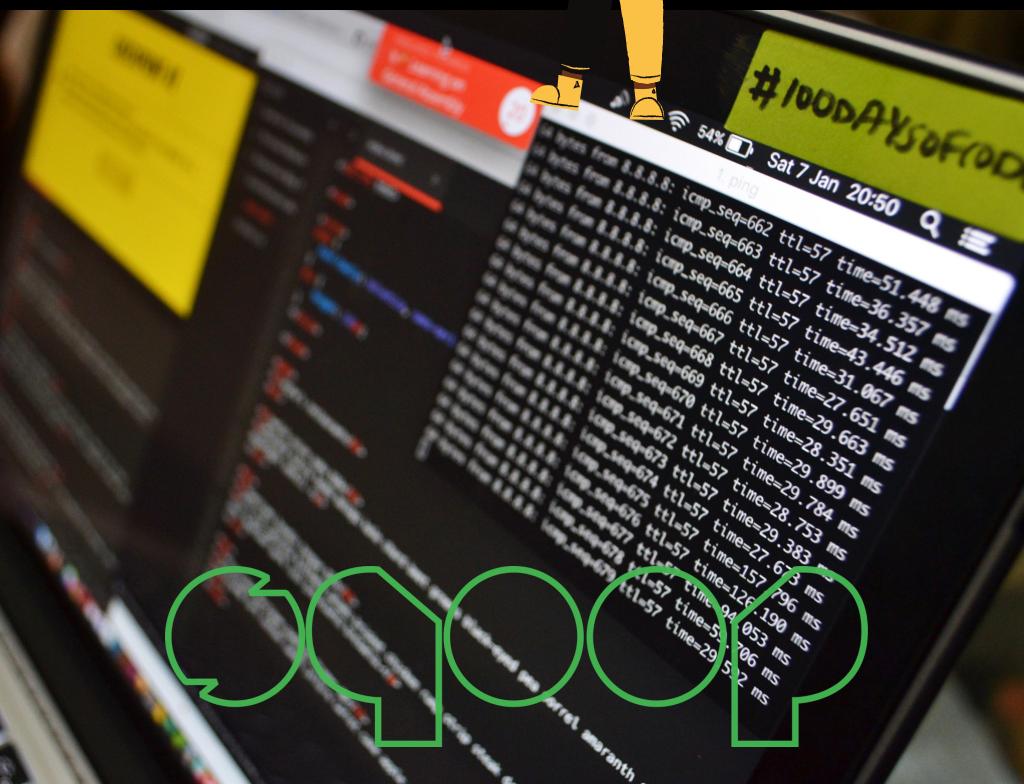


COMPREHENSIVE GUIDE TO INTERVIEWS FOR SQOOP FOR BIG DATA



ZEP ANALYTICS

Introduction

We've curated this series of interview which guides to accelerate your learning and your mastery of data science skills and tools.

From job-specific technical questions to tricky behavioral inquiries and unexpected brainteasers and guesstimates, we will prepare you for any job candidacy in the fields of data science, data analytics, or BI analytics and Big Data.

These guides are the result of our data analytics expertise, direct experience interviewing at companies, and countless conversations with job candidates. Its goal is to teach by example - not only by giving you a list of interview questions and their answers, but also by sharing the techniques and thought processes behind each question and the expected answer.

Become a global tech talent and unleash your next, best self with all the knowledge and tools to succeed in a data analytics interview with this series of guides.

COMPREHENSIVE GUIDE TO INTERVIEWS FOR DATA SCIENCE



Data Science interview questions cover a wide scope of multidisciplinary topics. That means you can never be quite sure what challenges the interviewer(s) might send your way.

That being said, being familiar with the type of questions you can encounter is an important aspect of your preparation process.

Below you'll find examples of real-life questions and answers. Reviewing those should help you assess the areas you're confident in and where you should invest additional efforts to improve.

Become a Tech Blogger at Zep!!

Why don't you start your journey as a blogger and enjoy unlimited free perks and cash prizes every month.

[Explore](#)



ZEP ANALYTICS

1. What is SQOOP..?

This is the short meaning of (SQL+HadOOP =SQOOP)
It is a tool designed to transfer data between Hadoop and relational databases or mainframes. You can use Sqoop to import data from a relational database management system (RDBMS) such as MySQL or Oracle or a mainframe into the Hadoop Distributed File System (HDFS), transform the data in Hadoop MapReduce, and then export the data back into an RDBMS.

