DLI Accelerated Data Science Teaching Kit

# Lecture 11.2 - Example Spark Programs

# Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")          [Base RDD]

errors = lines.filter(_.startsWith("ERROR"))

messages = errors.map(_.split('\t')(2))
cachedMsgs = messages.cache()                 [Transformed RDD]


cachedMsgs.filter(_.contains("foo")).count    [Action]

cachedMsgs.filter(_.contains("bar")).count

. . .
```
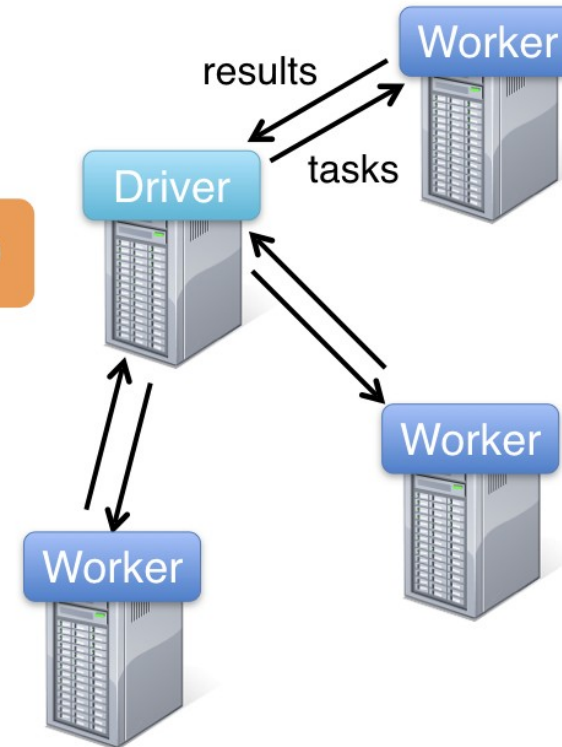
**Result:** scaled to 1 TB data in 5-7 sec
(vs 170 sec for on-disk data)

http://ananthakumaran.in/2010/03/29/scala-underscore-magic.html
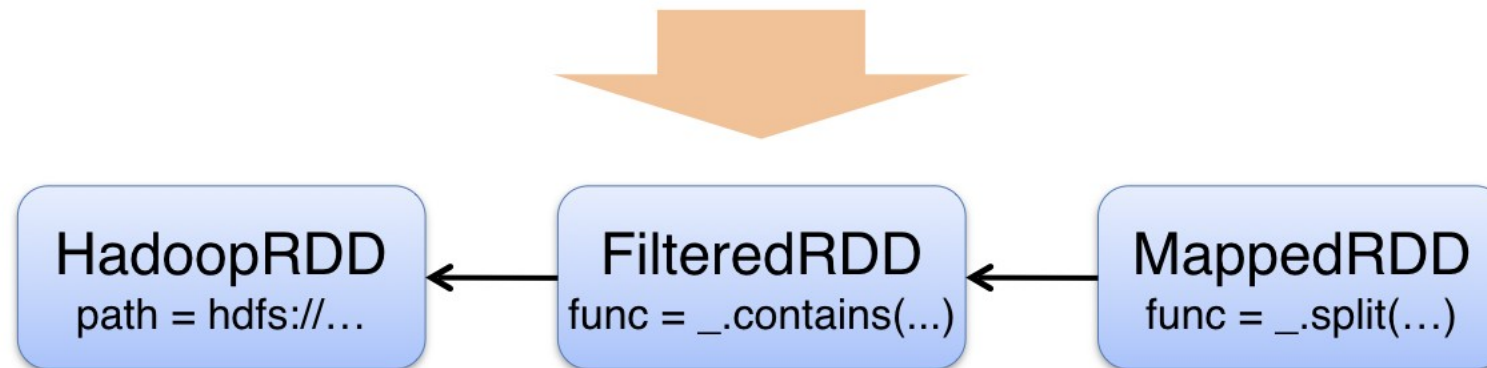http://www.slideshare.net/normation/scala-dreaded



3

# Fault Tolerance

RDDs track the series of transformations used to build them (their *lineage*) to recompute lost data

E.g: messages = textFile(...).filter(_.contains("error"))
                      .map(_.split('\t')(2))



| HadoopRDD | FilteredRDD | MappedRDD |
|---|---|---|
| path = hdfs://... | func = _.contains(...) | func = _.split(…) |

# Example: Logistic Regression

