# CLOUDERA
**Educational Services**

# Preparing with Cloudera Data Engineering
## Instructor Guide

# Agenda

- Why this new version?

- The environment

- The notebooks

- The content

- Timings

- Troubleshooting

- What's next?

- Questions

CLOUDERA
Educational Services

# Too Much RDDs

- The content is outdated, a lot of stuff predates 2016 (Spark 1.5). This translates to a lot of RDD based content which nobody in their right minds uses nowadays. (RDD appears in 160 slides vs 71 in the new version)
- The labs rely on stark interpreters which again nobody uses professionally.
- It does not foster HDFS  best practices: HDFS permissions are disabled (I guess for convenience).
- The Common Patterns in Spark Data Processing is a joke made of the description of the PageRank and kMeans algorithms that exist out of the box in recent versions and an awkward introduction to Machine Learning
- No content on Hive, although Hive and Spark have a long history of being closely integrated

**CLOUDERA**
Educational Services

# EDU-KEYCLOAK

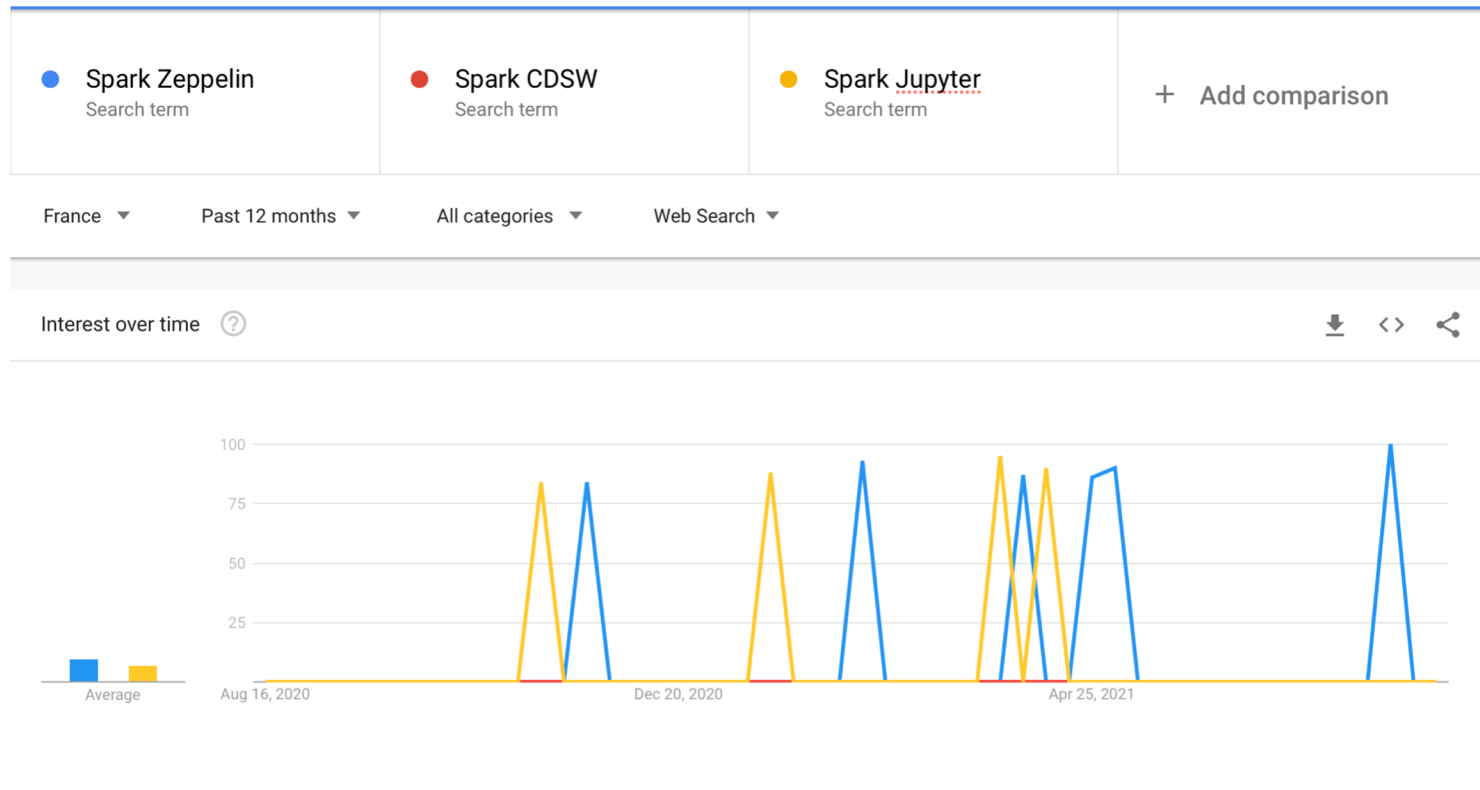## Log In

Username or email

Password

Log In

# No Commonality with Spark Performance

- Spark performance issues are much more likely to occur in Data Engineering than in Data Science stages
- This course should use the same environment so that both contents can be fungible
- Clients will likely request mixed custom courses that will leverage both courses

# Spark 3.0

- Spark 3.0 contains major performance related improvements that promise to perform up to 17* faster

- If we do not jump on this train fast our content will be quickly outdated

CLOUDERA
Educational Services

# Zeppelin is available in CDP and is part of the Spark ecosystem

CLOUDERA
Educational Services

# CLOUDERA
**Educational Services**

---

# The environment

# The Environment

- Started from the latest Spark Application Performance Tuning  template
- Created a /home/training/training-materials/dev home folder
- Updated to  **CDP PVC BE 7.1.7**,
- installed **Spark 3.1**,
