

Introduction to Apache Spark

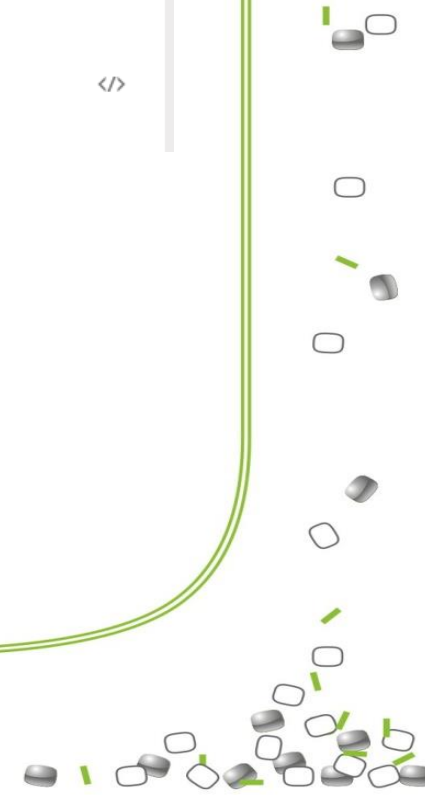
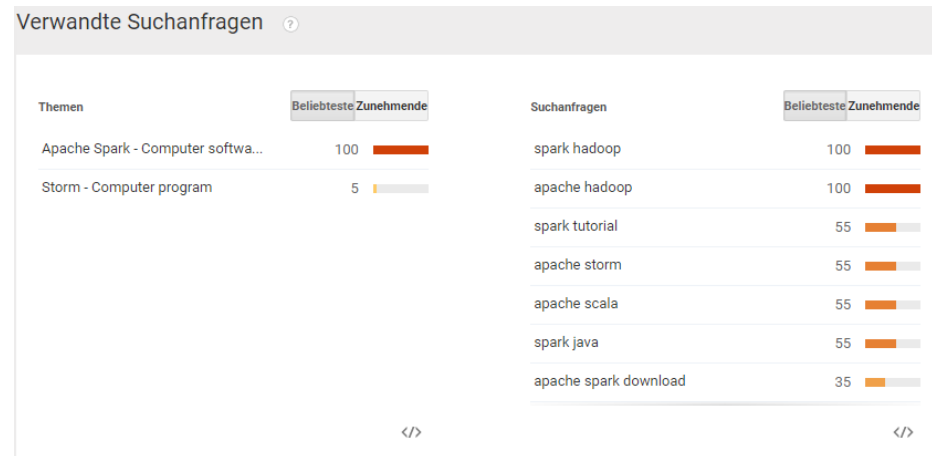
