

Steffen Jacobs

22. November 2019





Relevanz und Grundlagen

Was ist Apache Spark?

Wertversprechen

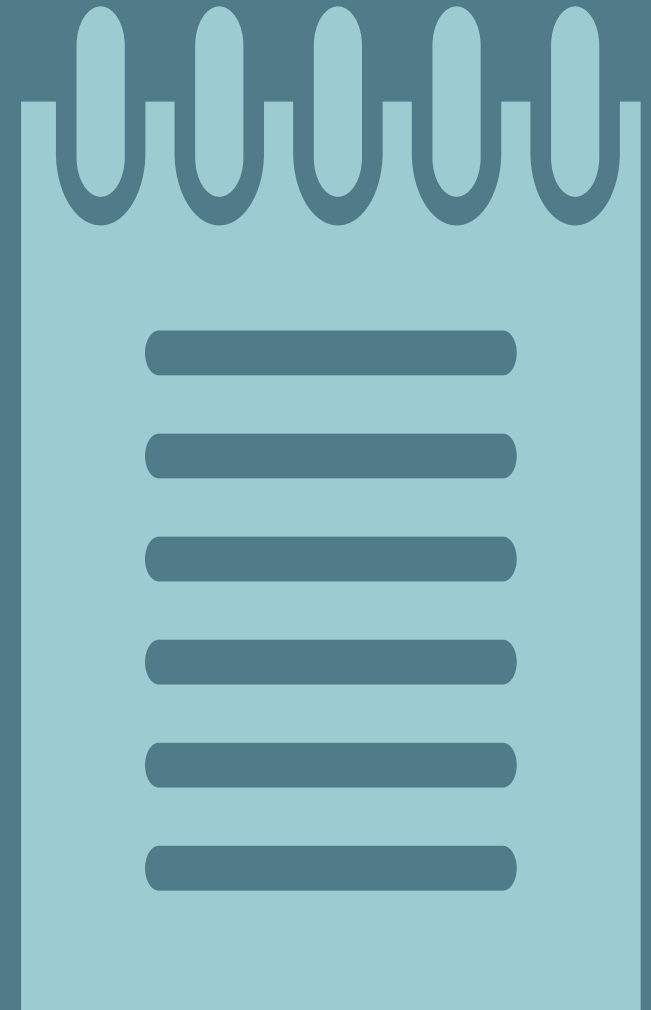
Was kann Apache Spark und wofür ist es zu gebrauchen?

Funktionsweise

Wie funktioniert Apache Spark?

