# Pepperdata does performance (for Big Data)

| 18 Thousand Production | 50 Million Jobs/Year | 200 Trillion |
|---|---|---|

pepperdata.

# Today's talk will cover…

- Why debugging performance problems is hard

- Data elements needed for a complete view of application performance from separate tools

- Bringing these elements together in a single tool

# 2 reasons why debugging performance problems is hard

# Reason # 1

Same external symptoms, but many possible causes

- code

- data

- configuration

- cluster weather

# Reason #2

Existing tools provide limited visibility

- Spark Web UI is the most popular
    - Execution plan with some aggregate performance data

- Ganglia, Ambari, CM etc
    - Time series data about cluster, not specific to Spark apps

- Code execution not connected to resource consumption

- Unhealthy hardware or load from other apps unaccounted

# 3 data elements form a complete picture of Spark application performance

1. Code execution plan

   – Indicates which block of code is being executed

2. Time series view

   – Visual of resource consumption of application

   – Outliers in resource use very easy to detect

3. Cluster weather

   – A view of all applications that runs on the cluster

   – A view of the health of all the nodes in the cluster

First half of the solution

## Spark Web UI

# Logical code execution plan from Spark: Jobs / Stages / DAG



Spark 1.6.0    Jobs    **Stages**    Storage    Environment    Executors      **ScalaPageRank** application UI

## Stages for All Jobs

**Completed Stages:** 6

### Completed Stages (6)

| Stage Id | Description | | Submitted | Duration | Tasks: Succeeded/Total | Input | Output | Shuffle Read | Shuffle Write |
|---|---|---|---|---|---|---|---|---|---|
| 5 | saveAsHadoopFile at IOCommon.scala:63 | +details | 2017/05/29 12:36:53 | 23 s | 8/8 | | 176.4 MB | 547.3 MB | |
| 4 | flatMap at SparkPageRank.scala:62 | +details | 2017/05/29 12:29:45 | 7.1 min | 8/8 | 10.4 GB | | 879.6 MB | 547.3 MB |
| 3 | flatMap at SparkPageRank.scala:62 | +details | 2017/05/29 12:25:32 | 4.2 min | 8/8 | 10.4 GB | | 829.8 MB | 547.3 MB |
| 2 | flatMap at SparkPageRank.scala:62 | +details | 2017/05/29 12:18:59 | 6.5 min | 8/8 | 10.4 GB | | 2.9 GB | 497.5 MB |
| 1 | distinct at SparkPageRank.scala:58 | +details | 2017/05/29 12:13:00 | 6.0 min | 24/24 | | | 2.7 GB | 2.6 GB |
| 0 | distinct at SparkPageRank.scala:58 | +details | 2017/05/29 12:06:55 | 6.1 min | 24/24 | 2.8 GB | | | 2.7 GB |

