

## CertKiller.CCA-500,55.Q&A

Number: CCA-500  
Passing Score: 800  
Time Limit: 120 min  
File Version: 14.15



<http://www.gratisexam.com/>



**Cloudera Certified Administrator for Apache Hadoop**

- Pretty much all the questions we study that may have multiple answers, no answer is in doubt, I got on the test.
- Got this vce from my friend who passed with 98% , each and every stuff in it. I am sharing with you guys.
- These are the most accurate study questions. Just focus on these and sit in your exam.
- Modified few questions, fixed few spelling mistakes and typos.
- Finally, I got right questions for this exam and share with you guys. Best Wishes.

## Exam A

### QUESTION 1

Your cluster's mapred-start.xml includes the following parameters

```
<name>mapreduce.map.memory.mb</name>  
<value>4096</value>  
<name>mapreduce.reduce.memory.mb</name>  
<value>8192</value>
```

And any cluster's yarn-site.xml includes the following parameters

```
<name>yarn.nodemanager.vmem-pmem-ratio</name>  
<value>2.1</value>
```

What is the maximum amount of virtual memory allocated for each map task before YARN will kill its Container?

- A. 4 GB
- B. 17.2 GB
- C. 8.9 GB
- D. 8.2 GB
- E. 24.6 GB

**Correct Answer:** D

**Section:** (none)

**Explanation**

**Explanation/Reference:**

In order to get maximum amount of virtual memory allocated for each map task, you have to multiply mapreduce.map.memory.mb with yarn.nodemanager.vmem-pmem-ratio. The result would be 8601.6 MB. So the nearest answer is 8.2 since 8.9 is more than 8601.6 MB.

Reference: <https://support.pivotal.io/hc/en-us/articles/201462036-Mapreduce-YARN-Memory-Parameters>

### QUESTION 2

You want to understand more about how users browse your public website. For example, you want to know which pages they visit prior to placing an order. You have a server farm of 200 web servers hosting your website. Which is the most efficient process to gather these web server across logs into your Hadoop cluster analysis?

- A. Sample the web server logs web servers and copy them into HDFS using curl
- B. Ingest the server web logs into HDFS using Flume

- C. Channel these clickstreams into Hadoop using Hadoop Streaming
- D. Import all user clicks from your OLTP databases into Hadoop using Sqoop
- E. Write a MapReduce job with the web servers for mappers and the Hadoop cluster nodes for reducers

**Correct Answer:** B

**Section:** (none)

**Explanation**

**Explanation/Reference:**

Apache Flume is a service for streaming logs into Hadoop.

Apache Flume is a distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of streaming data into the Hadoop Distributed File System (HDFS). It has a simple and flexible architecture based on streaming data flows; and is robust and fault tolerant with tunable reliability mechanisms for failover and recovery.

**QUESTION 3**

Which three basic configuration parameters must you set to migrate your cluster from MapReduce 1 (MRv1) to MapReduce V2 (MRv2)? (Choose three)

- A. Configure the NodeManager to enable MapReduce services on YARN by setting the following property in yarn-site.xml:  
<name>yarn.nodemanager.hostname</name>  
<value>your\_nodeManager\_shuffle</value>
- B. Configure the NodeManager hostname and enable node services on YARN by setting the following property in yarn-site.xml:  
<name>yarn.nodemanager.hostname</name>  
<value>your\_nodeManager\_hostname</value>
- C. Configure a default scheduler to run on YARN by setting the following property in mapred-site.xml:  
<name>mapreduce.jobtracker.taskScheduler</name>  
<value>org.apache.hadoop.mapred.JobQueueTaskScheduler</value>
- D. Configure the number of map tasks per job on YARN by setting the following property in mapred-site.xml:  
<name>mapreduce.job.maps</name>  
<value>2</value>
- E. Configure the ResourceManager hostname and enable node services on YARN by setting the following property in yarn-site.xml:  
<name>yarn.resourcemanager.hostname</name>  
<value>your\_resourceManager\_hostname</value>
- F. Configure MapReduce as a Framework running on YARN by setting the following property in mapred-site.xml:  
<name>mapreduce.framework.name</name>  
<value>yarn</value>

**Correct Answer:** AEF

**Section:** (none)

**Explanation**

**Explanation/Reference:**

**QUESTION 4**

You need to analyze 60,000,000 images stored in JPEG format, each of which is approximately 25 KB. Because your Hadoop cluster isn't optimized for storing and processing many small files, you decide to do the following actions:

1. Group the individual images into a set of larger files
2. Use the set of larger files as input for a MapReduce job that processes them directly with python using Hadoop streaming.

Which data serialization system gives the flexibility to do this?



<http://www.gratisexam.com/>

- A. CSV
- B. XML
- C. HTML
- D. Avro
- E. SequenceFiles
- F. JSON

**Correct Answer:** E

**Section:** (none)

**Explanation**

**Explanation/Reference:**

Sequence files are block-compressed and provide direct serialization and deserialization of several arbitrary data types (not just text). Sequence files can be generated as the output of other MapReduce tasks and are an efficient intermediate representation for data that is passing from one MapReduce job to another.