- Installed and built the data in HDFS
- Created the Hive tables
- Installed the required jar files
- Installed Arrow
- Installed jq
- Uploaded the notebooks in Zeppelin
- Installed the images for the notebooks

**CLOUDERA**
Educational Services

# Published Services

- ssh
- Cloudera Manager
- Zeppelin
- Hue
- Spark 2 History Server
- Spark 3 History Server

Published services: 6 ▼ Hide Published Services ✚ Add Published Service

This network adapter has the following published services: ✚ Add Published Service

| Internal port | Public address | Action |
|---|---|---|
| 22 | services-emea.skytap.com:16260 | ✖ Remove |
| 7180 | services-emea.skytap.com:16264 | ✖ Remove |
| 8885 | services-emea.skytap.com:16317 | ✖ Remove |
| 8889 | services-emea.skytap.com:16318 | ✖ Remove |
| 18088 | services-emea.skytap.com:16355 | ✖ Remove |
| 18089 | services-emea.skytap.com:9520 | ✖ Remove |

Public IP addresses: *None* ✚ Add Public IP with DNS ✚ Add Static Public IP ⓘ Learn more

Secondary IPs: *None* ▶ Manage Secondary IPs

# /home/training/training_materials/dev

- data
    - The datasets used in the legacy labs
- exercises
    - Spark-application
    - yarn
- scripts
    - Utility scripts
    - Don't trust the setup.sh script to build a new environment, it hasn't been tested. Start one from Skytap instead.
- notebooks
    - Zip file of the current version of notebooks

The entire dev folder is zipped into a single master file: `zip -r dev-20210816.zip dev/*`

CLOUDERA
Educational Services

# Available Services

- Removed unnecessary services to save on resources
  - Impala
  - Hue
  - Oozie
  - NiFi
  - NiFi Registry
- **0 warning starting point**
- HDFS permissions are enabled

---

✅ ClouderaDeveloperTraining ⋮

Cloudera Runtime 7.1.7 (Parcels)

| | | |
|---|---|---|
| ✅ | 🖥 1 Hosts | |
| ✅ | 🐘 HDFS | ⋮ |
| ✅ | 🐝 Hive | ⋮ |
| ✅ | 🐝 Hive on Tez | ⋮ |
| ✅ | Kafka | ⋮ |
| ✅ | Livy | ⋮ |
| ✅ | Livy for Spark 3 | ⋮ |
| ✅ | Spark | ⋮ |
| ✅ | Spark 3 | ⋮ |
| ⚪ | Tez | ⋮ |
| ✅ | YARN | ⋮ |
| ✅ | YARN Queue Manager | ⋮ |
| ✅ | Zeppelin | ⋮ |
| ✅ | ZooKeeper | ⋮ |