Sqoop automates most of this process, relying on the database to describe the schema for the data to be imported. Sqoop uses MapReduce to import and export the data, which provides parallel operation as well as fault tolerance.

The Sqoop main intended for:

- System and application programmers
- System administrators Database administrators Data analysts
- Data engineers

2. Why is the default maximum mappers are 4 in Sqoop?

As of my knowledge, the default number of mapper 4 is followed by minimum concurrent task for one machine. We will lead to set a higher number of concurrent tasks, which can result in faster job completion.

3. Is it possible set speculative execution in Sqoop ..?

In sqoop by default speculative execution is off, because if Multiple mappers run for single task, we get duplicates of data in HDFS. Hence to avoid this decrepency it is off. Also number of reducers for sqoop job is 0, since it is merely a job running a MAP only job that dumps data into HDFS. We are not aggregating anything.

4. What causes of hadoop throw ClassNotFoundException while sqoop integration ..?

The most causes of that the supporting library (like connectors) was not updated in sqoop's library path, so we need to update it on that specific path.

5. How to view all the databases and tables in RDBMS from SQOOP..?

Using below commands we can,

sqoop-list-databases

sqoop-list-tables

6. How to view table columns details in RDBMS from SQOOP..?

Unfortunately we don't have any commands like sqoop-list-columns, But we can achieve via free form query to check the information schema for the particular RDBMS tables.here is an example:

```
$ sqoop eval --connect 'jdbc:mysql://nameofmyserver;'  
database=nameofmydatabase; username=dineshkumar;  
password=dineshkumar --query "SELECT column_name,  
DATA_TYPE FROM  
INFORMATION_SCHEMA.Columns WHERE  
table_name='mytableofinterest'
```

7. I am getting FileAlreadyExists exception error in Sqoop while importing data from RDBMS to a hive table.? So How do we resolve it.?

you can specify the --hive-overwrite option to indicate that existing table in hive must be replaced. After your data is imported into HDFS or this step is omitted

8. What is the default file format to import data using Apache Sqoop?

Sqoop allows data to be imported using two file formats
i) Delimited Text File Format

This is the default file format to import data using Sqoop. This file format can be explicitly specified using the -as-textfile argument to the import command in Sqoop. Passing this as an argument to the command will produce the string based representation of all the records to the output files with the delimited characters between rows and columns.

ii) Sequence File Format

It is a binary file format where records are stored in custom record-specific data types which are shown as Java classes. Sqoop automatically creates these data types and manifests them as java classes.

9. How do I resolve a Communications Link Failure when connecting to MySQL?

Verify that you can connect to the database from the node where you are running Sqoop:

```
$ mysql --host=IP Address --database=test --  
user=username --password=password
```

Add the network port for the server to your my.cnf file
Set up a user account to connect via Sqoop. Grant permissions to the user to access the database over the network:

Log into MySQL as root mysql -u root -p
ThisIsMyPassword

Issue the following command: mysql> grant all privileges on test.* to 'testuser'@'%' identified by 'testpassword'

10. How do I resolve an IllegalArgumentException when connecting to Oracle?

This could be caused by a non-owner trying to connect to the table so prefix the table name with the schema, for example SchemaName.OracleTableName.

11. What's causing this Exception in thread main
java.lang.IncompatibleClassChangeError
when running non-CDH Hadoop with Sqoop?

Try building Sqoop 1.4.1-incubating with the command line property -Dhadoopversion=20.

12. I have around 300 tables in a database. I want to import all the tables from the database except the tables named Table298, Table 123, and Table299. How can I do this without having to import the tables one by one?

This can be accomplished using the import-all-tables import command in Sqoop and by specifying the exclude-tables option with it as follows-

```
sqoop import-all-tables --connect -username -password --exclude-tables Table298, Table 123, Table 299
```

13. Does Apache Sqoop have a default database?

Yes, MySQL is the default database.

bigdatascholars.blogspot.com/2018/08/sqoop-interview-question-and-answers.html

14. How can I import large objects (BLOB and CLOB objects) in Apache Sqoop?

Apache Sqoop import command does not support direct import of BLOB and CLOB large objects. To import large objects, JDBC based imports have to be used without the direct argument to the import utility.

15. How can you execute a free form SQL query in Sqoop to import the rows in a sequential manner?

This can be accomplished using the -m 1 option in the Sqoop import command. It will create only one MapReduce task which will then import rows serially.