# Physical execution plan from Spark: Executors / Tasks

Second half of solution

# Time series view

# Time series view of resource consumption for the App
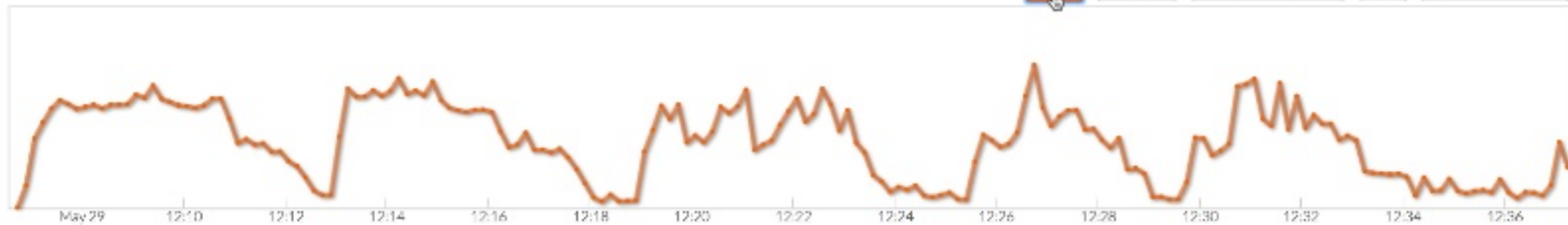
Best of both worlds

# Bring them together

# Bringing it all together: CPU across Stages

## ScalaPageRank

App Id: 1496028096235_0157

Cluster Load View

CPU | Memory | Non Heap Memory | GC | HDFS Bytes Read



| Stage | Job | File Name | Stage Start Time | Stage Duration ((hh:)mm:ss) | Average Executor CPU (%) | Stage Peak Heap Memory(MB) | Stage Peak Non Heap Memory(MB) | Total GC Time (seconds) | Parallel Stages |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | distinct at SparkPageRank.scala:58 | 2017/05/29-12:06 | 06:05 | 41.5 | 23,327 | 497 | 5.82 | none |
| 1 | 0 | distinct at SparkPageRank.scala:58 | 2017/05/29-12:13 | 06:58 | 37.8 | 20,106 | 507 | 6.6 | none |
| 2 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:18 | 06:32 | 32.1 | 24,256 | 514 | 22.2 | none |
| 3 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:25 | 04:12 | 31.6 | 26,676 | 517 | 15.5 | none |
| 4 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:29 | 07:08 | 26 | 27,414 | 518 | 31.3 | none |
| 5 | 0 | saveAsHadoopFile at IOCommon.scala:63 | 2017/05/29-12:36 | 00:22 | 26.7 | 23,639 | 524 | 0 | none |

SPARK SUMMIT 2017

# Memory across all Stages of App



**ScalaPageRank**

App Id: 1496028096235_0157

📈 Cluster Load View

| CPU | Memory | Non Heap Memory | GC | HDFS Bytes Read |

stage job
0-5   0

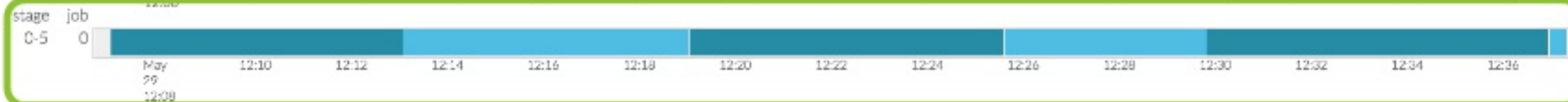| Stage | Job | File Name | Stage Start Time | Stage Duration ((hh):mm:ss) | Average Executor CPU (%) | Stage Peak Heap Memory(MB) | Stage Peak Non Heap Memory(MB) | Total GC Time (seconds) | Parallel Stages |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | distinct at SparkPageRank.scala:58 | 2017/05/29-12:06 | 06:05 | 41.5 | 23,327 | 497 | 5.82 | none |
| 1 | 0 | distinct at SparkPageRank.scala:58 | 2017/05/29-12:13 | 05:58 | 37.8 | 20,106 | 507 | 6.6 | none |
| 2 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:18 | 06:32 | 32.1 | 24,256 | 514 | 22.2 | none |
| 3 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:25 | 04:12 | 31.6 | 28,676 | 517 | 15.5 | none |
| 4 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:29 | 07:08 | 26 | 27,414 | 518 | 31.3 | none |
| 5 | 0 | saveAsHadoopFile at IOCommon.scala:63 | 2017/05/29-12:36 | 00:22 | 26.7 | 23,639 | 524 | 0 | none |