# Resource Management

- You should be able to run two applications simultaneously in the default queue

- As a good practice, stop the services you do not need

## Queue Properties
root.default

### Resource Allocation

Minimum User Limit  50.0 ⇅ %

User Limit Factor  0.5  ⇅

CLOUDERA
Educational Services

# CLOUDERA
**Educational Services**

## The notebooks

# Notebooks Design Principles

- A single artefact without internal duplication

  - to simplify maintenance

- That contains the solution, instructions yet that allows the student to type his code

  - for effective learning

- That runs without errors

  - to enable easy regression and performance testing

# Notebooks Structure

- All the Zeppelin notebooks share the same structure

    - About

    - Setup

    - Demo or Lesson

    - Lab *

    - Result

    - Solution

    - Tear down

    - Footer

    * This part is optional for Demo only notebooks

CLOUDERA
Educational Services

# About

- High level information that mainly helps decide whether this is the notebook you are looking for or not

  - Objective: <Short description of the notebook>

  - Files locations:

  - Successful outcome:

  - Before you begin: <Dependencies>

  - Related lessons:

  - Copyright

- Not always rigorously filled to be honest

**CLOUDERA**
Educational Services

# Setup

- This is the section in which all the preparation required for the lab should be carried out.

- It can contain code to retrieve the data required for the lab as well as catch up code to make the notebook independent from previous labs.

- In order for the notebook to always run without errors, special care should be taken when creating directories or files such as deleting them before recreating them.

CLOUDERA
Educational Services

# Demo or Lesson

- This is the section where the instructor walks the students through each paragraph to illustrate the topic of the notebook.

# Lab

- This is the section where the student will try to perform the lab steps.
  It should contain the lab instructions in markdown paragraphs interspersed with  empty code paragraphs with numbered titles

| | |
|---|---|
| Use a shell paragraph to list the content of the AdventureWorks home directory | FINISHED ▷ ⤢ 📖 ⚙ |

**1 - List the content of the AdventureWorks home directory**  FINISHED ▷ ⤢ 📖 ⚙

| | |
|---|---|
| The orders.csv file is the largest. Let's take a look at its content. | FINISHED ▷ ⤢ 📖 ⚙ |

**2 - Do a tail on orders.csv**  FINISHED ▷ ⤢ 📖 ⚙

**CLOUDERA**
Educational Services

# Result

■ This section summarizes what the student has just achieved.

## Result

FINISHED ▷ ⤢ 📖 ⚙

**You have now:** created an insightful dashboard with the data from this company using Spark DataFrames

CLOUDERA
Educational Services

# Solution

- This is the section where the student can look up solutions to the lab steps using the matching numbered titles. It contains only code paragraphs.

# Tear Down

- This section is specific to this new cluster that uses Livy as a broker between Zeppelin and Spark. It contains a single paragraph that deletes the Livy session that the notebook created at the beginning. This ensures that each notebook starts with a fresh Livy session thus avoiding accumulation phenomenons that eventually lead to random failures.
- The script relies on the jq framework that is installed on the cluster.

---

## Tear Down

FINISHED

Took 0 sec. Last updated by anonymous at October 02 2020, 2:04:24 PM.

---

**Delete the Livy session**

FINISHED

```sh
%sh

sessionId=$(curl -s localhost:8998/sessions | jq '.sessions[0].id')
curl -s localhost:8998/sessions/$sessionId -X DELETE
```

Took 4 sec. Last updated by anonymous at August 22 2020, 8:03:04 AM.

CLOUDERA
Educational Services

# Footer

- This footer of the Solution section provides links to additional resources as well as the Cloudera Educational Services home page.

## References

FINISHED

Handling Data Skew in Apache Spark

Took 0 sec. Last updated by anonymous at October 02 2020, 2:10:34 PM.

### Additional resources

FINISHED

We hope you've enjoyed this lab. Below are additional resources that you should find useful:

1. Cloudera Tutorials are your natural next step where you can explore Spark in more depth.
2. Cloudera Community is a great resource for questions and answers on Spark, Data Analytics/Science, and many more Big Data topics.
3. Apache Spark Documentation - official Spark documentation.
4. Apache Zeppelin Project Home Page - official Zeppelin web site.

Took 0 sec. Last updated by anonymous at October 02 2020, 2:10:30 PM.

FINISHED

CLOUDERA

Took 0 sec. Last updated by anonymous at October 02 2020, 2:10:37 PM.