16. What is the difference between Sqoop and DistCP command in Hadoop?

Both distCP (Distributed Copy in Hadoop) and Sqoop transfer data in parallel but the only difference is that distCP command can transfer any kind of data from one Hadoop cluster to another whereas Sqoop transfers data between RDBMS and other components in the Hadoop ecosystem like HBase, Hive, HDFS, etc.

17. What is Sqoop metastore?

Sqoop metastore is a shared metadata repository for remote users to define and execute saved jobs created using sqoop job defined in the metastore. The sqoop-site.xml should be configured to connect to the metastore.

18. You use -split-by clause but it still does not give optimal performance then how will you improve the performance further?

Using the -boundary-query clause. Generally, sqoop uses the SQL query select min (), max () from to find out the boundary values for creating splits. However, if this query is not optimal then using the -boundary-query argument any random query can be written to generate two numeric columns.

19. What is the significance of using -split-by clause for running parallel import tasks in Apache Sqoop?

--Split-by clause is used to specify the columns of the table that are used to generate splits for data imports. This clause specifies the columns that will be used for splitting when importing the data into the Hadoop cluster. – split-by clause helps achieve improved performance through greater parallelism. Apache Sqoop will create splits based on the values present in the columns specified in the -split-by clause of the import command. If the -split-by clause is not specified, then the primary key of the table is used to create the splits while data import. At times the primary key of the table might not have evenly distributed values between the minimum and maximum range. Under such circumstances -split-by clause can be used to specify some other column that has even distribution of data to create splits so that data import is efficient.

20. During sqoop import, you use the clause -m or --num-mappers to specify the number of mappers as 8 so that it can run eight parallel MapReduce tasks, however, sqoop runs only four parallel MapReduce tasks. Why?

Hadoop MapReduce cluster is configured to run a maximum of 4 parallel MapReduce tasks and the sqoop import can be configured with number of parallel tasks less than or equal to 4 but not more than 4.

21. You successfully imported a table using Apache Sqoop to HBase but when you query the table it is found that the number of rows is less than expected. What could be the likely reason?

If the imported records have rows that contain null values for all the columns, then probably those records might have been dropped off during import because HBase does not allow null values in all the columns of a record.

22. The incoming value from HDFS for a particular column is NULL. How will you load that row into RDBMS in which the columns are defined as NOT NULL?

Using the -input-null-string parameter, a default value can be specified so that the row gets inserted with the default value for the column that it has a NULL value in HDFS.

23. How will you synchronize the data in HDFS that is imported by Sqoop?

Data can be synchronised using incremental parameter with data import ---Incremental parameter can be used with one of the two options-
i) append-If the table is getting updated continuously with new rows and increasing row id values then incremental import with append option should be used where values of some of the columns are checked (columns to be checked are specified using -check-column) and if it discovers

any modified value for those columns then only a new row will be inserted.

ii) lastmodified - In this kind of incremental import, the source has a date column which is checked for. Any records that have been updated after the last import based on the lastmodified column in the source, the values would be updated.

24. What are the relational databases supported in Sqoop?

Below are the list of RDBMSs that are supported by Sqoop Currently.

MySQL

PostGreSQL

Oracle

Microsoft SQL

IBM's Netezza

Teradata

25. What are the destination types allowed in Sqoop Import command?

Currently Sqoop Supports data imported into below services.

HDFS

Hive

HBase

HCatalog

Accumulo

26. Is Sqoop similar to distcp in hadoop?

Partially yes, hadoop's distcp command is similar to Sqoop Import command. Both submits parallel map-only jobs.

But distcp is used to copy any type of files from Local FS/HDFS to HDFS and Sqoop is for transferring the data records only between RDMBS and Hadoop eco system services, HDFS, Hive and HBase.

27. What are the majorly used commands in Sqoop?

In Sqoop Majorly Import and export commands are used. But below commands are also useful some times.

codegen

eval

import-all-tables

job

list-databases

list-tables

merge

metastore

28. While loading tables from MySQL into HDFS, if we need to copy tables with maximum possible speed, what can you do ?

We need to use -direct argument in import command to use direct import fast path and this -direct can be used only with MySQL and PostGreSQL as of now.

29. While connecting to MySQL through Sqoop, I am getting Connection Failure exception what might be the root cause and fix for this error scenario?

This might be due to insufficient permissions to access your MySQL database over the network. To confirm this we can try the below command to connect to MySQL database from Sqoop's client machine.

```
$ mysql --host=MySQL node > --database=test --  
user= --password=
```

30. What is the importance of eval tool?

It allow users to run sample SQL queries against Database and preview the result on the console.

