# Scala/Spark

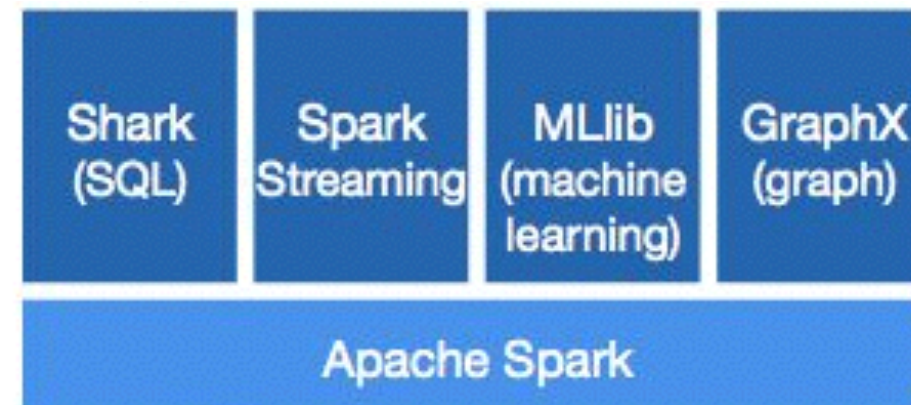## Wikipedia :

**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.