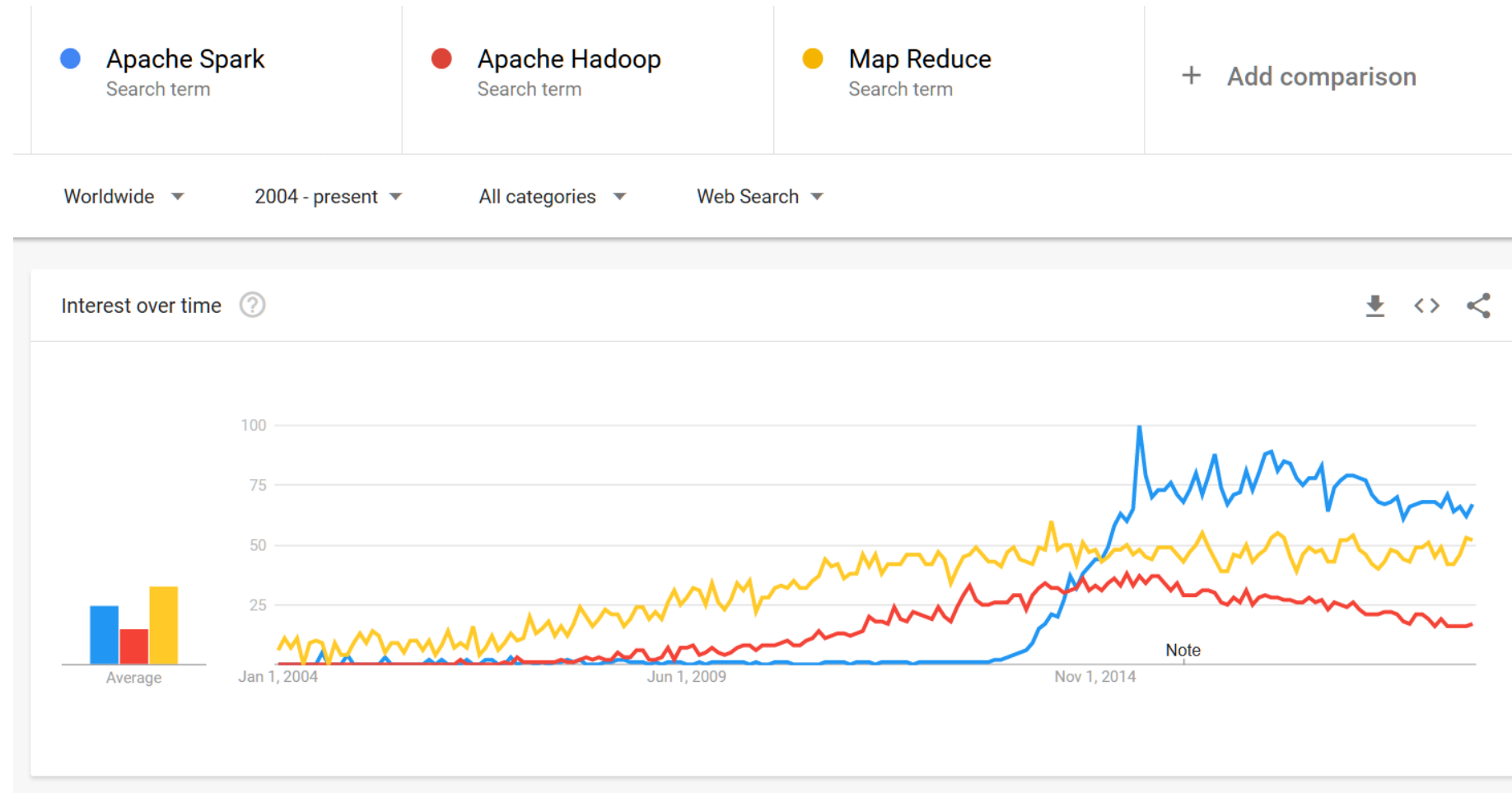
Demo

Wie funktioniert Apache Spark in der Praxis?



Relevanz

Findet Apache Spark überhaupt
Verwendung?



[Source: Google Trends](#)

Grundlagen

Was ist Apache Spark?



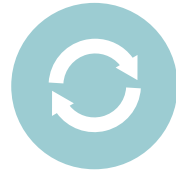
Datenverarbeitungseengine

Apache Spark ist eine clusterfähige Datenverarbeitungseengine. Alternativen sind Apache Flink, Apache Storm, Akka Reactive Streams, ggf. Hadoop.



In Memory Verarbeitung

Im Kontrast zu Map Reduce auf Hadoop verarbeitet Spark die Daten gleich im Speicher – Das Resultat ist ein significant höherer Durchsatz.



Transformation und Aggregation

Apache Spark soll dabei helfen, große Datenströme möglichst Echtzeit-nah zu verarbeiten.



