# Java-Success.com

Industrial strength Java/JEE Career Companion to open more doors

search here …      Go

**Home**   **Java FAQs**   **600+ Java Q&As**   **Career**   **Tutorials**   **Member**   **Why?**

Can u Debug?   Java 8 ready?   Top X   Productivity Tools   Judging Experience?

# 04: Q27 – Q36 Apache Spark interview questions & answers

Posted on May 14, 2016 by Arulkumaran Kumaraswamipillai

0
Like

Tweet

0

Share

G+1

Share

**Q27.** Where is Apache Spark used in the Hadoop eco system?

**A27.** Spark is essentially a data processing framework that is faster & more flexible than "Map Reduce". The Spark itself has grown into an eco system with Spark SQL, Spark streaming, Spark UI, GraphX, MLlib, and SparkR. Apache Spark can run on Hadoop clusters, as a standalone system or on the cloud. Spark can be used for fast processing (e.g. transforming sequence files into AVRO or Parquet file formats, reading from HBase, Hive, Cassandra, and any HDFS, etc), for sophisticated analytics (e.g. machine learning & graph algorithm), and for near real time (i.e. NRT)

## 600+ Full Stack Java/JEE Interview Q&As ♥Free ♦FAQs

open all | close all

streaming of "Discretized Streams or DStreams". DStreams are defined as sequences of RDD's.

Q28. Why is Apache Spark favoured over MapReduce as an open source big data processing framework
A28. Spark gives you a comprehensive and unified framework to manage big data processing requirements with near real time (i.e. NRT) latency.
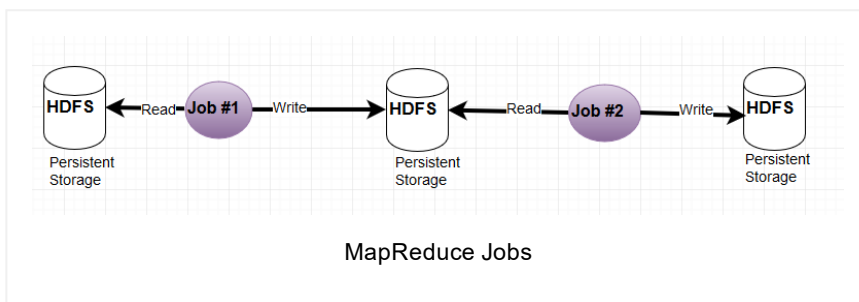
1) In MapReduce, the output data between each step has to be stored in the distributed file system before the next step can begin. This approach can be very slow for iterative tasks due to data replications across nodes & data storage I/O operations. Spark allows you for the steps to run

a) completely in memory for performance,
b) by writing everything to disk to handle large data sets and
c) by writing partially to disk & partially processing from the memory to get the best of both performance & ability to work with large data sets.

You have to look at your data and use cases to assess the memory requirements.

2) In MapReduce, if you want to perform complex processing, you need to string together a series of MapReduce jobs, and execute them in sequence.

1) read data from HDFS -> 2) apply map and reduce –> 3) write data back to HDFS –> 1) read data from HDFS –> 2) apply map and reduce –> 3) write data back to HDFS and so on……



MapReduce Jobs

Spark allows you to develop complex multi-step pipelines using DAG (i.e. Directed Acyclic Graph) pattern so that different jobs can work with the same data. This makes the development easier, but also makes Spark perform better even if you write everything to disk instead of processing from the memory.

Q29. What is a RDD in Spark?

A29. RDD stands for **R**esilient **D**istributed**D**atasets, which is a collection of fault-tolerant operational elements that run in parallel. The data is immutable & partitioned to run in a distributed manner. RDDs can be <u>cached</u> across computing nodes in a cluster.



RDDs for in-memory computations

JavaSparkContext's parallelize method on a list of integers are copied across the Hadoop cluster (i.e. Data nodes) to form a distributed datasets,

```
1
2  SparkConf conf = new SparkConf().setAppName("Sequ
3  JavaSparkContext sc = new JavaSparkContext(conf);
4
5  List<Integer> data = Arrays.asList(1, 2, 3, 4, 5)
6  JavaRDD<Integer> distData = sc.parallelize(data);
7
```

and can be operated on in parallel to sum up the elements.

```
1
2  distData.reduce((a, b) -> a + b)
3
```

You can create a RDD from any storage source like text files, sequence files, avro files, etc.

```
1
2  SparkConf conf = new SparkConf().setAppName("Sequ
3  JavaSparkContext sc = new JavaSparkContext(conf);
4
5  JavaRDD<String> distTextFileRDD = sc.textFile("da
6
```

Q30. What are the different types of RDD operations?
A30. RDD supports two types:

1) Transformations: Create a new dataset from an existing one. For example ".**map**" in the example below is a transformation that extracts values (i.e. _2) from key/value pairs.