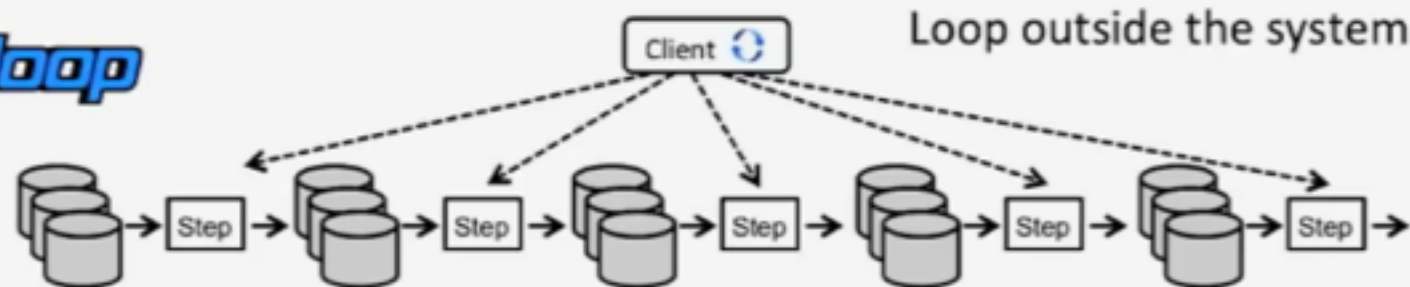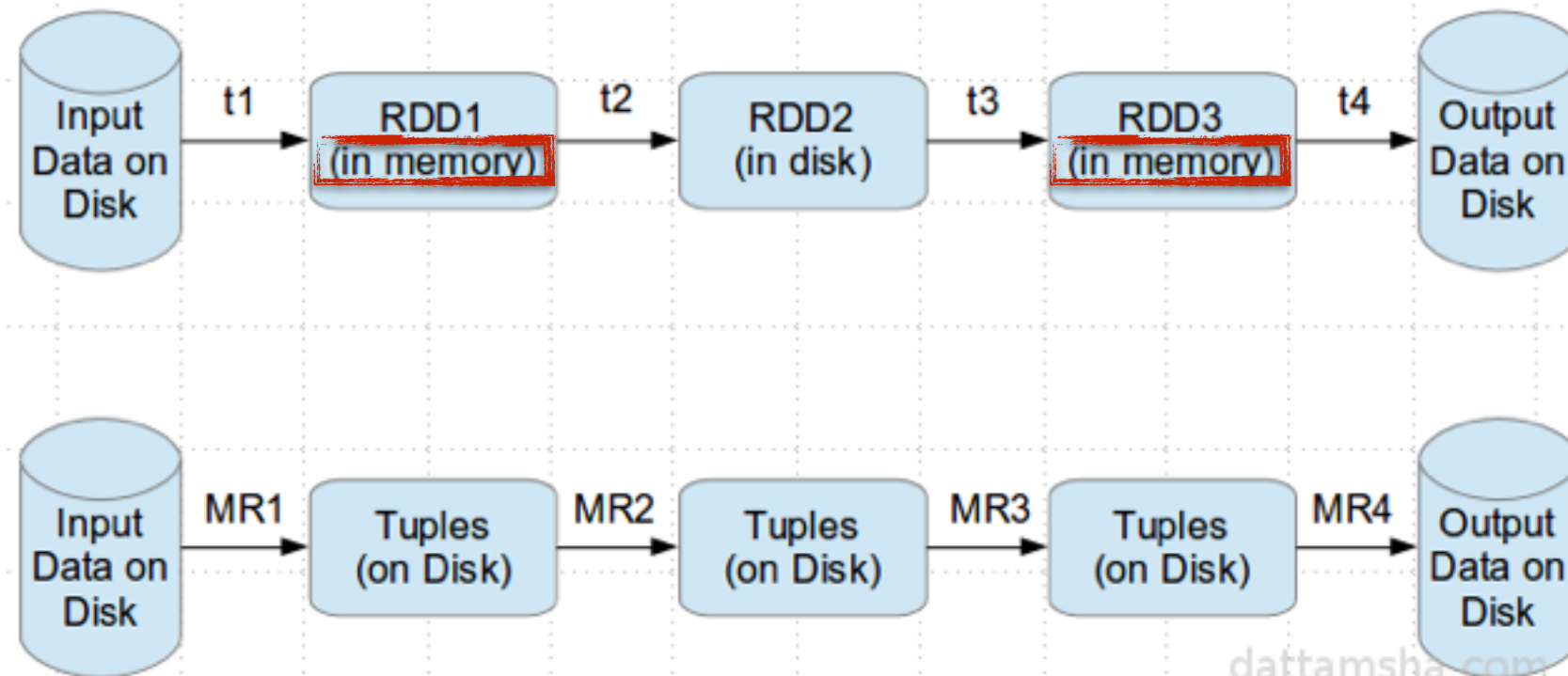
## Project Components

