# Talend Big Data Sandbox

**Big Data Insights Cookbook**

# Talend Big Data Sandbox

**Big Data Insights Cookbook**

**Overview**

**Pre-requisites**

**Setup & Configuration**

**Hadoop Distribution Download**

**Demo** *(Scenario)*

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|---|---|---|---|---|

## About this cookbook

### What is the Talend Cookbook?



Using the Talend Real-Time Big Data Platform, this Cookbook provides step-by-step instructions to build and run an end-to-end integration scenario.



The demos are built on real world use-casees and demonstrate how Talend, Spark, NoSQL and real-time messaging can be easily integrated into your daily business.



Whether batch, streaming or real-time integration, you will begin to understand how Talend can be used to address your big data challenges and move your business into the Data-Driven Age.

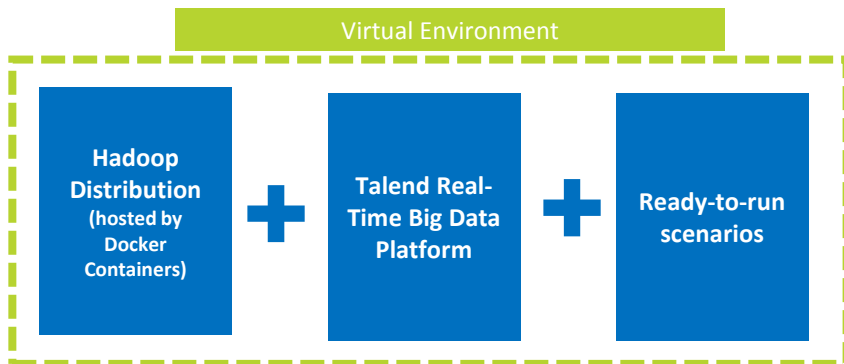# Talend Big Data Sandbox

**Big Data Insights Cookbook**
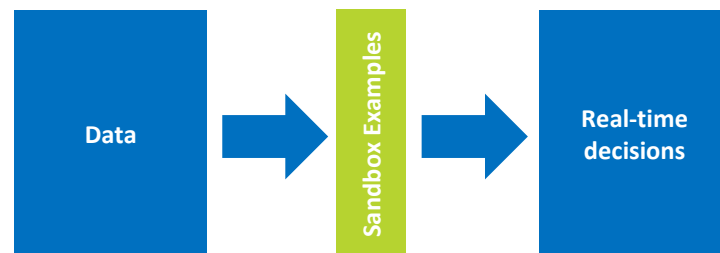
## What is the Big Data Sandbox?



The Talend Real-Time Big Data Sandbox is a virtual environment that combines the Talend Real-Time Big Data Platform with some sample scenarios pre-built and ready-to-run.

See how Talend can turn data into real-time decisions through sandbox examples that integrate Apache Kafka, Spark, Spark Streaming, Hadoop and NoSQL.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

## What Pre-requisites are required to run Sandbox?

Talend Platform for Big Data includes a graphical IDE (Talend Studio),
teamwork management, data quality, and advanced big data features.

## Internet connection required for the entire setup process

To see a full list of features please visit Talend's Website:
http://www.talend.com/products/real-time-big-data

You will need a Virtual Machine player such as VMWare
or Virtualbox, which can be downloaded here:
- VMware Player Site
- Virtualbox Site

**Follow the VM Player install instructions from the provider**

The recommended host machine should have:

**Memory 8-10GB**   **Disk Space 20GB**   (5GB is for the image download)

**Download the Sandbox Virtual Machine file**
https://info.talend.com/prodevaltpbdrealtimesandbox.html

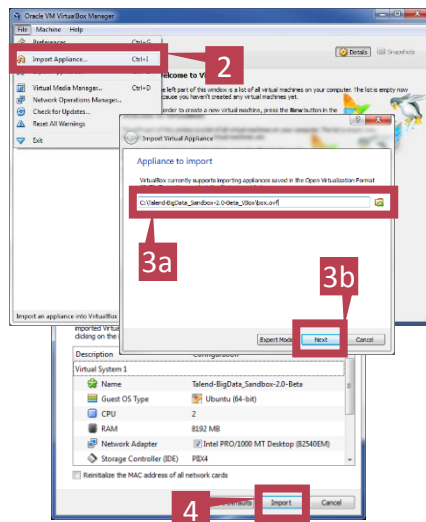# Talend Big Data Sandbox
## Big Data Insights Cookbook

## How do I set-up & configure Sandbox?

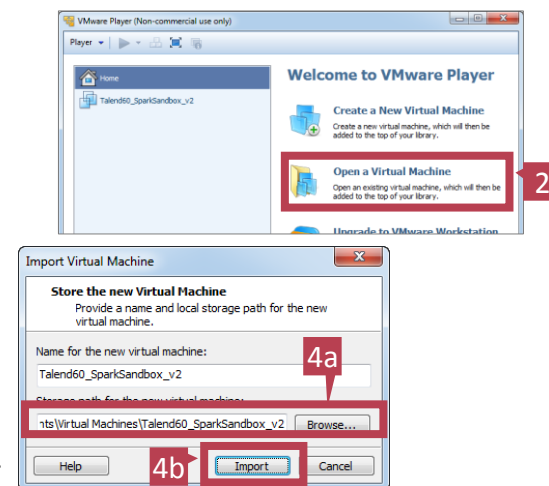**Follow the steps below to install and configure your Big Data Sandbox:**

- Save the downloaded Virtual Machine file to a location on your local PC that is easy to access (e.g. C:/TalendSandbox)
- Follow the instructions below based on the Virtual Machine Player and matching Sandbox file that you are using

### Virtualbox

1. Open Virtualbox.

2. From the menu bar, select **File > Import Appliance…**

3. Navigate to the **.ova** file that you downloaded. Select it and click **Next**.

4. Accept the default Appliance Settings by clicking **Import.**

### VMware Player

1. Open VMware Player.

2. Click on "**Open a Virtual Machine**"

3. Navigate to the **.ova** file that you downloaded. Select it and click **Open**.

4. Select the Storage path for the new Virtual Machine (e.g. C:/TalendSandbox/vmware) and then click **Import**.

**Note:** The Talend Big Data Sandbox Virtual Machines come pre-configured to run with 8GB RAM and 2 CPU's. You may need to adjust these settings based on your PC's capabilities. While not pre-configured, it is also recommended to enable a Sound Card/Devise before starting the VM to take advantage of Tutorial Videos within the Virtual Environment.
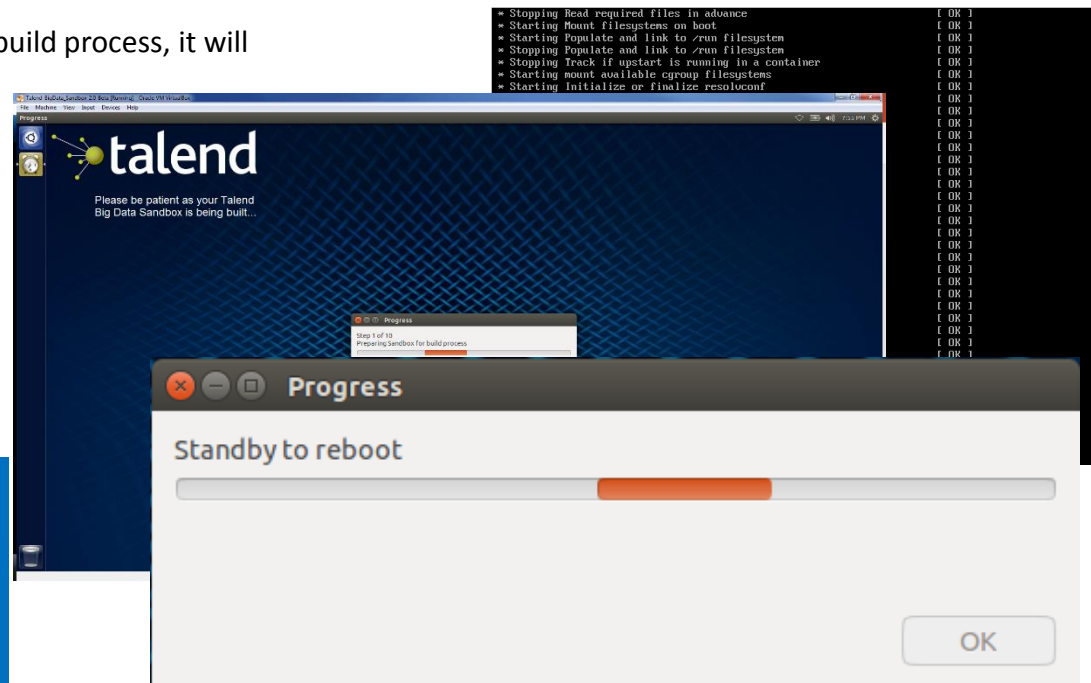
| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|---|---|---|---|---|

## Starting the VM for the first time…

- When you start the Talend Big Data Sandbox for the first time, the virtual machine will begin a **6-step process to build** the Virtual Environment.