CLOUDERA
Educational Services

# Adding the dev notebooks

To upload the notebooks in Zeppelin, the <span style="color:orange">students</span> need to run the command that is in the README file of their ~/training_materials/dev directory.

```
cd /home/training/training_materials/dev
sh scripts/uploadAllNotebooks.sh notebooks/dev-notebooks-20210830.zip
```
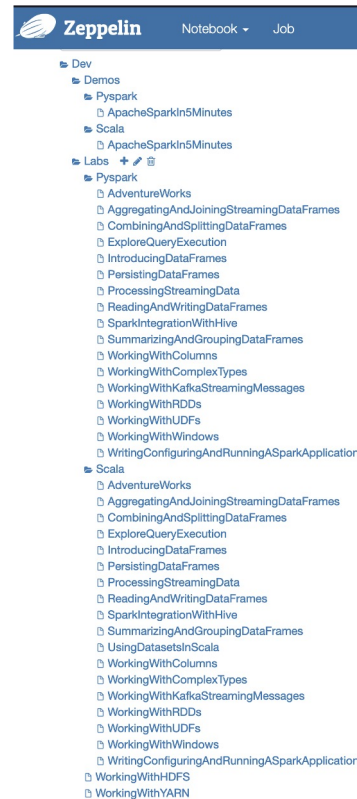
CLOUDERA
Educational Services

# Adding the perf notebooks

If you want to experiment with the Perf notebooks, you can install them using this script:

```
cd /home/training/training_materials/perf/install/scripts
sh uploadAllNotebooks.sh localhost.localdomain:8885 'userName=admin&password=admin'
/home/training/training_materials/perf/notebooks/Perf-notebooks20210810.zip
```

CLOUDERA
Educational Services

# Notebooks Folders Structure

- **Dev**
  - Labs
    - Pyspark
    - Scala
  - Demos
    - Pyspark
    - Scala
- **Perf**
  - Labs
    - Pyspark
  - Demos
    - Pyspark

CLOUDERA
Educational Services

# Notebooks

- There is **no Exercise Manual**, all the exercises are in notebooks

- Sometimes notebooks will contain only instructions and solutions in markdown cells
    - Spark Streaming
    - Submit a Spark Application

- Most of the times, the notebooks will run the code required to perform the tasks described in the instructions

- The notebooks come from several sources
    - Legacy DevSh
    - Data Science
    - HWX Spark course
    - Some I made specifically for this class

- **The notebooks coming from the Data Science class on day #2 include a Lesson part that instructors need to present**

CLOUDERA
Educational Services

# Working with Spark 3

- Spark3 has its own Application History Server listening on port 18089

- Zeppelin uses Livy for Spark3 to communicate with Spark 3

- Livy for Spark3 uses port 28998

- When Zeppelin restarts it defaults to the Livy for Spark2 port 8998

- **Always check in the interpreter page that the zeppelin.livy.url is set to 28998**

zeppelin.livy.url                                                    http://localhost.localdomain:28998

CLOUDERA
Educational Services

# Exhaustive use of setJobGroup

# The Paragraphs

- **There are no empty paragraphs in the notebooks**

- Just paragraphs waiting to be executed

CLOUDERA
Educational Services

# CLOUDERA
**Educational Services**

The content

# Course components

| Component | Version Control | Version |
|---|---|---|
| Students Presentation | Git | https://github.com/HortonworksUniversity/contentFactory |
| Zeppelin notebooks | Git | https://github.com/HortonworksUniversity/contentFactory |
| Skytap template | Template name includes date | |
| Instructor Guide | File name includes date | |
| Data Sheet | File name includes date | |

CLOUDERA
Educational Services

# Agenda

| Day 1 | Day 2 | Day 3 | Day 4 |
|---|---|---|---|
| Class Introduction | Spark DataFrames | Integration with Hive | Introduction to Streaming DataFrames |
| Zeppelin Introduction | Reading DataFrames | Visualization with Zeppelin | Kafka Introduction |
| HDFS Introduction | Working with Columns | Distributed Processing | Integration with Kafka |
| YARN Introduction | Transforming DataFrames | Distributed Persistence | Aggregating and Joining Streaming DataFrames |
| Distributed Processing History | Working with UDFs | Building Spark Applications | (*) Appendix: Scala Datasets |
| Spark RDDs | Working with Windows | | |

