓
AAI     5246821400   ✓
AAN     817567400
AAP     2802701500   ✓
AAR     49882000
AAV     834246600
AB      1125446300
ABA     11686500
ABB     4532301800
ABC     11439581700
ABD     469354400
ABG     458850900
ABK     11899868300
ABM     675519400
ABR     268351700
ABT     25664130200
ABV     1579314800
ABVT    49839000
ABX     16691172100
ACC     495415800
ACE     5224896200
ACF     5453798000
ACG     1481168200
ACH     1448279800
ACT     7219904300
```

*keys are sorted in ASc*

**≡ HUE** | Query ▾ | 🔍 Search saved documents... | Jobs 🔳 ↺ 👤 bigdatalab45644

📄 **File Browser**

‹ ▤ default

**Tables** (5000) + ↺

Filter...

- ▦ 1415_m6_dataset_credit_data_rx3_u95wx2jg
- ▦ 1a_bank3
- ▦ 1myemids
- ▦ 1stjuly
- ▦ 2008_data
- ▦ 5000sales
- ▦ 555555555555555555d
- ▦ 5july
- ▦ 5thjulyy
- ▦ 6thjuly
- ▦ 7jul2022
- ▦ 7th_july
- ▦ 7thjulaman
- ▦ 7thjuly
- ▦ 7thjuly007
- ▦ 7thjuly1
- ▦ 7thjuly11
- ▦ 7thjuly111
- ▦ 7thjuly2022
- ▦ 7thjuly2022apaxis
- ▦ 2022db
- ▦ y_fbn
- ▦ 7thjuly_silpa

↩ Back | 🏠 Home | Page 1 to 1 of 1 |◀ ◀◀ ▶▶ ▶|

✏ Edit file | / user / bigdatalab45644 / training / out1 / **part-r-00000**

↻ Refresh

▥ View as binary

⬇ Download

Last modified
05/18/2023 5:19 PM
+05:30

User
bigdatalab45644

Group
bigdatalab45644

Size
2.85 KB

Mode
100644

```
AA    42061448400
AAI   5246821400
AAN   817567400
AAP   2802701500
AAR   49882000
AAV   834246600
AB    1125446300
ABA   11686500
ABB   4532301800
ABC   11439581700
ABD   469354400
ABG   458850900
ABK   11899868300
ABM   675519400
ABR   268351700
ABT   25664130200
ABV   1579314800
ABVT  49839000
ABX   16691172100
ACC   495415800
ACE   5224896200
ACF   5453798000
ACG   1481168200
ACH   1448279800
ACT   72199043800
```



---

*Untitled - Notepad*

File Edit Format View Help

```
key = store_id
val = price*qty

m1
1,400
1,200
2,1600
2,600

m2
1,800
1,400
2,2400
2,800

1, [400,200,800,400]
2, [1600,600,2400,800]

Reducer
--------
1, [400+200+800+400]
2, [1600+600+2400+800]


output
1,1800
2,5400
```

Windows (CRLF) | Ln 83, Col 6 | 100%



---
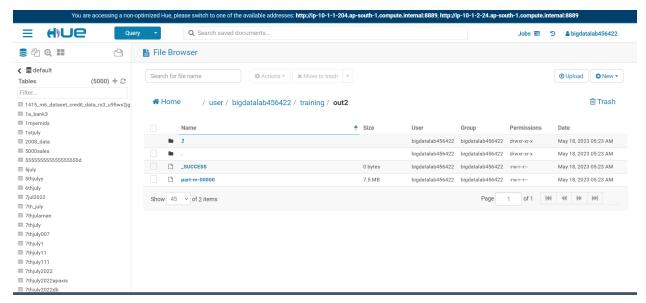
→

\# run only map task

\# change java code for reducerTask  as '0'

`job.setNumReduceTasks(0);`

\# generates part-r-0000 file, 'r'  means reducer output