- This process can take **10-20 mins depending on internet connection speeds** and network traffic. Popup messages will be present on screen to keep you informed of the progress.

- Once the Sandbox has completed it's build process, it will **automatically reboot**.



### Login Info

| | |
|---|---|
| **User:** | **talend** |
| **Password:** | **talend** |
| **Sudo Password:** | **talend** |

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|---|---|---|---|---|

## Starting the VM for the first time (cont.)

- Once the Virtual Machine reboots, the **Docker Components** that were installed during the build process will need to **initialize.**

- Additional **Popup messages** will appear to inform you of the progress.

- When complete, a message will show that the **System is Ready**!



**Login Info**

| | |
|---|---|
| **User:** | **talend** |
| **Password:** | **talend** |
| **Sudo Password:** | **talend** |

## Shutting Down the VM

- There are 2 methods to shut down the Virtual Machine - **Standard** and **Advanced**

- **Standard Shutdown**
  - This is the standard shutdown method for Ubuntu
  - It is available from the system menu at the Top-Right of the menu bar

- **Advanced Shutdown**
  - Found on the Desktop, Double-click the Power Icon to start a custom shutdown script that cleanly exits all running Docker Containers before shutting down the VM. It will take a few minutes to execute but will make Startup quicker.
  - You can choose to either <u>Shutdown</u> or <u>Reboot</u> your VM via the Advanced Shutdown method.
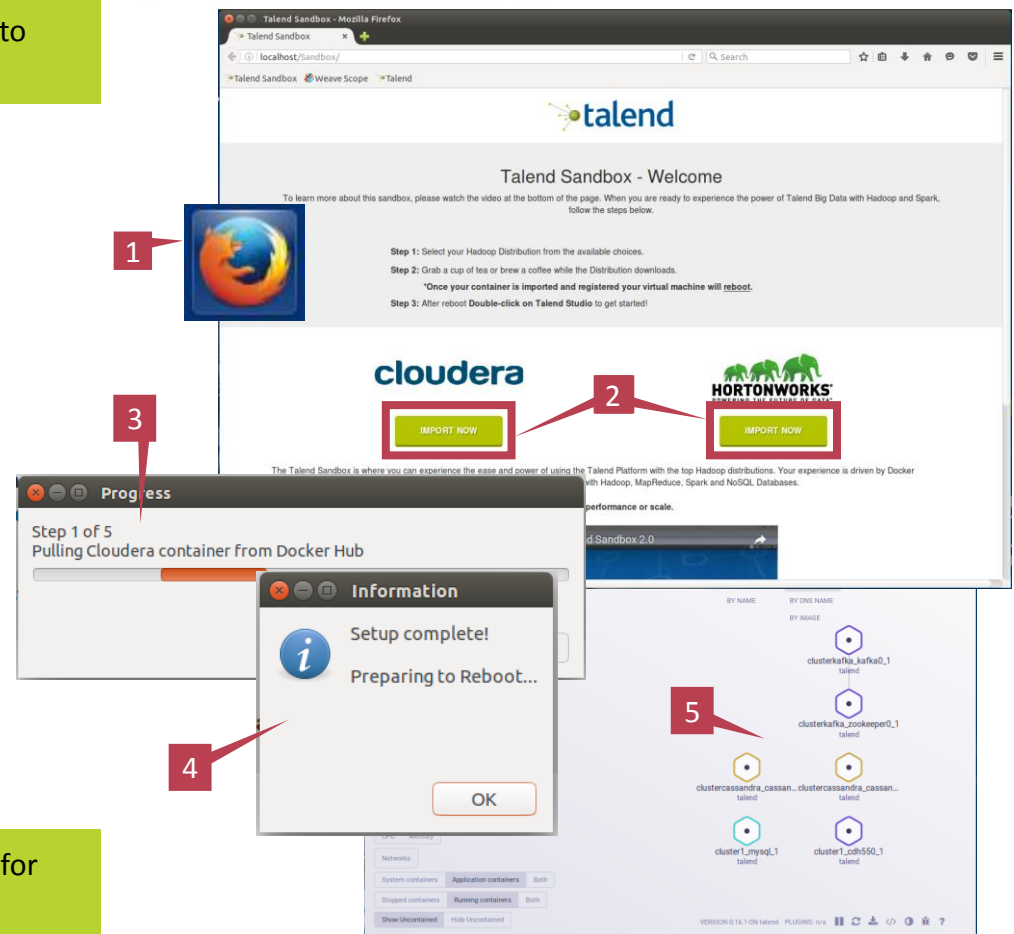
## Starting Talend Studio for the first time…

1. Start Talend Studio by double-clicking on the desktop Icon or single clicking on the Unity Bar Icon.

2. Click **I Accept** the End User License Agreement.

3. Click on **Manage Connections** and enter your email address, then Click **OK.**

4. Select the **Base_Project** – java project and click **Finish**.

5. Once Studio Loads, Close the Welcome Screen.

6. Install **Additional Talend Packages**.  Select *Required third-party libraries* and click **Finish**.

7. A popup will display all 3rd party licenses that need acceptance.  Click the "*I accept the terms of the selected license agreement*" radio button and click **Accept all.**

8. **Let the downloads complete before continuing.**

# Talend Big Data Sandbox
## Big Data Insights Cookbook

## Choosing a Distribution…

**Note:** It is not necessary to download a Distribution to evaluate the Sandbox.  Click Here to begin now!

**Follow the steps below to install a Hadoop Distribution in the Talend Big Data Sandbox:**

1. Start **Firefox**

2. **Choose the Distribution** you would like to evaluate with the Talend Platform.

3. **Be patient** as the Virtual Machine downloads and installs the selected Distributions.  Each distribution is approx. 2.5GB *(Notifications will indicate progress.)*

4. The **Virtual Machine will reboot** when the installation is complete.

5. Upon reboot, the system will initialize the installed Docker containers.  Once complete, you will see a popup stating "System is Ready!!"

**Note:** Be sure to watch the available Tutorial Videos for more information on the Sandbox

# Talend Big Data Sandbox
## Big Data Insights Cookbook

**Note:** This demo is available in Local Mode and Distribution Mode. In Local Mode it utilizes Talend's Local Spark Engine and Local File System. In Distribution Mode, it utilizes the selected Distro's Yarn Resource Manager and HDFS File System.

## Overview:

In this Demo you will see a simple version of making your website an Intelligent Application.

**You will experience:**

- Building a Spark Recommendation Model

- Setting up a new Kafka topic to help simulate live web traffic coming from Live web users browsing a retail web store.

- Most important you will see first-hand with Talend how you can take streaming data and turn it into real-time recommendations to help improve shopping cart sales.



**The following Demo will help you see the value that using Talend can bring to your big data projects:**
The Retail Recommendation Demo is designed to illustrate the simplicity and flexibility Talend brings to using Spark in your Big Data Architecture.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

**This Demo highlights:**

| Kafka | Machine Learning | Spark Streaming |
|-------|------------------|-----------------|
| Create a Kafka Topic to Produce and Consume real-time streaming data | Create a Spark recommendation model based on specific user actions | Stream live recommendations to a Cassandra NoSQL database for "Fast Data" access for a WebUI |

If you are familiar with the ALS model, you can update the ALS parameters to enhance the model or just leave the default values.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

## Demo Setup:

Before you can execute the Retail Recommendation Demo, you will need to generate the source data and pre-populate the Cassandra Lookup Tables.