(*) if time allows          Comes from the HWX Spark class          Comes
from the Data Science class

# Design Principles

- **#Death2RDDs**
  - Only one chapter and one notebook on RDDs
  - RDD appears in 71 slides vs 160 in the previous version
- **#Death2Slides**
  - 17 Pyspark notebooks
  - 19 Scala notebooks
  - 2 Shell notebooks
- **#Death2Monoliths**
  - The content is made of 59 components listed in a manifest file
  - Those components are managed in a github repository:
    https://github.com/HortonworksUniversity/contentFactory
- **#VivaKnowledgeChecks**
  - I adapted some of the Knowledge Checks we had at HWX and created one

# The Manifest (1 / 4)

Lesson|DEV/Spark/ClouderaDeveloperTrainingClassIntroduction
Lesson|DEV/Spark/ZeppelinIntroduction
Demo/Notebook|Dev_Demos_Pyspark_ApacheSparkIn5Minutes.json
Demo/Notebook|Dev_Demos_Scala_ApacheSparkIn5Minutes.json
Lesson|DEV/Spark/HDFSIntroduction
Lab/Notebook|Dev_Labs_WorkingWithHDFS.json
Lesson|DEV/Spark/YARNIntroduction
Lab/Notebook|Dev_Labs_WorkingWithYARN.json
Lesson|DEV/Spark/DistributedProcessingHistory
Lesson|DEV/Spark/WorkingWithRDDs
Lab/Notebook|Dev_Labs_Pyspark_WorkingWithRDDs.json
Lab/Notebook|Dev_Labs_Scala_WorkingWithRDDs.json
Lesson|DEV/Spark/WorkingWithDataFrames
Lab/Notebook|Dev_Labs_Pyspark_IntroducingDataFrames.json
Lab/Notebook|Dev_Labs_Scala_IntroducingDataFrames.json
Lab/Notebook|Dev_Labs_Pyspark_ReadingAndWritingDataFrames.json
Lab/Notebook|Dev_Labs_Scala_ReadingAndWritingDataFrames.json

# The Manifest (2 / 4)

```
Lab/Notebook|Dev_Labs_Pyspark_WorkingWithColumns.json
Lab/Notebook|Dev_Labs_Scala_WorkingWithColumns.json
Lab/Notebook|Dev_Labs_Pyspark_WorkingWithComplexTypes.json
Lab/Notebook|Dev_Labs_Scala_WorkingWithComplexTypes.json
Lab/Notebook|Dev_Labs_Pyspark_CombiningAndSplittingDataFrames.json
Lab/Notebook|Dev_Labs_Scala_CombiningAndSplittingDataFrames.json
Lab/Notebook|Dev_Labs_Pyspark_SummarizingAndGroupingDataFrames.json
Lab/Notebook|Dev_Labs_Scala_SummarizingAndGroupingDataFrames.json
Lab/Notebook|Dev_Labs_Pyspark_WorkingWithUDFs.json
Lab/Notebook|Dev_Labs_Scala_WorkingWithUDFs.json
Lab/Notebook|Dev_Labs_Pyspark_WorkingWithWindows.json
Lab/Notebook|Dev_Labs_Scala_WorkingWithWindows.json
Lesson|DEV/Spark/ApacheHiveIntroduction
Lesson|DEV/Spark/HiveSparkIntegration
Lab/Notebook|Dev_Labs_Pyspark_SparkIntegrationWithHive.json
Lab/Notebook|Dev_Labs_Scala_SparkIntegrationWithHive.json
Lesson|DEV/Spark/DataVisualizationWithZeppelin
```

