

Scala and the JVM for Big Data: Lessons from Spark

Scala World
Sept. 22, 2015



1

©Dean Wampler 2014-2015, All Rights Reserved

Tuesday, September 22, 15

Photos from Olympic National Park, Washington State, USA, Aug. 2015.
All photos are copyright (C) 2015, Dean Wampler, except where otherwise noted. All Rights Reserved.

polyglotprogramming.com/talks
dean.wampler@typesafe.com
[@deanwampler](https://twitter.com/deanwampler)



2

Tuesday, September 22, 15

You can find this and my other talks here.



spark

3

Tuesday, September 22, 15

Scala has become popular for Big Data tools and apps, such as Spark. Why is Spark itself so interesting?

Productivity?

Very concise, elegant, functional APIs.

- Scala, Java
- Python, R
- ... and SQL!

4

Tuesday, September 22, 15

We saw an example why this true.

While Spark was written in Scala, it has Java, Python, R, and even SQL APIs, too, which means it can be a single tool used across a Big Data organization, engineers and data scientists.

Productivity?

Interactive shell (REPL)

- Scala, Python, R, and SQL

Spark Notebook (iPython/Jupyter-like)

5

Tuesday, September 22, 15

This is especially useful for the SQL queries we'll discuss, but also handy once you know the API for experimenting with data and/or algorithms. In fact, I tend to experiment in the REPL or Spark Notebook, then copy the code to a more "permanent" form, when it's supposed to be compiled and run as a batch program.

Performance?

Lazy API, inspired by Scala collections.

- Dataflow DAG of processing nodes.
- Distributed and run on a cluster.
- MOAR optimizations coming...

6

Tuesday, September 22, 15

How is Spark more efficient? As we'll see, Spark programs are actually "lazy" dataflows definitions that are only evaluated on demand. Because Spark has this directed acyclic graph of steps, it knows what data to attempt to cache in memory between steps (with programmable tweaks) and it can combine many logical steps into one "stage" of computation, for efficient execution while still providing an intuitive API experience.

We'll come back to the optimizations.

Performance?

Lazy API, inspired by Scala collections.

- Combines node steps into “stages”.
- Can cache intermediate data.



8

Tuesday, September 22, 15



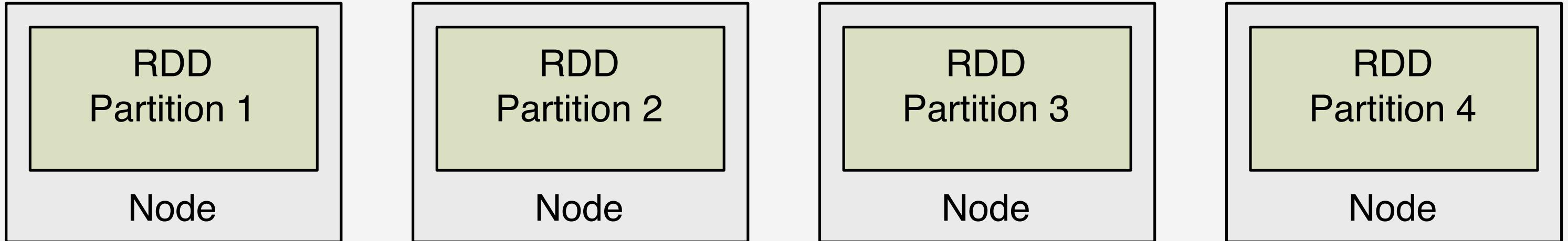
Resilient Distributed Dataset

9

Tuesday, September 22, 15

The core abstraction in the core of Spark.

Cluster



- Lazily constructed.
- Each DAG node has an RDD, but...
- Intermediate RDDs not materialized.

10

Tuesday, September 22, 15

The core concept is a Resilient Distributed Dataset, a partitioned collection. They are resilient because if one partition is lost, Spark knows the lineage and can reconstruct it. However, you can also cache RDDs to eliminate having to walk back too far. RDDs are immutable, but they are also an abstraction, so you don't instantiate a new RDD with wasteful data copies for each stage of the pipeline, but rather at the end of each stage..



Example: Inverted Index

11

Tuesday, September 22, 15

Let's look at a small, real actual Spark program, the Inverted Index.

wikipedia.org/hadoop

Hadoop provides
MapReduce and HDFS

...

wikipedia.org/hbase

HBase stores data in HDFS

...

inverse index

block

| | |
|--------|---|
| ... | ... |
| hadoop | (.../hadoop,1) |
| hbase | (.../hbase,1),(.../hive,1) |
| hdfs | (.../hadoop,1),(.../hbase,1),(.../hive,1) |
| hive | (.../hive,1) |
| ... | ... |

block

| | |
|-----|-----|
| ... | ... |
|-----|-----|

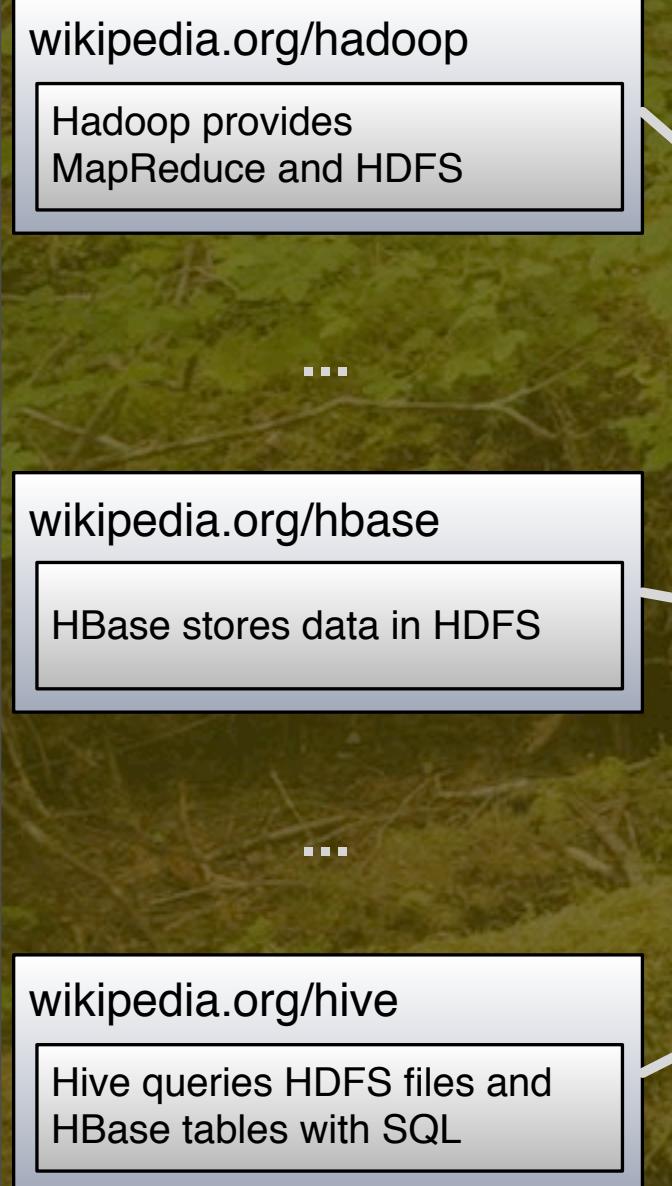
block

| | |
|-----|-----|
| ... | ... |
|-----|-----|

Tuesday, September 22, 15

Let's look at a small, real actual Spark program, the Inverted Index.

Web Crawl



Compute Inverted Index

| index | |
|----------------------|--------------------|
| block | |
| ... | ... |
| wikipedia.org/hadoop | Hadoop provides... |
| ... | ... |

| block | |
|---------------------|-----------------|
| ... | ... |
| wikipedia.org/hbase | HBase stores... |
| ... | ... |

| block | |
|--------------------|-----------------|
| ... | ... |
| wikipedia.org/hive | Hive queries... |
| ... | ... |

Miracle!!

| inverse index | |
|---------------|---|
| block | |
| ... | ... |
| hadoop | (.../hadoop,1) |
| hbase | (.../hbase,1),(.../hive,1) |
| hdfs | (.../hadoop,1),(.../hbase,1),(.../hive,1) |
| hive | (.../hive,1) |
| ... | ... |

| block | |
|-------|-----|
| ... | ... |

| block | |
|-------|-----|
| ... | ... |

| block | |
|-------|-----------------------------|
| ... | ... |
| and | (.../hadoop,1),(.../hive,1) |
| ... | ... |

```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._

val sparkContext = new
  SparkContext(master, "Inv. Index")
sparkContext.textFile("/path/to/input") .
map { line =>
  val array = line.split(",", 2)
  (array(0), array(1))
}.flatMap {
  case (id, contents) =>
    toWords(contents).map(w => ((w,id),1))
}.reduceByKey(_ + _).
map {
  case ((word,id),n) => (word,(id,n))
}.groupByKey.
mapValues {
  seq => sortByCount(seq)
}.saveAsTextFile("/path/to/output")
```

14

Tuesday, September 22, 15

Here is an actual Spark program, written as a script.

```
import org.apache.spark.SparkContext  
import org.apache.spark.SparkContext._  
  
val sparkContext = new  
  SparkContext(master, "Inv. Index")  
sparkContext.textFile("/path/to/input").  
map { line =>  
  val array = line.split(",", 2)  
  (array(0), array(1))  
}.flatMap {  
  case (id, contents) =>  
    toWords(contents).map(w => ((w, id), 1))  
}
```

15

Tuesday, September 22, 15

Zoom in to analyze...

Start with the imports. The SparkContext is the entry point for all programs.

```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._

val sparkContext = new
SparkContext(master, "Inv. Index")
sparkContext.textFile("/path/to/input").
map { line =>
  val array = line.split(",", 2)
  (array(0), array(1))
}.flatMap {
  case (id, contents) =>
  toWords(contents).map(w => ((w,id),1))
}
```

16

Tuesday, September 22, 15

Create a SparkContext, specifying the “master” (local for single core on the local machine, “local[*]” for all cores, “mesos://...” for Mesos, “yarn-*” for YARN, etc.)

```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._

val sparkContext = new
  SparkContext(master, "Inv. Index")
sparkContext.textFile("/path/to/input").
map { line =>
  val array = line.split(",", 2)
  (array(0), array(1))
}.flatMap {
  case (id, contents) =>
  toWords(contents).map(w => ((w,id),1))
```

17

Tuesday, September 22, 15

Using the sparkContext, read test data (the web crawl data). Assume it's comma-delimited and split it into the identifier (e.g., URL) and the contents. If the input is very large, multiple, parallel tasks will be spawned across the cluster to read and process each file block as a partition, all in parallel.

```
val array = line.split( ' ', 2 )
(array(0), array(1))
}.flatMap {
  case (id, contents) =>
    toWords(contents).map(w => ((w,id),1))
}.reduceByKey(_ + _).
map {
  case ((word,id),n) => (word,(id,n)))
}.groupByKey.
mapValues {
  seq => sortByCount(seq)
}.saveAsTextFile("/path/to/output")
```

18

Tuesday, September 22, 15

An idiom for counting...

Flat map to tokenize the contents into words (using a “toWords” function that isn’t shown), then map over those words to nested tuples, where the (word, id) is the key, and see count of 1 is the value.

Then use reduceByKey, which is an optimized groupBy where we don’t need the groups, we just want to apply a function to reduce the values for each unique key. Now the records are unique (word,id) pairs and counts ≥ 1 .

```
val array = line.split( ' ', 2 )
(array(0), array(1))
}.flatMap {
case (id, contents) =>
toWords(contents).map(w => ((w,id),1))
}.reduceByKey(_ + _).
map {
case ((word,id),n) => (word,(id,n)))
}.groupByKey.
mapvalues {
seq => sortByCount(seq)
}.saveAsTextFile("/path/to/output")
```

19

Tuesday, September 22, 15

I love how simple this line is anon. function is; by moving the nested parentheses, we setup the final data set, where the words alone are keys and the values are (path, count) pairs. Then we group over the words (the “keys”).

```
val array = line.split( ' ', 2 )
(array(0), array(1))
}.flatMap {
case (id, contents) =>
toWords(contents).map(w => ((w,id),1))
}.reduceByKey(_ + _).
map {
case ((word,id),n) => (word,(id,n)))
}.groupByKey.
mapValues {
seq => sortByCount(seq)
}.saveAsTextFile("/path/to/output")
```