# GC across all Stages of App



**ScalaPageRank**

App Id: 1496028096235_0157

Cluster Load View

Legend: CPU · Memory · Non Heap Memory · GC · HDFS Bytes Read
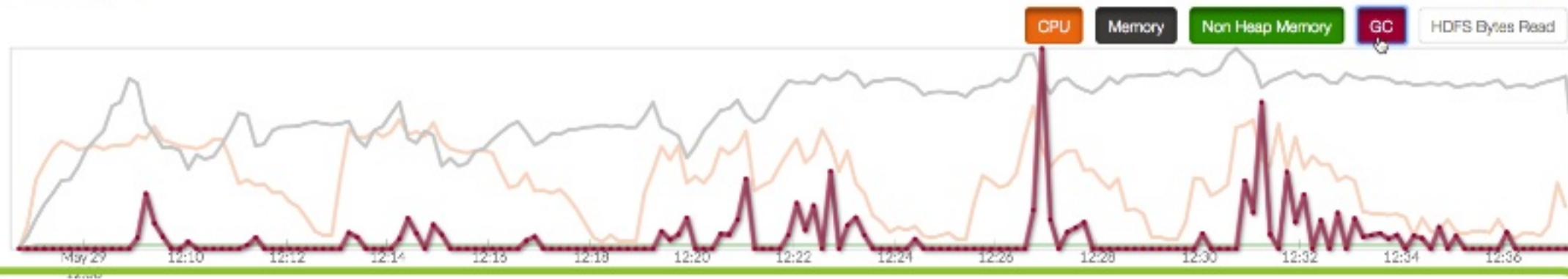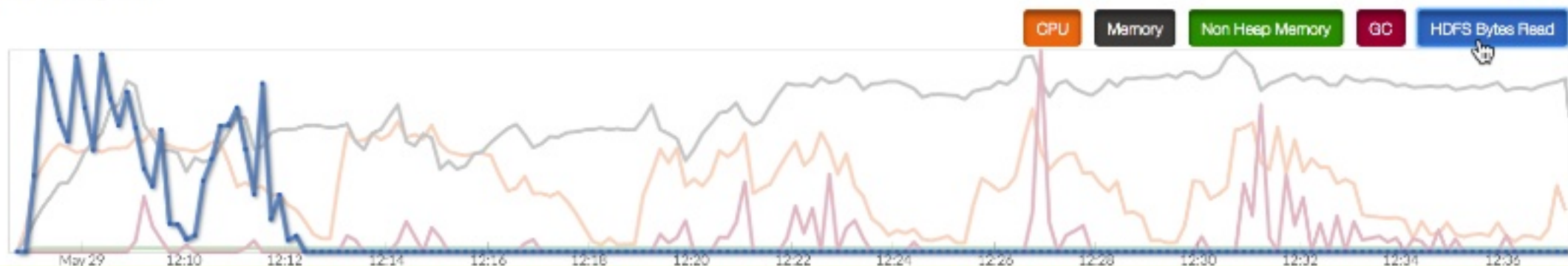
| Stage | Job | File Name | Stage Start Time | Stage Duration ((hh:)mm:ss) | Average Executor CPU (%) | Stage Peak Heap Memory(MB) | Stage Peak Non Heap Memory(MB) | Total GC Time (seconds) | Parallel Stages |
|-------|-----|-----------|------------------|------------------------------|---------------------------|-----------------------------|---------------------------------|--------------------------|------------------|
| 0 | 0 | distinct at SparkPageRank.scala:58 | 2017/05/29-12:06 | 06:05 | 41.5 | 23,327 | 497 | 5.82 | none |
| 1 | 0 | distinct at SparkPageRank.scala:58 | 2017/05/29-12:13 | 05:58 | 37.8 | 20,106 | 507 | 6.6 | none |
| 2 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:18 | 06:32 | 32.1 | 24,256 | 514 | 22.2 | none |
| 3 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:25 | 04:12 | 31.6 | 26,676 | 517 | 15.5 | none |
| 4 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:29 | 07:08 | 26 | 27,414 | 518 | 31.3 | none |
| 5 | 0 | saveAsHadoopFile at IOCommon.scala:63 | 2017/05/29-12:36 | 00:22 | 26.7 | 23,639 | 524 | 0 | none |