**QUESTION 5**

Identify two features/issues that YARN is designed to address: (Choose two)

- A. Standardize on a single MapReduce API
- B. Single point of failure in the NameNode

- C. Reduce complexity of the MapReduce APIs
- D. Resource pressure on the JobTracker
- E. Ability to run framework other than MapReduce, such as MPI
- F. HDFS latency

**Correct Answer:** DE

**Section:** (none)

**Explanation**

**Explanation/Reference:**

Reference: <http://www.revelytix.com/?q=content/hadoop-ecosystem> (YARN, first para)

#### **QUESTION 6**

Which YARN daemon or service monitors a Controller's per-application resource using (e.g., memory CPU)?

- A. ApplicationMaster
- B. NodeManager
- C. ApplicationManagerService
- D. ResourceManager

**Correct Answer:** A

**Section:** (none)

**Explanation**

**Explanation/Reference:**

#### **QUESTION 7**

Which is the default scheduler in YARN?

- A. YARN doesn't configure a default scheduler, you must first assign an appropriate scheduler class in yarn-site.xml
- B. Capacity Scheduler
- C. Fair Scheduler
- D. FIFO Scheduler

**Correct Answer:** B

**Section:** (none)

**Explanation**

**Explanation/Reference:**

Reference: <http://hadoop.apache.org/docs/r2.4.1/hadoop-yarn/hadoop-yarn-site/CapacityScheduler.html>

**QUESTION 8**

Which YARN process run as "container 0" of a submitted job and is responsible for resource requests?

- A. ApplicationManager
- B. JobTracker
- C. ApplicationMaster
- D. JobHistoryServer
- E. ResourceManager
- F. NodeManager

**Correct Answer: C**

**Section: (none)**

**Explanation**

**Explanation/Reference:**

**QUESTION 9**

Which scheduler would you deploy to ensure that your cluster allows short jobs to finish within a reasonable time without starting long-running jobs?

- A. Complexity Fair Scheduler (CFS)
- B. Capacity Scheduler
- C. Fair Scheduler
- D. FIFO Scheduler

**Correct Answer: C**

**Section: (none)**

**Explanation**

**Explanation/Reference:**

Reference: [http://hadoop.apache.org/docs/r1.2.1/fair\\_scheduler.html](http://hadoop.apache.org/docs/r1.2.1/fair_scheduler.html)

**QUESTION 10**

Your cluster is configured with HDFS and MapReduce version 2 (MRv2) on YARN. What is the result when you execute: `hadoop jar SampleJar MyClass` on a client machine?

- A. SampleJar.Jar is sent to the ApplicationMaster which allocates a container for SampleJar.Jar
- B. Sample.jar is placed in a temporary directory in HDFS
- C. SampleJar.jar is sent directly to the ResourceManager
- D. SampleJar.jar is serialized into an XML file which is submitted to the ApplicatoionMaster

**Correct Answer:** A

**Section:** (none)

**Explanation**

**Explanation/Reference:**

answer is modified.

#### **QUESTION 11**

You are working on a project where you need to chain together MapReduce, Pig jobs. You also need the ability to use forks, decision points, and path joins. Which ecosystem project should you use to perform these actions?

- A. Oozie
- B. ZooKeeper
- C. HBase
- D. Sqoop
- E. HUE

**Correct Answer:** A

**Section:** (none)

**Explanation**

**Explanation/Reference:**

#### **QUESTION 12**

Which YARN daemon or service negotiates map and reduce Containers from the Scheduler, tracking their status and monitoring progress?

- A. NodeManager
- B. ApplicationMaster
- C. ApplicationManager
- D. ResourceManager

**Correct Answer:** B

**Section:** (none)

**Explanation**

**Explanation/Reference:**

Reference: <http://www.devx.com/opensource/intro-to-apache-mapreduce-2-yarn.html> (See resource manager)

### QUESTION 13

During the execution of a MapReduce v2 (MRv2) job on YARN, where does the Mapper place the intermediate data of each Map Task?

- A. The Mapper stores the intermediate data on the node running the Job's ApplicationMaster so that it is available to YARN ShuffleService before the data is presented to the Reducer
- B. The Mapper stores the intermediate data in HDFS on the node where the Map tasks ran in the HDFS /usercache/&(user)/apache/application\_&(appid) directory for the user who ran the job
- C. The Mapper transfers the intermediate data immediately to the reducers as it is generated by the Map Task
- D. YARN holds the intermediate data in the NodeManager's memory (a container) until it is transferred to the Reducer
- E. The Mapper stores the intermediate data on the underlying filesystem of the local disk in the directories yarn.nodemanager.local-DIFS

**Correct Answer:** E

**Section:** (none)

**Explanation**

**Explanation/Reference:**

### QUESTION 14

You suspect that your NameNode is incorrectly configured, and is swapping memory to disk. Which Linux commands help you to identify whether swapping is occurring? (Select all that apply)

- A. free
- B. df
- C. memcat
- D. top
- E. jps
- F. vmstat
- G. swapinfo

**Correct Answer:** ADF

**Section:** (none)



## Explanation

### Explanation/Reference:

Reference: <http://www.cyberciti.biz/faq/linux-check-swap-usage-command/>

### QUESTION 15

On a cluster running CDH 5.0 or above, you use the `hadoop fs put` command to write a 300MB file into a previously empty directory using an HDFS block size of 64 MB. Just after this command has finished writing 200 MB of this file, what would another user see when they look in the directory?