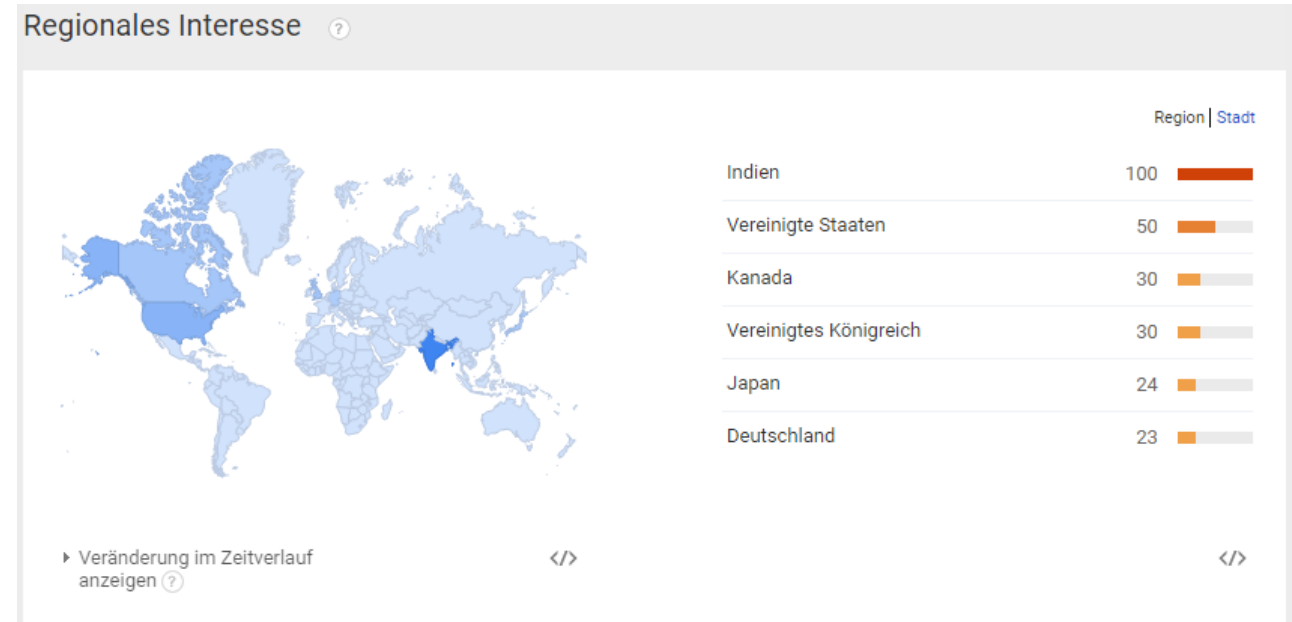
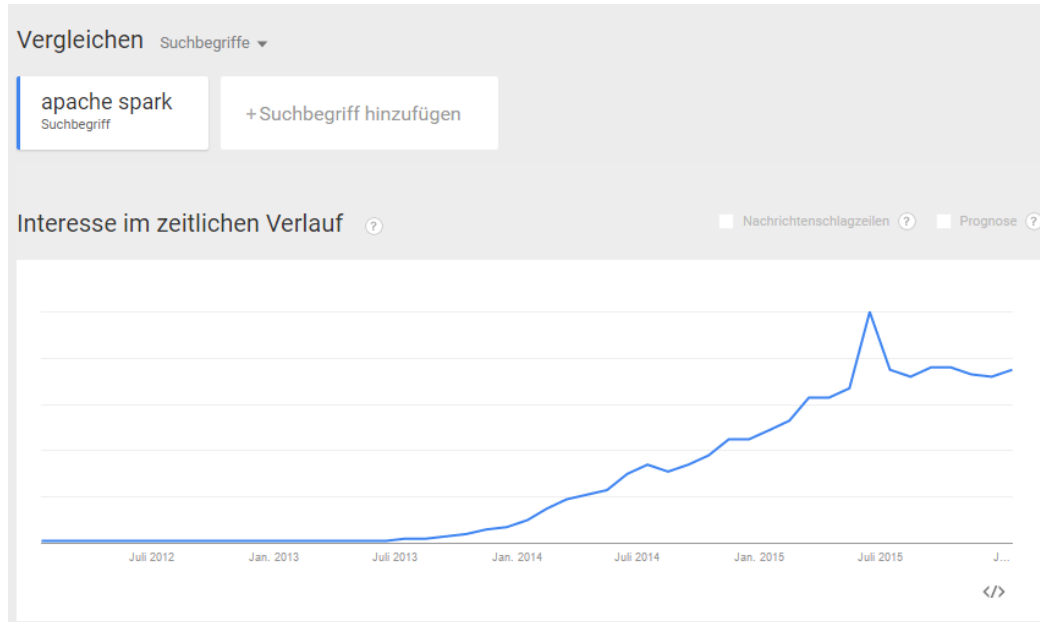