20

Tuesday, September 22, 15

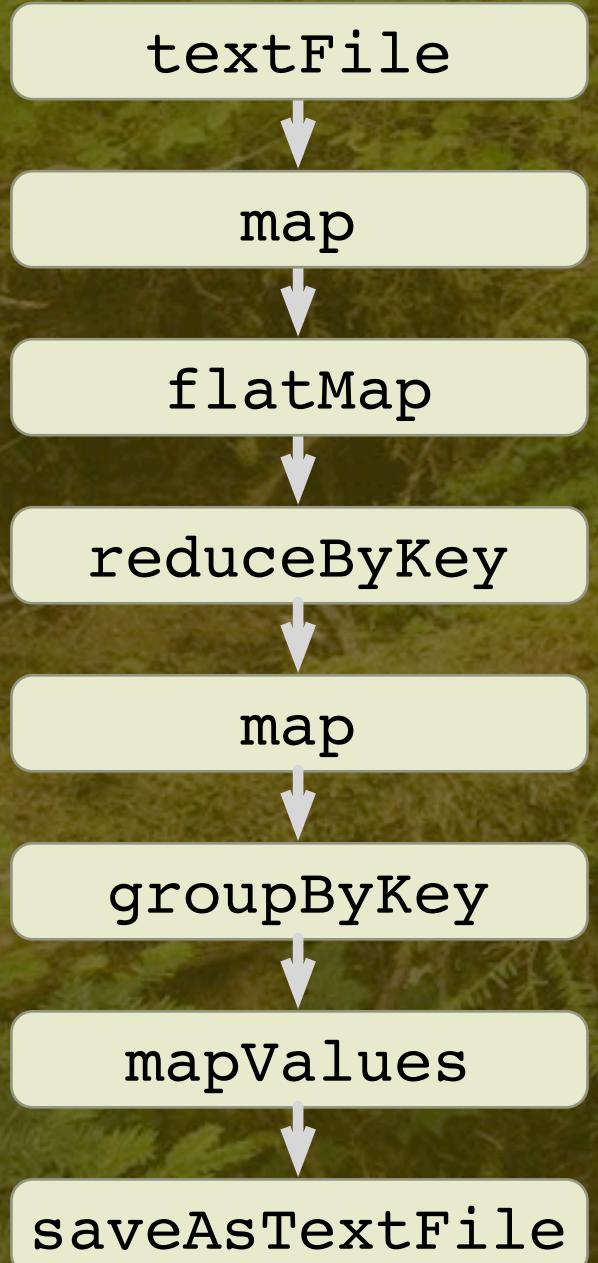
Finally, for each unique word, sort the nested collection of (path,count) pairs by count descending (this is easy with Scala's Seq.sortBy method). Then save the results as text files.

```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._

val sparkContext = new
  SparkContext(master, "Inv. Index")
sparkContext.textFile("/path/to/input").
map { line =>
  val array = line.split(",", 2)
  (array(0), array(1))
}.flatMap {
  case (id, contents) =>
  toWords(contents).map(w => ((w,id),1))
}.reduceByKey(_ + _).
map {
  case ((word,id),n) => (word,(id,n))
}.groupByKey.
mapValues {
  seq => sortByCount(seq)
}.saveAsTextFile("/path/to/output")
```

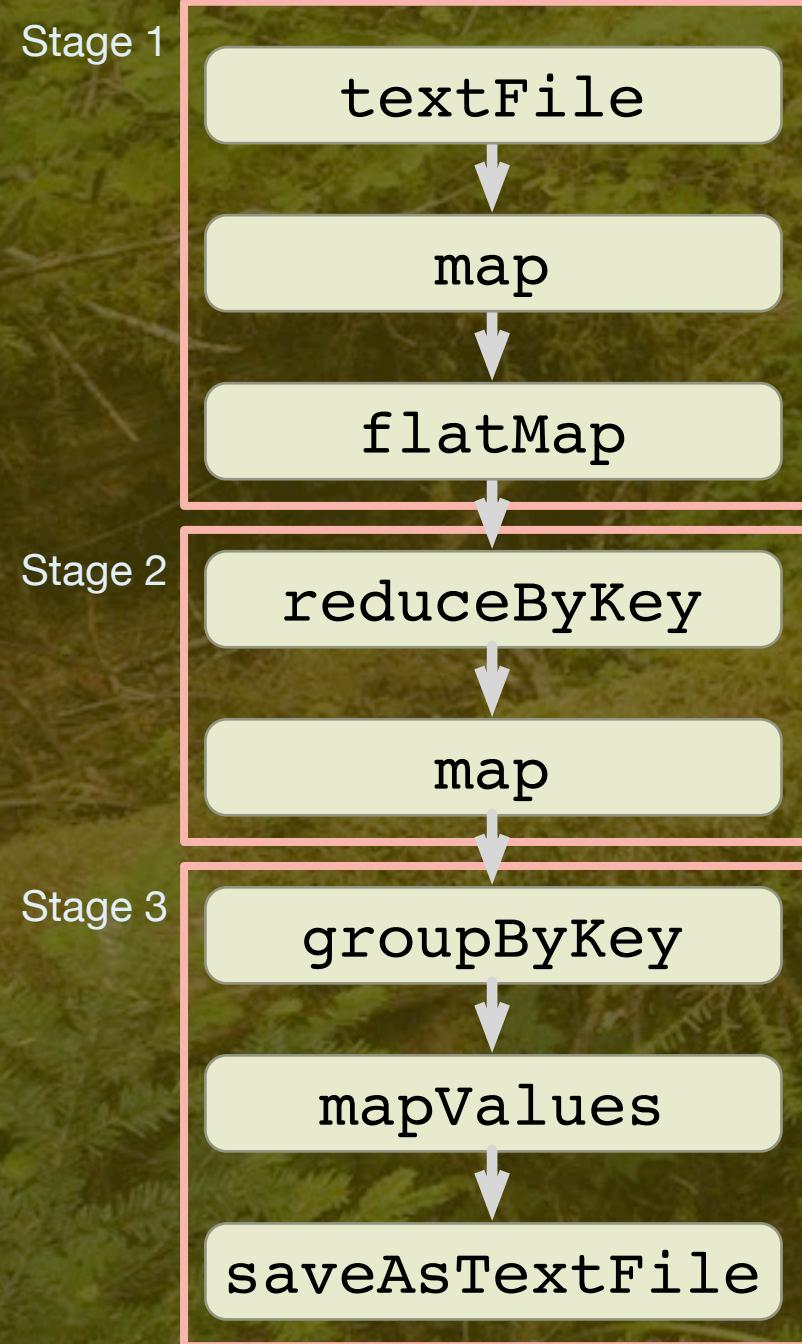
Altogether

Stages



```
import org.apache.spark.SparkContext  
import org.apache.spark.SparkContext._  
  
val sparkContext = new  
  SparkContext(master, "Inv. Index")  
sparkContext.textFile("/path/to/input").  
map { line =>  
  val array = line.split(",", 2)  
  (array(0), array(1))  
}.flatMap {  
  case (id, contents) =>  
    toWords(contents).map(w => ((w,id),1))  
}.reduceByKey(_ + _).  
map {  
  case ((word,id),n) => (word,(id,n))  
}.groupByKey.  
mapValues {  
  seq => sortByCount(seq)  
}.saveAsTextFile("/path/to/output")
```

Stages



```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._

val sparkContext = new
  SparkContext(master, "Inv. Index")
sparkContext.textFile("/path/to/input").
map { line =>
  val array = line.split(",", 2)
  (array(0), array(1))
}.flatMap {
  case (id, contents) =>
    toWords(contents).map(w => ((w,id),1))
}.reduceByKey(_ + _).
map {
  case ((word,id),n) => (word,(id,n))
}.groupByKey.
mapValues {
  seq => sortByCount(seq)
}.saveAsTextFile("/path/to/output")
```

23

Tuesday, September 22, 15

The transformation steps that don't require data from other partitions can be pipelined together into a single JVM process (per partition), called a Stage. When you do need to bring together data from different partitions, such as group-bys, joins, reduces, then data must be "shuffled" between partitions (i.e., all keys of a particular value must arrive at the same JVM instance for the next transformation step). That triggers a new stage, as shown. So, this algorithm requires three stages and the RDDs are materialized only for the last steps in each stage.



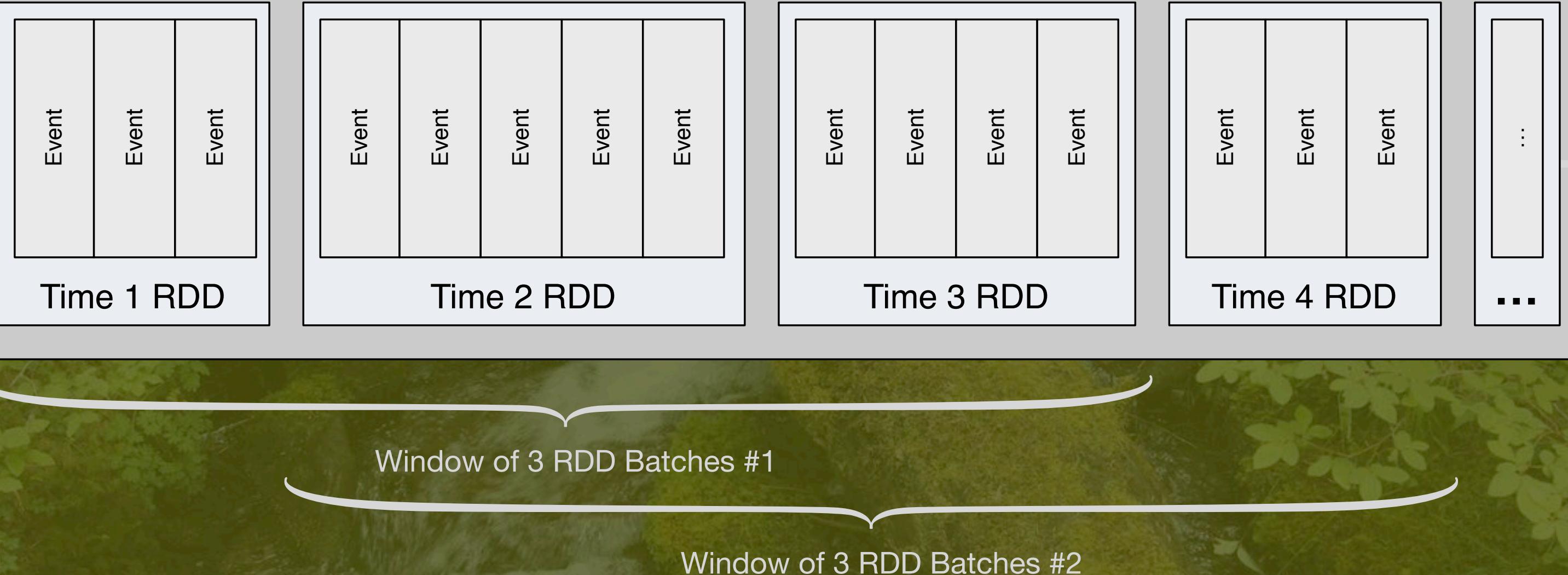
24

Tuesday, September 22, 15

DStreams: Spark Streaming

25

DStream (discretized stream)



26

Tuesday, September 22, 15

For streaming, one RDD is created per batch iteration, with a DStream (discretized stream) holding all of them, which supports window functions. Spark started life as a batch-mode system, just like MapReduce, but Spark's dataflow stages and in-memory, distributed collections (RDDs - resilient, distributed datasets) are lightweight enough that streams of data can be timesliced (down to ~1 second) and processed in small RDDs, in a “mini-batch” style. This gracefully reuses all the same RDD logic, including your code written for RDDs, while also adding useful extensions like functions applied over moving windows of these batches.

Flexibility

Composable operators, performance optimizations, and general flexibility are a foundation for higher-level APIs...

27

Tuesday, September 22, 15

A major step forward. Due to the lightweight nature of Spark processing, it can efficiently support a wider class of algorithms.

A scenic mountain landscape featuring a small lake in the distance, surrounded by dense green forests and rugged, rocky mountain peaks under a clear blue sky.

SQL/ DataFrames

28

Tuesday, September 22, 15

Like Hive for MapReduce, a subset of SQL (omitting transactions and the U in CRUD) is relatively easy to implement as a DSL on top of a general compute engine like Spark. Hence, the SQL API was born, but it's grown into a full-fledged programmatic API supporting both actual SQL queries and an API similar to Python's DataFrame API for working with structured data in a more type-safe way (errr, at least for Scala).

Example

```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.sql.SQLContext

val sparkContext = new
  SparkContext(master, "...")
val sqlContext = new
  SQLContext(sparkContext)
val flights =
  sqlContext.read.parquet(".../flights")
val planes =
  sqlContext.read.parquet(".../planes")
flights.registerTempTable("flights")
planes.registerTempTable("planes")
flights.cache(); planes.cache()

val planes_for_flights1 = sqlContext.sql("""
  SELECT * FROM flights f
  JOIN planes p ON f.tailNum = p.tailNum LIMIT 100""")

val planes_for_flights2 =
  flights.join(planes,
    flights("tailNum") ===
    planes ("tailNum")).limit(100)
```

Tuesday, September 22, 15

Example using SparkSQL to join two data sets (adapted from Typesafe's Spark Workshop training), data for flights and information about the planes.

```
import org.apache.spark.SparkContext  
import org.apache.spark.SparkContext._  
import org.apache.spark.sql.SQLContext  
  
val sparkContext = new  
  SparkContext(master, "...")  
val sqlContext = new  
  SQLContext(sparkContext)  
val flights =  
  sqlContext.read.parquet(".../flights")  
val planes =  
  salContext.read.parquet(".../planes")
```

30

Tuesday, September 22, 15

Now we also need a SQLContext.

```
+ SQLContext(sparkContext)  
val flights =  
    sqlContext.read.parquet(".../flights")  
val planes =  
    sqlContext.read.parquet(".../planes")  
flights.registerTempTable("flights")  
planes.registerTempTable("planes")  
flights.cache(); planes.cache()
```

```
val planes_for_flights1 =  
sqlContext.sql(  
    "SELECT * FROM flights f
```

31

Tuesday, September 22, 15

Read the data as Parquet files, which include the schemas. Create temporary tables (purely virtual) for SQL queries, and cache the tables for faster, repeated access.

```
val planes_for_flights1 =  
sqlContext.sql("""  
SELECT * FROM flights f  
JOIN planes p ON f.tailNum =  
p.tailNum LIMIT 100""")
```

```
val planes_for_flights2 =  
flights.join(planes,  
flights("tailNum") ===  
planes ("tailNum")).limit(100)
```

```
val planes_for_flights1 =  
sqlContext.sql("""  
SELECT * FROM flights f  
.JOIN planes p ON f.tailNum =  
p.tailNum LIMIT 100""")
```

```
val planes_for_flights2 =  
flights.join(planes,  
flights("tailNum") ===  
planes ("tailNum")).limit(100)
```

Altogether

```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.sql.SQLContext

val sparkContext = new
  SparkContext(master, "...")
val sqlContext = new
  SQLContext(sparkContext)
val flights =
  sqlContext.read.parquet(".../flights")
val planes =
  sqlContext.read.parquet(".../planes")
flights.registerTempTable("flights")
planes.registerTempTable("planes")
flights.cache(); planes.cache()

val planes_for_flights1 = sqlContext.sql("""
  SELECT * FROM flights f
  JOIN planes p ON f.tailNum = p.tailNum LIMIT 100""")

val planes_for_flights2 =
  flights.join(planes,
    flights("tailNum") ===
    planes ("tailNum")).limit(100)
```

34

Tuesday, September 22, 15

Example using SparkSQL to join two data sets (adapted from Typesafe's Spark Workshop training), data for flights and information about the planes.