SPARK SUMMIT 2017

# HDFS Reads across all Stages of App



## ScalaPageRank

App Id: 1496028096235_0157

Cluster Load View

Legend: CPU · Memory · Non Heap Memory · GC · HDFS Bytes Read
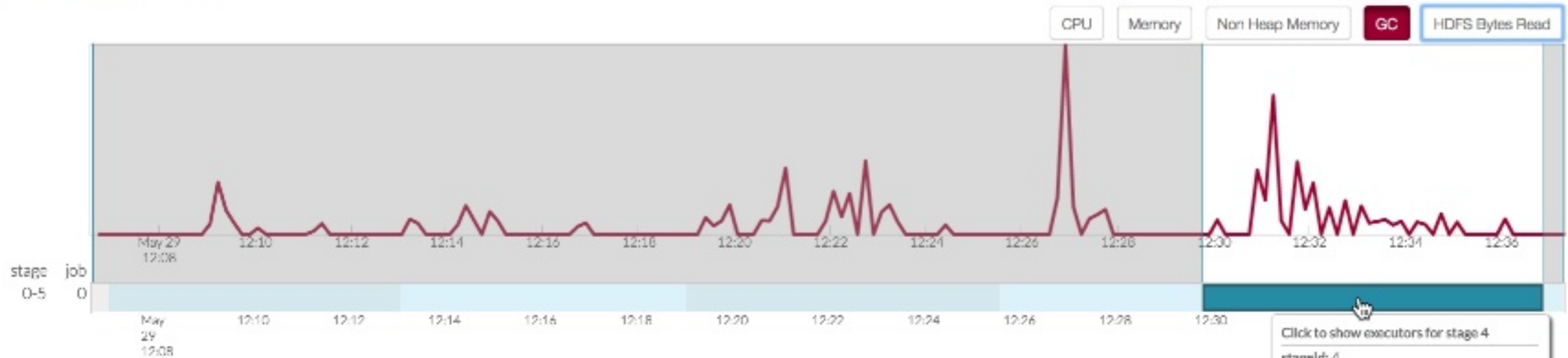
| Stage | Job | File Name | Stage Start Time | Stage Duration ((hh:)mm:ss) | Average Executor CPU (%) | Stage Peak Heap Memory(MB) | Stage Peak Non Heap Memory(MB) | Total GC Time (seconds) | Parallel Stages |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | distinct at SparkPageRank.scala:58 | 2017/05/29-12:06 | 06:05 | 41.5 | 23,327 | 497 | 5.82 | none |
| 1 | 0 | distinct at SparkPageRank.scala:58 | 2017/05/29-12:13 | 05:58 | 37.8 | 20,106 | 507 | 6.6 | none |
| 2 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:18 | 06:32 | 32.1 | 24,256 | 514 | 22.2 | none |
| 3 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:25 | 04:12 | 31.6 | 26,676 | 517 | 15.5 | none |
| 4 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:29 | 07:08 | 26 | 27,414 | 518 | 31.3 | none |
| 5 | 0 | saveAsHadoopFile at IOCommon.scala:63 | 2017/05/29-12:36 | 00:22 | 26.7 | 23,639 | 524 | 0 | none |

# Bringing it all together



**ScalaPageRank**

App Id: 1496028096235_0157

Cluster Load View

Legend: CPU | Memory | Non Heap Memory | GC | HDFS Bytes Read

| Stage | Job | File Name | Stage Start Time | Stage Duration ((hh:)mm:ss) | Average Executor CPU (%) | Stage Peak Heap Memory(MB) | Stage Peak Non Heap Memory(MB) | Total GC Time (seconds) | Parallel Stages |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | distinct at SparkPageRank.scala:58 | 2017/05/29-12:06 | 06:05 | 41.5 | 23,327 | 497 | 5.82 | none |
| 1 | 0 | distinct at SparkPageRank.scala:58 | 2017/05/29-12:13 | 05:58 | 37.8 | 20,106 | 507 | 6.6 | none |
| 2 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:18 | 06:32 | 32.1 | 24,256 | 514 | 22.2 | none |
| 3 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:25 | 04:12 | 31.8 | 26,676 | 517 | 15.5 | none |
| 4 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:29 | 07:08 | 26 | 27,414 | 518 | 31.3 | none |
| 5 | 0 | saveAsHadoopFile at IOCommon.scala:63 | 2017/05/29-12:36 | 00:22 | 26.7 | 23,639 | 524 | 0 | none |

# Let's examine GC activity in Stage 4

## ScalaPageRank

App Id: 1496028096235_0157

Cluster Load View

CPU | Memory | Non Heap Memory | **GC** | HDFS Bytes Read



Click to show executors for stage 4
stageId: 4
attemptId: 0
jobId: 0
file: flatMap at SparkPageRank.scala:62
start time: 2017/05/29-12:29
duration: 07:08
parallelStages: none

| Stage | Job | File Name | Stage Start Time | Stage Duration ((hh:)mm:ss) | Average Executor CPU (%) | Stage Peak Memory | GC Time (seconds) | Parallel Stages |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | distinct at SparkPageRank.scala:58 | 2017/05/29-12:06 | 06:05 | 41.5 | 20,347 | 497 | 6.82 | none |
| 1 | 0 | distinct at SparkPageRank.scala:58 | 2017/05/29-12:13 | 05:58 | 37.8 | 20,106 | 507 | 6.6 | none |
| 2 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:18 | 06:32 | 32.1 | 24,256 | 514 | 22.2 | none |
| 3 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:25 | 04:12 | 31.6 | 26,676 | 517 | 15.5 | none |
| 4 | 0 | flatMap at SparkPageRank.scala:62 | 2017/05/29-12:29 | 07:08 | 26 | 27,414 | 518 | 31.3 | none |
| 5 | 0 | saveAsHadoopFile at IOCommon.scala:83 | 2017/05/29-12:36 | 00:22 | 26.7 | 23,639 | 524 | 0 | none |