dsquare.de
Salzburger Straße 27
83073 Stephanskirchen

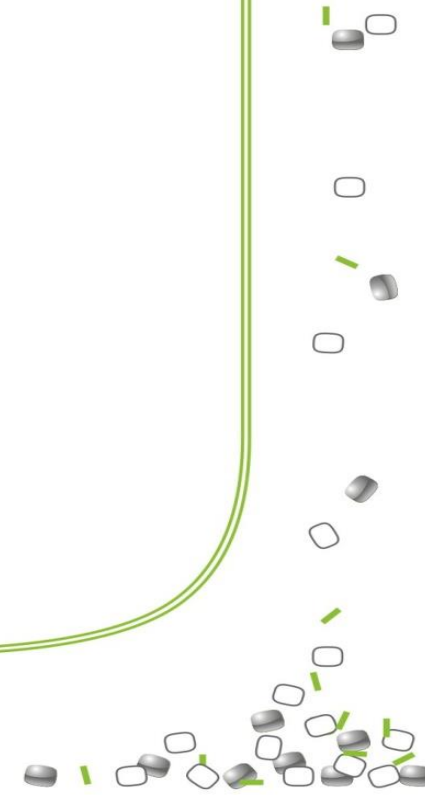
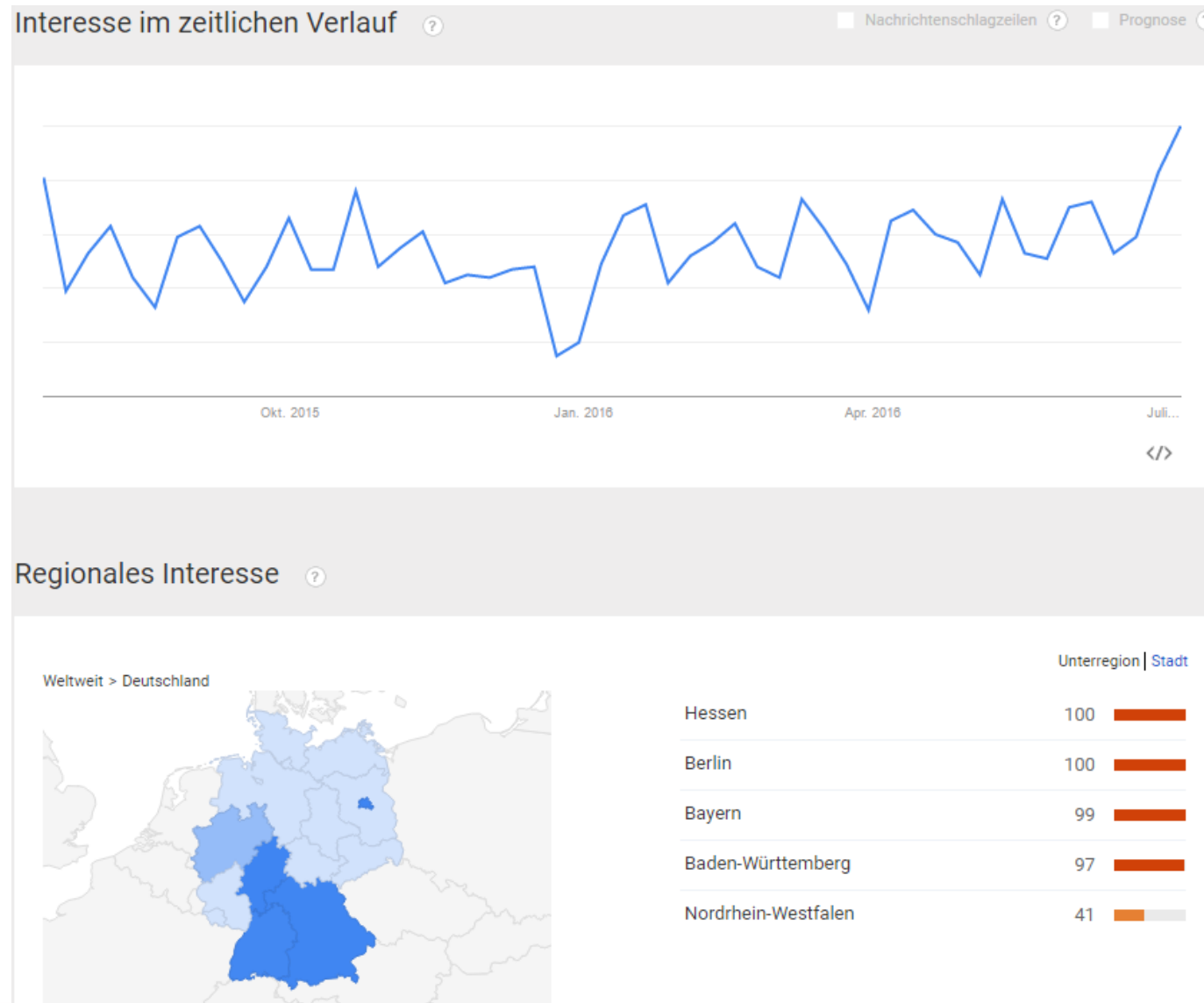
Tel.: 08031-234 1140
Mobil: 0172-1484 731
Email: info@dsquare.de
www.dsquare.de