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.
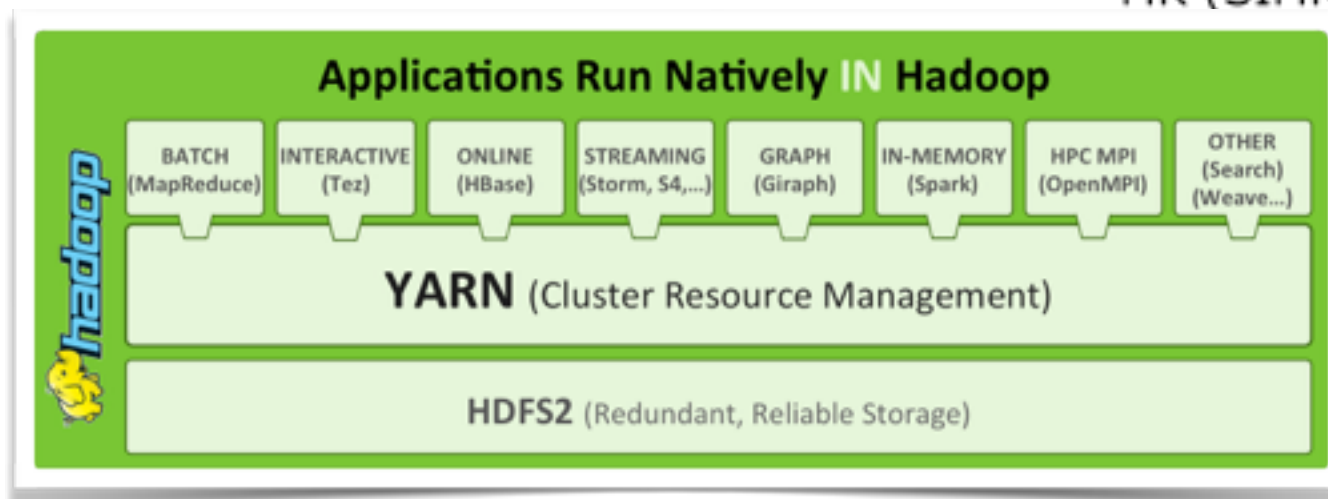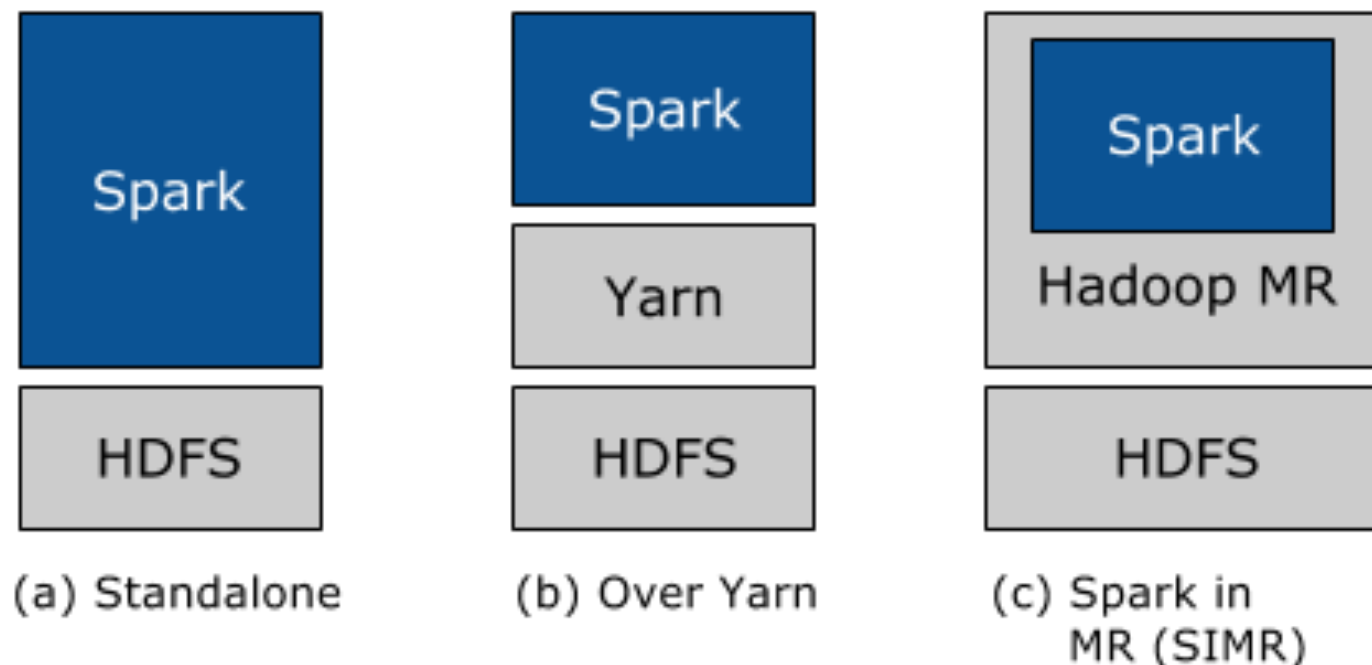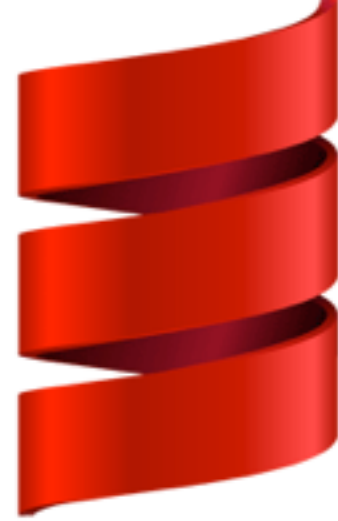
# Hadoop vs Spark