# The Manifest (3 / 4)

```
Lab/Notebook|Dev_Labs_Pyspark_AdventureWorks.json
Lab/Notebook|Dev_Labs_Scala_AdventureWorks.json
Lesson|DEV/Spark/DistributedProcessingChallenges
Lesson|DEV/Spark/SparkDistributedProcessing
Lab/Notebook|Dev_Labs_Pyspark_ExploreQueryExecution.json
Lab/Notebook|Dev_Labs_Scala_ExploreQueryExecution.json
Lesson|DEV/Spark/SparkDistributedPersistence
Lab/Notebook|Dev_Labs_Pyspark_PersistingDataFrames.json
Lab/Notebook|Dev_Labs_Scala_PersistingDataFrames.json
Lesson|DEV/Spark/WritingConfiguringAndRunningSparkApplications
Lab/Notebook|Dev_Labs_Pyspark_WritingConfiguringAndRunningASparkApplication.json
Lab/Notebook|Dev_Labs_Scala_WritingConfiguringAndRunningASparkApplication.json
Lesson|DEV/Spark/IntroductionToStructuredStreaming
Lab/Notebook|Dev_Labs_Pyspark_ProcessingStreamingData.json
Lab/Notebook|Dev_Labs_Scala_ProcessingStreamingData.json
Lesson|DEV/Spark/MessageProcessingWithApacheKafka
Lesson|DEV/Spark/StructuredStreamingWithApacheKafka
```

# The Manifest (4 / 4)

```
Lab/Notebook|Dev_Labs_Pyspark_WorkingWithKafkaStreamingMessages.json
Lab/Notebook|Dev_Labs_Scala_WorkingWithKafkaStreamingMessages.json
Lesson|DEV/Spark/AggregatingAndJoiningStreamingDataFrames
Lab/Notebook|Dev_Labs_Pyspark_AggregatingAndJoiningStreamingDataFrames.json
Lab/Notebook|Dev_Labs_Scala_AggregatingAndJoiningStreamingDataFrames.json
Lesson|DEV/Spark/ClouderaDeveloperTrainingClassConclusion
Lesson|DEV/Spark/WorkingWithDatasetsInScala
Lab/Notebook|Dev_Labs_Scala_UsingDatasetsInScala.json
```

**The orange notebooks include a lesson section that should be presented by instructors.**

**CLOUDERA**
Educational Services

# Optional content

- **Working with Datasets in Scala**
  - Scala specific

CLOUDERA
Educational Services

# About the Knowledge Checks

- Use them to engage with students and quickly check their understanding

- But don't get sidetracked into protracted conversations

- You don't have to ask all of them, use the ones

  - That you like or

  - That will validate an important point

- Some questions will be about topics that were not introduced in the slides, you can

  - Choose to skip them or

  - Use them to introduce the corresponding topic

CLOUDERA
Educational Services

# Day #1

| # | Item | Day | Duration |
|---|------|-----|----------|
| 1 | Lesson\|DEV/Spark/ClouderaDeveloperTrainingClassIntroduction | 1 | 45 |
| 2 | Lesson\|DEV/Spark/ZeppelinIntroduction | 1 | 20 |
| 3 | Demo/Notebook\|Dev_Demos_Pyspark_ApacheSparkIn5Minutes.json | 1 | 20 |
| 5 | Lesson\|DEV/Spark/HDFSIntroduction | 1 | 45 |
| 6 | Lab/Notebook\|Dev_Labs_WorkingWithHDFS.json | 1 | 30 |
| 7 | Lesson\|DEV/Spark/YARNIntroduction | 1 | 45 |
| 8 | Lab/Notebook\|Dev_Labs_WorkingWithYARN.json | 1 | 30 |
| 9 | Lesson\|DEV/Spark/DistributedProcessingHistory | 1 | 30 |
| 10 | Lesson\|DEV/Spark/WorkingWithRDDs | 1 | 45 |
| 11 | Lab/Notebook\|Dev_Labs_Pyspark_WorkingWithRDDs.json | 1 | 40 |

CLOUDERA
Educational Services

# Day #2

| # | Item | Day | Duration |
|---|---|---|---|
| 13 | Lesson\|DEV/Spark/WorkingWithDataFrames | 2 | 45 |
| 14 | Lab/Notebook\|Dev_Labs_Pyspark_IntroducingDataFrames.json | 2 | 50 |
| 16 | Lab/Notebook\|Dev_Labs_Pyspark_ReadingAndWritingDataFrames.json | 2 | 40 |
| 18 | Lab/Notebook\|Dev_Labs_Pyspark_WorkingWithColumns.json | 2 | 50 |
| 20 | Lab/Notebook\|Dev_Labs_Pyspark_WorkingWithComplexTypes.json | 2 | 30 |
| 22 | Lab/Notebook\|Dev_Labs_Pyspark_CombiningAndSplittingDataFrames.json | 2 | 40 |
| 24 | Lab/Notebook\|Dev_Labs_Pyspark_SummarizingAndGroupingDataFrames.json | 2 | 40 |
| 26 | Lab/Notebook\|Dev_Labs_Pyspark_WorkingWithUDFs.json | 2 | 30 |
| 28 | Lab/Notebook\|Dev_Labs_Pyspark_WorkingWithWindows.json | 2 | 30 |