Apache Spark according to Google Trends

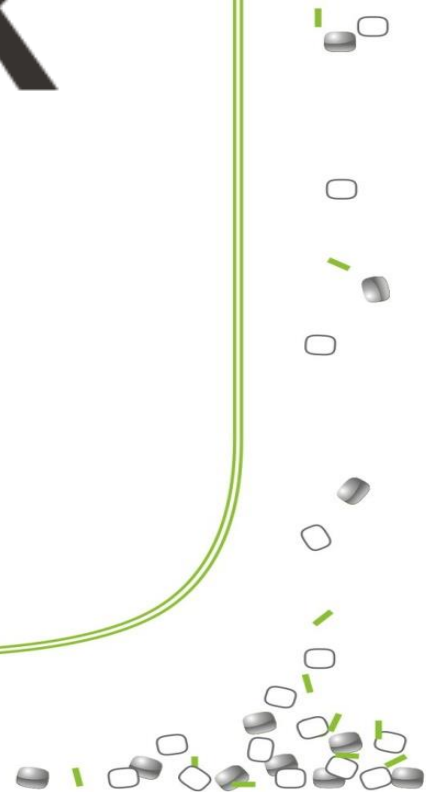


TDWI Conference sparked interest in Spark



What is Apache Spark?

- ❑ A cluster-based computing engine
- ❑ Developed since 2012
- ❑ Developed by students at UC Berkley
- ❑ APIs for
 - ❑ Python
 - ❑ Java
 - ❑ R
 - ❑ Scala
- ❑ Supports SQL, ML, Streaming Data, Graph processing
- ❑ Faster than Hadoops Map-Reduce



Timeline



Since late 1990s
APPLY Functions
In-memory
Single Process
Single Core

Not Scalable



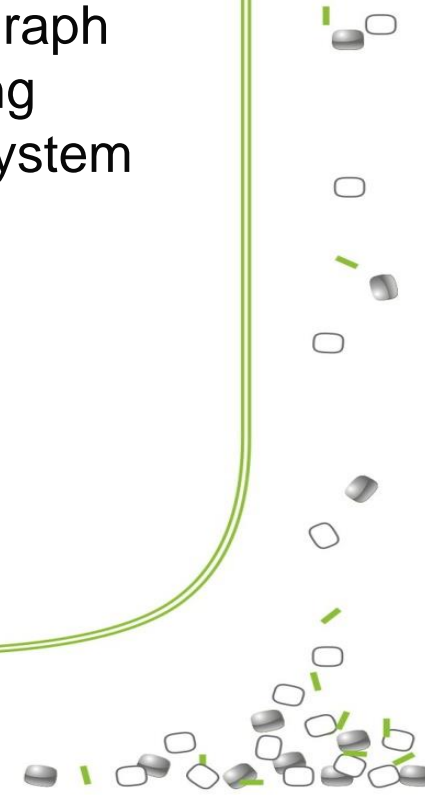
Since 2007
Map-Reduce
Parallel Computing
Distributed File System