MLlib

Tuesday, September 22, 15

Many machine learning libraries are being implemented on top of Spark. An important requirement is the ability to do linear algebra and iterative algorithms quickly and efficiently, which is used in the training algorithms for many ML models.

- Machine Learning requires:
 - Iterative training of models.
 - Good linear algebra perf.

MLlib

Tuesday, September 22, 15

Many machine learning libraries are being implemented on top of Spark. An important requirement is the ability to do linear algebra and iterative algorithms quickly and efficiently, which is used in the training algorithms for many ML models.



GraphX

Tuesday, September 22, 15

Similarly, efficient iteration makes graph traversal algorithms tractable, enabling “first-class” graph representations of data, like social networks.

- Graph algorithms require:
 - Incremental traversal.
 - Efficient edge and node reps.

GraphX

Tuesday, September 22, 15

Similarly, efficient iteration makes graph traversal algorithms tractable, enabling “first-class” graph representations of data, like social networks.

Foundation:

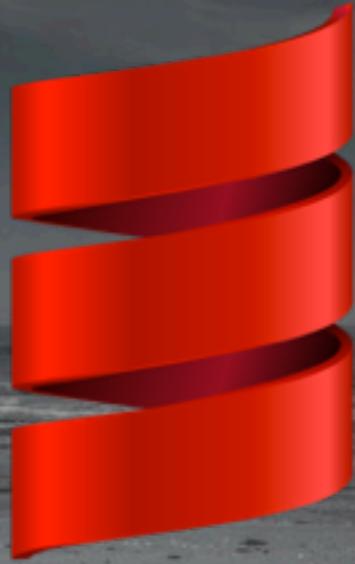
The JVM

39

Tuesday, September 22, 15

With that introduction, let's step back and look at the larger context, starting at the bottom; why is the JVM the platform of choice for Big Data?

Tools and Libraries



Akka
Breeze
Algebird
Spire & Cats
Axe

...

40

Tuesday, September 22, 15

We have a rich suite of mature(-ish) languages, development tools, and math-related libraries for building data applications...

<http://akka.io>

<https://github.com/scalanlp/breeze>

<https://github.com/twitter/algebird>

<https://github.com/non/spire>

<https://github.com/non/cats>

<http://axle-lang.org/>

20 Years of DevOps

and

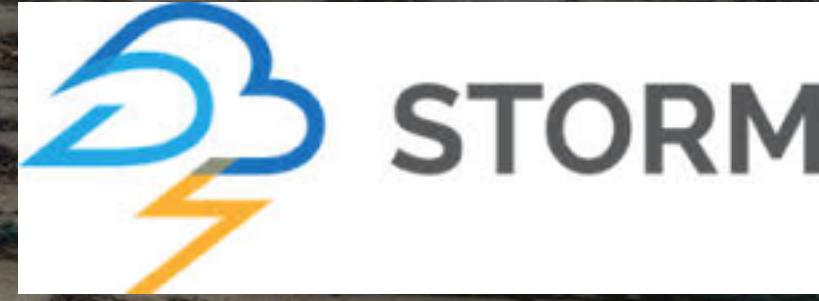
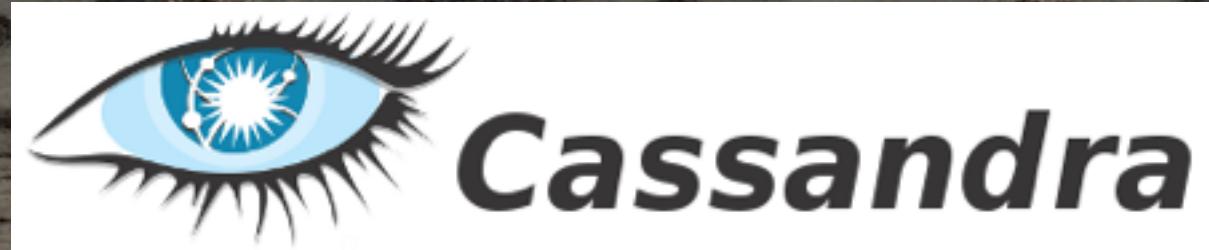
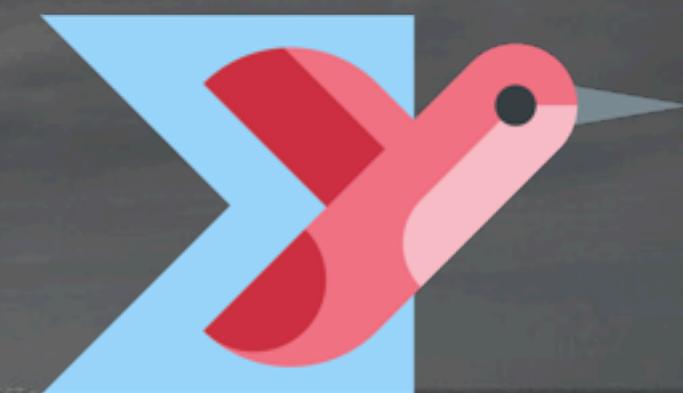
Lots of Java Devs

41

Tuesday, September 22, 15

We have 20 years of operations experience running JVM-based apps in production. We have many Java developers, some with 20 years of experience, writing JVM-based apps.

Big Data Ecosystem



Samza

42

Tuesday, September 22, 15

All this is why most Big Data tools in use have been built on the JVM, including Spark, especially those for OSS projects created and used by startups.

<http://spark.apache.org>

<https://github.com/twitter/scalding>

<https://github.com/twitter/summingbird>

<http://kafka.apache.org>

<http://hadoop.apache.org>

<http://cassandra.apache.org>

<http://storm.apache.org>

<http://samza.apache.org>

<http://lucene.apache.org/solr/>



But it's
not perfect....

43

Tuesday, September 22, 15

But no system is perfect. The JVM wasn't really designed with Big Data systems, as we call them now, in mind.

Garbage Collection

Issues

44

Tuesday, September 22, 15

Garbage collection is a wonder thing for removing a lot of tedium and potential errors from code, but the default settings in the JVM are not ideal for Big Data apps.

GC Challenges

- Typical heaps 10s-100s of GB.
- Uncommon in generic services.

45

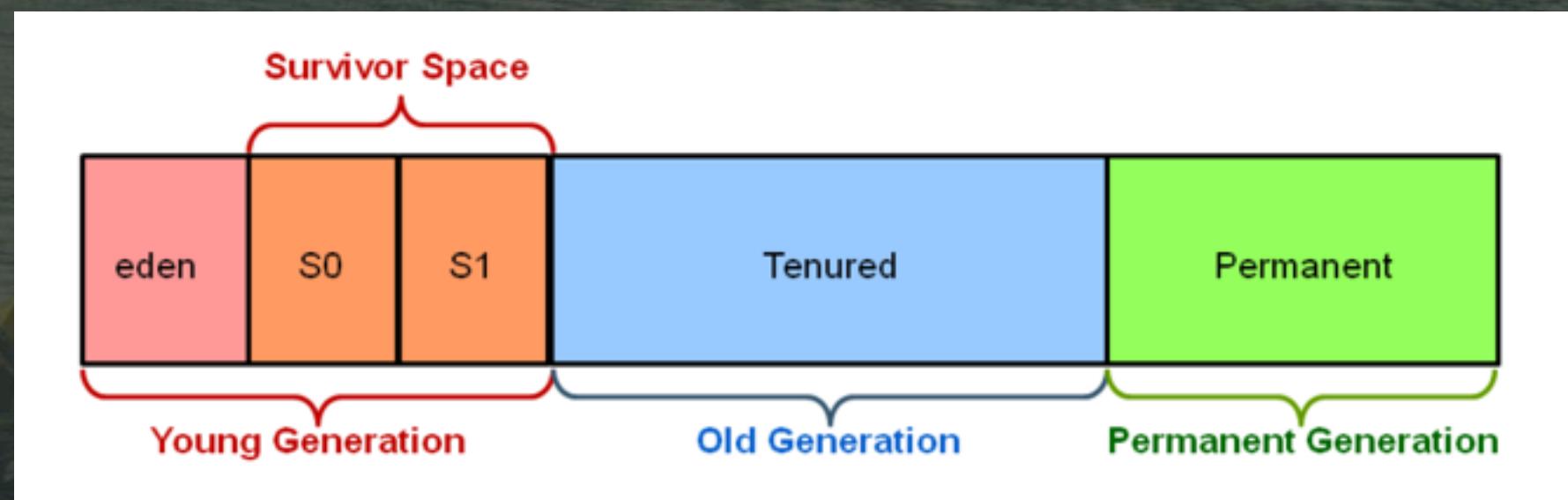
Tuesday, September 22, 15

We'll explain the details shortly, but the programmer controls which data sets are cached to optimize usage.

See <https://databricks.com/blog/2015/05/28/tuning-java-garbage-collection-for-spark-applications.html> for a detailed overview of GC challenges and tuning for Spark.

GC Challenges

- Too many cached RDDs leads to huge old generation garbage.
- Leads to long GC “stop the world” pauses.



46

Tuesday, September 22, 15

We'll explain the details shortly, but the programmer controls which data sets are cached to optimize usage.

See <https://databricks.com/blog/2015/05/28/tuning-java-garbage-collection-for-spark-applications.html> for a detailed overview of GC challenges and tuning for Spark.

Image from <http://www.oracle.com/webfolder/technetwork/tutorials/obe/java/gc01/index.html>.

GC Challenges

- Those cached RDDs can hold a lot of reference objects.
 - (Lots of cache misses.)

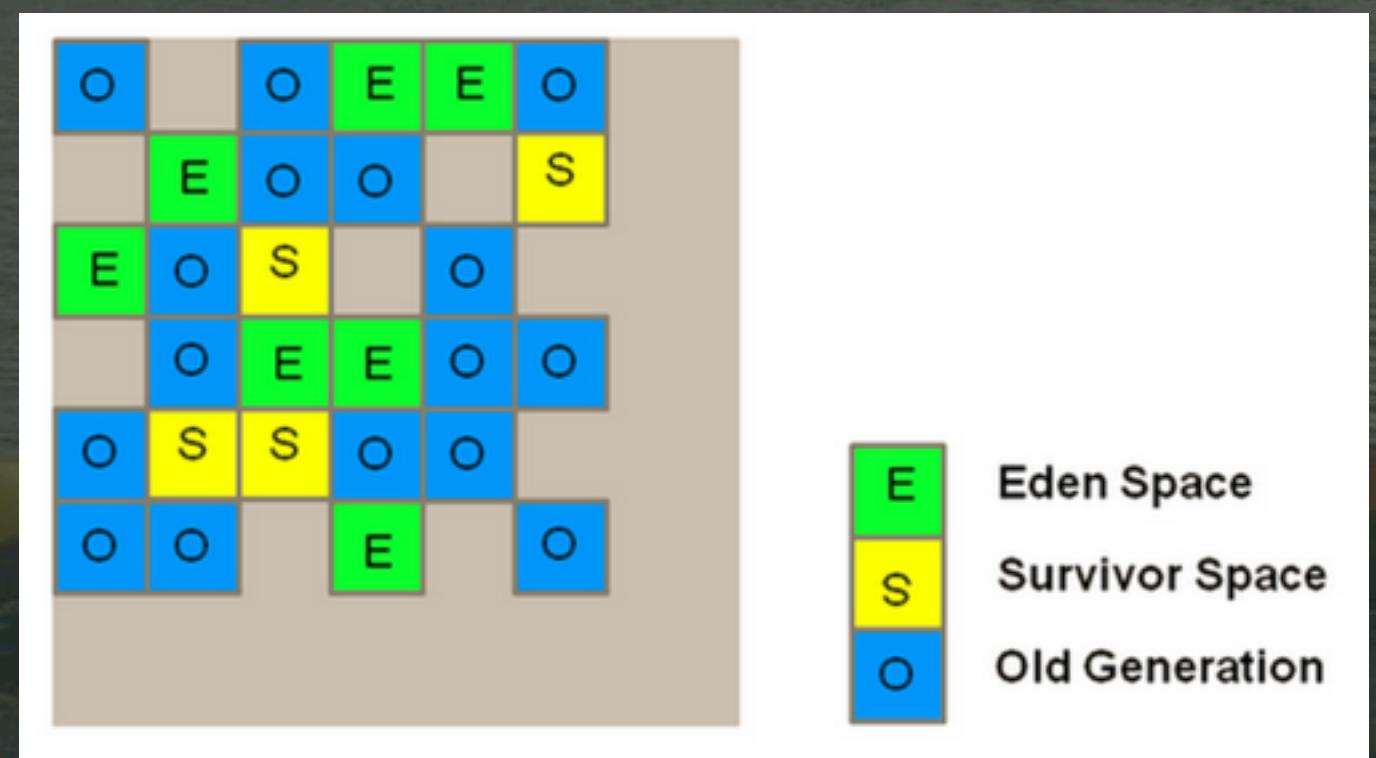
47

Tuesday, September 22, 15

The sheer number of objects can be huge, not just the size of the memory used.