# Executor skew increased Stage duration 2x



ScalaPageRank

App Id: 1496028096235_0157    Stage: 4:0

Cluster Load View

Executors for Stage 4 : Attempt 0

# Possible solution: increase number of partitions

Cluster weather

# What if it's not your fault?

# How does cluster weather impact your app ?

# No apparent reason for delay from Spark Web UI

| Stage Id | Description | | Submitted | Duration | Tasks: Succeeded/Total | Input | Output | Shuffle Read | Shuffle Write |
|---|---|---|---|---|---|---|---|---|---|
| 2 | saveAsNewAPIHadoopFile at ScalaTeraSort.scala:60 | +details | 2017/05/23 16:55:02 | 16 min | 8/8 | | 29.8 GB | 13.1 GB | |
| 1 | map at ScalaTeraSort.scala:49 | +details | 2017/05/23 16:51:46 | 3.3 min | 240/240 | 29.1 GB | | | 13.1 GB |
| 0 | BaseRangePartitioner at ScalaTeraSort.scala:56 | +details | 2017/05/23 16:51:30 | 15 s | 8/8 | | | | |

| Stage Id | Description | | Submitted | Duration | Tasks: Succeeded/Total | Input | Output | Shuffle Read | Shuffle Write |
|---|---|---|---|---|---|---|---|---|---|
| 2 | saveAsNewAPIHadoopFile at ScalaTeraSort.scala:60 | +details | 2017/05/23 15:12:33 | 41 min | 8/8 | | 29.8 GB | 13.1 GB | |
| 1 | map at ScalaTeraSort.scala:49 | +details | 2017/05/23 14:50:46 | 22 min | 240/240 | 29.8 GB | | | 13.1 GB |
| 0 | BaseRangePartitioner at ScalaTeraSort.scala:56 | +details | 2017/05/23 14:49:35 | 32 s | 8/8 | 768.0 MB | | | |

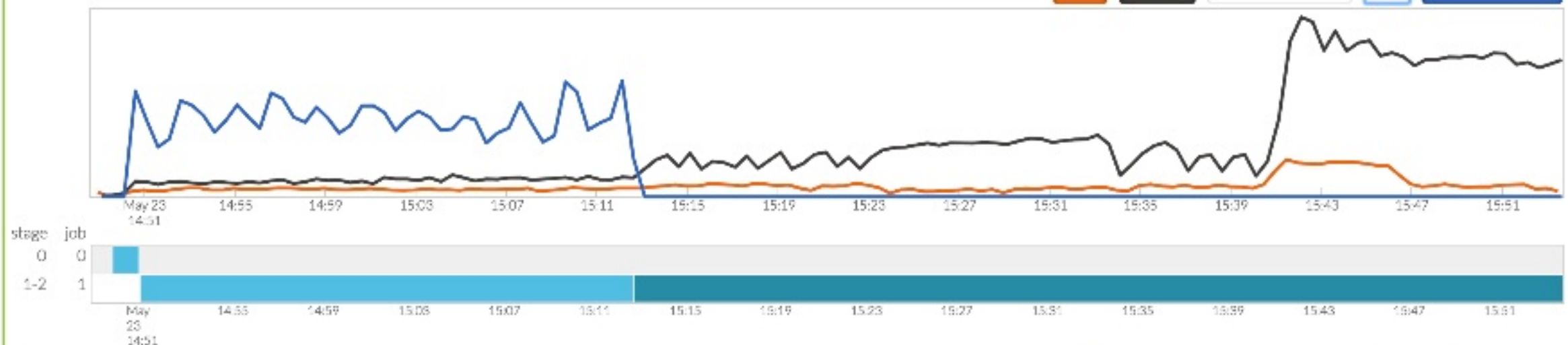# Time series shows slower run of App with much lower resources