Linear Scalability



Since 2009
Directed acyclic graph
Parallel Computing
Distributed File System

Linear Scalability



Map-reduce vs. Spark



Map-reduce

Directed acyclic graph

No writeback to HDFS necessary

Data passed to next processing step

Developer focused

Transformations available

Many APIs

In-Memory processing

RDD materialized in memory across cluster

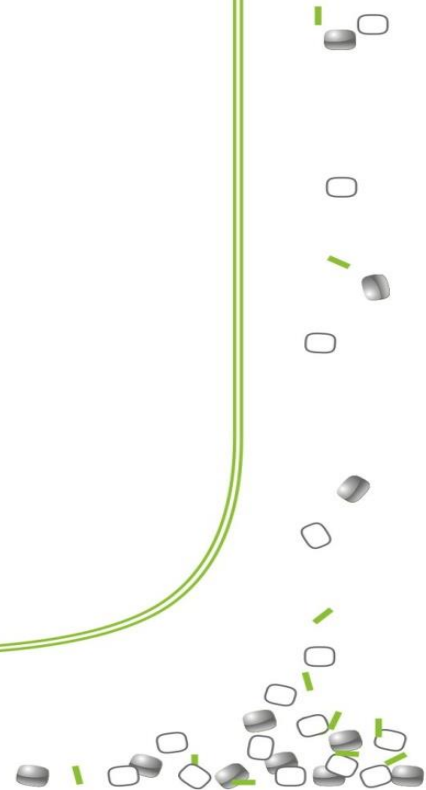
No need to reload from disc



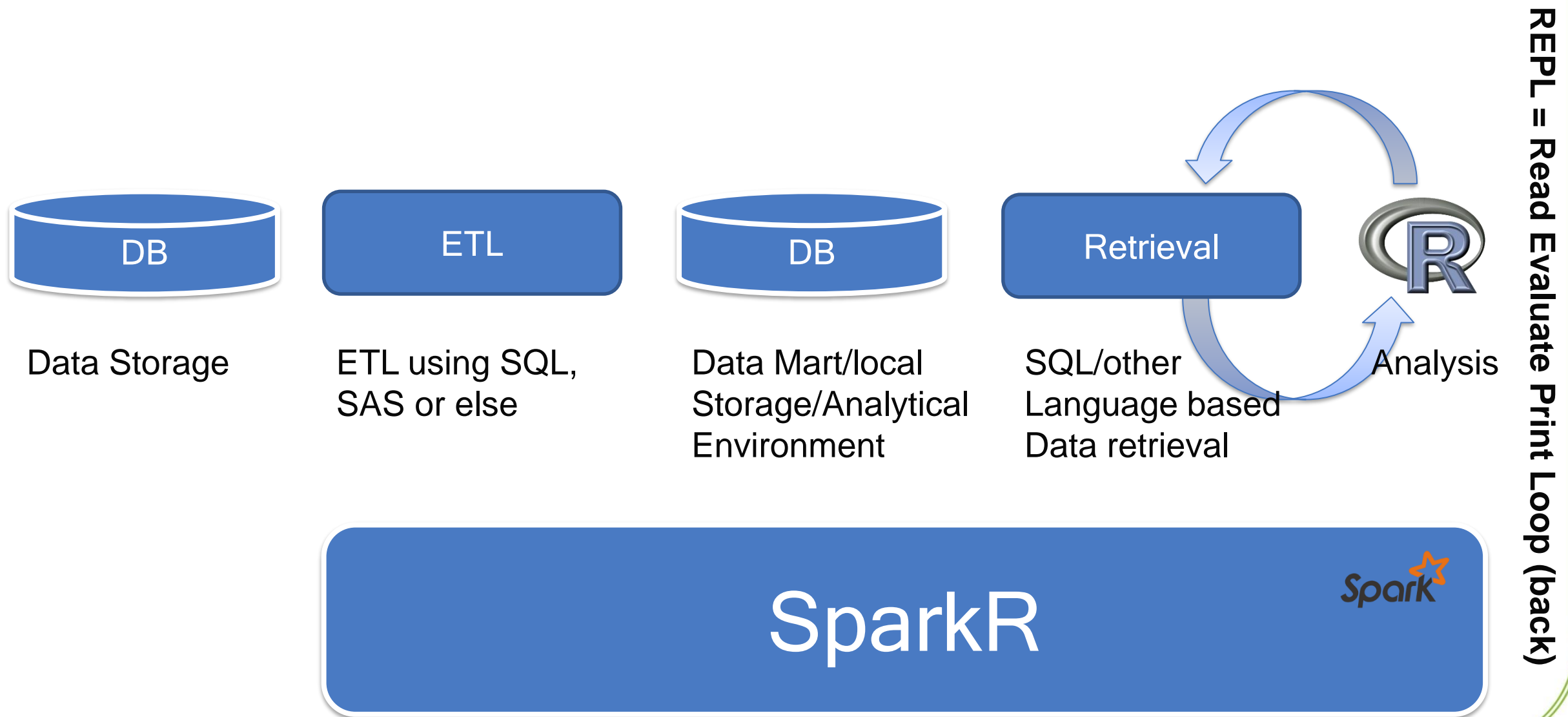
Spark is well suited for the needs of Data Scientists



Iterative application of algorithms
Multiple passes over data sets
Reactive applications