GC Challenges

- Garbage First GC (G1 - JVM 1.6).
 - Balance latency and throughput.
 - More flexible mem. region mgmt.



48

Tuesday, September 22, 15

Image: <http://www.oracle.com/webfolder/technetwork/tutorials/obe/java/G1GettingStarted/index.html>

Much better for Spark it turns out...

GC Challenges

- Best for Spark:
 - -XX:UseG1GC -XX:-ResizePLAB -Xms... -Xmx... -XX:InitiatingHeapOccupancyPercent=... -XX:ConcGCThreads=...

databricks.com/blog/2015/05/28/tuning-java-garbage-collection-for-spark-applications.html

49

Tuesday, September 22, 15

Summarizing a long, detailed blog post (<https://databricks.com/blog/2015/05/28/tuning-java-garbage-collection-for-spark-applications.html>) in one slide!
Their optimal settings for Spark for large data sets, before the “Tungsten” optimizations. The numbers elided (“...”) are scaled together.

GC Challenges

- Spark has two memory “fractions”:
 1. Cache RDDs (largest fraction).
 2. Other Heap usage.
- Configurable by:
`spark.storage.memoryFragment`

50

Tuesday, September 22, 15

Spark splits the heap you give it into two fractions, the largest (by default) is used to cache data that has already been computed, so you don't recompute it unnecessarily. The smaller fraction is for other heap usage, such as running transformations on the RDDs.

Understanding this separation is important for successful production use. Can be tuned depending on the job type.

Object

Overhead

51

Tuesday, September 22, 15

For general-purpose management of trees of objects, Java works well, but not for data orders of magnitude larger, where you tend to have many instances of the same “schema”.

Java Objects?

- “abcd”: 4 bytes for raw UTF8, right?
- 48 bytes for the Java object:
 - 12 byte header.
 - 20 bytes for array overhead.
 - 8 bytes for UTF16 chars.
 - 8 bytes for hash code.

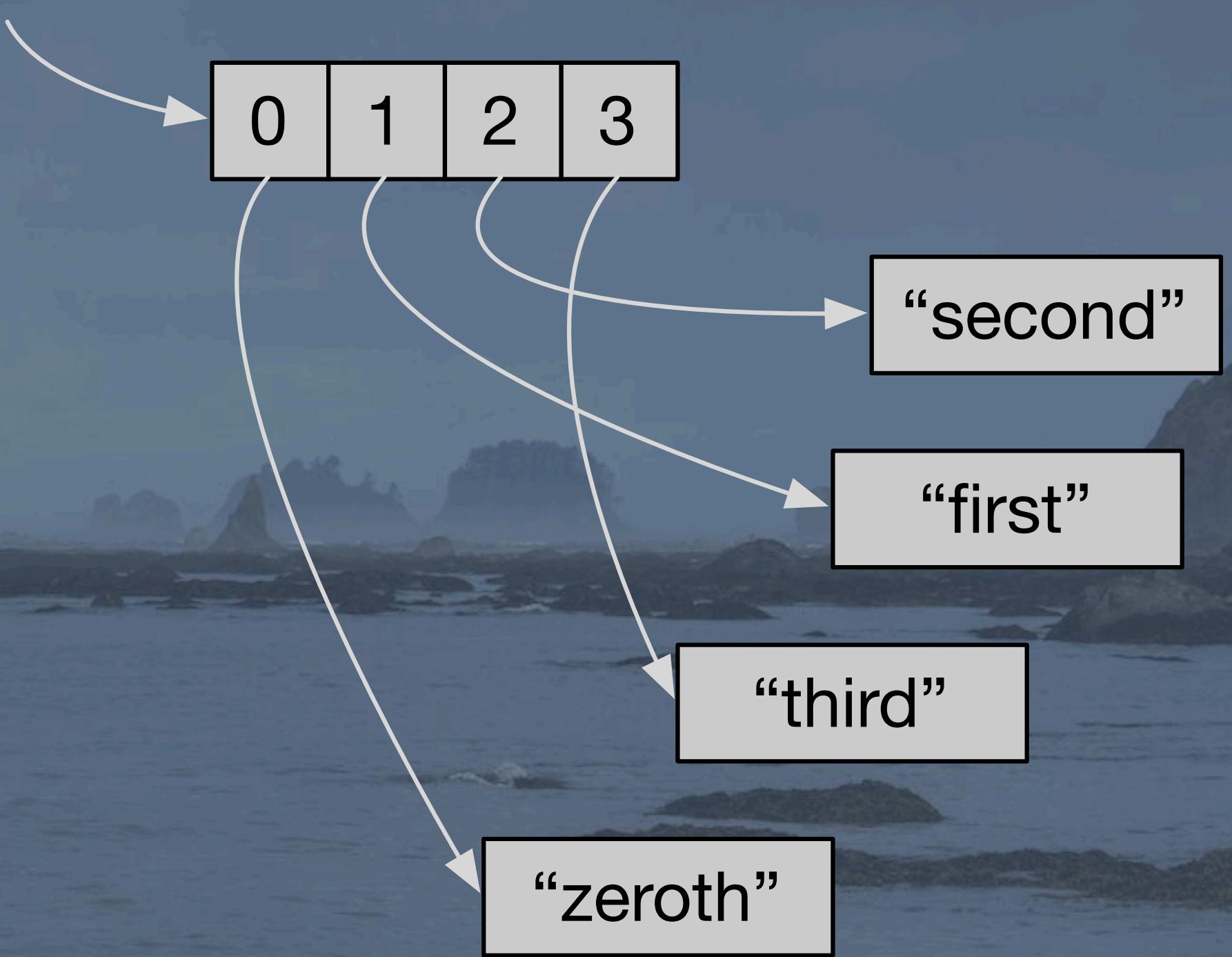
52

Tuesday, September 22, 15

From <http://www.slideshare.net/SparkSummit/deep-dive-into-project-tungsten-josh-rosen>

val myArray: Array[String]

Arrays



Java Objects?

- Case classes, tuples, & Spark's Row type used for "records", but their reference indirections add lots of overhead.
- ... and lead to cache misses.

54

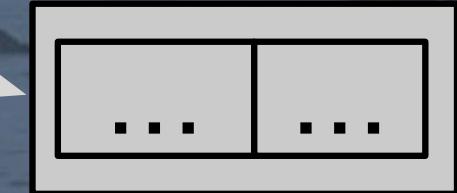
Tuesday, September 22, 15

From <http://www.slideshare.net/SparkSummit/deep-dive-into-project-tungsten-josh-rosen>

val person: Person

| | |
|---------------|----|
| name: String | |
| age: Int | 29 |
| addr: Address | |

“Buck Trends”

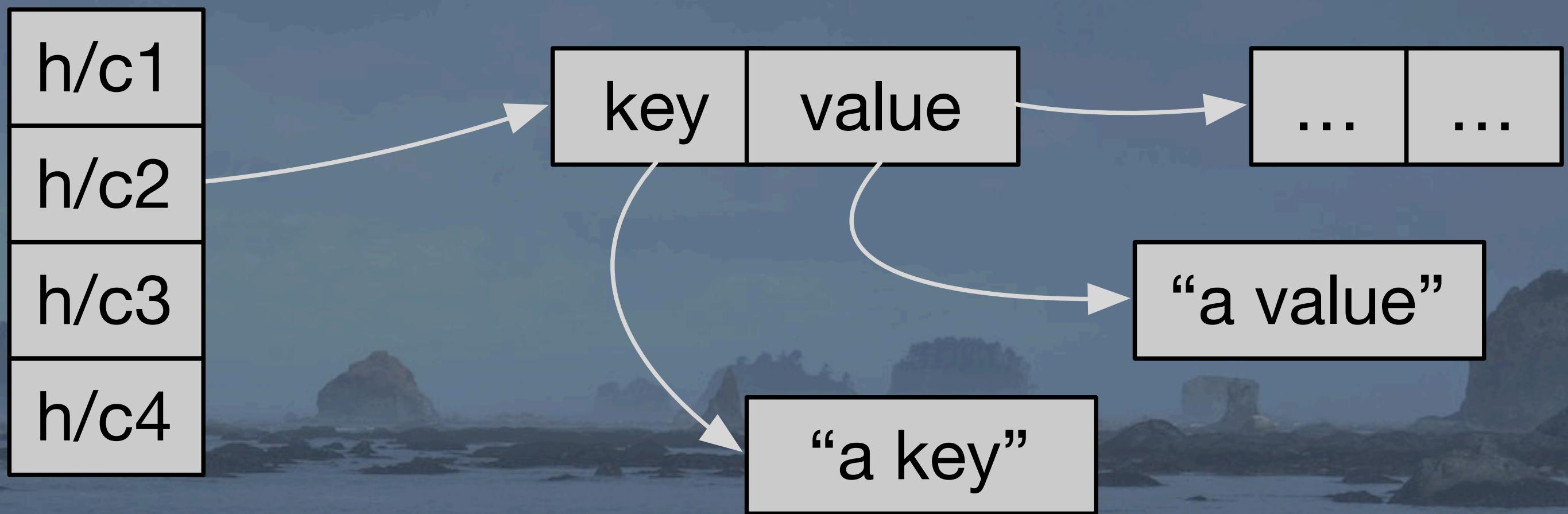


Case Classes

Tuesday, September 22, 15

There's the memory for the array, then the references to other objects for each element around the heap.

Hash Map



Hash Maps

56

Tuesday, September 22, 15

There's the memory for the array, then the references to other objects for each element around the heap.

No Unsigned Types

What's
factorial(-1)?

57

Tuesday, September 22, 15

Unsigned types are very useful for many applications and scenarios. Not all integers need to be signed!
Arrays are limited to $2^{(32-1)}$ elements, rather than $2^{(32)}$!



Richer data libs. in Python & R

Tuesday, September 22, 15

58

Python and R have a longer history as Data Science languages, so they have much richer and more mature libraries for Data Science, e.g., statistics, machine learning, natural language processing, etc.

Java OOP Thinking!

59

Tuesday, September 22, 15

Finally, while many developers have many years of Java experience, Java's historic object-oriented emphasis is largely a disadvantage when thinking about data problems. (However, it's useful for modularizing application code.)

Why Scala?

60

Tuesday, September 22, 15

Okay, all things considered, the JVM is a great platform for Big Data.
Scala has emerged as the de facto “data engineering” language, as opposed to data science, where Python, R, SAS, Matlab, and Mathematica still dominate (although Scala has its passionate advocates here, too.)

Pragmatic OOP + FP

61

Tuesday, September 22, 15

While some people hate the fact that Scala embraces OOP (and I know that you know that I know who you are...), this is pragmatic to make Scala accessible to Java developers and it provides a convenient, familiar modularity tool. However, data science is the killer app for FP, IMHO, so the core problem needs to use FP, which Scala enables nicely.

- Abstractions vs. implementations:
 - Pure functional abstractions?
 - Mutable implementations, when necessary.
 - Performance is critical.

62

Tuesday, September 22, 15

Raw performance is extremely important for competitive data computation systems. Providing ~pure, high-level abstractions with well-defined semantics, enabling flexibility in the implementations to exploit mutability for performance (but hopefully not prematurely applied).

Scala Big Data Sandwich



63

Tuesday, September 22, 15

Image: https://en.wikipedia.org/wiki/Hamburger#/media/File:NCI_Visuals_Food_Hamburger.jpg

Scala Big Data Sandwich

scopes

Objects as Modules



Functional APIs

Optimized (Mutable?) Code

64

Tuesday, September 22, 15

Image: https://en.wikipedia.org/wiki/Hamburger#/media/File:NCI_Visuals_Food_Hamburger.jpg



REPL & Notebooks

65

Tuesday, September 22, 15

Interactive exploration is a powerful tool for exploring data, algorithms, and analysis approaches. Data scientists are familiar with notebook metaphors, such as iPython. Now Spark Notebook brings this way of using the REPL to Scala and Spark, providing a richer experience than the REPL or IDE alone.