- A. The directory will appear to be empty until the entire file write is completed on the cluster
- B. They will see the file with a `._COPYING_` extension on its name. If they view the file, they will see contents of the file up to the last completed block (as each 64MB block is written, that block becomes available)
- C. They will see the file with a `._COPYING_` extension on its name. If they attempt to view the file, they will get a `ConcurrentFileAccessException` until the entire file write is completed on the cluster
- D. They will see the file with its original name. If they attempt to view the file, they will get a `ConcurrentFileAccessException` until the entire file write is completed on the cluster

**Correct Answer:** B

**Section:** (none)

### Explanation

### Explanation/Reference:

answer is corrected.

### QUESTION 16

Which command does Hadoop offer to discover missing or corrupt HDFS data?

- A. `Hdfs fs du`
- B. `Hdfs fsck`
- C. `Dskchk`
- D. The map-only checksum
- E. Hadoop does not provide any tools to discover missing or corrupt data; there is no need because three replicas are kept for each data block

**Correct Answer:** B

**Section:** (none)

### Explanation

### Explanation/Reference:

Reference: <https://twiki.grid.iu.edu/bin/view/Storage/HadoopRecovery>

### QUESTION 17

You are planning a Hadoop cluster and considering implementing 10 Gigabit Ethernet as the network fabric. Which workloads benefit the most from faster network fabric?

- A. When your workload generates a large amount of output data, significantly larger than the amount of intermediate data
- B. When your workload consumes a large amount of input data, relative to the entire capacity of HDFS
- C. When your workload consists of processor-intensive tasks
- D. When your workload generates a large amount of intermediate data, on the order of the input data itself

**Correct Answer:** A

**Section:** (none)

**Explanation**

**Explanation/Reference:**

### QUESTION 18

Your cluster is running MapReduce version 2 (MRv2) on YARN. Your ResourceManager is configured to use the FairScheduler. Now you want to configure your scheduler such that a new user on the cluster can submit jobs into their own queue application submission. Which configuration should you set?



<http://www.gratisexam.com/>

- A. You can specify new queue name when user submits a job and new queue can be created dynamically if the property `yarn.scheduler.fair.allow-undecleared-pools = true`
- B. `Yarn.scheduler.fair.user.fair-as-default-queue = false` and `yarn.scheduler.fair.allow-undecleared-pools = true`
- C. You can specify new queue name when user submits a job and new queue can be created dynamically if `yarn.scheduler.fair.user-as-default-queue = false`
- D. You can specify new queue name per application in `allocations.xml` file and have new jobs automatically assigned to the application queue

**Correct Answer:** A

**Section:** (none)

**Explanation**

**Explanation/Reference:**

**QUESTION 19**

A slave node in your cluster has 4 TB hard drives installed (4 x 2TB). The DataNode is configured to store HDFS blocks on all disks. You set the value of the `dfs.datanode.du.reserved` parameter to 100 GB. How does this alter HDFS block storage?

- A. 25GB on each hard drive may not be used to store HDFS blocks
- B. 100GB on each hard drive may not be used to store HDFS blocks
- C. All hard drives may be used to store HDFS blocks as long as at least 100 GB in total is available on the node
- D. A maximum of 100 GB on each hard drive may be used to store HDFS blocks

**Correct Answer: C**

**Section: (none)**

**Explanation**

**Explanation/Reference:**

**QUESTION 20**

What two processes must you do if you are running a Hadoop cluster with a single NameNode and six DataNodes, and you want to change a configuration parameter so that it affects all six DataNodes. (Choose two)

- A. You must modify the configuration files on the NameNode only. DataNodes read their configuration from the master nodes
- B. You must modify the configuration files on each of the DataNodes machines
- C. You don't need to restart any daemon, as they will pick up changes automatically
- D. You must restart the NameNode daemon to apply the changes to the cluster
- E. You must restart all six DataNode daemons to apply the changes to the cluster

**Correct Answer: AD**

**Section: (none)**

**Explanation**

**Explanation/Reference:**

Typically one machine in the cluster is designated as the NameNode and another machine as the JobTracker, exclusively. These are the masters. The rest of the machines in the cluster act as both DataNode and TaskTracker. These are the slaves

**QUESTION 21**

You have installed a cluster HDFS and MapReduce version 2 (MRv2) on YARN. You have no dfs.hosts entry(ies) in your hdfs-site.xml configuration file. You configure a new worker node by setting fs.default.name in its configuration files to point to the NameNode on your cluster, and you start the DataNode daemon on that worker node. What do you have to do on the cluster to allow the worker node to join, and start sorting HDFS blocks?

- A. Without creating a dfs.hosts file or making any entries, run the commands `hadoop.dfsadmin- refreshModes` on the NameNode
- B. Restart the NameNode
- C. Creating a dfs.hosts file on the NameNode, add the worker Node's name to it, then issue the command `hadoop dfsadmin refresh Nodes =` on the Namenode
- D. Nothing; the worker node will automatically join the cluster when NameNode daemon is started

**Correct Answer:** A

**Section:** (none)

**Explanation**

**Explanation/Reference:**

#### QUESTION 22

You use the `hadoop fs put` command to add a file "sales.txt" to HDFS. This file is small enough that it fits into a single block, which is replicated to three nodes in your cluster (with a replication factor of 3). One of the nodes holding this file (a single block) fails. How will the cluster handle the replication of file in this situation?