```
[bigdatalab456422@ip-10-1-1-204 ~]$ hadoop jar myjar.jar StockVolume
training/NYSE.csv training/out2
WARNING: Use "yarn jar" to launch YARN applications.
23/05/18 12:23:09 INFO client.RMProxy: Connecting to ResourceManager
at ip-10-1-1-204.ap-south-1.compute.internal/10.1.1.204:8032
23/05/18 12:23:10 WARN mapreduce.JobResourceUploader: Hadoop
command-line option parsing not performed. Implement the Tool
interface and execute your application with T
oolRunner to remedy this.
23/05/18 12:23:10 INFO mapreduce.JobResourceUploader: Disabling
Erasure Coding for path:
/user/bigdatalab456422/.staging/job_1684298513961_0137
23/05/18 12:23:10 INFO input.FileInputFormat: Total input files to
process : 1
23/05/18 12:23:10 INFO mapreduce.JobSubmitter: number of splits:1
23/05/18 12:23:10 INFO Configuration.deprecation:
yarn.resourcemanager.system-metrics-publisher.enabled is deprecated.
Instead, use yarn.system-metrics-publisher.enable
d
23/05/18 12:23:10 INFO mapreduce.JobSubmitter: Submitting tokens for
job: job_1684298513961_0137
23/05/18 12:23:10 INFO mapreduce.JobSubmitter: Executing with tokens:
[]
23/05/18 12:23:11 INFO conf.Configuration: resource-types.xml not
found
23/05/18 12:23:11 INFO resource.ResourceUtils: Unable to find
'resource-types.xml'.
23/05/18 12:23:11 INFO impl.YarnClientImpl: Submitted application
application_1684298513961_0137
23/05/18 12:23:11 INFO mapreduce.Job: The url to track the job:
http://ip-10-1-1-204.ap-south-1.compute.internal:6066/proxy/applicati
on_1684298513961_0137/
23/05/18 12:23:11 INFO mapreduce.Job: Running job:
job_1684298513961_0137
23/05/18 12:23:20 INFO mapreduce.Job: Job job_1684298513961_0137
running in uber mode : false
23/05/18 12:23:20 INFO mapreduce.Job:  map 0% reduce 0%
23/05/18 12:23:26 INFO mapreduce.Job:  map 100% reduce 0%
23/05/18 12:23:26 INFO mapreduce.Job: Job job_1684298513961_0137
completed successfully
23/05/18 12:23:26 INFO mapreduce.Job: Counters: 33
        File System Counters
                FILE: Number of bytes read=0
                FILE: Number of bytes written=222430
                FILE: Number of read operations=0
```

```
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=40990986
                HDFS: Number of bytes written=7842509
                HDFS: Number of read operations=7
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
                HDFS: Number of bytes read erasure-coded=0
        Job Counters
                Launched map tasks=1
                Data-local map tasks=1
                Total time spent by all maps in occupied slots
(ms)=4102
                Total time spent by all reduces in occupied slots
(ms)=0
                Total time spent by all map tasks (ms)=4102
                Total vcore-milliseconds taken by all map tasks=4102
                Total megabyte-milliseconds taken by all map
tasks=4200448
        Map-Reduce Framework
                Map input records=735026
                Map output records=735026
                Input split bytes=124
                Spilled Records=0
                Failed Shuffles=0
                Merged Map outputs=0
                GC time elapsed (ms)=69
                CPU time spent (ms)=3170
                Physical memory (bytes) snapshot=361504768
                Virtual memory (bytes) snapshot=2584989696
                Total committed heap usage (bytes)=480772096
                Peak Map Physical memory (bytes)=361504768
                Peak Map Virtual memory (bytes)=2584989696
        File Input Format Counters
                Bytes Read=40990862
        File Output Format Counters
                Bytes Written=7842509
```

Query

Search saved documents...

Jobs   bigdatalab456422

## File Browser

Search for file name          Actions ▾    ✖ Move to trash                          ⊕ Upload    ⊕ New ▾

🏠 Home      /  user /  bigdatalab456422 /  training /  **out2**                                    🗑 Trash

| | Name | Size | User | Group | Permissions | Date |
|---|---|---|---|---|---|---|
| ☐ | 📁 ⬆ | | bigdatalab456422 | bigdatalab456422 | drwxr-xr-x | May 18, 2023 05:23 AM |
| ☐ | 📁 . | | bigdatalab456422 | bigdatalab456422 | drwxr-xr-x | May 18, 2023 05:23 AM |
| ☐ | 📄 _SUCCESS | 0 bytes | bigdatalab456422 | bigdatalab456422 | -rw-r--r-- | May 18, 2023 05:23 AM |
| ☐ | 📄 part-m-00000 | 7.5 MB | bigdatalab456422 | bigdatalab456422 | -rw-r--r-- | May 18, 2023 05:23 AM |

Show  45  of 2 items                                              Page  1  of 1   |◀  ◀◀  ▶▶  ▶|

### default

Tables                    (5000) ✛ ⟳

Filter...

⊞ 1415_m6_dataset_credit_data_rx3_u95wx2jg
⊞ 1a_bank3
⊞ 1myemids
⊞ 1stjuly
⊞ 2008_data
⊞ 5000sales
⊞ 555555555555555555d
⊞ 5july
⊞ 5thjulyy
⊞ 6thjuly
⊞ 7jul2022
⊞ 7th_july
⊞ 7thjulaman
⊞ 7thjuly
⊞ 7thjuly007
⊞ 7thjuly1
⊞ 7thjuly11
⊞ 7thjuly111
⊞ 7thjuly2022
⊞ 7thjuly2022apaxis
⊞ 7thjuly2022db

here , in part-m-0000, 'm' means output from mapper