2) Actions: Return a value to the driver program after running a computation on the dataset. For example, ".**collect**" in the example below is an action that returns a collection of values.

```
1
2  JavaPairRDD<IntWritable, BytesWritable> distSeqFi
3         sc.sequenceFile(inputFile.getPath(), IntW
4  List<String> valuesXml = distSeqFileRDD.map(x ->
5
```

Q31. What is a "RDD Lineage"?
A31. Spark does not support data replication in memory, hence in an event of any data loss, it is rebuilt using the "RDD Lineage". It is a process of reconstructing lost data partitions.

Q32. How does Spark support development of complex multi-step pipelines?
A32. Spark allows you to develop complex multi-step pipelines using DAG (i.e. Directed Acyclic Graph) pattern so that different jobs can work with the same data. This makes the development easier, but also makes Spark perform better

even if you write everything to disk instead of processing from the memory.

In Spark, a job is associated with a chain of RDD dependencies organized in a direct acyclic graph (DAG) that looks like the following:

**Q33.** What is a partition in a Spark job?
**A33.** Partitioning is the process that logically divides units of data to be processed in parallel to speed up data processing. RDDs created in 2 partitions

```
1
2  List<Integer> data = Arrays.asList(1, 2, 3, 4, 5,
3  JavaRDD<Integer> distData = sc.parallelize(data,
4
```

**Q34.** How are Spark variables shared across nodes?
**A34.** When a map or reduce operator is executed on a remote node, it works on separate copies of all the variables used within the operation at a particular node, and any updates to these variables are not propagated back to the driver program. Spark provides 2 approaches to share variables across nodes in a cluster.

**1) Accumulators**: Variables that can be used to aggregate values from worker nodes back to the driver program.
**2) Broadcast variables**: Shared variable to efficiently distribute large read-only values to all the worker nodes.

# Accumulators

Counting the number of blank lines in a given text input.

```
1
2   JavaRDD<String> lines = sc.textFile("data.txt");
3   final Accumulator<Integer> blankLines = sc.accum
4   JavaPairRDD<String, Integer> counts = lines.flatl
5      {
6          if ("".equals(line)) {
7              blankLines.add(1); // increment the
8          }
9          return Arrays.asList(line.split(" "));
10     }).mapToPair(word -> new Tuple2<String, Integ
```

```
11        .reduceByKey((x, y) -> x + y);
12
13 System.out.println("Blank lines count: " + blank
14
```

# Broadcast variables

Broadcast the list of words to ignore to all the nodes in a
cluster.

```
1
2 JavaRDD<String> lines = sc.textFile("data.txt");
3 final Broadcast<List<String>> wordsToIgnore = sc.
4
```

Q35. What is a SparkContext?

A35. A "SparkContext" is the main entry point for a Spark job.
A "SparkContext" represents the connection to a Spark
cluster, and can be used to create RDDs, accumulators and
broadcast variables on that cluster.

# Create a SparkContext

```
1
2  import org.apache.spark.SparkConf;
3  import org.apache.spark.api.java.JavaPairRDD;
4  import org.apache.spark.api.java.JavaRDD;
5  import org.apache.spark.api.java.JavaSparkContex
6  import org.slf4j.Logger;
7
8  //...
9
10    final static JavaSparkContext sc;
11
12    static {
13      SparkConf conf =
14        new SparkConf().setAppName("Sequence To
15            .set("spark.executor.memory", "1g")
16            .set("spark.serializer", "org.apache
17      sc = new JavaSparkContext(conf);
18    }
19
20  //.....
21
```

# Create RDDs

```
1
2  JavaRDD<String> lines = sc.textFile("data.txt");
3
```

# Shared variables: Accumulators & Broadcast variables

```
1
2  JavaRDD<String> lines = sc.textFile("data.txt");
3  final Accumulator<Integer> blankLines = sc.accumu
4  final Broadcast<List<String>> wordsToIgnore = sc.
5
```

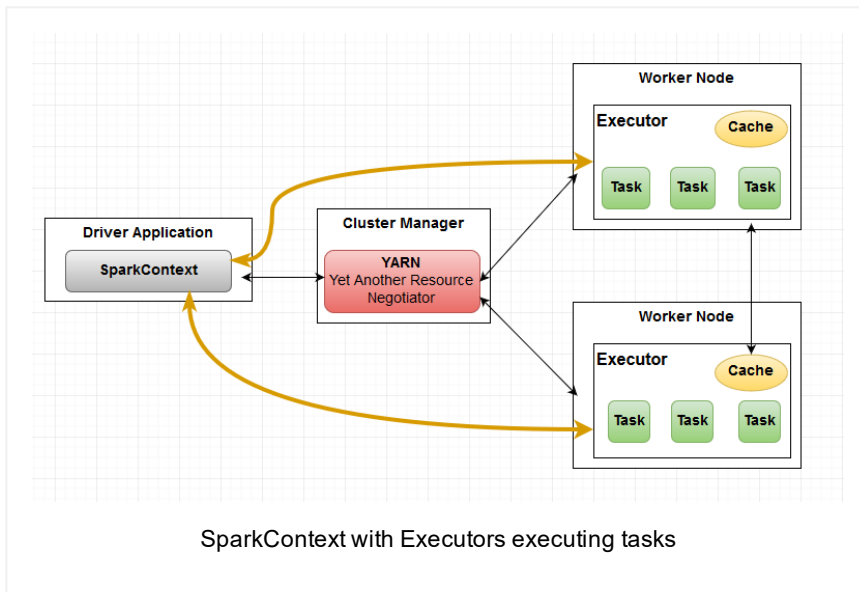Q36. What is a "Spark streaming"?