- A. The file will remain under-replicated until the administrator brings that node back online
- B. The cluster will re-replicate the file the next time the system administrator reboots the NameNode daemon (as long as the file's replication factor doesn't fall below)
- C. This will be immediately re-replicated and all other HDFS operations on the cluster will halt until the cluster's replication values are resorted
- D. The file will be re-replicated automatically after the NameNode determines it is under- replicated based on the block reports it receives from the NameNodes

**Correct Answer:** D

**Section:** (none)

**Explanation**

**Explanation/Reference:**

The NameNode marks all blocks stored on the dead DataNode as under-replicated and orchestrates their rereplication. If the NameNode has marked a DataNode as dead, it will not include that DataNode in the list of machines it returns to clients as containing blocks they have requested. The NameNode will begin the process of re-replication, contacting a DataNode which contains a copy of each block which was on the now-dead DataNode and telling it to re-replicate that block to another DataNode.

#### QUESTION 23

Given:

```
[user1@host1 ~] yarn application -list
```

Total Applications: 3

Application ID	Application-Name	Application-Type	User	Queue	State	Final-State	Progress	Tracking
Application_1374638600275_0109	Sleep Job	MAPREDUCE	user1	KILLED	KILLED	KILLED	100%	host1:54059
Application_1374638600275_0121	Sleep Job	MAPREDUCE	user1	FINISHED	SUCCEEDED	SUCCEEDED	100%	host1:19888/Jobhistory/Job_1374638600275_0121
Application_1374638600275_0020	Sleep Job	MAPREDUCE	user1	FINISHED	SUCCEEDED	SUCCEEDED	100%	host1:19888/Jobhistory/Job_1374638600275_0020

You want to clean up this list by removing jobs where the State is KILLED. What command you enter?

- A. Yarn application refreshJobHistory
- B. Yarn application kill application\_1374638600275\_0109
- C. Yarn radmin refreshQueue
- D. Yarn radmin kill application\_1374638600275\_0109

**Correct Answer: B**

**Section: (none)**

**Explanation**

**Explanation/Reference:**

Reference: [http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.1-latest/bk\\_using-apache-hadoop/content/common\\_mrv2\\_commands.html](http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.1-latest/bk_using-apache-hadoop/content/common_mrv2_commands.html)

#### QUESTION 24

Assume you have a file named foo.txt in your local directory. You issue the following three commands:

```
Hadoop fs mkdir input
Hadoop fs put foo.txt input/foo.txt
Hadoop fs put foo.txt input
```

What happens when you issue the third command?

- A. The write succeeds, overwriting foo.txt in HDFS with no warning
- B. The file is uploaded and stored as a plain file named input
- C. You get a warning that foo.txt is being overwritten
- D. You get an error message telling you that foo.txt already exists, and asking you if you would like to overwrite it.
- E. You get a error message telling you that foo.txt already exists. The file is not written to HDFS
- F. You get an error message telling you that input is not a directory

G. The write silently fails

**Correct Answer:** CE

**Section:** (none)

**Explanation**

**Explanation/Reference:**

#### **QUESTION 25**

You are configuring a server running HDFS, MapReduce version 2 (MRv2) on YARN running Linux. How must you format underlying file system of each DataNode?

- A. They must be formatted as HDFS
- B. They must be formatted as either ext3 or ext4
- C. They may be formatted in any Linux file system
- D. They must not be formatted - - HDFS will format the file system automatically

**Correct Answer:** B

**Section:** (none)

**Explanation**

**Explanation/Reference:**

#### **QUESTION 26**

You are migrating a cluster from MapReduce version 1 (MRv1) to MapReduce version 2 (MRv2) on YARN. You want to maintain your MRv1 TaskTracker slot capacities when you migrate. What should you do?

- A. Configure `yarn.applicationmaster.resource.memory-mb` and `yarn.applicationmaster.resource.cpu-vcores` so that ApplicationMaster container allocations match the capacity you require.
- B. You don't need to configure or balance these properties in YARN as YARN dynamically balances resource management capabilities on your cluster
- C. Configure `mapred.tasktracker.map.tasks.maximum` and `mapred.tasktracker.reduce.tasks.maximum` in `yarn-site.xml` to match your cluster's capacity set by the `yarn-scheduler.minimum-allocation`
- D. Configure `yarn.nodemanager.resource.memory-mb` and `yarn.nodemanager.resource.cpu-vcores` to match the capacity you require under YARN for each NodeManager

**Correct Answer:** D

**Section:** (none)

**Explanation**

**Explanation/Reference:**

**QUESTION 27**

You are running Hadoop cluster with all monitoring facilities properly configured.

Which scenario will go undeselected?

- A. HDFS is almost full
- B. The NameNode goes down
- C. A DataNode is disconnected from the cluster
- D. Map or reduce tasks that are stuck in an infinite loop
- E. MapReduce jobs are causing excessive memory swaps

**Correct Answer:** B

**Section:** (none)

**Explanation**

**Explanation/Reference:**

**QUESTION 28**

You decide to create a cluster which runs HDFS in High Availability mode with automatic failover, using Quorum Storage. What is the purpose of ZooKeeper in such a configuration?