Dataflow Programming

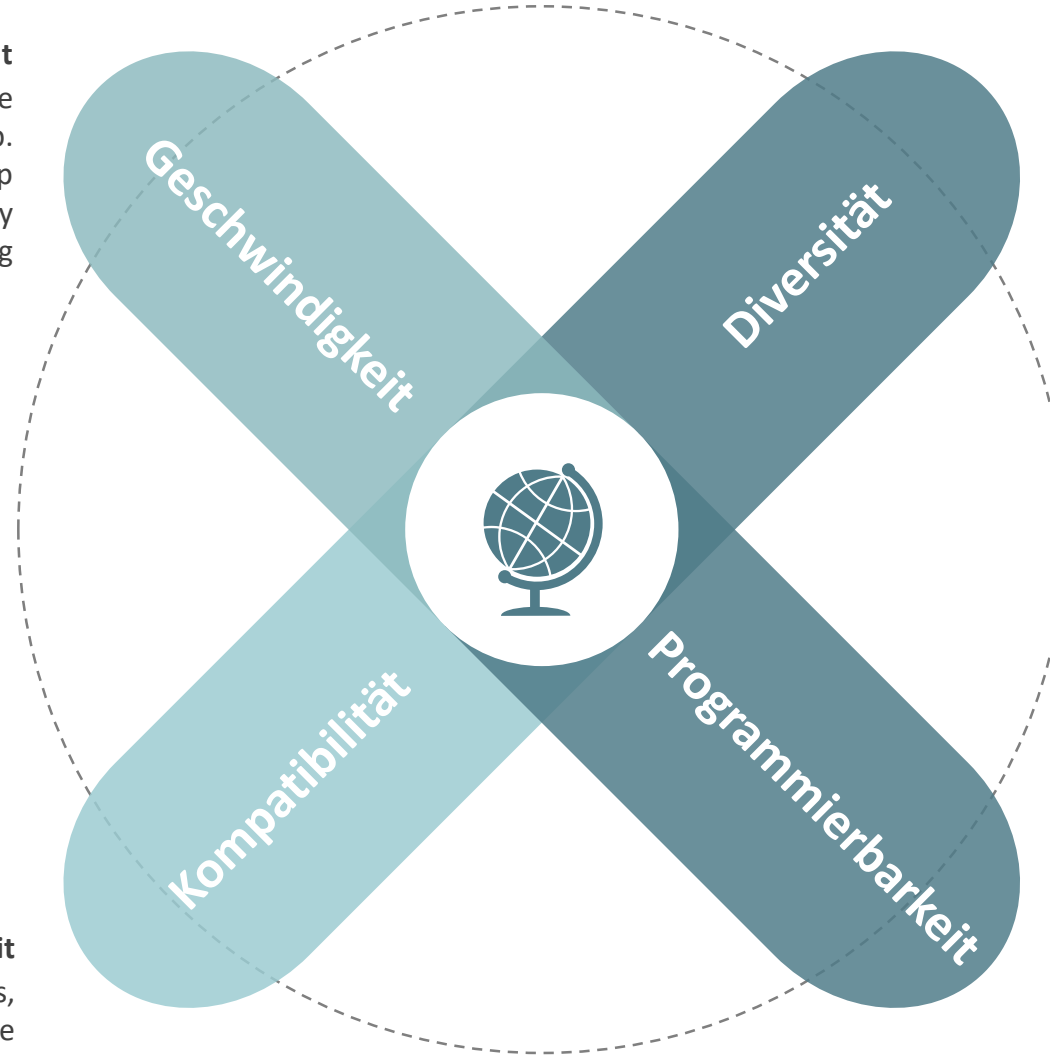
Apache Spark arbeitet direkt mit den Datenitems

Wert- versprechen

Was kann Apache Spark und
wofür ist es zu gebrauchen?

Geschwindigkeit
Signifikant bessere
Performance als zb.
Map-Reduce in Hadoop
dank In-Memory
Verarbeitung

Kompatibel mit
Hadoop, Apache Mesos,
Kubernetes, Standalone



All-In-One
Spark SQL, Spark Streaming,
Mlib (ML), GraphX

Polymorph
Java, Scala, Python, R, SQL, ...

Wert- versprechen

Was kann Apache Spark und
wofür ist es zu gebrauchen?

Geschwindigkeit
100TB sortieren

2013 Record:
Hadoop

2100 machines



72 minutes



2014 Record:
Spark

207 machines



23 minutes



Also sorted 1PB in 4 hours

[Source: Reynold Xin \(Stanford\)](#)

Wert- versprechen

Was kann Apache Spark und
wofür ist es zu gebrauchen?

Programmierbarkeit

```
1 package org.myorg;
2
3 import java.io.IOException;
4 import java.util.*;
5
6 import org.apache.hadoop.fs.Path;
7 import org.apache.hadoop.conf.*;
8 import org.apache.hadoop.io.*;
9 import org.apache.hadoop.mapreduce.*;
10 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
11 import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
12 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
13 import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
14
15 public class WordCount {
16
17     public static class Map extends Mapper<LongWritable, Text, Text, IntWritable> {
18         private final static IntWritable one = new IntWritable(1);
19         private Text word = new Text();
20
21         public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
22             String line = value.toString();
23             StringTokenizer tokenizer = new StringTokenizer(line);
24             while (tokenizer.hasMoreTokens()) {
25                 word.set(tokenizer.nextToken());
26                 context.write(word, one);
27             }
28         }
29     }
30
31     public static class Reduce extends Reducer<Text, IntWritable, Text, IntWritable> {
32
33         public void reduce(Text key, Iterable<IntWritable> values, Context context)
34             throws IOException, InterruptedException {
35             int sum = 0;
36             for (IntWritable val : values) {
37                 sum += val.get();
38             }
39             context.write(key, new IntWritable(sum));
40         }
41     }
42
43     public static void main(String[] args) throws Exception {
44         Configuration conf = new Configuration();
45
46         Job job = new Job(conf, "wordcount");
47
48         job.setOutputKeyClass(Text.class);
49         job.setOutputValueClass(IntWritable.class);
50
51         job.setMapperClass(Map.class);
52         job.setReducerClass(Reduce.class);
53
54         job.setInputFormatClass(TextInputFormat.class);
55         job.setOutputFormatClass(TextOutputFormat.class);
56
57         FileInputFormat.addInputPath(job, new Path(args[0]));
58         FileOutputFormat.setOutputPath(job, new Path(args[1]));
59
60         job.waitForCompletion(true);
61     }
62
63 }
```

```
val textFile = sc.textFile("hdfs://...")
val counts = textFile.flatMap(line => line.split(" "))
                        .map(word => (word, 1))
                        .reduceByKey(_ + _)
counts.saveAsTextFile("hdfs://...")
```

