

Scala/Spark



Wikipedia :

Apache Spark is an [open-source](#) cluster computing framework originally developed in the AMPLab at [UC Berkeley](#). In contrast to [Hadoop](#)'s two-stage disk-based [MapReduce](#) paradigm, Spark's in-memory primitives provide performance up to 100 times faster for certain applications.^[1] By allowing user programs to load data into a cluster's memory and query it repeatedly, Spark is well-suited to machine learning algorithms.^[2]

Spark requires a cluster manager and a distributed storage system. For cluster management, Spark supports standalone (native Spark cluster), [Hadoop YARN](#), or [Apache Mesos](#).^[3] For distributed storage, Spark can interface with a wide variety, including [Hadoop Distributed File System \(HDFS\)](#),^[4] [Cassandra](#),^[5] [OpenStack Swift](#), and [Amazon S3](#). Spark also supports a pseudo-distributed local mode, usually used only for development or testing purposes, where distributed storage is not required and the local file system can be used instead; in this scenario, Spark is running on a single machine with one executor per CPU core.

Project Components

1. Spark Core and Resilient Distributed Datasets (RDDs)
2. Spark SQL
3. Spark Streaming
4. MLlib Machine Learning Library
5. GraphX

Features

- Java, Scala, Python, and R APIs.
- Scalability to over 8000 nodes in production.^[11]
- Ability to cache datasets in memory for interactive data analysis: extract a working set, cache it, query it repeatedly.
- Interactive command line interface (in Scala or Python) for low-latency [horizontally scalable](#) data exploration.
- Higher level library for stream processing^[12], through Spark Streaming.
- Support for structured and relational query processing (SQL), through Spark SQL.
- Higher level libraries for machine learning and graph processing.

Examples Scala

- **Declare a list of integers as a variable called “myNumbers”:**

```
scala> val myNumbers = List(1, 2, 5, 4, 7, 3)
myNumbers: List[Int] = List(1, 2, 5, 4, 7, 3)
```

- **Declare a function, cube, that computes the cube (third power) of an Int.**

```
scala> def cube(a: Int): Int = a * a * a
cube: (a: Int)Int
```

- **Apply the function to myNumbers using the map function.**

```
scala> myNumbers.map(x => cube(x))
res: List[Int] = List(1, 8, 125, 64, 343, 27)
// Scala also provides some shorthand ways of writing this:
// myNumbers.map(cube(_))
// myNumbers.map(cube)
```

- **Then also try writing the function inline in a map call, using closure notation.**