Spark can unify an analytical environment



RDD

Data

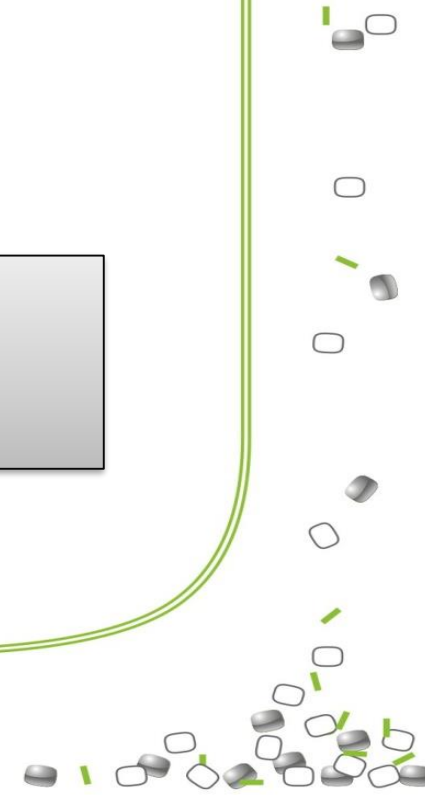
Col1	Col2	Col3
Item 1	Item 4	Item 7
Item 2	Item 5	Item 8
Item 3	Item 6	Item 9

This could be an RDD = Resilient Distributed Dataset



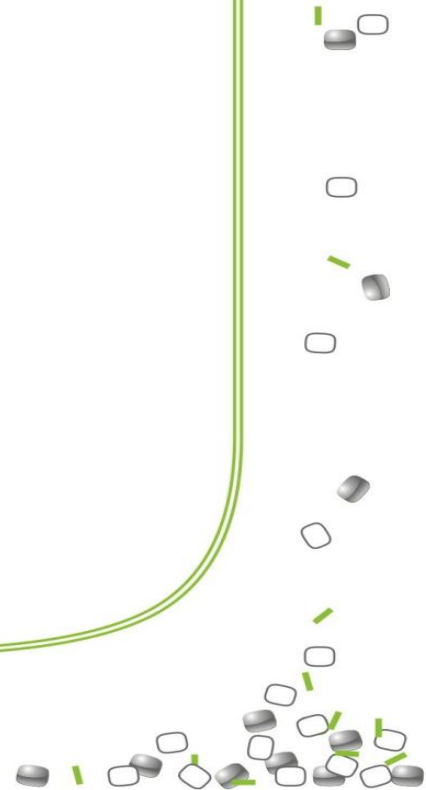
Worker Nodes

Worker nodes:
They Cache the data and do the (lazy) evaluation.



Preliminaries using Spark in R-Studio

```
.libPaths(c(.libPaths(), '/opt/spark-1.6.1-bin-hadoop2.6/R/lib')) Sys.setenv(SPARK_HOME =  
'/opt/spark-1.6.1-bin-hadoop2.6')  
Sys.setenv(PATH = paste(Sys.getenv('PATH'), '/opt/spark-1.6.1-bin-hadoop2.6/bin', sep =  
' : '))  
library(SparkR)  
  
d.csv <- "com.databricks:spark-csv_2.11:1.4.0,"  
d.pg <- "org.postgresql:postgresql94-jdbc-9.4:1207-1"  
  
sc <- sparkR.init(sparkPackages=c(d.csv))  
sqlContext <- sparkRSQL.init(sc)
```