# Spark and Hadoop : Working Together



(a) Standalone
(b) Over Yarn
(c) Spark in MR (SIMR)

**Applications Run Natively IN Hadoop**

| BATCH (MapReduce) | INTERACTIVE (Tez) | ONLINE (HBase) | STREAMING (Storm, S4,...) | GRAPH (Giraph) | IN-MEMORY (Spark) | HPC MPI (OpenMPI) | OTHER (Search) (Weave...) |

**YARN** (Cluster Resource Management)

**HDFS2** (Redundant, Reliable Storage)

Spark In MapReduce (SIMR): For the Hadoop users that are not running YARN yet, another option, in addition to the standalone deployment, is to use SIMR to launch Spark jobs inside MapReduce. With SIMR, users can start experimenting with Spark and use its shell within a couple of minutes after downloading it! This tremendously lowers the barrier of deployment, and lets virtually everyone play with Spark.

## History

The design of Scala started in 2001 at the École Polytechnique Fédérale de Lausanne (EPFL) by Martin Odersky, following on from work on Funnel, a programming language combining ideas from functional programming and Petri nets.[10] Odersky had previously worked on Generic Java and javac, Sun's Java compiler.[10]

After an internal release in late 2003, Scala was released publicly in early 2004 on the Java platform,[11] and on the .NET platform in June 2004.[6][10][12] A second version (v2.0) followed in March 2006.[6] The .NET support was officially dropped in 2012.[13]

Although Scala had extensive support for functional programming from the beginning, Java remained a purely object oriented language until the introduction of lambda expressions with Java 8 in 2014.



* In Scala, the compiler is incredibly smart, so this avoids the developer needing to specify explicitly those things that the compiler can infer.
* Scala is ported on Spark (java libraries can be used)

## Build Scala project with SBT



http://www.scala-sbt.org/0.13/tutorial/index.html

## Scala crash course

https://www.sics.se/~amir/files/download/dic/scala.pdf

# 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 * … * 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
```
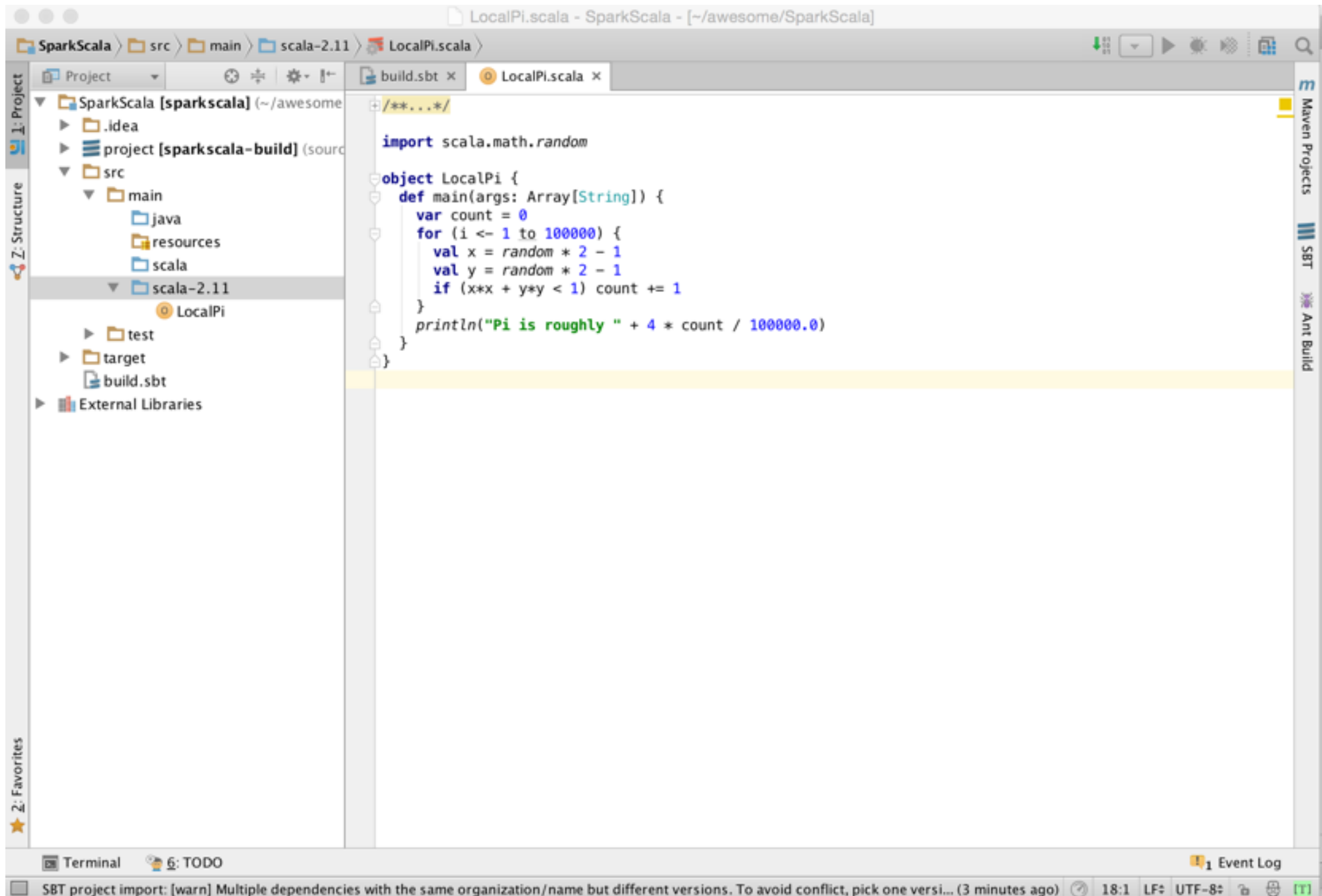
# Setup a new Project

# Add dependencies

# Pi / Local