Spark: 5 LoC

Map Reduce: 63 LoC

[Source: Rainer Gemulla \(Mannheim\)](#)

Funktionsweise

Wie funktioniert Apache Spark?



Apache Spark mit Hadoop

Wie funktioniert Apache Spark mit Hadoop?

Apache Spark vs. Map Reduce

Apache Spark oder Map Reduce?

Resilient Distributed Datasets

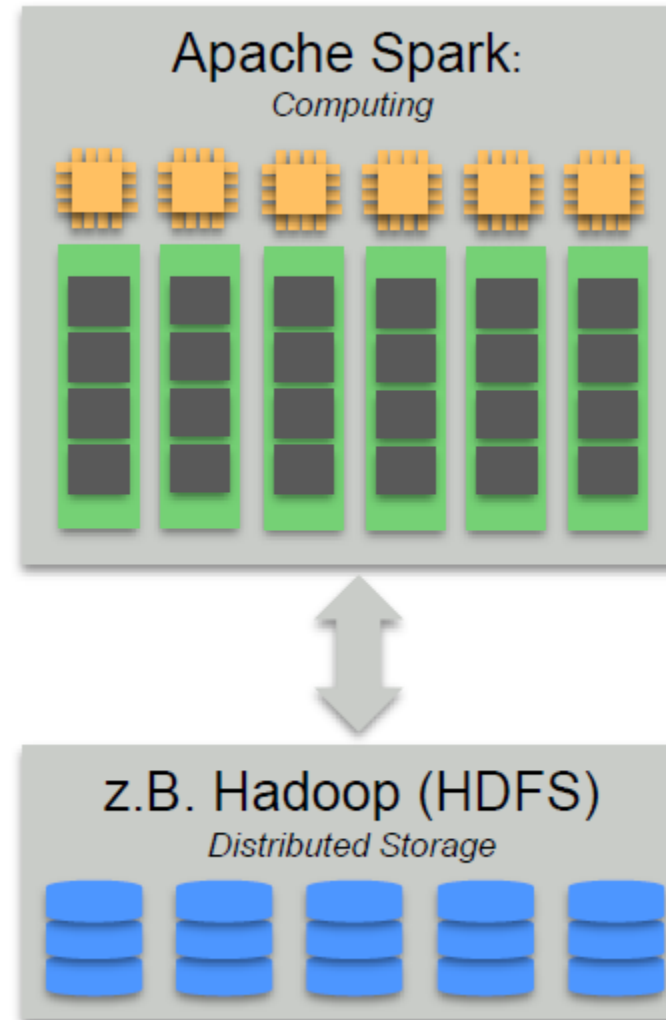
Die kleinste Dateneinheit in Apache Spark



Funktionsweise

Wie funktioniert Apache Spark?