1. Navigate to the **Job Designs** folder.

2. Click on **Standard Jobs > Realtime_Recommendation_Demo**

3. Double click on **Step_1a_Recommendation_DemoSetup 0.1** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.

5. When the job is finished, repeat steps 1-4 for **Step_1b_Recommendation_DemoSetup 0.1**

**Execute the Retail Recommendation Demo:**

Create a Kafka Topic:

1. Navigate to the **Job Designs** folder:

2. Click on **Standard Jobs > Realtime_Recommendation_Demo**

3. Double click on **Step_2_Recommendation_Create_KafkaTopic 0.1** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.

Now you can generate the recommendation model by loading the product ratings data into the Alternating Least Squares (ALS) Algorithm. Rather than coding a complex algorithm with Scala, a single Spark component available in Talend Studio simplifies the model creation process. **The resultant model can be stored in HDFS or in this case, locally.**

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|----------|----------------|------------------------|---------------------|-------------------|

| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |

**Execute the Retail Recommendation Demo:**

Generate a Recommendation Model using Spark.

1. Navigate to the **Job Designs** folder.

2. Click on **Big Data Batch** >
   **Realtime_Recommendations_Demo**

3. Double click on
   **Step_3_Recommendation_Build_Model_Spark 0.1.**
   This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute

With the Recommendation model created, your lookup tables populated and your Kafka topic ready to consume data, you can now stream your Clickstream data into your Recommendation model and put the results into your Cassandra tables for reference from a WebUI.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|---|---|---|---|---|

| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |

**Execute the Retail Recommendation Demo:**

1. Navigate to the **Job Designs** folder.

2. Click on **Standard Jobs > Realtime_Recommendations_Demo**

3. Double click on **Step_4a_Recommendation_Push_to_Kafka 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.

This job is setup to simulate real-time streaming of web traffic and clickstream data into a Kafka topic that will then be consumed by our recommendation engine to produce our recommendations.

**After starting the Push to Kafka, continue to the next steps of the demo.**

## Execute the Retail Recommendation Demo:

In this job:
- A Kafka Consumer reads in Clickstream Data.
- The data is fed into the Recommendation Engine, producing Real-time "offers" based on the current user's activity.
- The tWindow component controls how often recommendations are generated.
- The recommendations are sent to 3 output streams
  - ✓ **Execution window** for viewing purposes
  - ✓ **File System** for later processing in your Big Data Analytics environment
  - ✓ **Cassandra** for use in a "Fast Data" layer by a WebUI

**With data streaming to the Kafka Topic…**
Start the recommendation pipeline.

1. Navigate to the **Job Designs** folder.

2. Click on **Big Data Streaming > Realtime_Recommendation_Demo**

3. Double click on **Step_4b_Recommendation_Realtime_Engine _Pipeline 0.1.** This opens the job in the designer window.

4. Click on **Run** to Start Recommendation Engine

# Talend Big Data Sandbox
**Big Data Insights Cookbook**

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|----------|----------------|-----------------------|---------------------|-------------------|

| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |
|-----------------------|-------------|-------------|--------------|---------------|

## Execute the Retail Recommendation Demo:

**Watch the execution output window**. You will now see your real-time data coming through with recommended products based on your Recommendation Model.

**Recommendations** are also written to a Cassandra database so they can be referenced by a WebUI to offer, for instance, last minute product suggestions when a customer is about to check-out.

➢ Once you have seen the results, you can **Kill** the Recommendation Engine and the Push to Kafka jobs to stop the streaming recommendations.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|---|---|---|---|---|

| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |
|---|---|---|---|---|

**Note:**  Execution of this demo requires a Hadoop distribution.  If a distro hasn't been selected, <u>click here</u>.

## Overview:

In this example we will utilize real-time streaming of data through a Kafka queue to track on-field player movements at a sporting event.

**You will experience:**

- Creating and populating a Kafka queue with real-time streaming data from an IoT device (i.e. field camera sensors).

- Using Spark Streaming technology to calculate speed and distance traveled by individual players.

- Charting player speed and distance in a real-time web-based dashboard.

# Talend Big Data Sandbox
**Big Data Insights Cookbook**

**This Demo highlights:**

| IoT data to Kafka | Spark Streaming | REST Service to Live Dashboard |
| --- | --- | --- |
| Capture IoT data in XML files, then load that data to a Kafka Queue for real-time processing. | Use Spark Streaming Technology to quickly calculate player distance and speed as their positions change on the playing field. | Use a restful web service to track player movements in a web-based dashboard. |

**Execute the Sport Stats Demo:**

Create a Kafka topic from which live data will stream

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > Realtime_SportStats_Demo**

3. Double click on **Step_1_SportStats_Create_KafkaTopic 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|----------|----------------|----------------------|---------------------|-------------------|

| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |
|-----------------------|-------------|-------------|--------------|---------------|

**Execute the Sport Stats Demo:**

Read data from an XML file (generated by sensor readings, for example) and populate the Kafka topic

1. Navigate to the **Job Designs** folder.
2. Click on **Standard > Realtime_SportStats_Demo**
3. Double click on **Step_2_SportStats_Read_Dataset 0.1.** This opens the job in the designer window.
4. From the Run tab, click on **Run** to execute.

> **This step simulates live player-tracking data being fed to a Kafka topic.**

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|---|---|---|---|---|

| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |
|---|---|---|---|---|

## Execute the Sport Stats Demo:

In this job:
- A Kafka Consumer reads the sensor data.
- A tWindow component controls how often data is read from the Kafka topic – in this case, 10 seconds worth of data is read every 10 seconds.
- The data is normalized for easier processing.
- Using the tCache components the process calculates distance and speed based on current and previous player positions.
- The resultant data is sent to 2 output streams
  - ✓ **Execution window** for viewing purposes
  - ✓ **MySQL Database** where it will be read by a web service to generate dashboard graphics. (MySQL is running on a Docker container)

1. Navigate to the **Job Designs** folder.

2. Click on **Big Data Streaming > Realtime_SportStats_Demo**

3. Double click on **Step_3_SportStats_LiveStream 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.

> **After starting the SportStats Live Stream, continue to the next steps of the demo.**

## Execute the Sport Stats Demo:

Start the Web Service to populate the Sport Stats web-based dashboard

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > Realtime_SportStats_Demo**

3. Double click on **Step_4_SportStats_WebService 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.

> **With the Web Service running, continue to the next step in this demo.**

# Talend Big Data Sandbox
## Big Data Insights Cookbook

**Execute the Sport Stats Demo:**

Watch the Live Dashboard reflect player movements with real-time updates

1. Open Firefox Web Browser.

2. On the Bookmarks toolbar, click on **Demos > SportStats Demo**



➢ Once you have seen the results, back in Talend Studio, you can **Kill** both the Web Service job and the Live Streaming job.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
| --- | --- | --- | --- | --- |

| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |
| --- | --- | --- | --- | --- |

**Note:** Execution of this demo requires a Hadoop distribution. If a distro hasn't been selected, <u>click here</u>.

## Overview:

In this example we demonstrate using native Map Reduce to enrich a dataset and aggregate the results for different web-based dashboards.

**You will experience:**

- Data loading to HDFS.

- Using MapReduce to enrich and aggregate data within the Hadoop Environment.

- Use of 3rd party graphing tools to generate a web-based dashboard of the calculated results.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

**This Demo will highlight:**

### HDFS

Read and Write data to HDFS with simple components from Talend

### Native MapReduce

Use Talend's MapReduce components to enrich and analyze data, natively, in Hadoop

### Insights

Feed your analysis data to a graphing tool such as Microsoft Excel or Tableau for stunning displays of the results.

**Demo Setup:**

Load data to HDFS.

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > Clickstream_Scenario > Pre_Requirements**

3. Double click on **LoadWeblogs 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute

When this job completes, look-up files will be uploaded to HDFS for use by the MapReduce jobs.

## Execute the Clickstream Demo:

The result of this process is aggregated data indicating the product interests of different areas across the United States for visualization within a Google Chart.