```scala
/**...*/

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)
  }
}
```
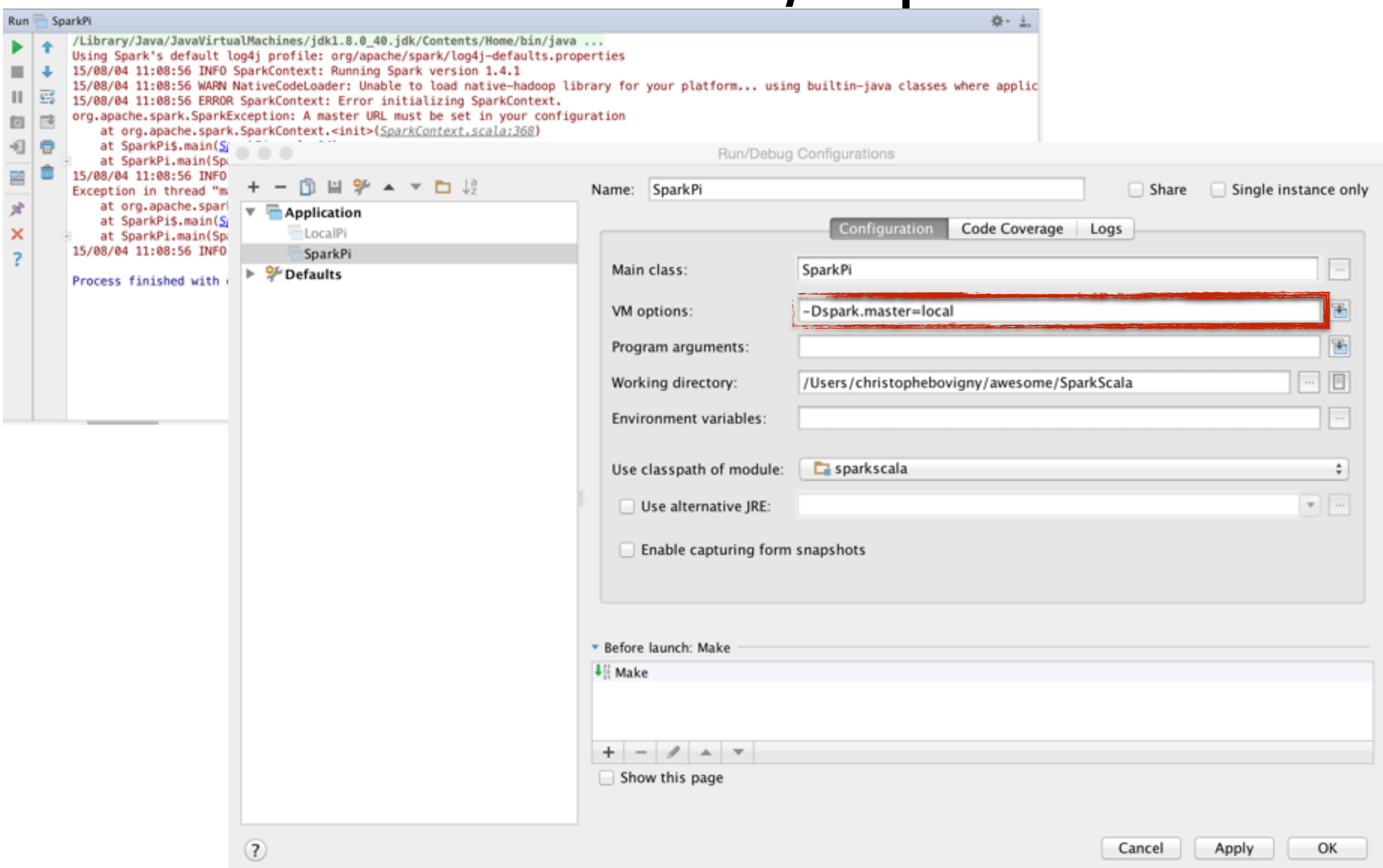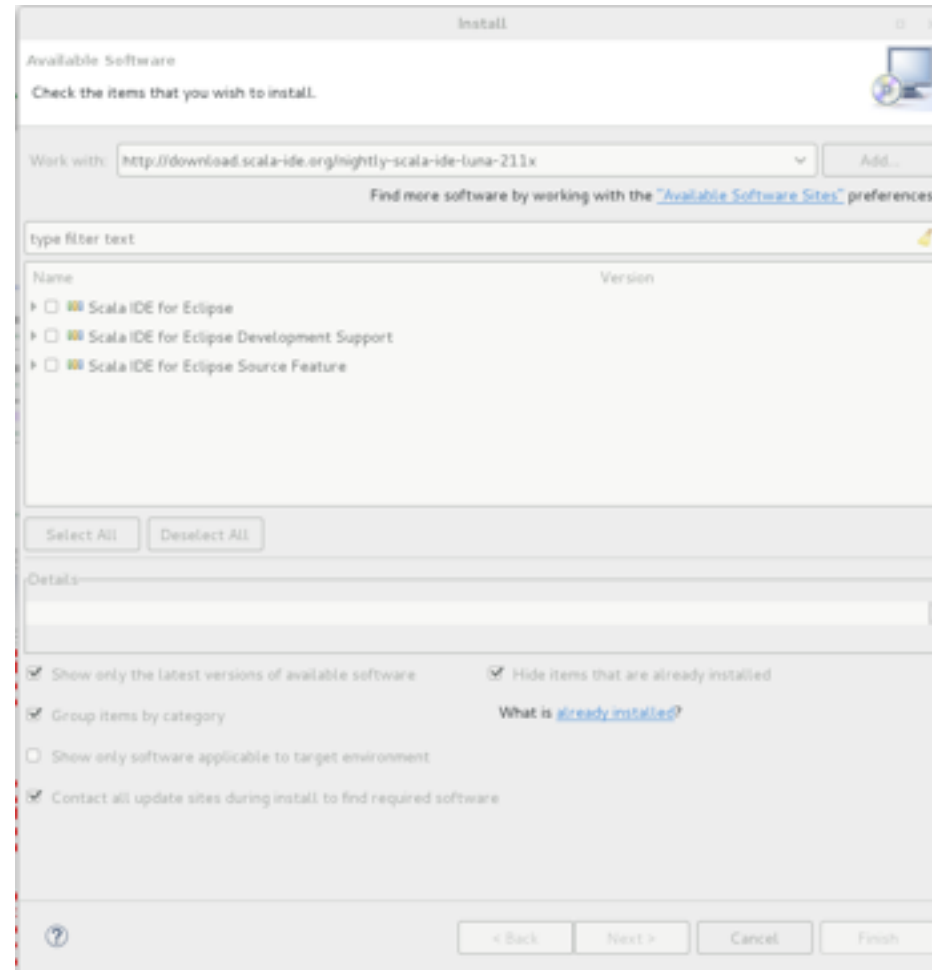
# Pi / Spark

# Run Pi /Spark

# eclipse Setup a new Project

- Install sbt for your machine

  www.scala-sbt.org/release/tutorial/Setup.html

- In eclipse Help -> Install New Software

# Setup a new Project

- Create the project folder structure  : *mkdir -p src/main/scala*