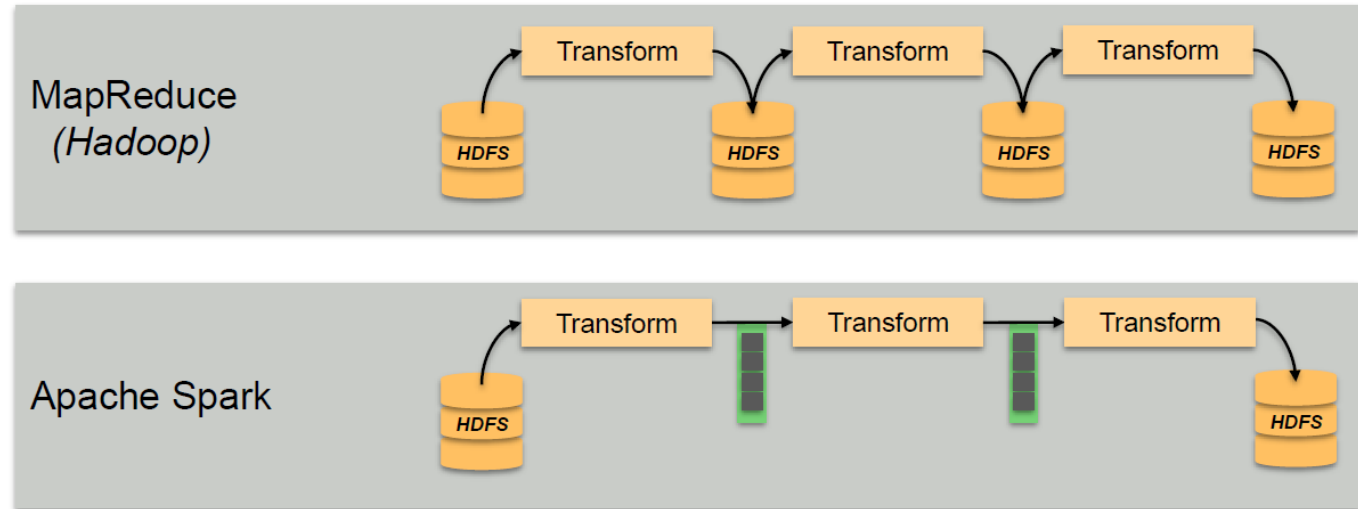
Apache Spark mit
Hadoop



Funktionsweise

Wie funktioniert Apache Spark?

Apache Spark
Vs
Map Reduce



Funktionsweise

Wie funktioniert Apache Spark?

Resilient Distributed Datasets



Konzeptionelle Objektmenge

Vergleichbar mit Reihen in einer Datenbanktabelle



Verteilbar

RDDs werden über das Cluster verteilt



Immutable und Disposable

RDDs sind unveränderbar

RDDs speichern keine Daten sondern halten sie nur temporär



Unterstützen Operationen

Transformationen

Actions

Demo

Wie funktioniert das in der Praxis?



Fragen

Gibt es Fragen?