- A. It only keeps track of which NameNode is Active at any given time
- B. It monitors an NFS mount point and reports if the mount point disappears
- C. It both keeps track of which NameNode is Active at any given time, and manages the Edits file. Which is a log of changes to the HDFS filesystem
- D. It only manages the Edits file, which is log of changes to the HDFS filesystem
- E. Clients connect to ZooKeeper to determine which NameNode is Active

**Correct Answer:** A

**Section:** (none)

**Explanation**

**Explanation/Reference:**

Reference: [http://www.cloudera.com/content/cloudera-content/cloudera- \(page 15\) docs/CDH4/latest/PDF/CDH4-High-Availability-Guide.pdf](http://www.cloudera.com/content/cloudera-content/cloudera- (page 15) docs/CDH4/latest/PDF/CDH4-High-Availability-Guide.pdf)

**QUESTION 29**

Choose three reasons why should you run the HDFS balancer periodically? (Choose three)

- A. To ensure that there is capacity in HDFS for additional data
- B. To ensure that all blocks in the cluster are 128MB in size
- C. To help HDFS deliver consistent performance under heavy loads
- D. To ensure that there is consistent disk utilization across the DataNodes
- E. To improve data locality MapReduce

**Correct Answer:** CDE

**Section:** (none)

**Explanation**

**Explanation/Reference:**

<http://www.quora.com/Apache-Hadoop/It-is-recommended-that-you-run-the-HDFS-balancer-periodically-Why-Choose-3>

**QUESTION 30**

Your cluster implements HDFS High Availability (HA). Your two NameNodes are named nn01 and nn02. What occurs when you execute the command: `hdfs haadmin failover nn01 nn02`?

- A. nn02 is fenced, and nn01 becomes the active NameNode
- B. nn01 is fenced, and nn02 becomes the active NameNode
- C. nn01 becomes the standby NameNode and nn02 becomes the active NameNode
- D. nn02 becomes the standby NameNode and nn01 becomes the active NameNode

**Correct Answer:** B

**Section:** (none)

**Explanation**

**Explanation/Reference:**

Explanation:

`failover` initiate a failover between two NameNodes

This subcommand causes a failover from the first provided NameNode to the second. If the first NameNode is in the Standby state, this command simply transitions the second to the Active state without error. If the first NameNode is in the Active state, an attempt will be made to gracefully transition it to the Standby state. If this fails, the fencing methods (as configured by `dfs.ha.fencing.methods`) will be attempted in order until one of the methods succeeds. Only after this process will the second NameNode be transitioned to the Active state. If no fencing method succeeds, the second NameNode will not be transitioned to the Active state, and an error will be returned.



**QUESTION 31**

You have a Hadoop cluster HDFS, and a gateway machine external to the cluster from which clients submit jobs. What do you need to do in order to run Impala on the cluster and submit jobs from the command line of the gateway machine?

- A. Install the impalad daemon, statestored daemon, and daemon on each machine in the cluster, and the impala shell on your gateway machine
- B. Install the impalad daemon, the statestored daemon, the catalogd daemon, and the impala shell on your gateway machine
- C. Install the impalad daemon and the impala shell on your gateway machine, and the statestored daemon and catalogd daemon on one of the nodes in the cluster
- D. Install the impalad daemon on each machine in the cluster, the statestored daemon and catalogd daemon on one machine in the cluster, and the impala shell on your gateway machine
- E. Install the impalad daemon, statestored daemon, and catalogd daemon on each machine in the cluster and on the gateway node

**Correct Answer:** D

**Section:** (none)

**Explanation**

**Explanation/Reference:**

**QUESTION 32**

You have just run a MapReduce job to filter user messages to only those of a selected geographical region. The output for this job is in a directory named westUsers, located just below your home directory in HDFS. Which command gathers these into a single file on your local file system?

- A. Hadoop fs getmerge R westUsers.txt
- B. Hadoop fs getemerge westUsers westUsers.txt
- C. Hadoop fs cp westUsers/\* westUsers.txt
- D. Hadoop fs get westUsers westUsers.txt

**Correct Answer:** B

**Section:** (none)

**Explanation**

**Explanation/Reference:**

**QUESTION 33**

In CDH4 and later, which file contains a serialized form of all the directory and files inodes in the filesystem, giving the NameNode a persistent checkpoint of the filesystem metadata?



<http://www.gratisexam.com/>

- A. fstime
- B. VERSION
- C. Fsimage\_N (where N reflects transactions up to transaction ID N)
- D. Edits\_N-M (where N-M transactions between transaction ID N and transaction ID N)

**Correct Answer:** C

**Section:** (none)

**Explanation**

**Explanation/Reference:**

Reference: <http://mikepluta.com/tag/namenode/>

#### QUESTION 34

You are running a Hadoop cluster with a NameNode on host mynamenode. What are two ways to determine available HDFS space in your cluster?

- A. Run `hdfs fs du /` and locate the DFS Remaining value
- B. Run `hdfs dfsadmin report` and locate the DFS Remaining value
- C. Run `hdfs dfs /` and subtract NDFS Used from configured Capacity
- D. Connect to `http://mynamenode:50070/dfshealth.jsp` and locate the DFS remaining value

**Correct Answer:** B

**Section:** (none)

**Explanation**

**Explanation/Reference:**

#### QUESTION 35

You have recently converted your Hadoop cluster from a MapReduce 1 (MRv1) architecture to MapReduce 2 (MRv2) on YARN architecture. Your developers are accustomed to specifying map and reduce tasks (resource allocation) tasks when they run jobs: A developer wants to know how specify to reduce tasks when a

specific job runs. Which method should you tell that developers to implement?