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > Clickstream_Scenario**

3. Double click on **Step_1_Clickstream_MasterJob 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.

**Note:** If asked to Download and Install additional Jar files, Click on **Download and Install**.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|----------|----------------|------------------------|---------------------|--------------------|

| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |

**Execute the Apache Weblog Demo:**

View the data in HDFS.

1. Open Firefox

2. Click on the bookmarked link titled **HDFS Browser**

3. In the **Utilities** Dropdown, select **Browse the File System** and navigate to */user/talend/clickstream_demo/output/results*

4. To view the data file, you must download it from HDFS.  This can be done right within the web browser by clicking on *part-00000* and choosing **download**

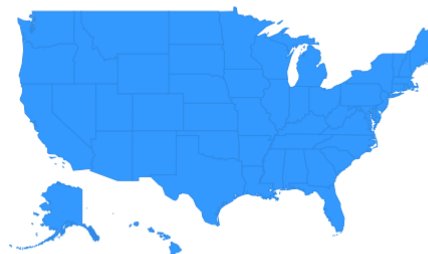# Talend Big Data Sandbox
## Big Data Insights Cookbook

**Execute the Clickstream Demo:**

View the analytical analysis dashboard.

1. Open **Firefox** Web Browser.

2. On the Bookmarks toolbar, click on **Demos > Clickstream Demo**

3. **Mouse over** the states to see the counts**.**

Each color is a category based on Category ID. From this we can see the top and bottom categories in different states/regions of the US based on the clickstream data gathered from the visitors to our retailers website. This gives us a picture as to where to focus our energy -- take care of the top categories and put the bottom ones for more sales.
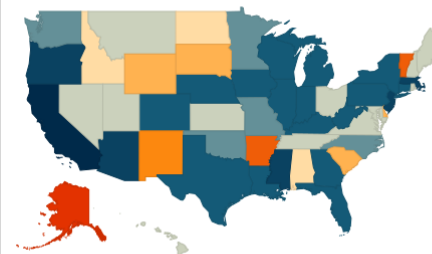
**Top Categories**

**Bottom Categories**

Table below is sortable - click on header to sort.

| State | Category | Clicks |
|---|---|---|
| US-MA | accessories | 1 |
| US-GA | accessories | 8 |
| US-MI | accessories | 2 |
| US-AZ | accessories | 1 |
| US-OR | accessories | 1 |
| US-MS | accessories | 1 |
| US-NJ | accessories | 3 |
| US-NC | automotive | 14 |
| US-NH | automotive | 6 |

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|---|---|---|---|---|

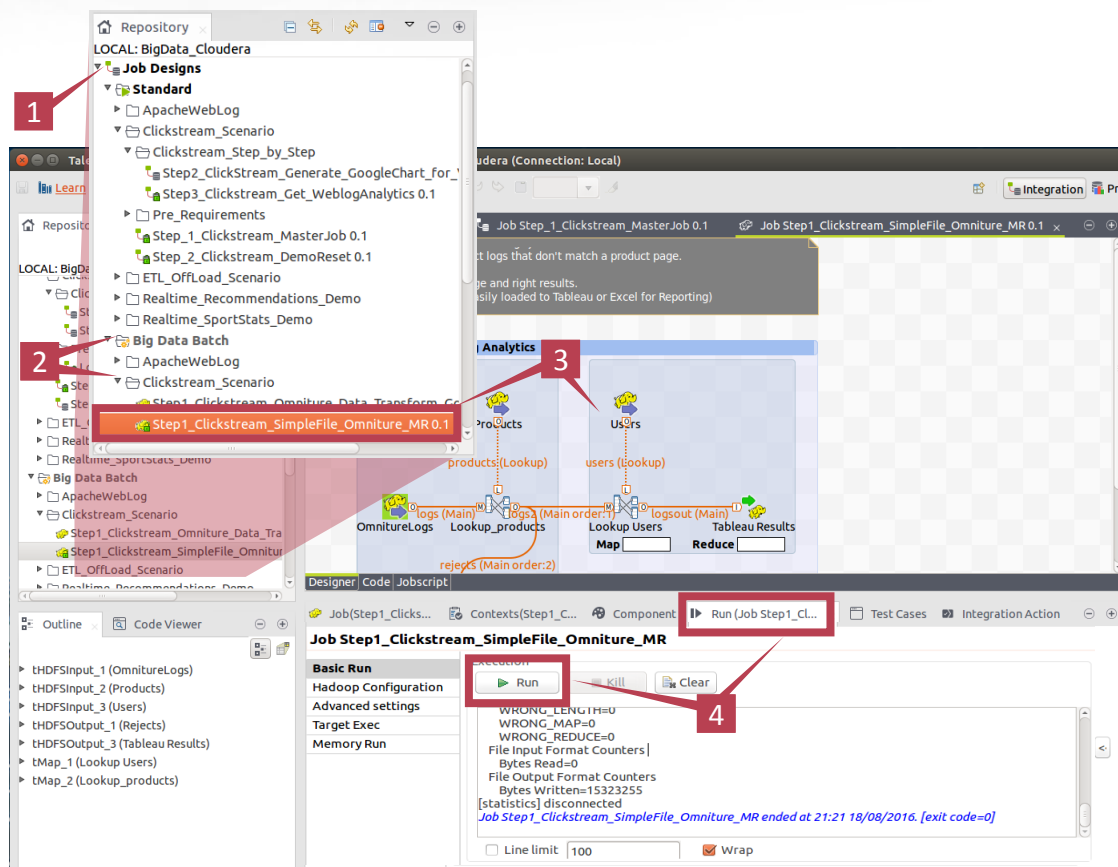| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |
|---|---|---|---|---|

## Execute the Clickstream Demo:

Additional analysis can be done to calculate the age and gender of users accessing specific links.

1. Navigate to the **Job Designs** folder.

2. Click on **Big Data Batch > Clickstream_Scenario**

3. Double click on **Step_1_Clickstream_SimpleFile_Omniture_MR 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute

The results of this job can be found in the HDFS File Browser:
*/user/talend/clickstream_demo/output/results*

With our analysis data in HDFS, we can load into a Hive Table for further querying or import to a visualization tool. Continue to the next steps of the demo to see how this can be done.

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|----------|----------------|----------------------|---------------------|-------------------|

| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |
|-----------------------|-------------|-------------|--------------|---------------|

## Execute the Clickstream Demo:

With our analysis complete, we can pull the raw file from HDFS or even put it into a Hive table for further querying.

1. Navigate to the **Job Designs** folder.

2. Navigate to **Standard > Clickstream_Scenario > Clickstream_Step_by_Step**

3. Double click on **Step3_Clickstream_Get_WeblogAnalytics 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute

The results of this job can be found on the local VM file system:
*/home/talend/Documents/Clickstream/webloganalytics.csv*

➢ This file could be imported to MS Excel or other BI tools like Tableau (not included in the Big Data Sandbox) to generate additional dashboards.