Cell Toolbar: None

```
wiugers.GraphicalUI  
    nodes.zipWithIndex.map{case (v,i) => Node(i, v, color=color) } :::  
    links.zipWithIndex.map { case ((i,j,v), k) => Edge(k, (i,j), v, color=color)},  
    maxPoints=nodes.size+links.size,  
    linkStrength=0.1  
)  
  
colorAtoms: List[Char] = List(0, 1, 2, 3, 4, 5, 6, 7, 8, 9, :, ;, <, =, >, ?, @, A, B, C, D, E,  
F, A, B, C, D, E, F)  
color: String  
res4: notebook.front.widgets.GraphChart[List[Product with Serializable with notebook.front.widgets.magic.Graph[Int]]] = <GraphChart widget>
```

Out[5]: 627 items



Collections API



67

Tuesday, September 22, 15

Collections API a natural foundation for data flows.

- * Powerful, yet concise.
- * Abstracts over implementation – can be parallelized, distributed.
- * Generalizes to streaming or a fixed collection.

Inspired Spark's API

But less complex

68

Tuesday, September 22, 15

Spark's API looks very similar to the "conceptual" abstractions of Scala's collections API, i.e., what you see in the simplified method signatures shown in the Scaladoc, as opposed to the actual signatures with the extra type parameters, CanBuildFrom implicit, etc.

Arguably, Spark's API offers a cleaner separation between interface and implementation.

```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._

val sparkContext = new SparkContext(master, "Inv. Index")
sparkContext.textFile("/path/to/input")
.map { line =>
  val array = line.split(",", 2)
  (array(0), array(1))
}.flatMap {
  case (id, contents) => toWords(contents).map(w => ((w,id),1))
}.reduceByKey {
  (count1, count2) => count1 + count2
}.map {
  case ((word, path), n) => (word, (path, n)))
}.groupByKey
.map {
  case (word, list) => (word, sortByCount(list))
}.saveAsTextFile("/path/to/output")
```

Recall...

69

Tuesday, September 22, 15

Recall the inverted index Spark script.

```
val array = line.split("", 2)
```

```
(array(0), array(1))
```

```
.flatMap {
```

```
  case (id, contents) => toWords(contents)
```

```
.reduceByKey {
```

```
  (count1, count2) => count1 + count2
```

```
.map {
```

```
  case ((word, path), n) => (word, (path, n))
```

```
.groupByKey
```

```
.map {
```

```
  case (word, list) => (word, sortByCount(1))
```

```
.saveAsTextFile("/path/to/output")
```

70

Tuesday, September 22, 15

Beautiful “combinators”.

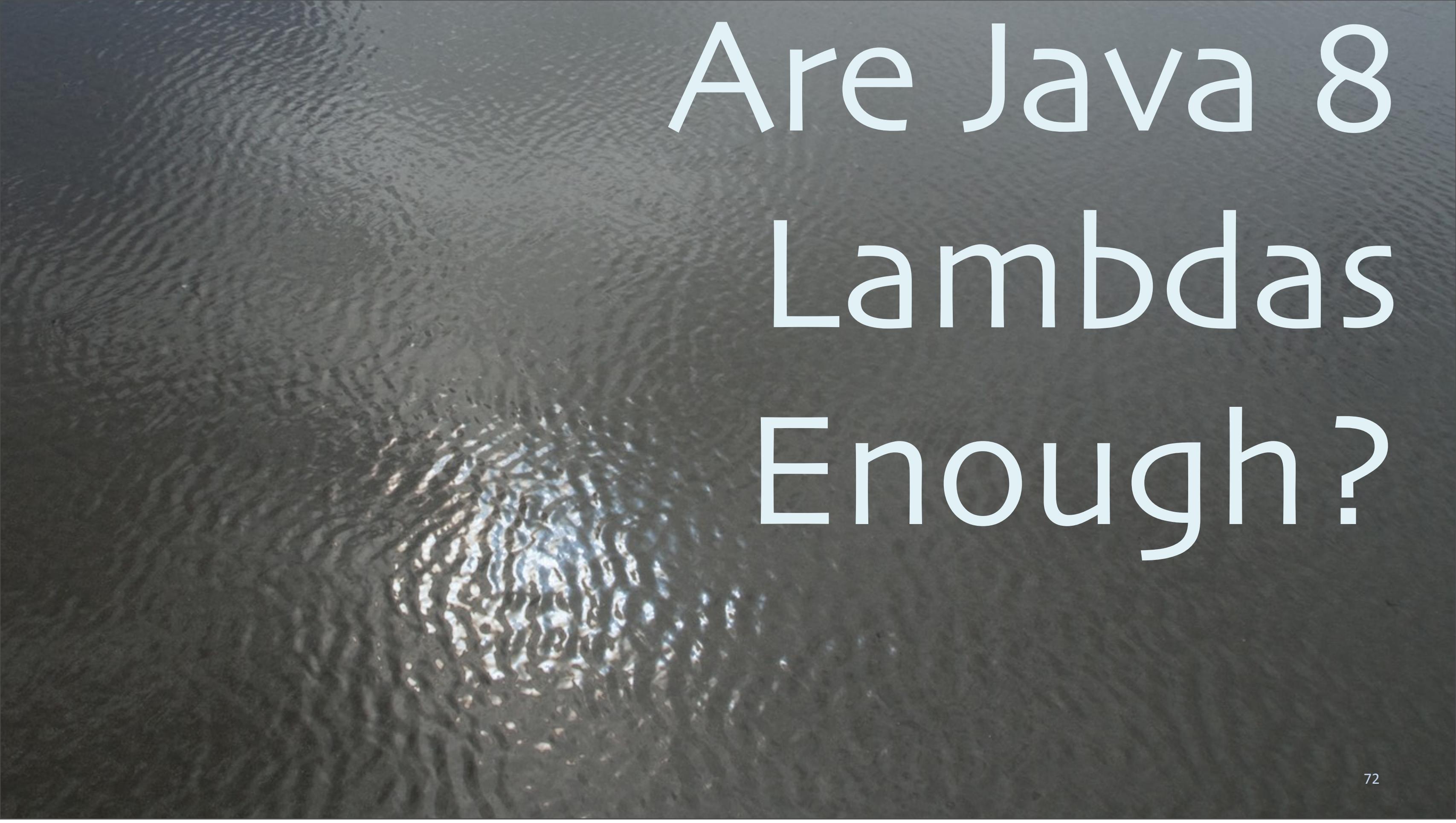
Claim

Spark does fine without
the Reader Monad.

71

Tuesday, September 22, 15

While the purists amongst us might prefer that all state mutation be explicitly encapsulated in monads, in fact the “informal” approach of Spark makes it approachable by a wide class of developers and the boundaries between lazy “transformations” and side-affecting “actions”, while not explicitly obvious, become so once you gain some experience with the API.



Are Java 8 Lambdas Enough?

72

Tuesday, September 22, 15

Spark is untenable in Java 7. The verbosity of anonymous inner classes started a flight to Scala. Then Java 8 added lambdas. Do they make Java good enough? No...

Tuples

73

Tuesday, September 22, 15

First, Tuples are very useful for writing concise code, as we saw in the previous example.

A photograph of a person walking along a wet, sandy beach. The water is shallow and reflects the surrounding environment. The person is wearing a dark jacket and shorts, and is carrying a backpack. The beach is bordered by a dense forest of evergreen trees. In the foreground, large white letters spell out "Pattern Matching".

Pattern Matching

74

Tuesday, September 22, 15

Pattern matching makes it easy to extract fields in records, match strings with regular expressions, test types, etc.



Type Inference

75

Tuesday, September 22, 15

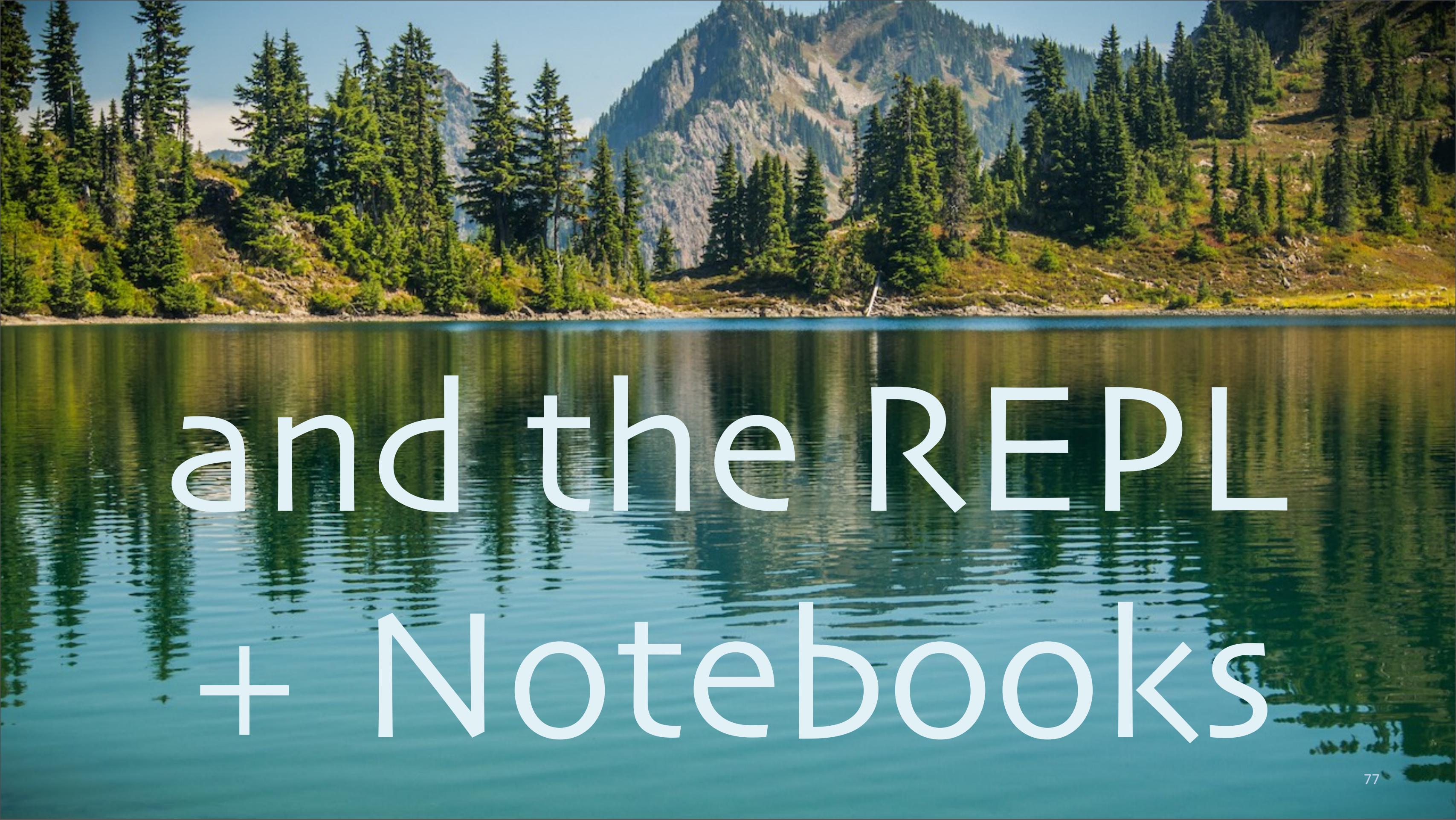
Type inference removes the tedium of explicit typing, as in Java, while providing type safety and useful feedback in interactive sessions.

DSLs

76

Tuesday, September 22, 15

Expressive DSLs are relatively easy to create in Scala. The collections API itself is a DSL of sorts for data flows. Another simple example is the ability to define your own mathematical types, like vectors, matrices, complex numbers, etc. and operators for them, like +, -, *, and / that you use as if the types were built in. We'll another Spark DSL shortly.



and the REPL
+ Notebooks

77

Tuesday, September 22, 15

Finally, the REPL and Notebook UIs for it are essential for exploring data and experimenting with algorithms.

Back to Spark



78

Tuesday, September 22, 15

Let's revisit Spark.

Performance

DataFrame API in SparkSQL:

- Records with known schemas.
- Relational-like operations.
- Inspired by Python/R DataFrames.