31. What is the process to perform an incremental data load in Sqoop?

The process to perform incremental data load in Sqoop is to synchronize the modified or updated data (often referred as delta data) from RDBMS to Hadoop. The delta data can be facilitated through the incremental load command in Sqoop.

Incremental load can be performed by using Sqoop import command or by loading the data into hive without overwriting it. The different attributes that need to be specified during incremental load in Sqoop are-

- 1) Mode (incremental) -The mode defines how Sqoop will determine what the new rows are. The mode can have value as Append or Last Modified.
- 2) Col (Check-column) -This attribute specifies the column that should be examined to find out the rows to be imported.
- 3) Value (last-value) -This denotes the maximum value of the check column from the previous import operation.

32. What is the significance of using -compress-codec parameter?

To get the out file of a sqoop import in formats other than .gz like .bz2 compressions when we use the -compress -code parameter.

33. Can free form SQL queries be used with Sqoop import command? If yes, then how can they be used?

Sqoop allows us to use free form SQL queries with the import command. The import command should be used with the -e and -query options to execute free form SQL queries. When using the -e and -query options with the import command the -target-dir value must be specified.

34. What is the purpose of sqoop-merge?

The merge tool combines two datasets where entries in one dataset should overwrite entries of an older dataset preserving only the newest version of the records between both the data sets.

35. How do you clear the data in a staging table before loading it by Sqoop?

By specifying the -clear-staging-table option we can clear the staging table before it is loaded. This can be done again and again till we get proper data in staging.

36. How will you update the rows that are already exported?

The parameter -update-key can be used to update existing rows. In a comma-separated list of columns is used which uniquely identifies a row. All of these columns is used in the WHERE clause of the generated UPDATE query. All other table columns will be used in the SET part of the query.

37. What is the role of JDBC driver in a Sqoop set up?

To connect to different relational databases sqoop needs a connector. Almost every DB vendor makes this connector available as a JDBC driver which is specific to that DB. So Sqoop needs the JDBC driver of each of the database it needs to interact with.

38. When to use --target-dir and --warehouse-dir while importing data?

To specify a particular directory in HDFS use --target-dir but to specify the parent directory of all the sqoop jobs use --warehouse-dir. In this case under the parent directory sqoop will create a directory with the same name as the table.

39. When the source data keeps getting updated frequently, what is the approach to keep it in sync with the data in HDFS imported by sqoop?

sqoop can have 2 approaches.

To use the --incremental parameter with append option where value of some columns are checked and only in case of modified values the row is imported as a new row.

To use the --incremental parameter with lastmodified option where a date column in the source is checked for records which have been updated after the last import.

40. sqoop takes a long time to retrieve the minimum and maximum values of columns mentioned in --split-by parameter. How can we make it efficient?

We can use the --boundary-query parameter in which we specify the min and max value for the column based on which the split can happen into multiple mapreduce tasks. This makes it faster as the query inside the --boundary-query parameter is executed first and the job is ready with the information on how many mapreduce tasks to create before executing the main query.

41. Is it possible to add a parameter while running a saved job?

Yes, we can add an argument to a saved job at runtime by using the --exec option sqoop job --exec jobname -- -- newparameter.

42. How will you implement all-or-nothing load using sqoop?

Using the staging-table option we first load the data into a staging table and then load it to the final target table only if the staging load is successful.

43. How will you update the rows that are already exported ?

The parameter --update-key can be used to update existing rows. In it a comma-separated list of columns is used which uniquely identifies a row. All of these columns is used in the WHERE clause of the generated UPDATE query. All other table columns will be used in the SET part of the query.

44. How can you sync a exported table with HDFS data in which some rows are deleted?

Truncate the target table and load it again.

45. How can we load to a column in a relational table which is not null but the incoming value from HDFS has a null value?

By using the -input-null-string parameter we can specify a default value and that will allow the row to be inserted into the target table.

46. How can you schedule a sqoop job using Oozie?

Oozie has in-built sqoop actions inside which we can mention the sqoop commands to be executed.

47. Sqoop imported a table successfully to HBase but it is found that the number of rows is fewer than expected. What can be the cause?

Some of the imported records might have null values in all the columns. As Hbase does not allow all null values in a row, those rows get dropped.

48. How can you force sqoop to execute a free form Sql query only once and import the rows serially ?

By using the -m 1 clause in the import command, sqoop creates only one mapreduce task which will import the rows sequentially.