# Talend Big Data Sandbox
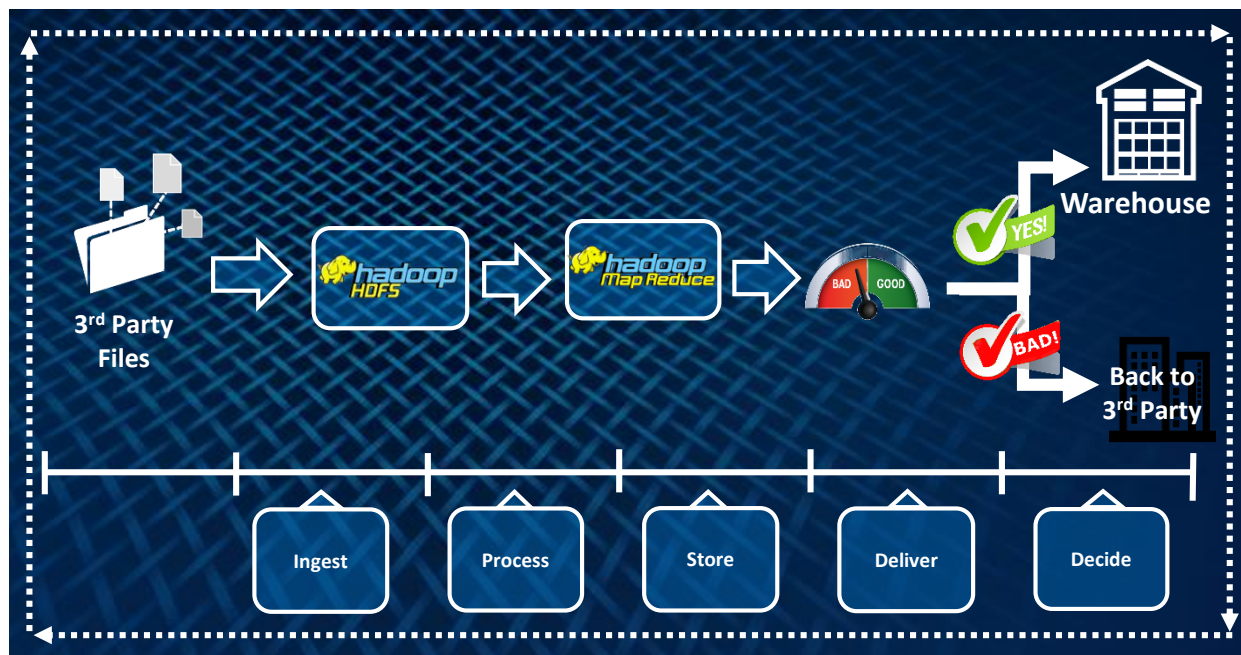## Big Data Insights Cookbook

**Note:** Execution of this demo requires a Hadoop distribution. If a distro hasn't been selected, <u>click here</u>.

## Overview:

In this example we demonstrate how using Talend with Hadoop can speed up and simplify processing large volumes of 3rd Party Data. The sample data is simulating a Life Sciences Prescribing habits data file from a 3rd Party vendor.

**You will experience:**

- optimizing your data warehouse by off-loading the ETL overhead to Hadoop and HDFS.

- Fast, Pre-load analytics on large volume datasets.

- Multiple Reports from same datasets to make informed and intelligent business decision that could decrease spend or increase revenue.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

**This Demo will highlight:**

### Large volume processing

With Talend and Hadoop, you can process Gigabytes and Terabytes of data in a fraction of the time.

### Pre-load Analytics

By analyzing large volumes of data BEFORE loading it to your Data Warehouse, you eliminate the overhead of costly data anomalies in the Data Warehouse.

### ETL Off-loading

Utilizing Talend with a Hadoop Cluster, you can optimize your Data Warehouse by removing the costly overhead of data processing.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|----------|----------------|----------------------|---------------------|-------------------|

| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |
|-----------------------|-------------|-------------|--------------|---------------|

**Note:** To quickly and easily see the value of the ETL Off-Load Demo, proceed with the below steps. If you would like a more in-depth experience and more control over the source data, [Click here…](#)

### Demo Setup:

To Execute this demo, you must first generate the source files for processing within Hadoop

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > ETL_OffLoad_Scenario**

3. Double click on **Step_1_ProductAnalysis_DemoSetup 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute

When this job completes, Source files will reside on the Virtual Machine to be processed by the demo. Additionally, within HDFS an initial report will have been generated by which the demo will compare for analysis.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|----------|----------------|----------------------|---------------------|-------------------|

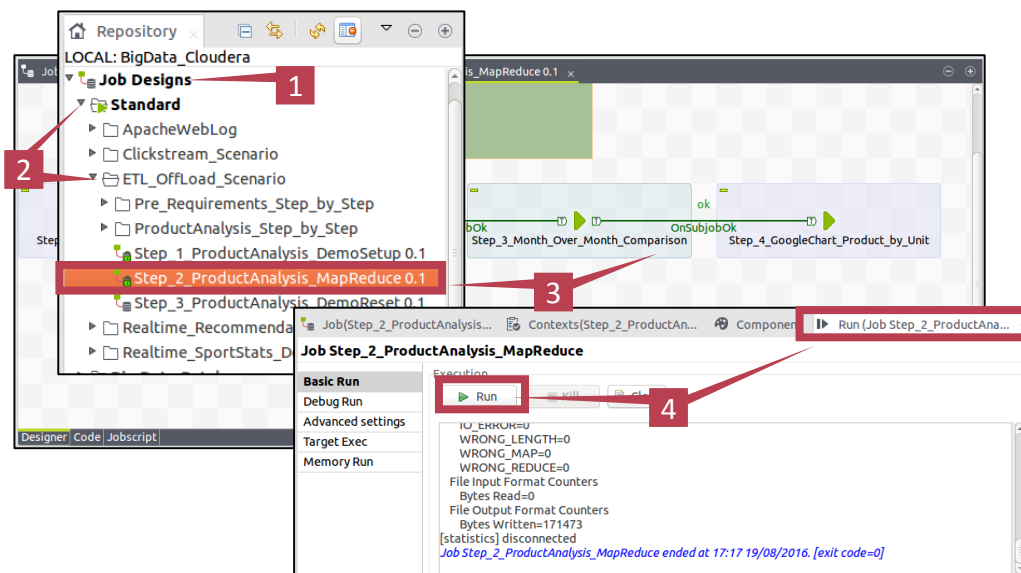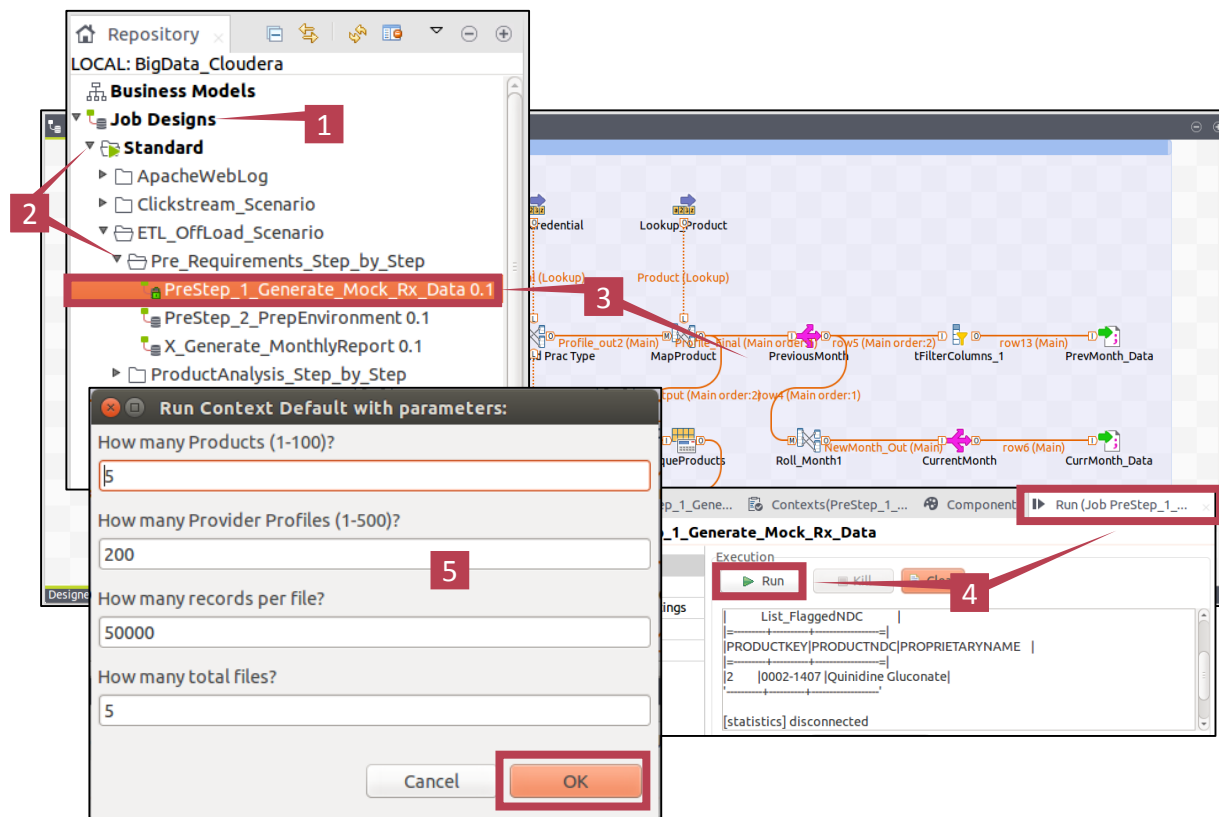| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |

### Execute the ETL Off-Load "One-Click" Demo:

In this "One-Click" version of the demo:
- Source files are placed in HDFS.
- MapReduce is used to collectively analyze all the compressed files.
- The resultant analysis is then compared to the previous months results and reports are generated.
- The generated reports are then sent to the Google Charts API for a graphical representation of the data.
- The resultant reports can be viewed in a web browser:
  - ✓ **Product by Physician** shows the number of prescriptions a physician has written for a particular drug
  - ✓ **Net Change** shows the total number of prescriptions for a particular drug across all physicians

1. Navigate to the **Job Designs** folder.

2. Click on **Standard >ETL_OffLoad_Scenario**

3. Double click on
   **Step_2_ProductAnalysis_MapReduce 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.



[Click Here](#) to finish this demo…

# Talend Big Data Sandbox
## Big Data Insights Cookbook

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|----------|----------------|----------------------|--------------------|--------------------|

| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |

## Demo Setup:

In this step-by-step version of the demo, you will see just how simple it is to work with Talend and Hadoop. You will also have more control over the source data used within the demo for a more personalized experience.

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > ETL_OffLoad_Scenario > Pre_Requirements_Step_by_Step**

3. Double click on **PreStep_1_Generate_Mock_Rx_Data 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.

5. When the job starts, **edit the values** as you wish *(staying within the suggested parameters and keeping in mind you are working in a virtual environment with limited space)* or leave the default values. Click **OK** when done.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

## Demo Setup (cont.):

Once you have generated your Mock Rx data, you will need to initialize the Hadoop environment with comparison data – in this case, it would be the "Previous Month" analysis.

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > ETL_OffLoad_Scenario > Pre_Requirements_Step_by_Step**

3. Double click on **PreStep_2_PrepEnvironment 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.

**Note:** Once this job completes, you are ready to execute the step-by-step ETL Off-Load Demo.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|----------|----------------|-----------------------|---------------------|-------------------|

| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |
|-----------------------|-------------|-------------|--------------|---------------|

**Execute the ETL Off-Load "Step-by-Step" Demo:**

With the demo environment setup complete, we can begin examining the ETL Off-load Process.

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > ETL_OffLoad_Scenario > ProductAnalysis_Step_by_Step**

3. Double click on **Step_ 1_PutFiles_on_HDFS 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.


When this job is complete, you will have your custom-generated source files on HDFS. To view the files:

5. Open **Firefox**

6. Click on the **HDFS Browser** link on the Bookmarks Toolbar

7. Select **Browse the file system** from the Utilities dropdown.

8. Navigate to */user/talend/ Product_demo/Input*

# Talend Big Data Sandbox
## Big Data Insights Cookbook

## Execute the ETL Off-Load "Step-by-Step" Demo:

Now that your source data is in HDFS, we can use the power of Hadoop and MapReduce to analyze the large dataset.

1. Navigate to the **Job Designs** folder.

2. Click on **Big Data Batch > ETL_OffLoad_Scenario**

3. Double click on
   **Step_2_Generate_MonthlyReport_mr 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.

When this job is complete, you can again navigate to the Hadoop file system to view the generated file:

5. Open **Firefox**

6. Click on the **HDFS Browser** link on the Bookmarks Toolbar

7. Select **Browse the file system** from the Utilities dropdown.

8. Navigate to */user/talend/Product_demo/Output*

**Execute the ETL Off-Load "Step-by-Step" Demo:**

With the Previous Month analysis as our baseline, we can now compare our Current Month analysis and track any anomalies.

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > ETL_OffLoad_Scenario > ProductAnalysis_Step_by_Step**

3. Double click on **Step_3_Month_Over_Month_Comparison 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.

When this job is complete, you can again navigate to the Hadoop file system to view the generated files:

5. Open **Firefox**

6. Click on the **HDFS Browser** link on the Bookmarks Toolbar

7. Select **Browse the file system** from the Utilities dropdown.

8. Navigate to */user/talend/Product_demo/Output*



## Directory

/user/talend/Product_demo/Output    Go!

| Permission | Owner | Group | Size | Replication | Block Size | Name |
|---|---|---|---|---|---|---|
| -rw-r--r-- | talend | talend | 167.45 KB | 3 | 128 MB | Current_Month_Report.txt |
| -rw-r--r-- | talend | talend | 50.64 KB | 3 | 128 MB | Previous_Month_Report.gz |
| -rw-r--r-- | talend | talend | 105 B | 3 | 128 MB | Product_NetChange_Report.csv |
| -rw-r--r-- | talend | talend | 10.27 KB | 3 | 128 MB | Product_Threshold_Report.gz |
| -rw-r--r-- | talend | talend | 7.56 KB | 3 | 128 MB | Product_by_Physician_Report.csv |
| drwxr-xr-x | talend | talend | 0 B | 0 | 0 B | Working |

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|---|---|---|---|---|

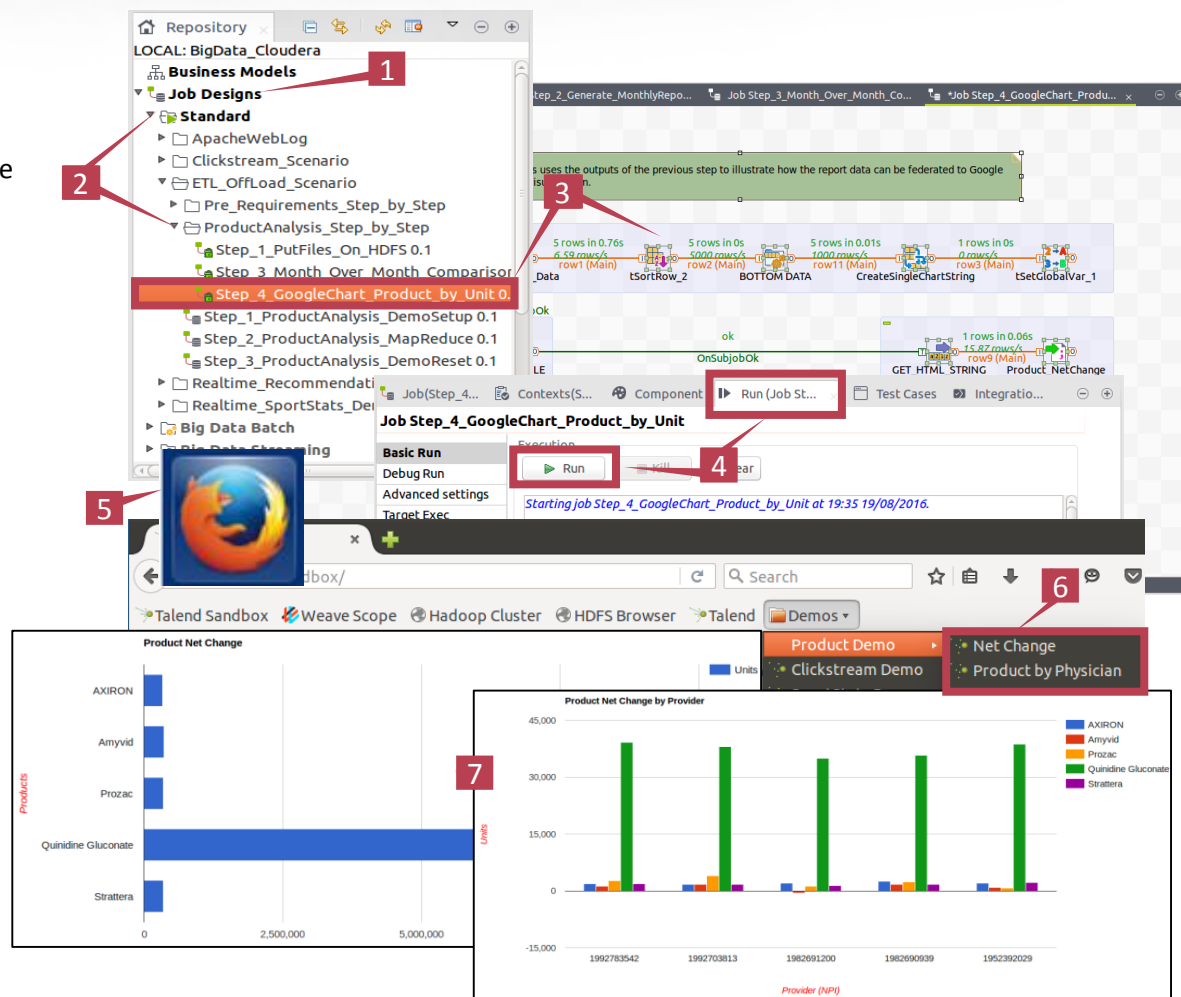| Retail Recommendation | Sport Stats | Clickstream | ETL Off-Load | Apache Weblog |
|---|---|---|---|---|