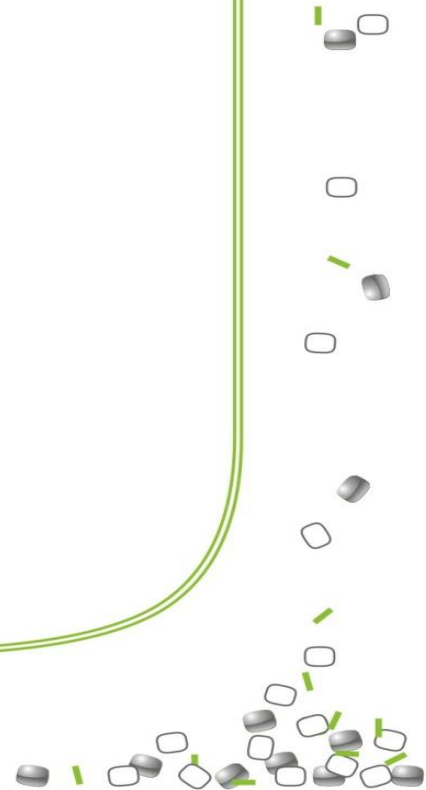


Get data from Spark

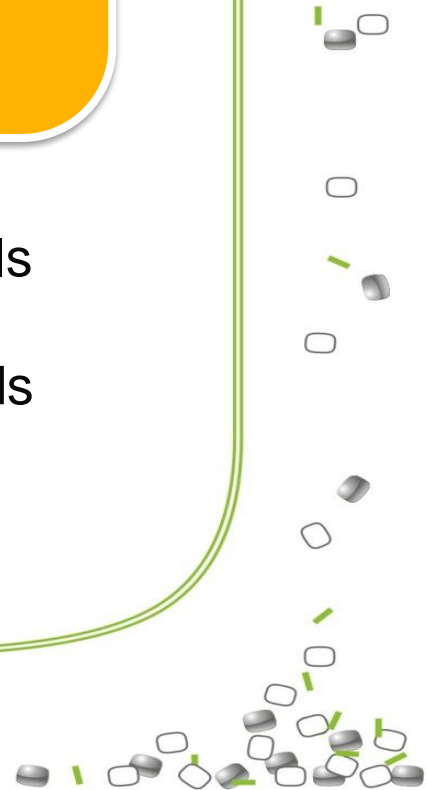
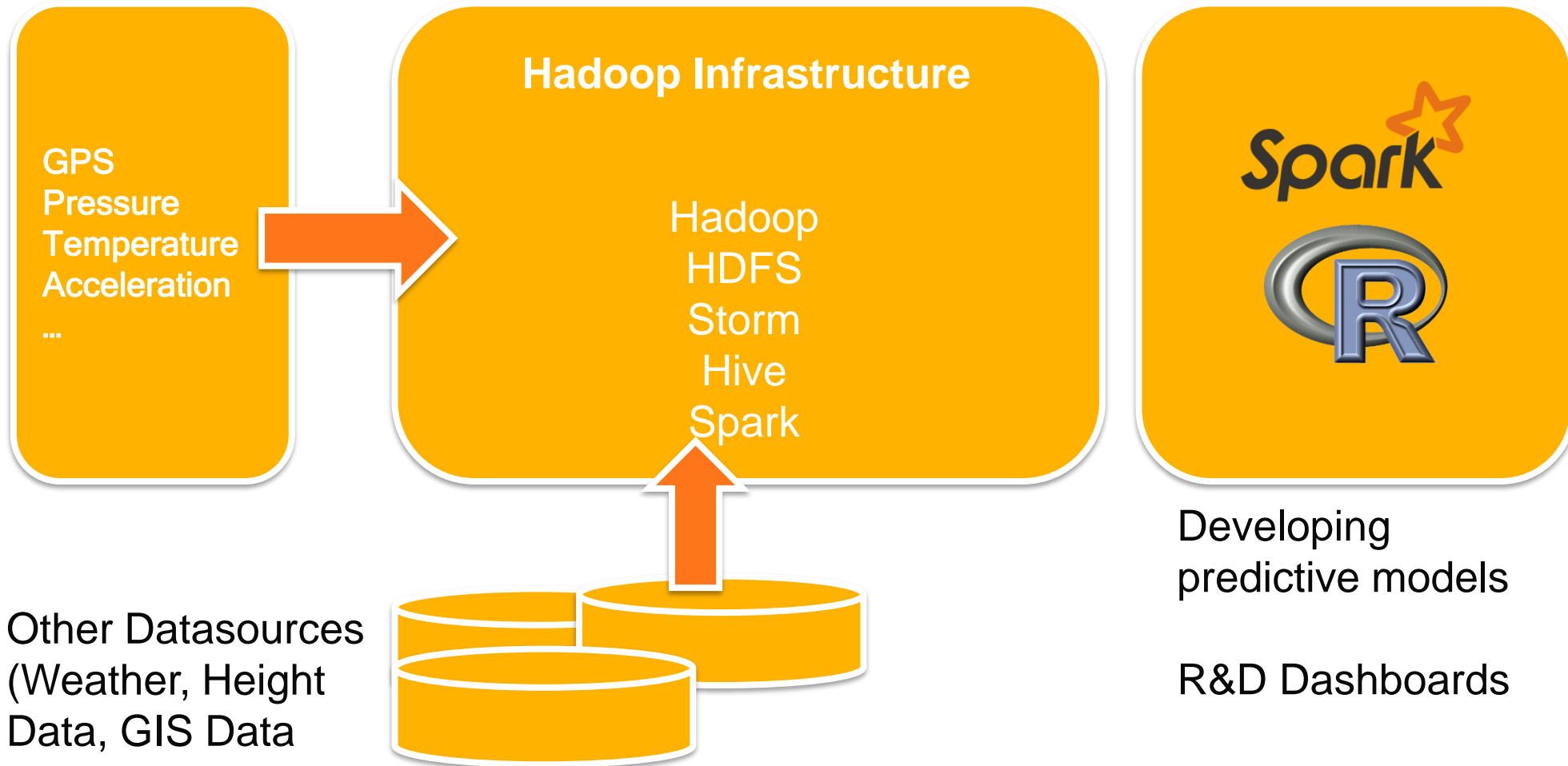
```
s.df <- read.df(sqlContext,  
  source = "com.databricks.spark.csv",  
  path = "/var/data/server-sample.log",  
  delimiter = " ", header = "false")
```

```
cache(s.df) # Bring Spark data.frame to R  
registerTempTable(s.df, "logs")
```

```
rc <- sql(sqlContext, "SELECT C0 AS ip, COUNT(*) AS n FROM logs GROUP BY C0 ORDER BY  
COUNT(*) DESC")
```