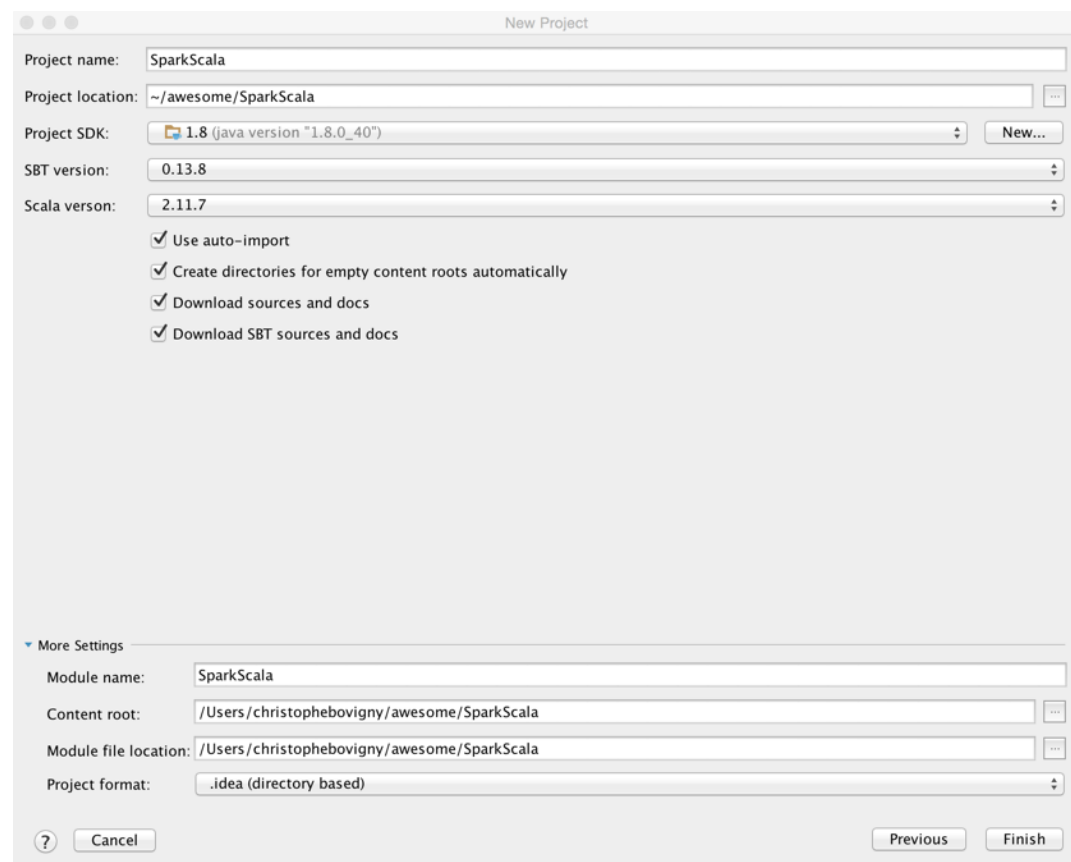
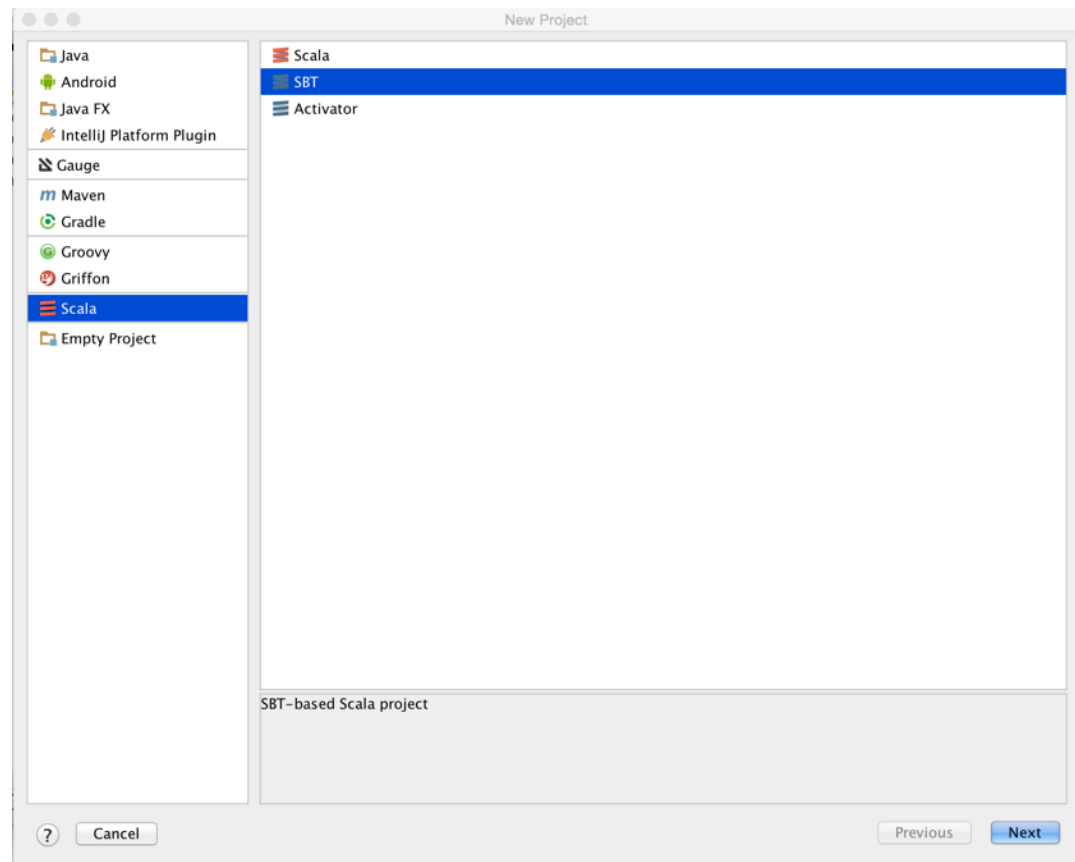
```
scala> myNumbers.map{x => x * x * x}
res: List[Int] = List(1, 8, 125, 64, 343, 27)
```

- **Define a factorial function that computes $n! = 1 * 2 * \dots * n$ given input n . You can use either a loop or recursion, in our solution we use recursion (see steps 5-7 of [First Steps to Scala](#)). Then compute the sum of factorials in myNumbers.**

```
scala> def factorial(n: Int): Int = if (n == 0) 1 else n * factorial(n-1) // From http://bit.ly/b2sVKI
factorial: (Int)Int
scala> myNumbers.map(factorial).sum
res: Int = 5193
```



Setup a new Project





Add dependencies

The screenshot shows an IDE window titled "build.sbt - SparkScala - [~/awesome/SparkScala]". The left sidebar displays the project structure for "SparkScala [sparkscala] (~/awesome/SparkScala)", including folders like ".idea", "project [sparkscala-build] (source)", "src", "target", and "build.sbt". The main editor area shows the content of "build.sbt":

```
name := "SparkScala"
version := "1.0"
scalaVersion := "2.11.7"

libraryDependencies ++= Seq(
  "org.apache.spark" %% "spark-streaming" % "1.4.1",
  "org.apache.spark" %% "spark-streaming-kafka" % "1.4.1",
  "javax.servlet" % "javax.servlet-api" % "3.0.1"
)
```

The right sidebar shows build tool options: "Maven Projects", "SBT", and "Ant Build". The bottom status bar indicates "SBT: [info] Loading project definition from /Users/christophebovigny/awesome/SparkSc..." and includes icons for Terminal, TODO, and Event Log.