**Execute the ETL Off-Load "Step-by-Step" Demo:**

The final step is to generate the charts using Google Charts API

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > ETL_OffLoad_Scenario > ProductAnalysis_Step_by_Step**

3. Double click on **Step_4_GoogleChart_Product_by_Unit 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.

When this job is complete, you can view the generated reports from the web browser:

5. Open **Firefox**

6. From a new tab, Click on **Demos > Product Demo > Net Change** to view the report.

7. Repeat Step 2 above to open the **Product by Physician** Report.

## Execute the ETL Off-Load "Step-by-Step" Demo:

Reset the demo and run it again! You can run this demo over and over and get different results by changing the Source Files.

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > ETL_OffLoad_Scenario**

3. Double click on
**Step_3_ProductAnalysis_DemoReset 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute.

# Run through this demo again!

**One Click**

**Step-by-step**

# Talend Big Data Sandbox
## Big Data Insights Cookbook

**Note:** Execution of this demo requires a Hadoop distribution. If a distro hasn't been selected, <u>click here</u>.

## Overview:

In this example we demonstrate using different Big Data Methods to aggregate and analyze large volumes of weblog data.

**You will experience:**

- Using Hive to store and access Data in a Hadoop Distributed File System

- Using standard MapReduce to analyze and count IP addresses in an Apache log file.

- Performing the same Analysis (count of IP addresses in an Apache log file) using Pig.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

**This Demo will highlight:**

| **Hive Components** | **Native MapReduce** | **Pig Components** |
|---|---|---|
| Connect, Create, Read and Write with Hive Components to access data in HDFS | Use Talend's MapReduce components to access and analyze data, natively, from HDFS | Understand the flexibility of Talend's capabilities to perform the same operations with multiple technologies |

# Talend Big Data Sandbox
## Big Data Insights Cookbook

**Execute the Apache Weblog Demo:**

Create the Hive Tables in HDFS and clear out old datasets.

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > ApacheWebLog**

3. Double click on **Step_1_ApacheWebLog_HIVE_Create 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute

When this job completes, old datasets from previous executions will have been cleaned up and a fresh Hive Table will be generated in HDFS.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

**Execute the Apache Weblog Demo:**

Filter the Apache Weblog files and load them to HDFS.

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > ApacheWebLog**

3. Double click on **Step_2_ApacheWeblog_Load 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute

This job filters "301" codes from the Weblog and loads the data to HDFS where it can be viewed by both the HDFS file browser or a Hive Query.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

**Execute the Apache Weblog Demo:**

View the data in HDFS.

1. Open **Firefox**

2. Click on the bookmarked link titled **HDFS Browser**

3. In the **Utilities** Dropdown, select **Browse the File System** and navigate to */user/talend/weblog*

> **Note:** While the data is loaded into HDFS, it is also saved in a location where the created Hive table is expecting data. So now you can view the data both through a Hive query or HDFS file browsing.
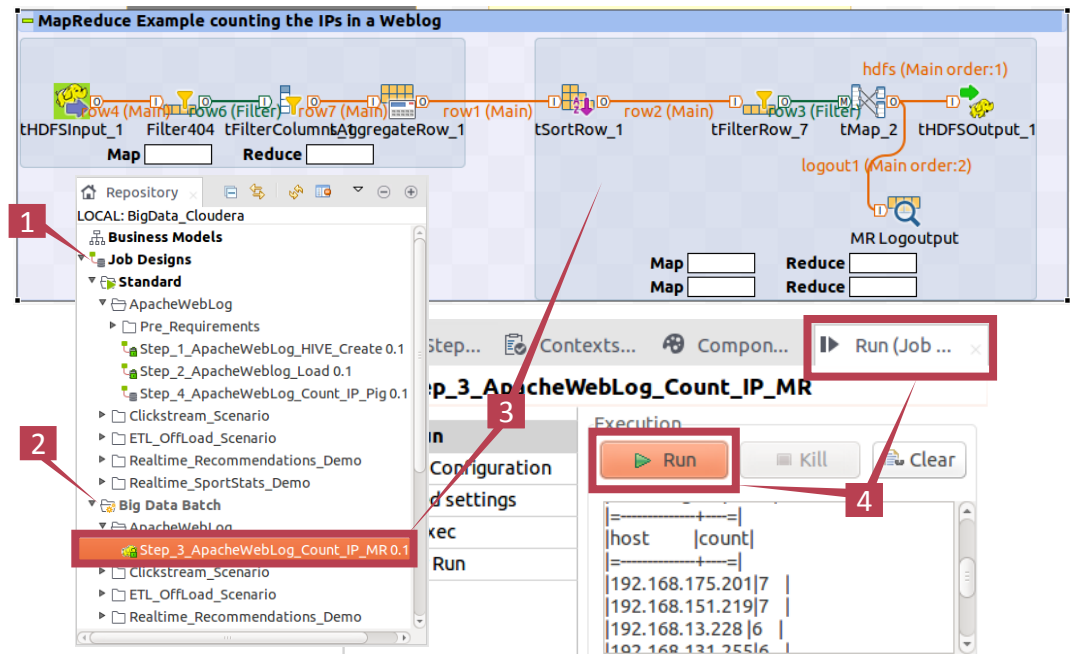
# Talend Big Data Sandbox
## Big Data Insights Cookbook

**Execute the Apache Weblog Demo:**

Use MapReduce to analyze and calculate distinct IP count.

1. Navigate to the **Job Designs** folder.

2. Click on **Big Data Batch > ApacheWebLog**

3. Double click on
   **Step_3_ApacheWeblog_Count_IP_MR 0.1.** This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute

5. When the job completes, You can view the results that are output to the Job Execution Window.



> ➢ The data from this job is also saved to HDFS. In the HDFS File Browser navigate to */user/talend/weblogMR/mr_apache_ip_out* to see the new files.

# Talend Big Data Sandbox

## Big Data Insights Cookbook

**Execute the Apache Weblog Demo:**

Use MapReduce to analyze and calculate distinct IP count.

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > ApacheWebLog**

3. Double click on
   **Step_4_ApacheWeblog_Count_IP_Pig 0.1.**
   This opens the job in the designer window.

4. From the Run tab, click on **Run** to execute

> ➤ The data from this job is also saved to HDFS. In the HDFS File Browser navigate to */user/talend/weblogPIG/apache_ip_cnt* to see the new files.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|----------|----------------|------------------------|----------------------|---------------------|

**Simple Conversion**

**This Bonus Demo will highlight:**

| Talend Flexibility | Simple Configuration | Future-Proof Technology |
|--------------------|----------------------|--------------------------|
| Switch between MapReduce and Spark Frameworks with just a few clicks of the mouse. | Once a job has been converted to a new Framework, configuration is quick and simple in Talend Studio. | Talend's Code-generation Architecture makes it easy to abreast of the latest Technology Trends. |

# Talend Big Data Sandbox
## Big Data Insights Cookbook

## Simple Conversion

This Demo takes a text version of a Talend Blog Post – Getting Started With Big Data – and does a simple word count. The Word Count example is a basic teaching tool for understanding how Map Reduce Technology works.