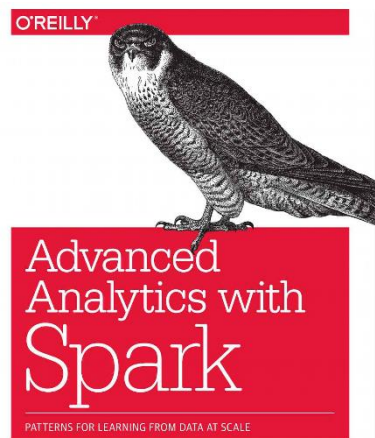
Analysing Sensorial Data



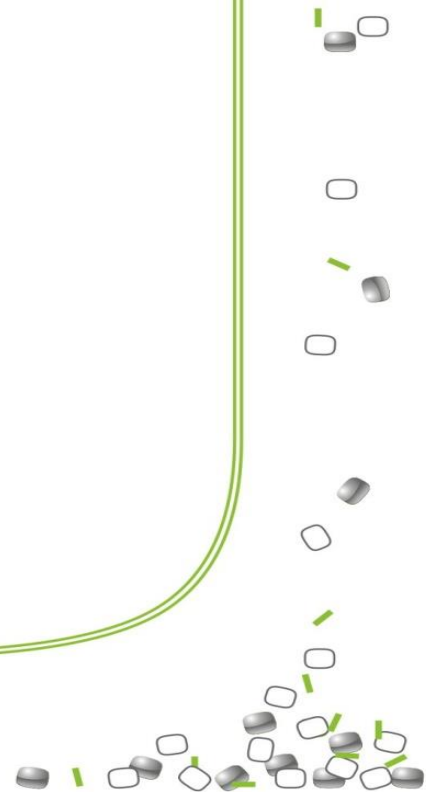
Further Sources

<https://spark.apache.org/docs/latest/api/R/index.html>

<https://spark.apache.org/docs/latest/sparkr.html>



Sandy Ryza, Uri Laserson,
Sean Owen & Josh Wills





Für Fragen stehen wir Ihnen gerne zur Verfügung!

© dsquare.de (2007-2015):

Diese Präsentation ist urheberrechtlich geschützt. Alle Nutzungs- und Verwertungsrechte liegen exklusiv bei der dsquare.de. Jede urheberrechtlich relevante Nutzung oder Verwertung dieser Präsentation oder von Teilen dieser Präsentation ist nur mit ausdrücklicher schriftlicher Zustimmung von dsquare.de zulässig. Dies gilt auch für die Weitergabe dieser Präsentation oder von Teilen dieser Präsentation an Dritte, für die diese Präsentation nicht bestimmt ist.