# Day #3

| # | Item | Day | Duration |
|---|------|-----|----------|
| 30 | Lesson|DEV/Spark/ApacheHiveIntroduction | 3 | 20 |
| 31 | Lesson|DEV/Spark/HiveSparkIntegration | 3 | 20 |
| 32 | Lab/Notebook|Dev_Labs_Pyspark_SparkIntegrationWithHive.json | 3 | 40 |
| 34 | Lesson|DEV/Spark/DataVisualizationWithZeppelin | 3 | 20 |
| 35 | Lab/Notebook|Dev_Labs_Pyspark_AdventureWorks.json | 3 | 40 |
| 37 | Lesson|DEV/Spark/DistributedProcessingChallenges | 3 | 20 |
| 38 | Lesson|DEV/Spark/SparkDistributedProcessing | 3 | 30 |
| 39 | Lab/Notebook|Dev_Labs_Pyspark_ExploreQueryExecution.json | 3 | 40 |
| 41 | Lesson|DEV/Spark/SparkDistributedPersistence | 3 | 30 |
| 42 | Lab/Notebook|Dev_Labs_Pyspark_PersistingDataFrames.json | 3 | 40 |
| 44 | Lesson|DEV/Spark/WritingConfiguringAndRunningSparkApplications | 3 | 30 |
| 45 | Lab/Notebook|Dev_Labs_Pyspark_WritingConfiguringAndRunningASparkApplication.json | 3 | 40 |

# Day #4

| # | Item | Day | Duration |
|---|------|-----|----------|
| 47 | Lesson\|DEV/Spark/IntroductionToStructuredStreaming | 4 | 40 |
| 48 | Lab/Notebook\|Dev_Labs_Pyspark_ProcessingStreamingData.json | 4 | 60 |
| 50 | Lesson\|DEV/Spark/MessageProcessingWithApacheKafka | 4 | 40 |
| 51 | Lesson\|DEV/Spark/StructuredStreamingWithApacheKafka | 4 | 30 |
| 52 | Lab/Notebook\|Dev_Labs_Pyspark_WorkingWithKafkaStreamingMessages.json | 4 | 60 |
| 54 | Lesson\|DEV/Spark/AggregatingAndJoiningStreamingDataFrames | 4 | 30 |
| 55 | Lab/Notebook\|Dev_Labs_Pyspark_AggregatingAndJoiningStreamingDataFrames.json | 4 | 60 |
| 57 | Lesson\|DEV/Spark/ClouderaDeveloperTrainingClassConclusion | 4 | 10 |
| 58 | Lesson\|DEV/Spark/WorkingWithDatasetsInScala | 4 | 30 |
| 59 | Lab/Notebook\|Dev_Labs_Scala_UsingDatasetsInScala.json | 4 | 30 |

# CLOUDERA
**Educational Services**

## Troubleshooting

# The Spark History Server becomes weak in the knees

- Just click on the 'Refresh History Server' in the Spark Actions menu:

CLOUDERA
Educational Services

# Your Livy session is dead or smells funny

■ This is unlikely to occur because I delete the Livy session at the bottom of each notebook so that you get a new one with each notebook but if it does, open the interpreter binding of the notebook and click on the blue loop icon next to the livy button.

**Settings**

**Interpreter binding**

Bind interpreter for this note. Click to Bind/Unbind interpreter. Drag and drop to reorder interpreters. The first interpreter on the list becomes default. To create/remove interpreters, go to Interpreter menu.

↻ livy %livy (default), %sql, %pyspark, %sparkr, %shared

↻ md %md

↻ angular %angular

↻ sh %sh

Save    Cancel


Jazz is not dead,
it just smells funny.