A36. Spark is a batch processing platform like Apache Hadoop, and Spark Streaming is a real-time processing tool that runs on top of the Spark engine. Spark streaming is related to Apache Storm, which is the most popular real-time processing platform for Big Data.

The primitive data type for Spark streaming is still RDD's encapsulated by a continuous stream of data known as "Discretized Streams" or DStreams. DStreams are defined as sequences of RDD's. A "DStream" is created from an input source, such as Apache Kafka, or from the transformation of another DStream.

```
1
2  SparkConf conf = new SparkConf().setMaster("loca
3  JavaStreamingContext jsc = new JavaStreamingCont
4  JavaReceiverInputDStream<String> lines = jsc.soc
5  JavaDStream<String> words = lines.flatMap(
6    new FlatMapFunction<String, String>() {
7      @Override public Iterable<String> call(Strin
8        return Arrays.asList(x.split(" "));
9      }
10  });
11
```

Q37. What is a Spark Executor?

A37. The "Driver Application" creates tasks & schedule them to be run on the "Spark Executors".

SparkContext with Executors executing tasks

Executors are worker nodes' processes in charge of running individual tasks in a given Spark job. Spark Executors are launched at the beginning of a Spark application and typically run for the entire lifetime of an application. Once they have finished running the tasks they send the results to the "Driver Application". "**Spark Executors**" also provide in-memory storage for RDDs that are cached.

```
1
2   rdd4.cache()
3
```

# Popular Posts

♦ 11 Spring boot interview questions & answers

**828 views**

♦ Q11-Q23: Top 50+ Core on Java OOP Interview Questions & Answers

**768 views**

18 Java scenarios based interview Questions and Answers

**400 views**

001A: ♦ 7+ Java integration styles & patterns interview questions & answers

**389 views**

01b: ♦ 13 Spring basics Q8 – Q13 interview questions & answers

**296 views**

♦ 7 Java debugging interview questions & answers

**293 views**

01: ♦ 15 Ice breaker questions asked 90% of the time
in Java job interviews with hints

**286 views**

♦ 10 ERD (Entity-Relationship Diagrams) Interview
Questions and Answers

**280 views**

♦ Q24-Q36: Top 50+ Core on Java classes, interfaces
and generics interview questions & answers

**240 views**

001B: ♦ Java architecture & design concepts
interview questions & answers

**202 views**

| Bio | **Latest Posts** |

## Arulkumaran Kumaraswamipillai

Mechanical Eng to freelance Java
developer in 3 yrs. Contracting since 2003,
and attended 150+ Java job interviews, and
often got 4 - 7 job offers to choose from. It
pays to prepare. So, published Java
interview Q&A books via Amazon.com in
2005, and sold 35,000+ copies. Books are
outdated and replaced with this subscription
based site.**945+** paid members. join my
LinkedIn Group. **Reviews**

**About** Arulkumaran Kumaraswamipillai

Mechanical Eng to freelance Java
developer in 3 yrs. Contracting since
2003, and attended 150+ Java job
interviews, and often got 4 - 7 job offers
to choose from. It pays to prepare. So, published Java
interview Q&A books via Amazon.com in 2005, and sold
35,000+ copies. Books are outdated and replaced with

this subscription based site.**945+** paid members. join my
LinkedIn Group. **Reviews**

**Posted in** Hadoop & BigData Interview Q&A**,** member-paid

# Empowers you to open more doors, and fast-track

### Technical Know Hows

☀ Java generics in no time ☀ Top 6 tips to transforming your thinking from OOP to FP ☀ How does a HashMap internally work? What is a hashing function?
☀ 10+ Java String class interview Q&As ☀ Java auto un/boxing benefits & caveats ☀ Top 11 slacknesses that can come back and bite you as an experienced Java developer or architect

### Non-Technical Know Hows

☀ 6 Aspects that can motivate you to fast-track your career & go places ☀ Are you reinventing yourself as a Java developer? ☀ 8 tips to safeguard your Java career against offshoring ☀ My top 5 career mistakes

# Prepare to succeed

☀ Turn readers of your Java CV go from "Blah blah" to "Wow"? ☀ How to prepare for Java job interviews? ☀ 16 Technical Key Areas ☀ How to choose from multiple Java job offers?

Select Category ▼

# © Disclaimer

or usefulness of content, though I do make an effort to be accurate. Links to external sites do not imply endorsement of the linked-to sites.