- A. MapReduce version 2 (MRv2) on YARN abstracts resource allocation away from the idea of "tasks" into memory and virtual cores, thus eliminating the need for a developer to specify the number of reduce tasks, and indeed preventing the developer from specifying the number of reduce tasks.
- B. In YARN, resource allocations is a function of megabytes of memory in multiples of 1024mb.  
Thus, they should specify the amount of memory resource they need by executing `D mapreduce-reduces.memory-mb-2048`
- C. In YARN, the ApplicationMaster is responsible for requesting the resource required for a specific launch. Thus, executing `D yarn.applicationmaster.reduce.tasks=2` will specify that the ApplicationMaster launch two task contains on the worker nodes.
- D. Developers specify reduce tasks in the exact same way for both MapReduce version 1 (MRv1) and MapReduce version 2 (MRv2) on YARN. Thus, executing `D mapreduce.job.reduces-2` will specify reduce tasks.
- E. In YARN, resource allocation is function of virtual cores specified by the ApplicationManager making requests to the NodeManager where a reduce task is handled by a single container (and thus a single virtual core). Thus, the developer needs to specify the number of virtual cores to the NodeManager by executing `p yarn.nodemanager.cpu-vcORES=2`

**Correct Answer:** D

**Section:** (none)

**Explanation**

**Explanation/Reference:**

### QUESTION 36

Your Hadoop cluster contains nodes in three racks. You have not configured the `dfs.hosts` property in the NameNode's configuration file. What results?

- A. The NameNode will update the `dfs.hosts` property to include machines running the DataNode daemon on the next NameNode reboot or with the command `dfsadmin refreshNodes`
- B. No new nodes can be added to the cluster until you specify them in the `dfs.hosts` file
- C. Any machine running the DataNode daemon can immediately join the cluster
- D. Presented with a blank `dfs.hosts` property, the NameNode will permit DataNodes specified in `mapred.hosts` to join the cluster

**Correct Answer:** C

**Section:** (none)

**Explanation**

**Explanation/Reference:**

### QUESTION 37

You are running a Hadoop cluster with MapReduce version 2 (MRv2) on YARN. You consistently see that MapReduce map tasks on your cluster are running slowly

because of excessive garbage collection of JVM, how do you increase JVM heap size property to 3GB to optimize performance?

- A. yarn.application.child.java.opts=-Xsx3072m
- B. yarn.application.child.java.opts=-Xmx3072m
- C. mapreduce.map.java.opts=-Xms3072m
- D. mapreduce.map.java.opts=-Xmx3072m

**Correct Answer:** C

**Section:** (none)

**Explanation**

**Explanation/Reference:**

Reference: <http://hortonworks.com/blog/how-to-plan-and-configure-yarn-in-hdp-2-0/>

### QUESTION 38

You have a cluster running with a FIFO scheduler enabled. You submit a large job A to the cluster, which you expect to run for one hour. Then, you submit job B to the cluster, which you expect to run a couple of minutes only.

You submit both jobs with the same priority.

Which two best describes how FIFO Scheduler arbitrates the cluster resources for job and its tasks? (Choose two)

- A. Because there is a more than a single job on the cluster, the FIFO Scheduler will enforce a limit on the percentage of resources allocated to a particular job at any given time
- B. Tasks are scheduled on the order of their job submission
- C. The order of execution of job may vary
- D. Given job A and submitted in that order, all tasks from job A are guaranteed to finish before all tasks from job B
- E. The FIFO Scheduler will give, on average, and equal share of the cluster resources over the job lifecycle
- F. The FIFO Scheduler will pass an exception back to the client when Job B is submitted, since all slots on the cluster are use

**Correct Answer:** AD

**Section:** (none)

**Explanation**

**Explanation/Reference:**

### QUESTION 39

A user comes to you, complaining that when she attempts to submit a Hadoop job, it fails. There is a Directory in HDFS named /data/input. The Jar is named j.jar, and the driver class is named DriverClass.

She runs the command:

```
Hadoop jar j.jar DriverClass /data/input/data/output
```

The error message returned includes the line:  
PrivilegedActionException as:training (auth:SIMPLE)  
cause:org.apache.hadoop.mapreduce.lib.input.InvalidInputException:  
Input path does not exist: file:/data/input

What is the cause of the error?

- A. The user is not authorized to run the job on the cluster
- B. The output directory already exists
- C. The name of the driver has been spelled incorrectly on the command line
- D. The directory name is misspelled in HDFS
- E. The Hadoop configuration files on the client do not point to the cluster

**Correct Answer:** A

**Section:** (none)

**Explanation**

**Explanation/Reference:**

answer is up-to-date.

#### **QUESTION 40**

Your company stores user profile records in an OLTP databases. You want to join these records with web server logs you have already ingested into the Hadoop file system. What is the best way to obtain and ingest these user records?