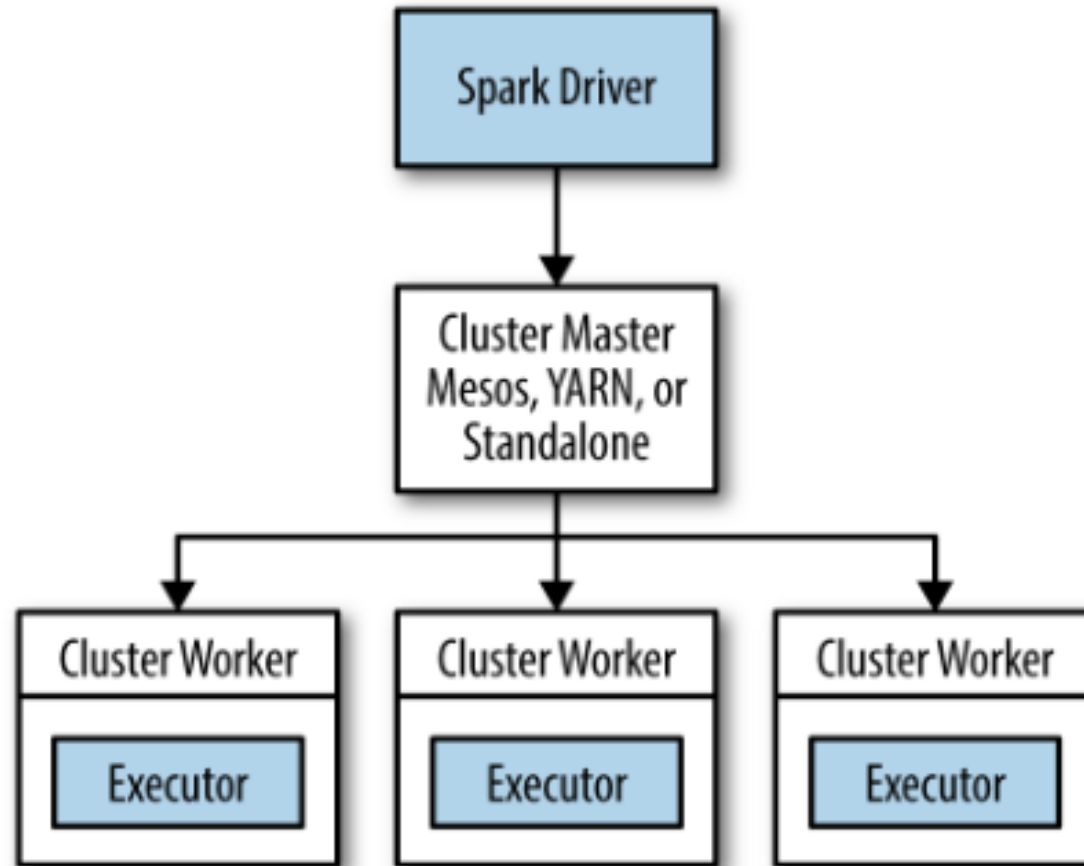


Vielen Dank



Backup

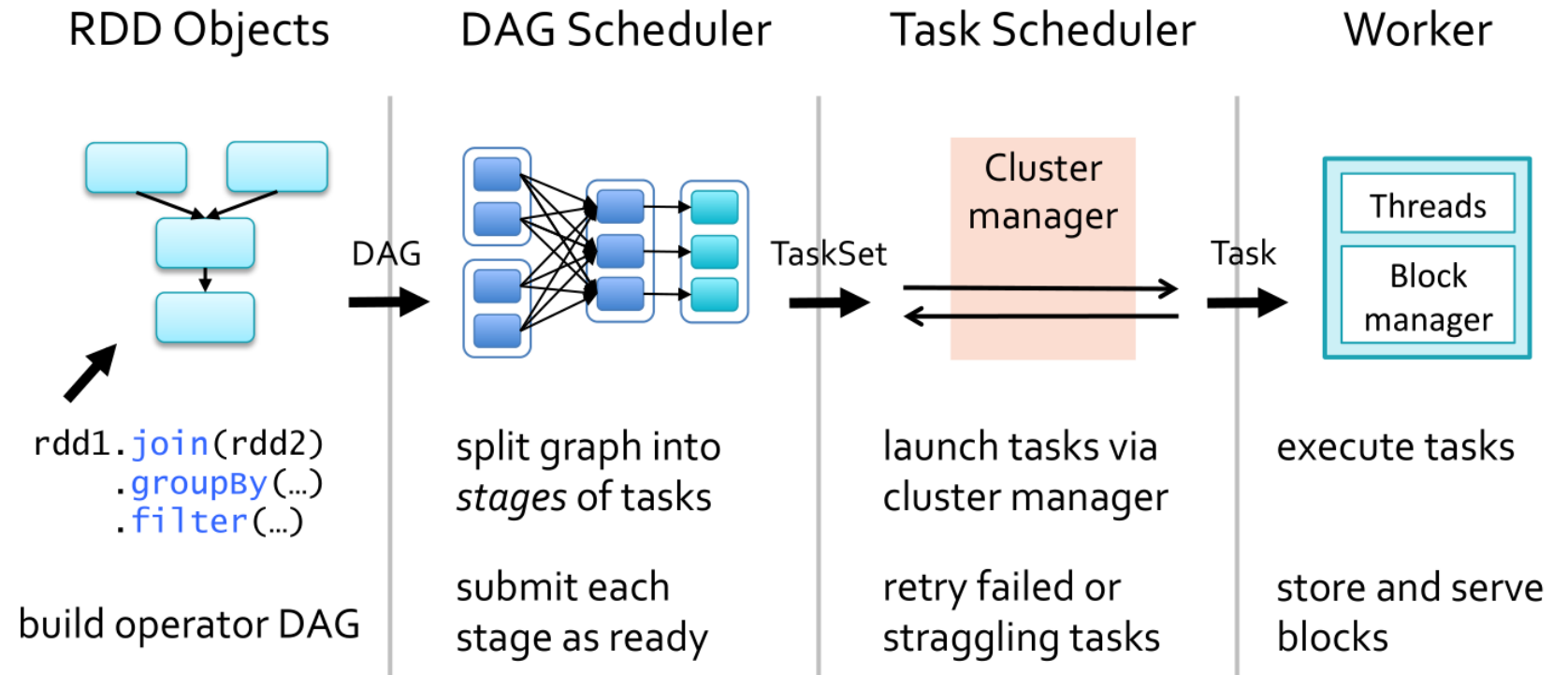
Spark Architektur



[Source: Rainer Gemulla \(Mannheim\)](#)

Backup

Spark Ausführung



Source: Rainer Gemulla (Mannheim)