π

3.141592653589793

Pi / Local

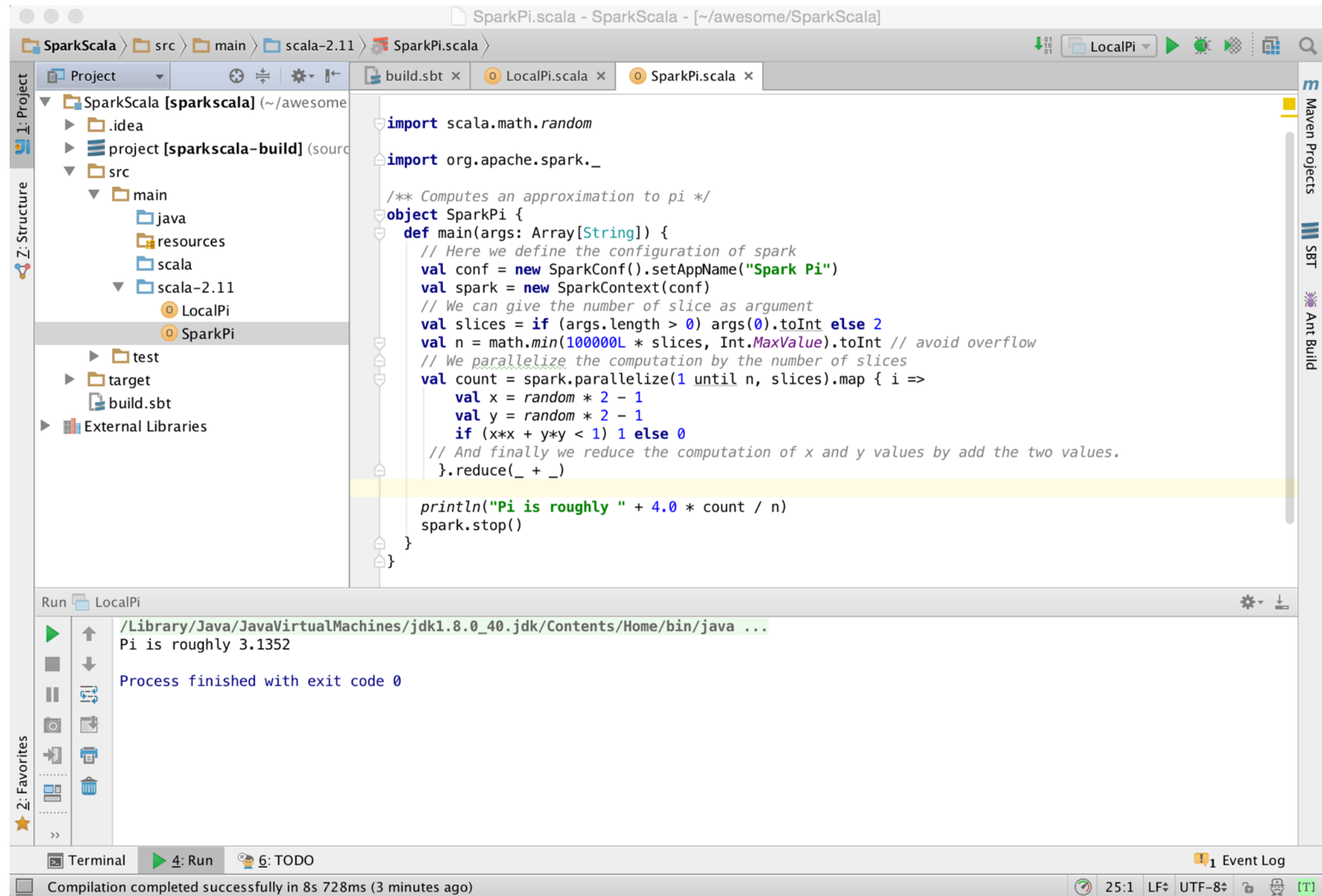
The screenshot shows an IDE window titled "LocalPi.scala - SparkScala - [~/awesome/SparkScala]". The project structure on the left includes "SparkScala [sparkscala] (~/awesome)" with subdirectories ".idea", "project [sparkscala-build] (source)", "src" (containing "main" with "java", "resources", "scala", and "scala-2.11" subdirectories), "test", "target", "build.sbt", and "External Libraries". The "LocalPi" file is selected under "src/main/scala-2.11".

The code in "LocalPi.scala" is as follows:

```
/**...*/  
import scala.math.random  
  
object LocalPi {  
  def main(args: Array[String]) {  
    var count = 0  
    for (i <- 1 to 100000) {  
      val x = random * 2 - 1  
      val y = random * 2 - 1  
      if (x*x + y*y < 1) count += 1  
    }  
    println("Pi is roughly " + 4 * count / 100000.0)  
  }  
}
```

The status bar at the bottom shows "Terminal", "6: TODO", "1 Event Log", and a warning message: "SBT project import: [warn] Multiple dependencies with the same organization/name but different versions. To avoid conflict, pick one versi... (3 minutes ago)". The status bar also displays "18:1", "LF", "UTF-8", and icons for a lock, a printer, and a terminal.

Pi / Spark



Run Pi /Spark