- A. Ingest with Hadoop streaming
- B. Ingest using Hive's IQAD DATA command
- C. Ingest with sqoop import
- D. Ingest with Pig's LOAD command
- E. Ingest using the HDFS put command

**Correct Answer:** C

**Section:** (none)

**Explanation**

**Explanation/Reference:**

**QUESTION 41**

Which two are features of Hadoop's rack topology? (Choose two)

- A. Configuration of rack awareness is accomplished using a configuration file. You cannot use a rack topology script.
- B. Hadoop gives preference to intra-rack data transfer in order to conserve bandwidth
- C. Rack location is considered in the HDFS block placement policy
- D. HDFS is rack aware but MapReduce daemon are not
- E. Even for small clusters on a single rack, configuring rack awareness will improve performance

**Correct Answer:** BC

**Section:** (none)

**Explanation**

**Explanation/Reference:**

**QUESTION 42**

Each node in your Hadoop cluster, running YARN, has 64GB memory and 24 cores. Your yarn.site.xml has the following configuration:

```
<property>
<name>yarn.nodemanager.resource.memory-mb</name>
<value>32768</value>
</property>
<property>
<name>yarn.nodemanager.resource.cpu-vcores</name>
<value>12</value>
</property>
```

You want YARN to launch no more than 16 containers per node. What should you do?

- A. Modify yarn-site.xml with the following property:  
<name>yarn.scheduler.minimum-allocation-mb</name>  
<value>2048</value>
- B. Modify yarn-sites.xml with the following property:  
<name>yarn.scheduler.minimum-allocation-mb</name>  
<value>4096</value>
- C. Modify yarn-site.xml with the following property:

<name>yarn.nodemanager.resource.cpu-vccores</name>

D. No action is needed: YARN's dynamic resource allocation automatically optimizes the node memory and cores

**Correct Answer:** A

**Section:** (none)

**Explanation**

**Explanation/Reference:**

#### **QUESTION 43**

You want to node to only swap Hadoop daemon data from RAM to disk when absolutely necessary. What should you do?

- A. Delete the /dev/vmswap file on the node
- B. Delete the /etc/swap file on the node
- C. Set the ram.swap parameter to 0 in core-site.xml
- D. Set vm.swapfile file on the node
- E. Delete the /swapfile file on the node

**Correct Answer:** D

**Section:** (none)

**Explanation**

**Explanation/Reference:**

#### **QUESTION 44**

You are configuring your cluster to run HDFS and MapReducer v2 (MRv2) on YARN. Which two daemons needs to be installed on your cluster's master nodes? (Choose two)

- A. HMaster
- B. ResourceManager
- C. TaskManager
- D. JobTracker
- E. NameNode
- F. DataNode

**Correct Answer:** BE

**Section: (none)**

**Explanation**

**Explanation/Reference:**

answer is perfected.

**QUESTION 45**

You observed that the number of spilled records from Map tasks far exceeds the number of map output records. Your child heap size is 1GB and your io.sort.mb value is set to 1000MB. How would you tune your io.sort.mb value to achieve maximum memory to disk I/O ratio?

- A. For a 1GB child heap size an io.sort.mb of 128 MB will always maximize memory to disk I/O
- B. Increase the io.sort.mb to 1GB
- C. Decrease the io.sort.mb value to 0
- D. Tune the io.sort.mb value until you observe that the number of spilled records equals (or is as close to equals) the number of map output records.

**Correct Answer: D**

**Section: (none)**

**Explanation**

**Explanation/Reference:**

**QUESTION 46**

You are running a Hadoop cluster with a NameNode on host mynamenode, a secondary NameNode on host mysecondarynamenode and several DataNodes.

Which best describes how you determine when the last checkpoint happened?

- A. Execute hdfs namenode report on the command line and look at the Last Checkpoint information
- B. Execute hdfs dfsadmin saveNamespace on the command line which returns to you the last checkpoint value in fstime file
- C. Connect to the web UI of the Secondary NameNode (<http://mysecondary:50090/>) and look at the "Last Checkpoint" information
- D. Connect to the web UI of the NameNode (<http://mynamenode:50070/>) and look at the "Last Checkpoint" information

**Correct Answer: C**

**Section: (none)**

**Explanation**

**Explanation/Reference:**

Reference: <https://www.inkling.com/read/hadoop-definitive-guide-tom-white-3rd/chapter-10/hdfs>



**QUESTION 47**

What does CDH packaging do on install to facilitate Kerberos security setup?

- A. Automatically configures permissions for log files at & MAPRED\_LOG\_DIR/userlogs
- B. Creates users for hdfs and mapreduce to facilitate role assignment
- C. Creates directories for temp, hdfs, and mapreduce with the correct permissions
- D. Creates a set of pre-configured Kerberos keytab files and their permissions
- E. Creates and configures your kdc with default cluster values

**Correct Answer:** B

**Section:** (none)

**Explanation**

**Explanation/Reference:**

**QUESTION 48**

Which process instantiates user code, and executes map and reduce tasks on a cluster running MapReduce v2 (MRv2) on YARN?