49) In a sqoop import command you have mentioned to run 8 parallel Mapreduce task but sqoop runs only 4. What can be the reason?

The Mapreduce cluster is configured to run 4 parallel tasks. So the sqoop command must have number of parallel tasks less or equal to that of the MapReduce cluster.

50. What happens when a table is imported into a HDFS directory which already exists using the -append parameter?

Using the --append argument, Sqoop will import data to a temporary directory and then rename the files into the normal target directory in a manner

that does not conflict with existing filenames in that directory.

51. How to import only the updated rows from a table into HDFS using sqoop assuming the source has last update timestamp details for each row?

By using the lastmodified mode. Rows where the check column holds a timestamp more recent than the timestamp specified with --last-value are imported.

52. Give a Sqoop command to import all the records from employee table divided into groups of records by the values in the column department_id.

```
$ sqoop import --connect jdbc:mysql://DineshDB  
-table EMPLOYEES --split-by dept_id -m2
```

53. What does the following query do?

```
$ sqoop import --connect jdbc:mysql://DineshDB  
-table sometable --where "id > 1000" --target-dir  
"/home/dinesh/sqoopincremental" --append
```

It performs an incremental import of new data, after having already imported the first 1000 rows of a table

54. can sqoop run without a hadoop cluster?

To run Sqoop commands, Hadoop is a mandatory prerequisite. You cannot run sqoop commands without the Hadoop libraries.

55. What is the importance of \$CONDITIONS in Sqoop?

Sqoop performs highly efficient data transfers by inheriting Hadoop's parallelism.

To help Sqoop split your query into multiple chunks that can be transferred in parallel, you need to include the \$CONDITIONS placeholder in the where clause of your query.

Sqoop will automatically substitute this placeholder with the generated conditions specifying which slice of data should be transferred by each individual task.

While you could skip \$CONDITIONS by forcing Sqoop to run only one job using the --num-mappers 1 parameter, such a limitation would have a severe performance impact.

For example:-

If you run a parallel import, the map tasks will execute your query with different values substituted in for \$CONDITIONS. one mapper may execute "select * from TblDinesh WHERE (salary >= 0 AND salary < 10000)", and the next mapper may execute "select * from TblDinesh WHERE (salary >= 10000 AND salary < 20000)" and so on.

56. How to use Sqoop validation?

You can use this parameter (--validate) to validate the counts between what's imported/exported between RDBMS and HDFS

57. Is it possible to import a file in fixed column length from the database using sqoop import?

Importing column of a fixed length from any database you can use free form query like below

```
sqoop import --connect jdbc:oracle:*
--username Dinesh
--password pwd
-e "select substr(COL1,1,4000),substr(COL2,1,4000)
from table where \$CONDITIONS"
--target-dir /user/dineshkumar/table_name
--as-textfile -m 1
```

58. How to pass Sqoop command as file arguments in Sqoop?

specify an options file, simply create an options file in a convenient location and pass it to the command line via -

-options-file argument.

eg: sqoop --options-file

```
/users/homer/work/import.txt --table TEST
```

59. Is it possible to import data apart from HDFS and Hive?

Sqoop supports additional import targets beyond HDFS and Hive. Sqoop can also import records into a table in HBase and Accumulo.

60) Is it possible to use sqoop --direct command in Hbase ?

This function is incompatible with direct import. But Sqoop can do bulk loading as opposed to direct writes. To use bulk loading, enable it using --hbase-bulkload.

61. Can I configure two sqoop command so that they are dependent on each other? Like if the first sqoop job is successful, second gets triggered. If first fails, second should not run?

No, using sqoop commands it is not possible, but You can use oozie for this. Create an oozie workflow. Execute the second action only if the first action succeeds.

62. What is UBER mode and where is the settings to enable in Hadoop?

Normally mappers and reducers will run by ResourceManager (RM), RM will create separate container for mapper and reducer. Uber configuration, will allow to run mapper and reducers in the same process as the ApplicationMaster (AM).

This brings our list of 60+ SQOOP for Big Data interview questions to an end.

We believe this concise guide will help you “expect the unexpected” and enter your first data analytics interview with confidence

We, at Zep provide a platform for Education, where your demand gets fulfilled. You demand we fulfil all your learning needs without costing you extra.



Ready to take the next steps?

Zep offers a platform for education to learn, grow & earn.

Become a Tech Blogger at Zep!!

Why don't you start your journey as a blogger and enjoy unlimited free perks and cash prizes every month.

[Explore](#)