```
val data = spark.textFile(...).map(readPoint).cache()

var w = Vector.random(D)

for (i <- 1 to ITERATIONS) {
  val gradient = data.map(p =>
    (1 / (1 + exp(-p.y*(w dot p.x))) - 1) * p.y * p.x
  ).reduce(_ + _)
  w -= gradient
}

println("Final w: " + w)
```
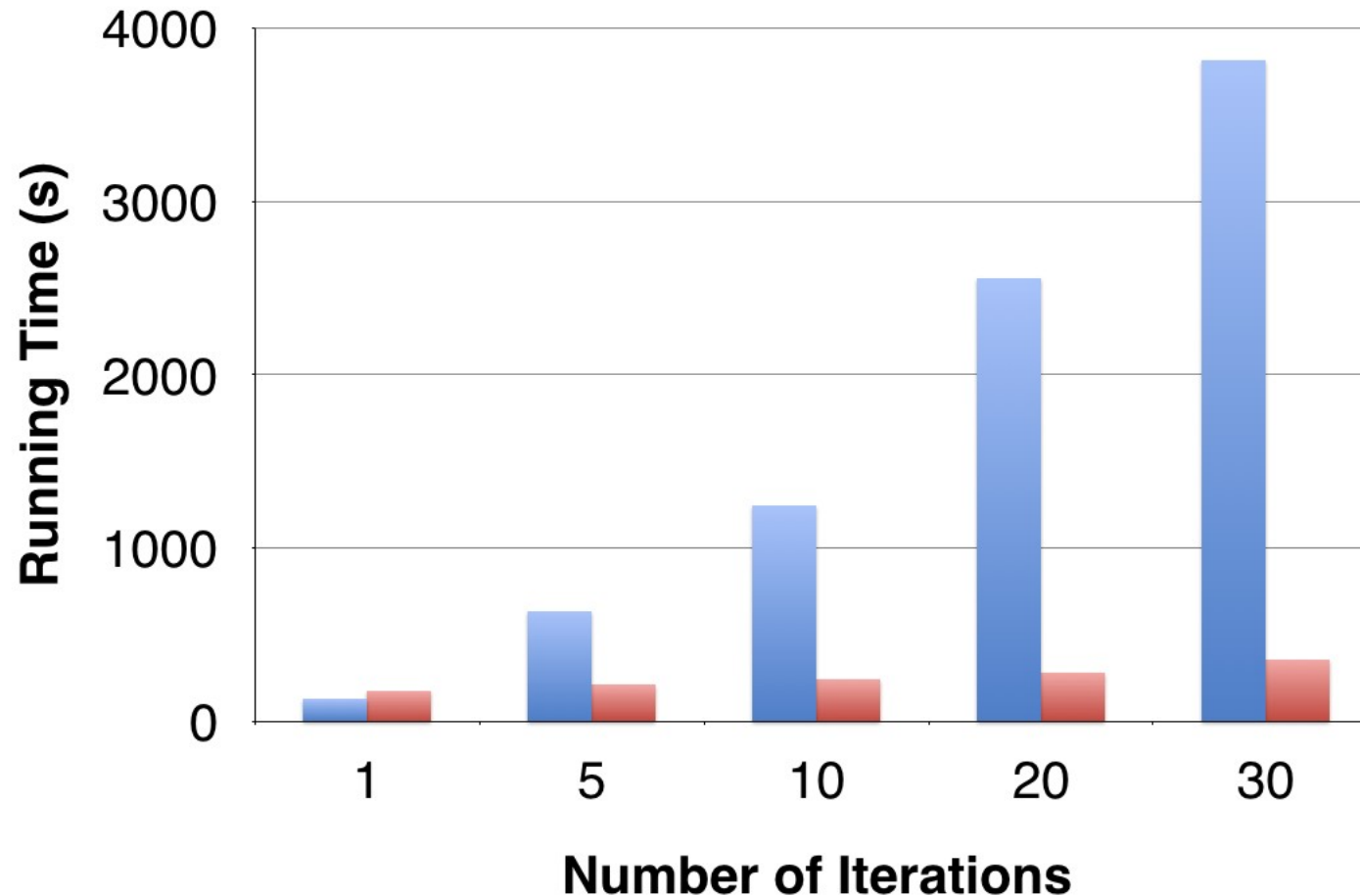
Load data in memory once

Initial parameter vector

Repeated MapReduce steps to do gradient descent

# Logistic Regression Performance

DLI Accelerated Data Science Teaching Kit

# Thank You

We thank Dr. Matei Zaharia for sharing teaching materials for Spark.