79

Tuesday, September 22, 15

This is a major step forward. Previously for Hadoop, Data Scientists often developed models in Python or R, then an engineering team ported them to Java MapReduce. Previously with Spark, you got good performance from Python code, but about 1/2 the efficiency of corresponding Scala code. Now, the performance is the same.

```
def join(right: DataFrame, ...): DataFrame = {  
  def groupBy(cols: Column*): GroupedData = {  
    def orderBy(sortExprs: Column*): DataFrame = {  
      def select(cols: Column*): DataFrame = {  
        def where(condition: Column): DataFrame = {  
          def limit(n: Int): DataFrame = {  
            def unionAll(other: DataFrame): DataFrame = {  
              def intersect(other: DataFrame): DataFrame = {  
                def sample(withReplacement: Boolean, fraction, seed): DataFrame = {  
                  def drop(col: Column): DataFrame = {  
                    def map[R: ClassTag](f: Row => R): RDD[R] = {  
                      def flatMap[R: ClassTag](f: Row => Traversable[R]): RDD[R] = {  
                        def foreach(f: Row => Unit): Unit = {  
                          def take(n: Int): Array[Row] = {  
                            def count(): Long = {  
                              def distinct(): DataFrame = {  
                                def agg(exprs: Map[String, String]): DataFrame = {
```

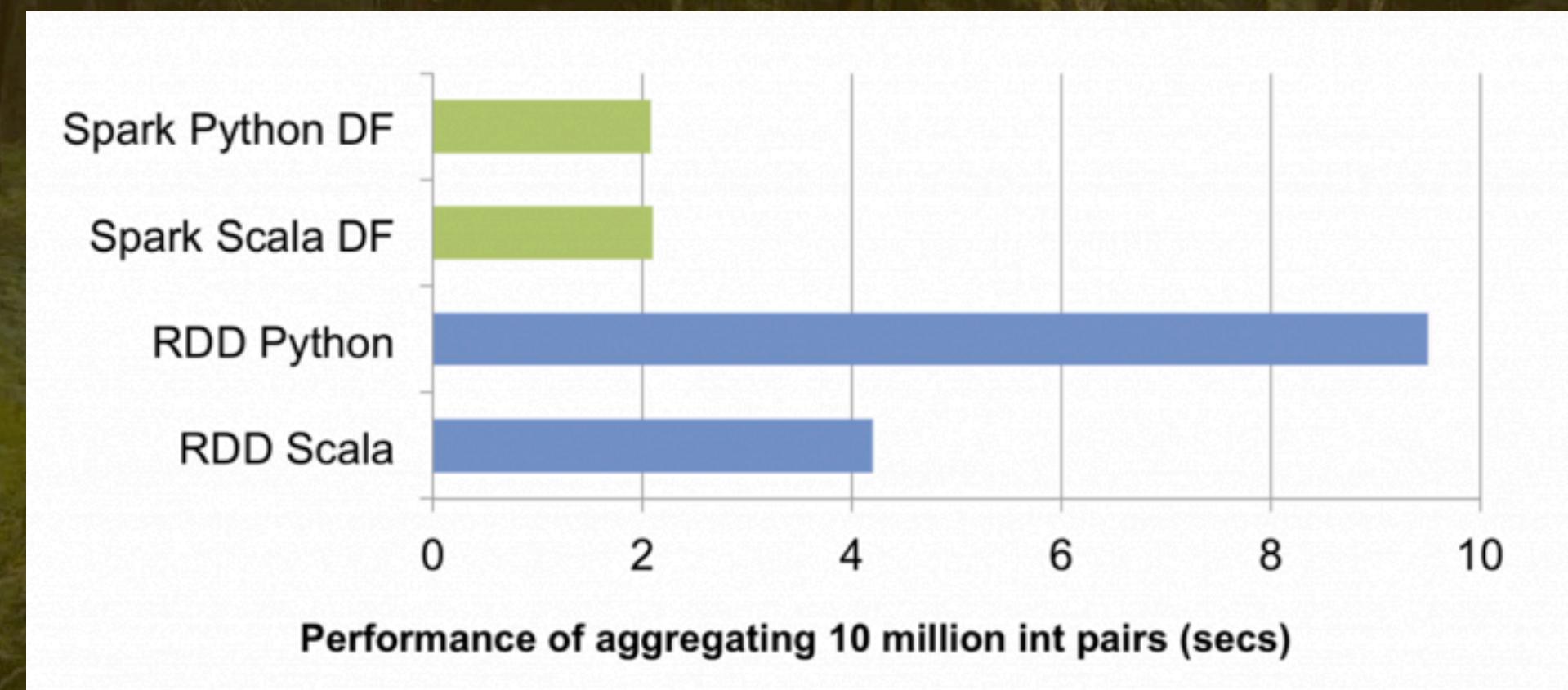
80

Tuesday, September 22, 15

So, the DataFrame exposes a more limited set of abstractions than the more general RDD API. This narrower interface enables broader (more aggressive) optimizations and other implementation choices underneath.

Performance

The DataFrame API has the same performance for all languages:
Scala, Java,
Python, R,
and SQL!



Tuesday, September 22, 15

This is a major step forward. Previously for Hadoop, Data Scientists often developed models in Python or R, then an engineering team ported them to Java MapReduce. Previously with Spark, you got good performance from Python code, but about 1/2 the efficiency of corresponding Scala code. Now, the performance is the same.

Graph from: <https://databricks.com/blog/2015/04/24/recent-performance-improvements-in-apache-spark-sql-python-dataframes-and-more.html>



Project Tungsten

82

Tuesday, September 22, 15

Project Tungsten is a multi-release initiative to improve Spark performance, focused mostly on the DataFrame implementation. References:

<http://www.slideshare.net/databricks/2015-0616-spark-summit>

<http://www.slideshare.net/SparkSummit/deep-dive-into-project-tungsten-josh-rosen>

<https://databricks.com/blog/2015/04/28/project-tungsten-bringing-spark-closer-to-bare-metal.html>



Initiative to greatly improve
performance of DataFrames.

Project Tungsten

83

Tuesday, September 22, 15

Project Tungsten is a multi-release initiative to improve Spark performance, focused mostly on the DataFrame implementation. References:

<http://www.slideshare.net/databricks/2015-0616-spark-summit>

<http://www.slideshare.net/SparkSummit/deep-dive-into-project-tungsten-josh-rosen>

Observations

Spark jobs are CPU bound:

- Improve network I/O? ~2% better.
- Improve disk I/O? ~20% better.

So, we need to look elsewhere...

84

Tuesday, September 22, 15

MapReduce jobs tend to be CPU bound, but research by Kay Ousterhout (http://www.eecs.berkeley.edu/~keo/talks/2015_06_15_SparkSummit_MakingSense.pdf) and other observers indicate that for Spark, optimizing I/O only improve performance ~5% for network improvements and ~20% for disk I/O. So, Spark jobs tend to be CPU bound.

Why?

- HW: 10Gbs networks, SSDs.
- Smart pruning of unneeded data.
- Effective use of caching in Spark.
- Better data formats, like Parquet.
- Serialization, compression, and hashing are CPU intensive.

85

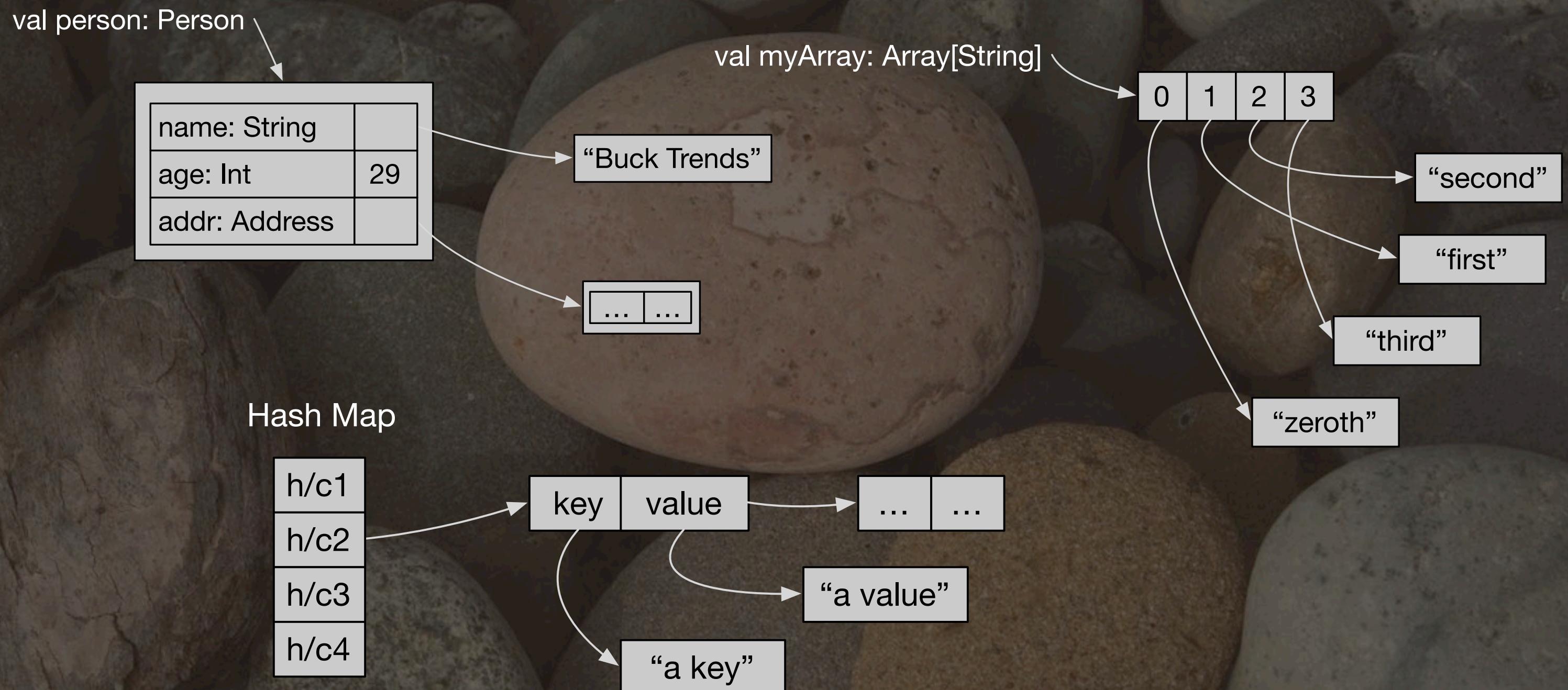
Tuesday, September 22, 15

First, the general-purpose object structure and memory management in the JVM is suboptimal for data-centric applications with huge datasets. So, Spark

Three-pronged Initiative

- Java object model overhead suboptimal for data.
- GC expensive.
- Cache awareness important.

Java Objects?



87

Tuesday, September 22, 15

From <http://www.slideshare.net/SparkSummit/deep-dive-into-project-tungsten-josh-rosen>

We discussed this overhead earlier. Here's a recap ;)

GC Challenges

- Recall:
- Typical heaps 10s-100s of GB.
 - Uncommon in generic services.
 - Too many cached RDDs leads to huge old generation sections.
 - Puts extreme pressure on GC.

88

Tuesday, September 22, 15

See <https://databricks.com/blog/2015/05/28/tuning-java-garbage-collection-for-spark-applications.html> for a detailed overview of GC challenges and tuning for Spark.

Expensive Logic

- Overhead evaluating expressions, like in SQL queries:
 - Virtual function calls.
 - Object creation from boxing.
 - Branching (if statements, etc.)

```
sqlContext.sql("""SELECT a + b FROM table""")
```

89



Tungsten

90

Tuesday, September 22, 15

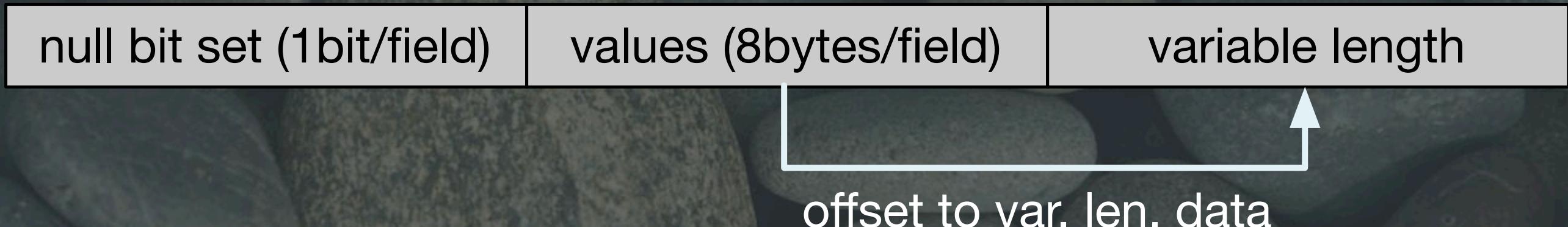
Project Tungsten tries to improve in all these areas using three lines of attack...

Tungsten

- Custom memory management.
 - Rather than Java objects.
- Cache-aware algorithms and data structures.
- Code gen for optimal performance.

sun.misc.Unsafe

- Off and On Heap.
- Compact Row:



- Hash code and equals on raw bytes.

92

Tuesday, September 22, 15

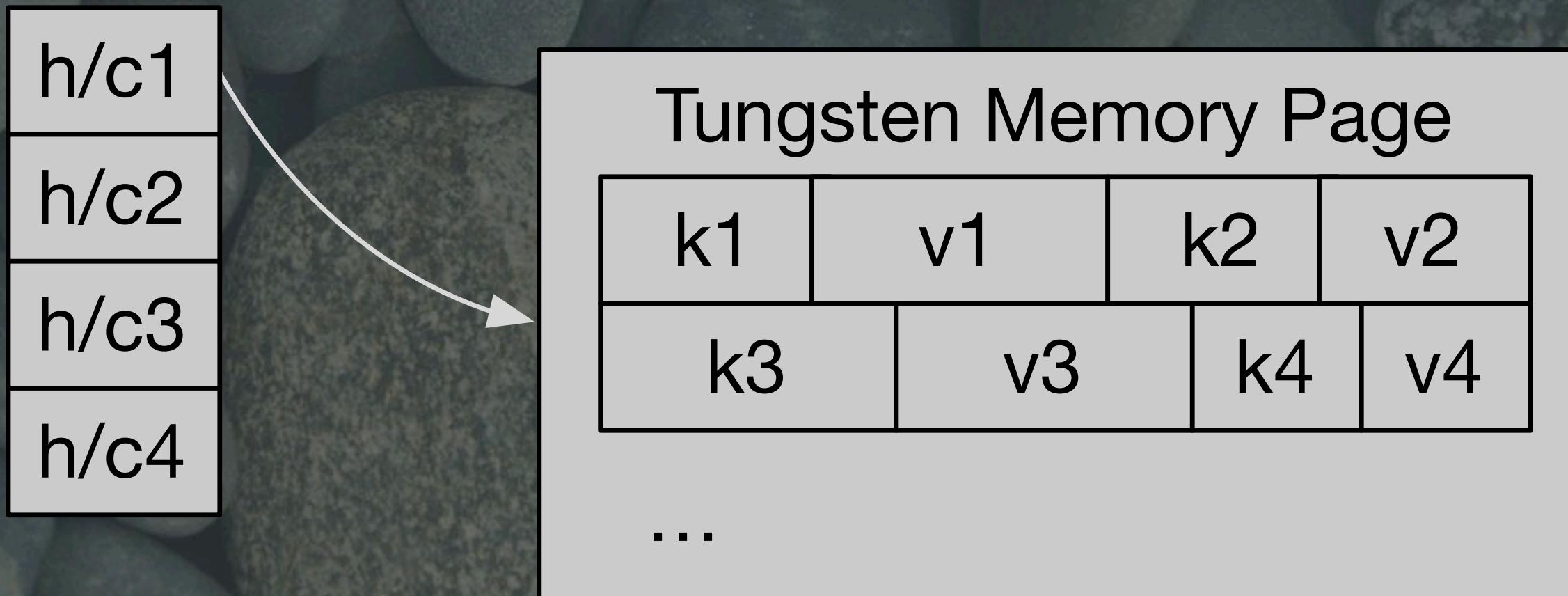
Unsafe allows you to manipulate memory directory, like in C/C++.