Fast Run
High peaks

Slow Run
Low peaks

CPU    Memory    Non Heap Memory    GC    HDFS Bytes Read

# View cluster weather conditions for slower run of App
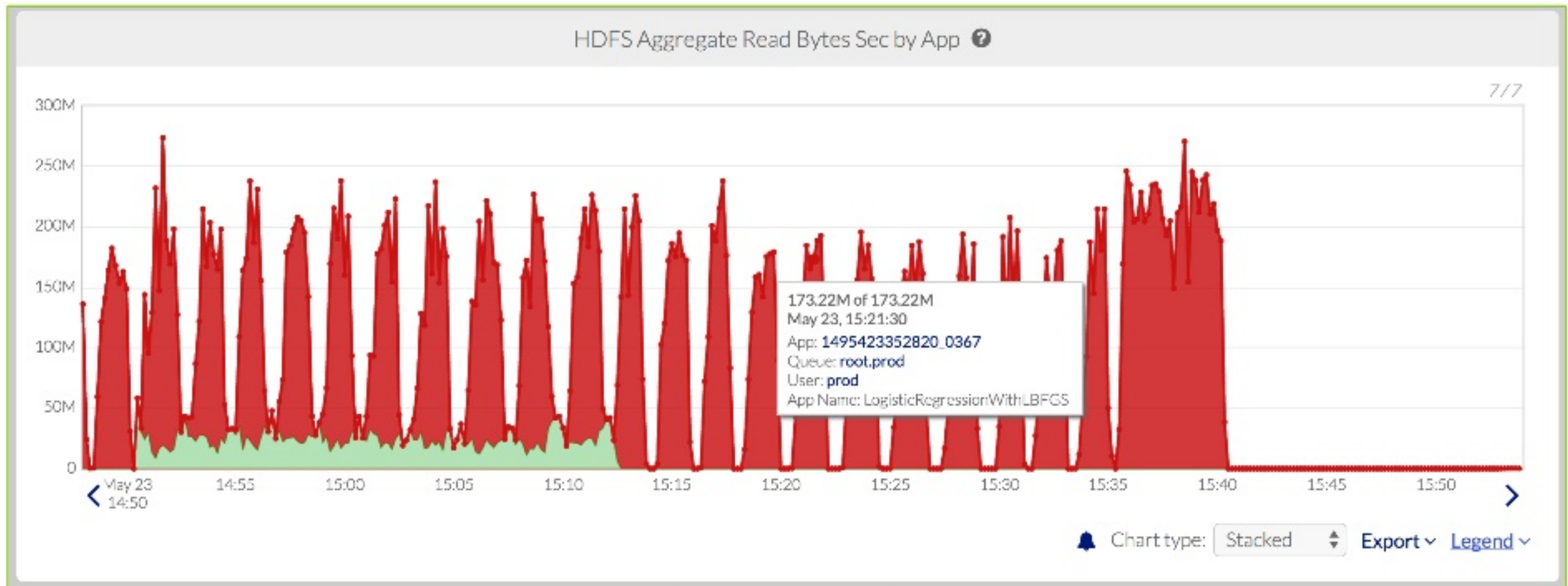
# Cluster weather reveals reason for CPU constraints on slower app

# Cluster weather reveals reason for memory constraints on slower app



Task JVM Used Heap Memory by App

16.68G of 23.7G
May 23, 14:59:20
App: 1495423352820_0367
Queue: root.prod
User: prod
App Name: LogisticRegressionWithLBFGS

# Cluster weather reveals reason for HDFS constraints on slower app

# To recap

- Execution plan integrated with time series shines a spot light on problems

- Stage and code section integrated with resources consumed enables focus on most impactful areas for optimization

- Knowing cluster weather can prevent time wasted debugging non-existent performance issues

# Code Analyzer for Apache Spark

- Free during Early Access starting June 5th, 2017

- Early Access is for development teams

- To learn more visit **pepperdata** booth #101

- My contact vinod@pepperdata.com

**pepperdata.com/products/code-analyzer**

# Thank You.

www.pepperdata.com/products/code-analyzer/

SPARK
SUMMIT
2017