- Frank Zappa -

CLOUDERA
Educational Services

# You get strange error messages about properties being unknown

- You restarted Zeppelin
- It defaulted to the Livy for Spark2 port
- You need to change the url of zeppelin.livy.url in the livy interpreter settings back to 28998

zeppelin.livy.url                                                    http://localhost.localdomain:28998

CLOUDERA
Educational Services

# Your first %jdbc cell reports an error

- If your first %jdbc cell reports an error just restart the jdbc interpreter
  - Open the interpreter binding panel by clicking on the small cog icon at the top right of the screen
  - Click on the looping arrows next to the jdbc blue button
  - Click Save

**Settings**

**Interpreter binding**

Bind interpreter for this note. Click to Bind/Unbind interpreter. Dr
The first interpreter on the list becomes default. To create/remo

⟳  livy %livy (default), %sql, %pyspark, %sparkr, %shared

⟳  md %md

⟳  angular %angular

⟳  sh %sh

⟳  jdbc %jdbc

[ Save ]  [ Cancel ]

CLOUDERA
Educational Services

# CLOUDERA
**Educational Services**

# Knowledge Checks
Cheat sheet

# HDFS - Answers

1 - Blocks, replicas/copies, reliability/robustness, data locality
2 - NameNode
3 - DataNode
4 - True.  This allows the NameNode to be so highly-available
5 - False.  Clients write to the first DataNode in the list created by the NameNode and the DN's then "pipeline" the data writing to the additional DN's

# YARN - Answers

1 – master = ResourceManager and worker = NodeManager
2 - Memory
3 – False, they run on worker nodes and control a specific job
4 – The ApplicationMaster must decide what to do which is not always to spin up a replacement
5 – False, queues are given a certain percentage of the whole.  Bonus answer is that "node labels" could be assigned to particular worker nodes which then can be configured for a given queue, but we didn't discuss this in the preso.

# Distributed Processing History - Answers

1 – Map and Reduce

2 – The number of blocks the input data is persisted to on HDFS.  User can supply the number of reducers.

3 – Mappers are called with a single KVP and can return 0..m.  Reducers receive a single KVP (the value is a list of values) and can also return 0..m.

4 – False, but you can have a Map-only job which would be the quickest possible job due to the lack of shuffle/sort and Reducer phases.

5 – Because they are "easier" for most developers and analysts

6 – False, Spark took many concepts from MapReduce and implemented them in a new way.

7 – False, no code change is required the GPU library mirrors the existing CPU one

# Working with RDDs - Answers

1. Resilient Distributed Dataset
2. sc.parallelize() and sc.textfile()
3. False. Transformations result in new RDDs being created. In Spark, data is immutable.
4. flatmap()
5. True
6. Lazy evaluation
7. False. The intersection function performs this task. The distinct function would remove duplicates elements, so that each element is only listed once regardless of how many times it appeared in the original RDD.
8. True

# Working with DataFrames - Answers

1) An API that allows to SQL to generate Spark jobs
2) Using spark.read, available formats are csv, json,parquet, orc, text
3) A DataFrame is a Dataset of type row, Datasets are strongly typed objects
4) False

# Hive Introduction - Answers

1 – HiveServer2 (HS2)

2 – MR, Spark & Tez.

3 – When a managed table is dropped, all of its underlying files will be deleted.

5 – Tez, ORC, Vectorization, ACID transactions

# Data Visualization with Zeppelin - Answers

1. Enable humans to make inferences and draw conclusions about large sets of data that would be impossible to make by looking at the data in tabular format.
2. Five
3. Export to JSON format, then they can import it.
4. Give them the note URL.
5. The Report view.
6. Link the paragraph
7. Dynamic forms

# Spark Distributed Processing - Answers

1 - False
2 - False:  an Application is a sequence of jobs
3 - Catalyst and Tungsten
4 - False, Lambdas impede Catalyst
5 - False
6 - False
7 - False, they are equivalent

CLOUDΞRA
Educational Services

# For You

- Read through the presentation
    - If you have comments/questions use the Google Slides internal commenting feature
    - I am the only one to be able to edit the presentation
- Launch an instance of the Skytap template
- Run through the notebooks

CLOUDERA
Educational Services

# Ideas for future evolution

- Consider storing the datasets in S3 for additional modularity

- Consider adding Flink modules to offer an alternative to the Spark Streaming content

CLOUDERA
Educational Services

# CLOUDERA
## Educational Services

THANK YOU