The new Compact Row format (for each record) uses a bit vector to mark values that are null, then packs the values together, each in 8 bytes. If a value is a fixed-size item and fits in 8 bytes, it's inlined.

Otherwise, the 8 bytes holds a reference to a location in variable-length segment (e.g., strings).

Rows are 8-byte aligned. Hashing and equality can be done on the raw bytes.

- BytesToBytesMap



Tungsten

- These data structures:
 - Improve cache hits.
 - Reduce GC drastically.

Tungsten

- Byte code generation for expressions to eliminate:
 - Boxing/Unboxing.
 - Virtual method call overhead.
 - Branching

95

Tuesday, September 22, 15

Their test results indicate ~3x improvement in performance, very close to hand-written, optimized code!



Reactive Streams

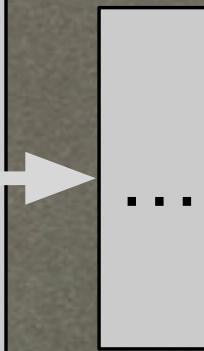
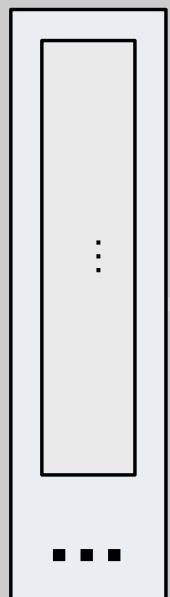
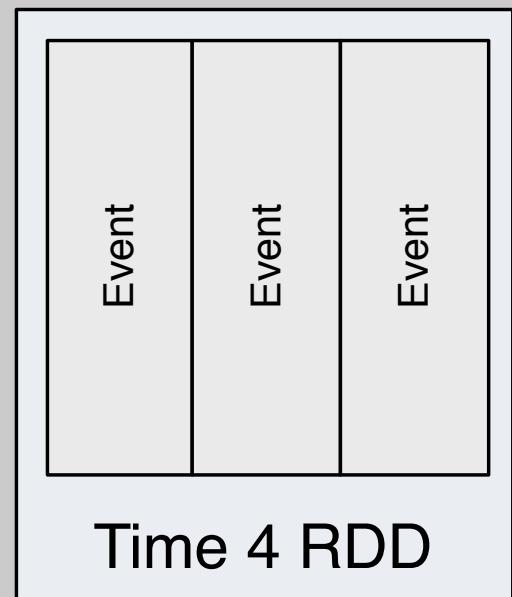
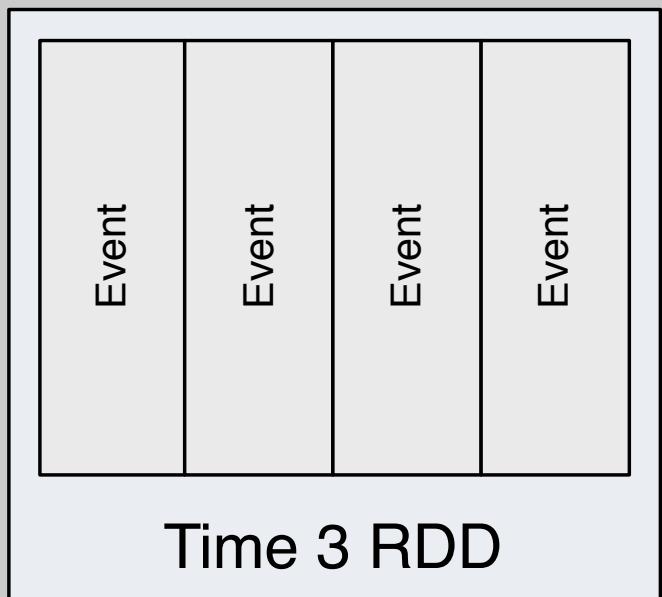
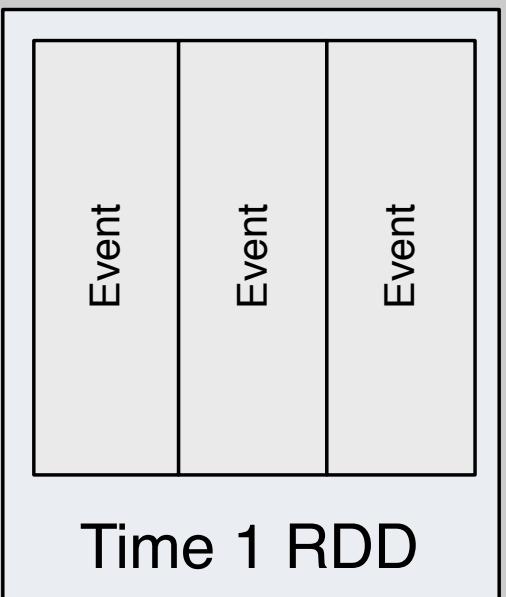
96

Tuesday, September 22, 15

Lastly, let's talk about streams again.

Recall:

DStream (discretized stream)



Window of 3 RDD Batches #1

Window of 3 RDD Batches #2

A photograph of a pile of dark green and yellowish seaweed resting on a light-colored, textured sand surface. The sand has distinct wavy patterns. The seaweed is piled in the upper right quadrant of the frame.

Robustness?

98

Tuesday, September 22, 15

Streams can run a very long time. Lots of problems can arise, such as consumers getting backed up or producers producing more and more stuff...

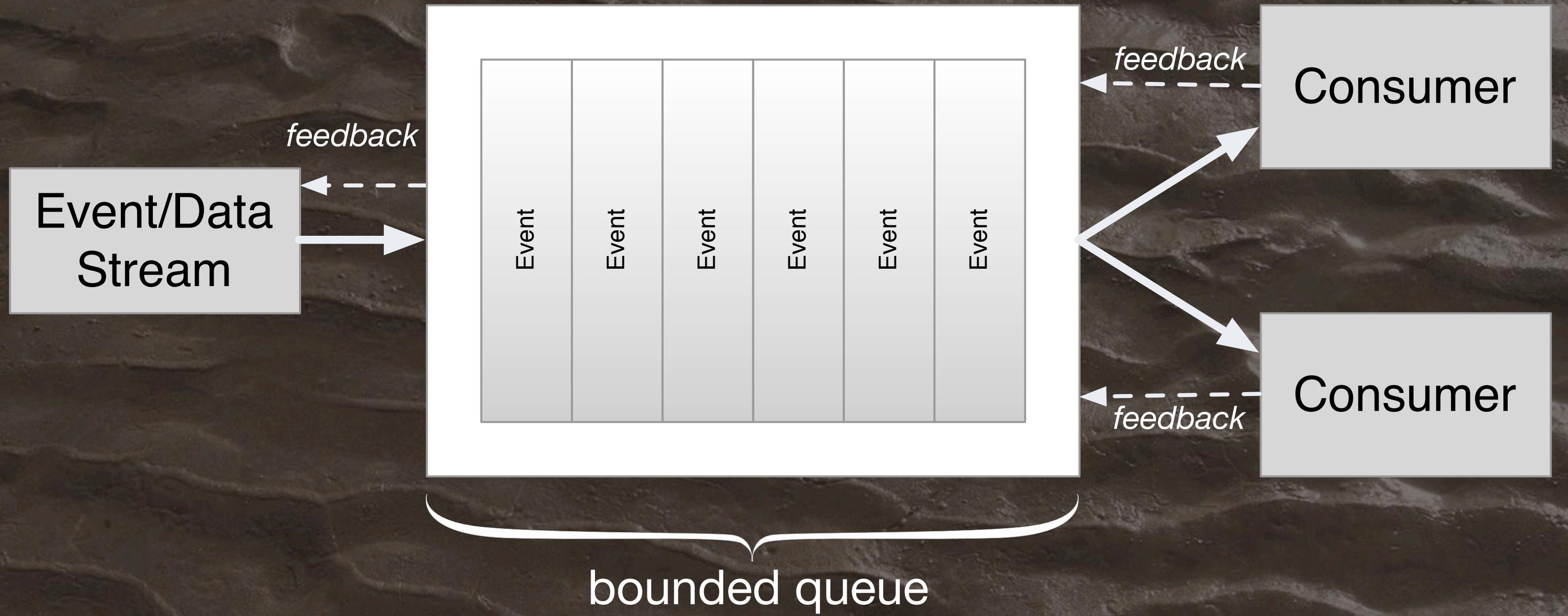


Backpressure

99

Tuesday, September 22, 15

Necessary for flow control. In Spark 1.5, Typesafe contributed a dynamic backpressure mechanism. (Blog post coming...) Going forward, we hope to contribute a Reactive Streams-compliant consumer.



100

Tuesday, September 22, 15

Necessary for flow control. In Spark 1.5, Typesafe contributed a dynamic backpressure mechanism. (Blog post coming...) Going forward, we hope to contribute a Reactive Streams-compliant consumer.