- Create a file in src/main/scala : *echo 'object Hi { def main(args: Array[String]) = println(« Hi !»)' > hw.scala*

- To allow managing dependencies, project name, Scala Version etc… create a file named *build.sbt* in the project root.
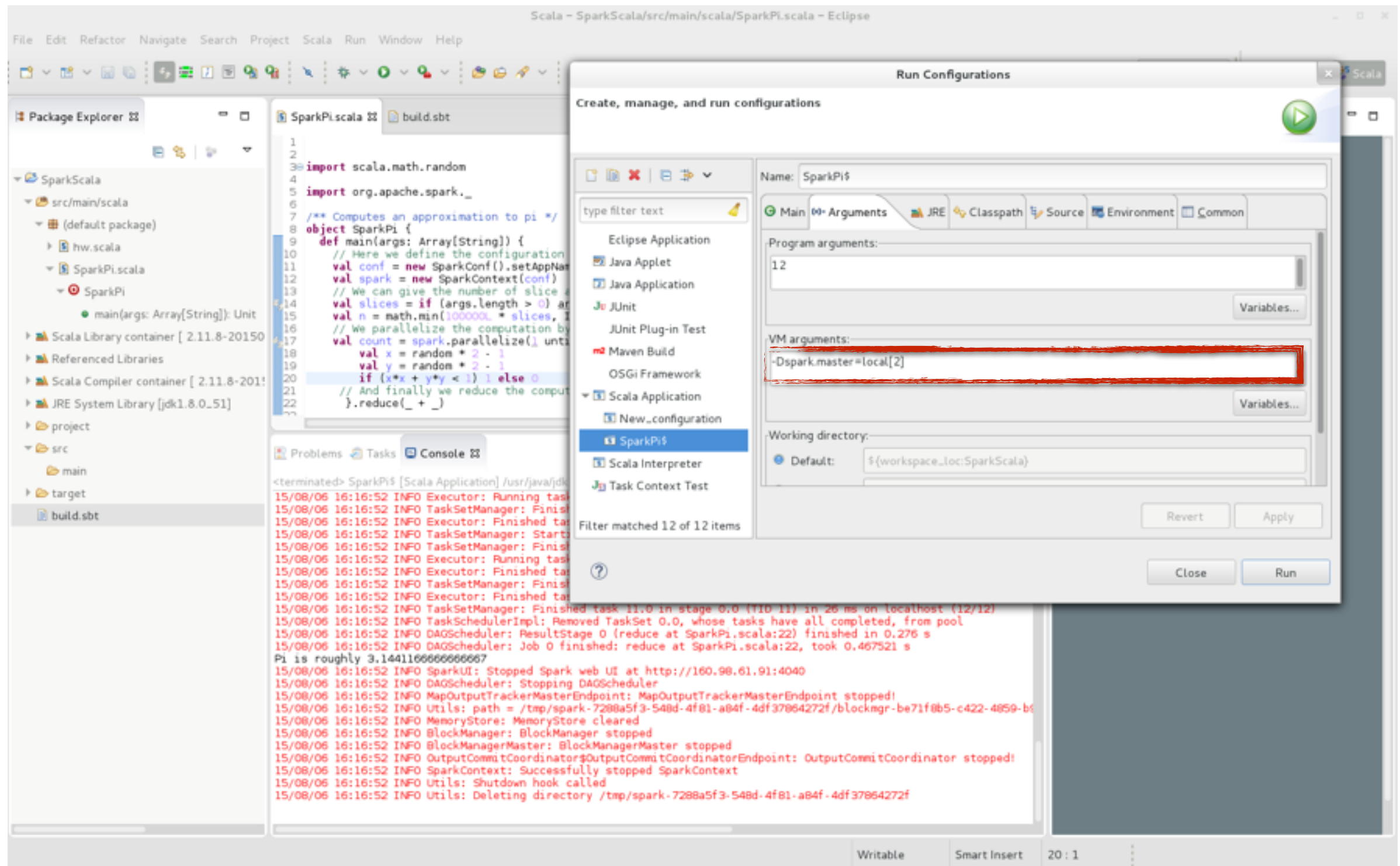
```
SparkPi.scala     build.sbt

1 name := "SparkScala"
2
3 version := "1.0"
4
5 scalaVersion := "2.11.7"
6
7
8 libraryDependencies ++= Seq(
9    "org.apache.spark" %% "spark-streaming" % "1.4.1",
10   "org.apache.spark" %% "spark-streaming-kafka" % "1.4.1",
11   "javax.servlet" % "javax.servlet-api" % "3.0.1"
12 )
13
14
```

# Setup a new Project

- Open $ ~/.sbt/0.13/plugins/build.sbt *addSbtPlugin(«com.typesafe.sbteclipse » % « sbteclipse-plugin » % 2.2.0)*

- In the project root folder : *sbt eclipse with-source=true*

- File->Import-> Existing Projects into Workspace

# Run Spark Project

# Scala Spark Examples

- https://github.com/apache/spark/tree/master/examples/src/main/scala/org/apache/spark/examples

- http://spark.apache.org