- A. NodeManager
- B. ApplicationMaster
- C. TaskTracker
- D. JobTracker
- E. NameNode
- F. DataNode
- G. ResourceManager

**Correct Answer:** A

**Section:** (none)

**Explanation**

**Explanation/Reference:**

**QUESTION 49**

Cluster Summary:

45 files and directories, 12 blocks = 57 total. Heap size is 15.31 MB/193.38MB(7%)

<b>Configured capacity</b>	<b>:</b>	<b>17.33GB</b>
<b>DFS Used</b>	<b>:</b>	<b>144KB</b>
<b>Non DFS Used</b>	<b>:</b>	<b>5.49GB</b>
<b>DFS Remaining</b>	<b>:</b>	<b>11.84GB</b>
<b>DFS Used %</b>	<b>:</b>	<b>0%</b>
<b>DFS Remaining %</b>	<b>:</b>	<b>68.32GB</b>
<b>Live Nodes</b>	<b>:</b>	<b>6</b>
<b>Dead Nodes</b>	<b>:</b>	<b>1</b>
<b>Decommissioning Nodes</b>	<b>:</b>	<b>0</b>
<b>Number of Under-Replicated Blocks</b>	<b>:</b>	<b>6</b>

Refer to the above screenshot.

You configure a Hadoop cluster with seven DataNodes and one of your monitoring UIs displays the details shown in the exhibit.

What does this tell you?

- A. The DataNode JVM on one host is not active
- B. Because your under-replicated blocks count matches the Live Nodes, one node is dead, and your DFS Used % equals 0%, you can't be certain that your cluster has all the data you've written it.
- C. Your cluster has lost all HDFS data which had blocks stored on the dead DataNode
- D. The HDFS cluster is in safe mode

**Correct Answer: A**

**Section: (none)**

**Explanation**

**Explanation/Reference:**

#### **QUESTION 50**

Which two features does Kerberos security add to a Hadoop cluster? (Choose two)

- A. User authentication on all remote procedure calls (RPCs)
- B. Encryption for data during transfer between the Mappers and Reducers

- C. Encryption for data on disk ("at rest")
- D. Authentication for user access to the cluster against a central server
- E. Root access to the cluster for users hdfs and mapred but non-root access for clients

**Correct Answer:** AD

**Section:** (none)

**Explanation**

**Explanation/Reference:**

#### **QUESTION 51**

Assuming a cluster running HDFS, MapReduce version 2 (MRv2) on YARN with all settings at their default, what do you need to do when adding a new slave node to cluster?

- A. Nothing, other than ensuring that the DNS (or/etc/hosts files on all machines) contains any entry for the new node.
- B. Restart the NameNode and ResourceManager daemons and resubmit any running jobs.
- C. Add a new entry to /etc/nodes on the NameNode host.
- D. Restart the NameNode of dfs.number.of.nodes in hdfs-site.xml

**Correct Answer:** A

**Section:** (none)

**Explanation**

**Explanation/Reference:**

[http://wiki.apache.org/hadoop/FAQ#I\\_have\\_a\\_new\\_node\\_I\\_want\\_to\\_add\\_to\\_a\\_running\\_Hadoop\\_cluster.3B\\_how\\_do\\_I\\_start\\_services\\_on\\_just\\_one\\_node.3F](http://wiki.apache.org/hadoop/FAQ#I_have_a_new_node_I_want_to_add_to_a_running_Hadoop_cluster.3B_how_do_I_start_services_on_just_one_node.3F)

#### **QUESTION 52**

Assuming you're not running HDFS Federation, what is the maximum number of NameNode daemons you should run on your cluster in order to avoid a "split-brain" scenario with your NameNode when running HDFS High Availability (HA) using Quorum-based storage?

- A. Two active NameNodes and two Standby NameNodes
- B. One active NameNode and one Standby NameNode
- C. Two active NameNodes and one Standby NameNode
- D. Unlimited. HDFS High Availability (HA) is designed to overcome limitations on the number of NameNodes you can deploy

**Correct Answer:** B

**Section:** (none)

## Explanation

### Explanation/Reference:

#### QUESTION 53

Table schemas in Hive are:

- A. Stored as metadata on the NameNode
- B. Stored along with the data in HDFS
- C. Stored in the Metadata
- D. Stored in ZooKeeper

**Correct Answer:** B

**Section:** (none)

### Explanation

### Explanation/Reference:

<http://stackoverflow.com/questions/22989592/how-to-get-hive-table-name-based-on-hdfs-location-path-without-connecting-to-m>

#### QUESTION 54

For each YARN job, the Hadoop framework generates task log file. Where are Hadoop task log files stored?

- A. Cached by the NodeManager managing the job containers, then written to a log directory on the NameNode
- B. Cached in the YARN container running the task, then copied into HDFS on job completion
- C. In HDFS, in the directory of the user who generates the job
- D. On the local disk of the slave node running the task

**Correct Answer:** D

**Section:** (none)

### Explanation

### Explanation/Reference:

#### QUESTION 55

You have a cluster running with the fair Scheduler enabled. There are currently no jobs running on the cluster, and you submit a job A, so that only job A is running on the cluster. A while later, you submit Job B. now Job A and Job B are running on the cluster at the same time. How will the Fair Scheduler handle these two jobs? (Choose two)

- A. When Job B gets submitted, it will get assigned tasks, while job A continues to run with fewer tasks.
- B. When Job B gets submitted, Job A has to finish first, before job B can get scheduled.
- C. When Job A gets submitted, it doesn't consume all the task slots.
- D. When Job A gets submitted, it consumes all the task slots.

**Correct Answer:** B

**Section:** (none)

**Explanation**

**Explanation/Reference:**

corrected.



<http://www.gratisexam.com/>