Final Thoughts



101

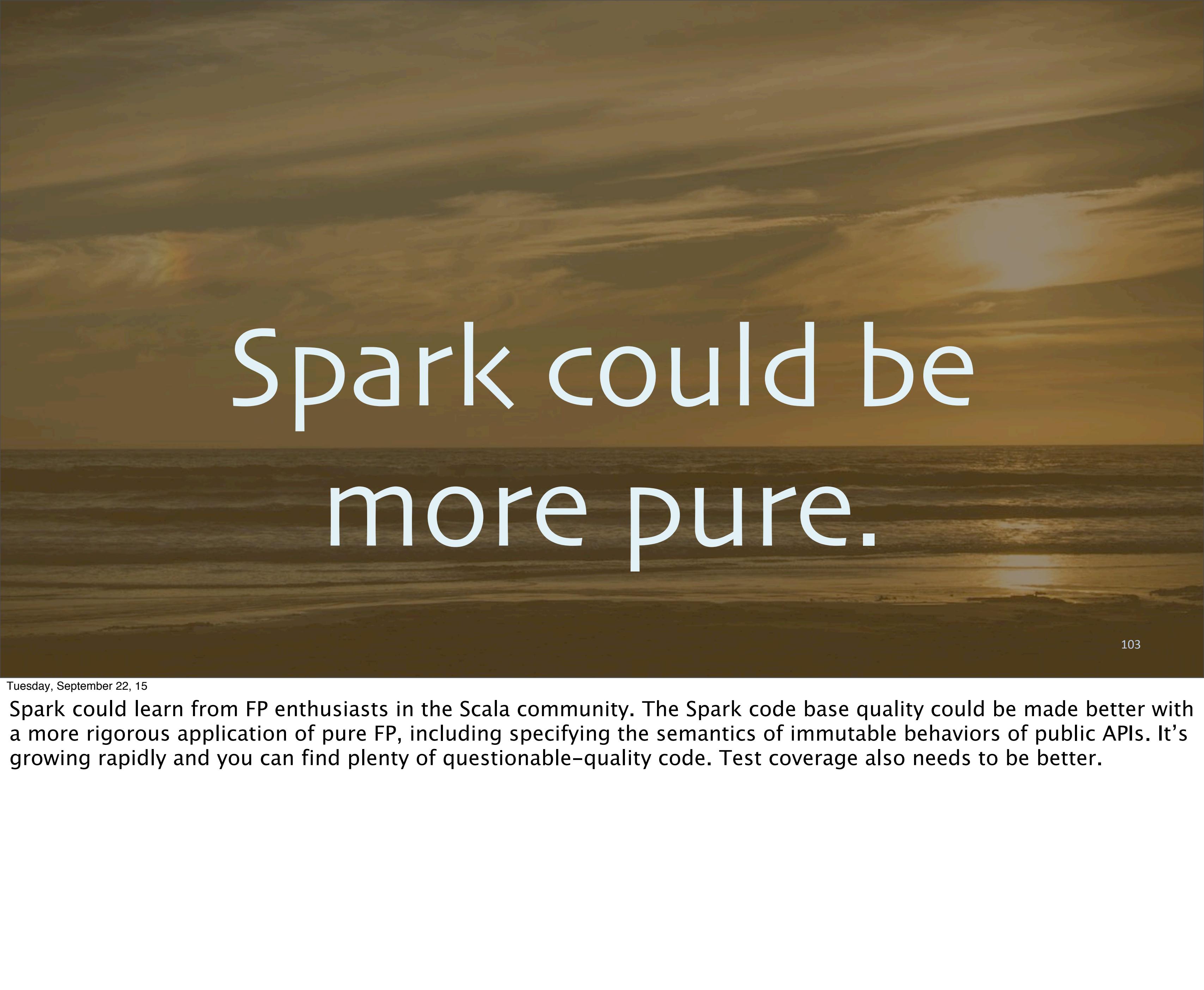
Tuesday, September 22, 15



Spark Is Driving Scala Adoption

102

Tuesday, September 22, 15

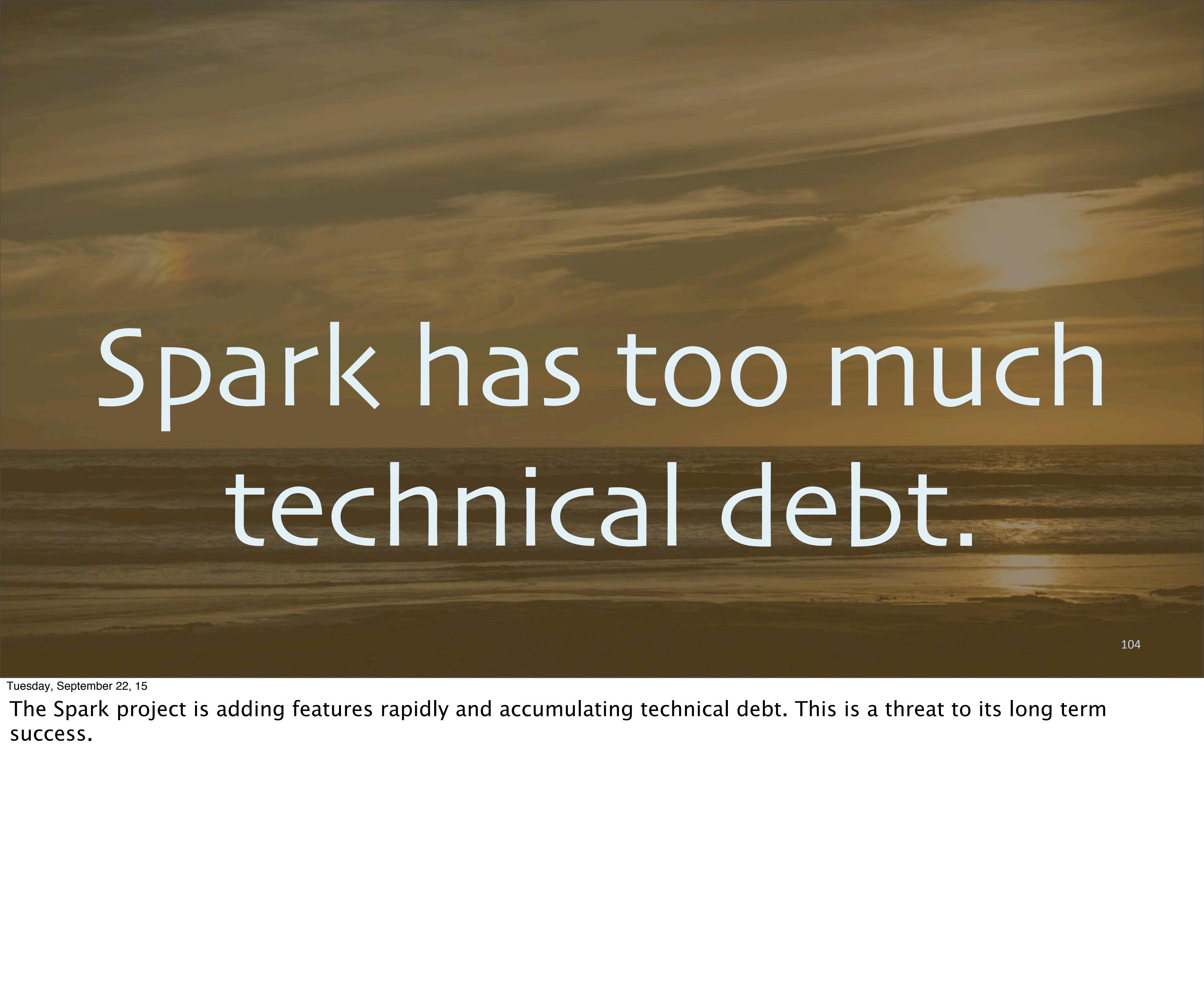


Spark could be
more pure.

103

Tuesday, September 22, 15

Spark could learn from FP enthusiasts in the Scala community. The Spark code base quality could be made better with a more rigorous application of pure FP, including specifying the semantics of immutable behaviors of public APIs. It's growing rapidly and you can find plenty of questionable-quality code. Test coverage also needs to be better.

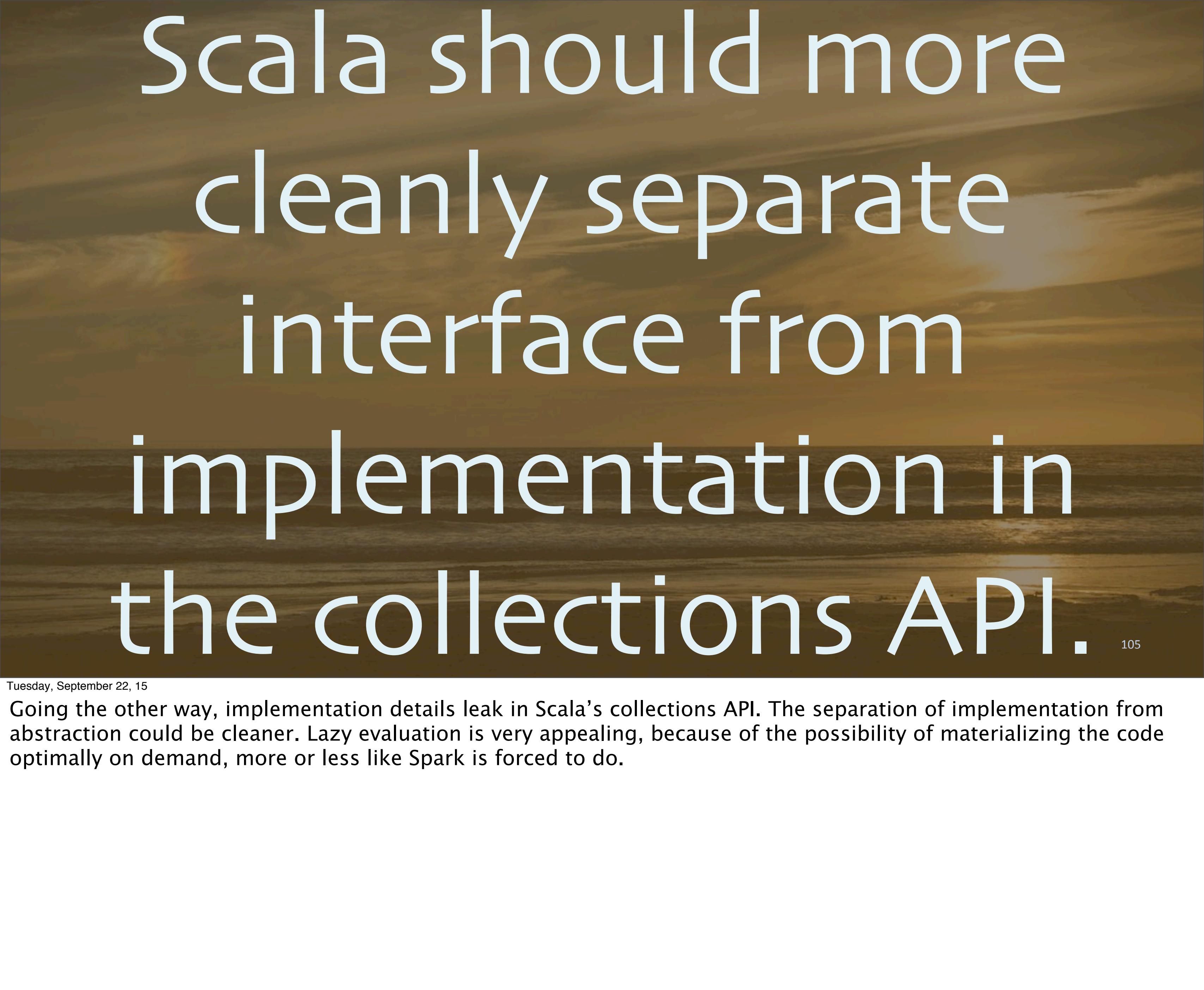


Spark has too much technical debt.

104

Tuesday, September 22, 15

The Spark project is adding features rapidly and accumulating technical debt. This is a threat to its long term success.

A photograph of a sunset or sunrise over a body of water. The sky is filled with warm, orange, and yellow clouds, with darker blues at the bottom. The water reflects these colors. The overall mood is peaceful and contemplative.

Scala should more
cleanly separate
interface from
implementation in
the collections API.

105

Tuesday, September 22, 15

Going the other way, implementation details leak in Scala's collections API. The separation of implementation from abstraction could be cleaner. Lazy evaluation is very appealing, because of the possibility of materializing the code optimally on demand, more or less like Spark is forced to do.

Scala could adopt
Spark's extensions to
the collections API.

spark-summit.org/eu-2015/events/spark-the-ultimate-scala-collections/

106

Tuesday, September 22, 15

Going the other way, Martin has even said publicly that Scala's collections API could adopt extensions in Spark.

```
[22.11.20] dearmwmpic@dearmwmpic-OptiPlex-5090:~/Documents/training/sparkWorkshop/spark-workshop-exercises
✓ grep 'def.*ByKey' ~/projects/spark/spark-git/core/src/main/scala/org/apache/spark/rdd/PairRDD
def combineByKey[C](createCombiner: V => C,
def combineByKey[C](createCombiner: V => C,
def aggregateByKey[U: ClassTag](zeroValue: U, partitioner: Partitioner)(seqOp: (U, V) => U,
def aggregateByKey[U: ClassTag](zeroValue: U, numPartitions: Int)(seqOp: (U, V) => U,
def aggregateByKey[U: ClassTag](zeroValue: U)(seqOp: (U, V) => U,
def foldByKey(
def foldByKey(zeroValue: V, numPartitions: Int)(func: (V, V) => V): RDD[(K, V)] = self.withScope
def foldByKey(zeroValue: V)(func: (V, V) => V): RDD[(K, V)] = self.withScope {
def sampleByKey(withReplacement: Boolean,
def sampleByKeyExact(
def reduceByKey(partitioner: Partitioner, func: (V, V) => V): RDD[(K, V)] = self.withScope {
def reduceByKey(func: (V, V) => V, numPartitions: Int): RDD[(K, V)] = self.withScope {
def reduceByKey(func: (V, V) => V): RDD[(K, V)] = self.withScope {
def reduceByKeyLocally(func: (V, V) => V): Map[K, V] = self.withScope {
def reduceByKeyToDriver(func: (V, V) => V): Map[K, V] = self.withScope {
def countByKey(): Map[K, Long] = self.withScope {
def countByKeyApprox(timeout: Long, confidence: Double = 0.95)
def countApproxDistinctByKey(
def countApproxDistinctByKey(
def countApproxDistinctByKey(
def countApproxDistinctByKey(relativeSD: Double = 0.05): RDD[(K, Long)] = self.withScope {
def groupByKey(partitioner: Partitioner): RDD[(K, Iterable[V])] = self.withScope {
def groupByKey(numPartitions: Int): RDD[(K, Iterable[V])] = self.withScope {
def combineByKey[C](createCombiner: V => C, mergeValue: (C, V) => C, mergeCombiners: (C, C) =>
def groupByKey(): RDD[(K, Iterable[V])] = self.withScope {
def subtractByKey[W: ClassTag](other: RDD[(K, W)]): RDD[(K, V)] = self.withScope {
def subtractByKey[W: ClassTag](
def subtractByKey[W: ClassTag](other: RDD[(K, W)], p: Partitioner): RDD[(K, V)] = self.withScope
```

Tuesday, September 22, 15

Here's a subset of possibilities, "by key" functions that assume a schema of (key,value).

Scala should adopt Tungsten optimizations?

108

Tuesday, September 22, 15

Going the other way, Scala could also adopt some Tungsten optimizations, such as the optimized Record and HashMap types.

Spark + Mesos

A photograph of a sunset over the ocean. The sky is filled with warm orange and yellow hues, transitioning into a darker blue at the top. The sun is a bright orange circle near the horizon. In the foreground, there are several small, dark rock formations or low-lying islands partially submerged in the water. The ocean waves are visible in the background, creating a sense of depth and motion.

typesafe.com/reactive-big-data

Tuesday, September 22, 15

We at Typesafe think the future of Spark is bright. We also think the next generation cluster management framework, Mesos, is the future of Big Data platforms, while also supporting cluster management for all your applications, databases, web services, etc., etc. That's why Typesafe is investing in this infrastructure and now offering commercial support for it.

typesafe.com/reactive-big-data
polyglotprogramming.com/talks
dean.wampler@typesafe.com



110

Tuesday, September 22, 15