**Execute the Simple Conversion Demo:**

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > Simple_Conversion**

3. Double click on **Blog_WordCount 0.1.** This opens the job in the designer window.

4. From the **Run** tab, click on **Run** to execute

The output appears in the execution window, displaying the number of occurrences of each word within the blog post.

> ➢ Now lets see how quickly we can convert this job to a Map Reduce Job and then a Spark Job.
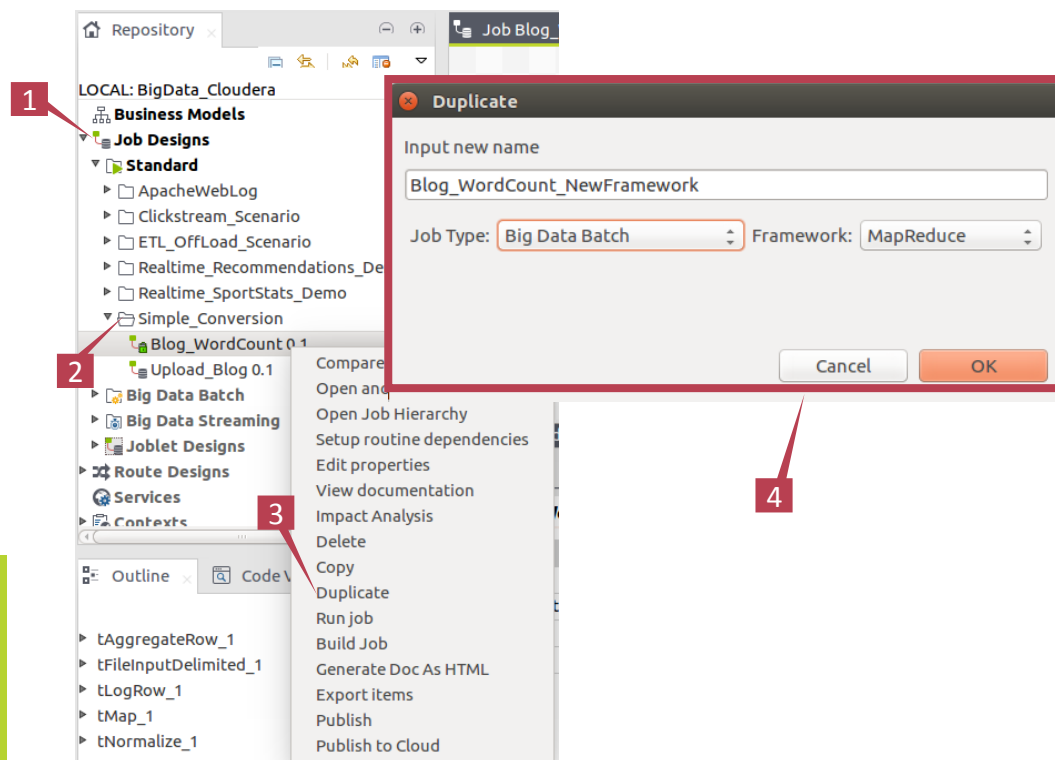
## Simple Conversion

**Converting to a Map Reduce Batch Job:**

1. Navigate to the **Job Designs** folder.

2. Click on **Standard > Simple_Conversion**

3. Right-click on **Blog_WordCount 0.1.** From the drop-down menu, select **Duplicate**.

4. Rename the job **Blog_WordCount_NewFramework** then in the Job Type dropdown, choose **Big Data Batch**. Finally, in the Framework dropdown, choose **MapReduce** and click **OK**.

   **Just that quickly, the standard job has been converted to a Big Data Batch job using the MapReduce Framework**

➢ Before you can run the job in MapReduce or Spark, you need to execute the **Upload_Blog** job in **Job Designs > Standard > Simple_Conversion**. Right-click on the job and select **Run Job** from the dropdown window.

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|---|---|---|---|---|

**Simple Conversion**

**Configure the new job and Execute:**

1. Navigate to the **Job Designs** folder.

2. Click on **Big Data Batch > Simple_Conversion**

3. Your newly converted job exists here. Double click on **Blog_WordCount_NewFramework 0.1** to open in the Design Window.

4. Next, In the Designer Window, click on the **tFileInputDelimited** component and then click on the **Component Tab**. Here, change the **Folder/File** property to *context.**HDFS_Dir** + context.File_Name*

# Talend Big Data Sandbox
## Big Data Insights Cookbook

| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|----------|----------------|-----------------------|---------------------|-------------------|

## Simple Conversion

### Configure the new job and Execute:

1. Now click on the **Run** tab and choose **Hadoop Configuration**.

2. Set the Property Type to **Repository** then click the **ellipsis** to pull up the Repository Content.

3. Click on the Hadoop Cluster Dropdown and then select **BigData_Hadoop** and click **OK**. This will import the preset Hadoop configuration into the job configuration.

4. Now Click on **Basic Run** and then **Run**. The job now executes on the Hadoop Cluster as a MapReduce Job.

# Talend Big Data Sandbox
## Big Data Insights Cookbook

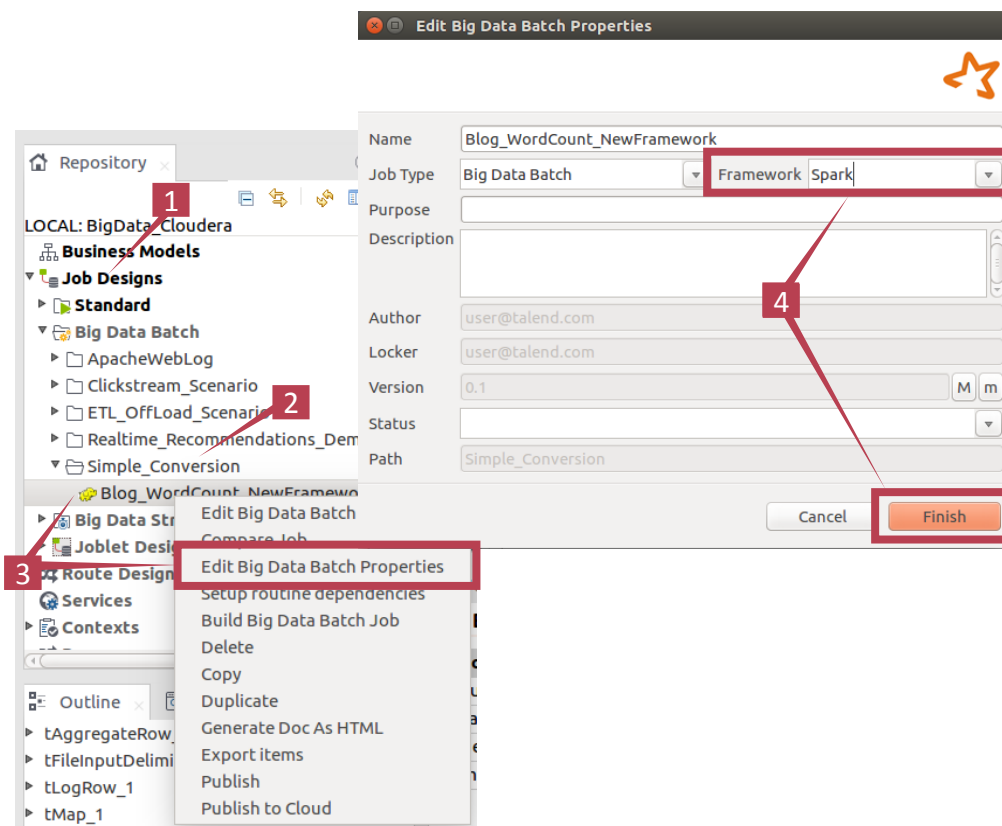| Overview | Pre-requisites | Setup & Configuration | Hadoop Distribution | Demo *(Scenario)* |
|----------|----------------|----------------------|---------------------|-------------------|

## Simple Conversion

> ➤ We can now take this same MapReduce job and convert it to execute on the Spark Framework.

**Convert to Spark Job:**

1. In the **Job Designs** folder.

2. Click on **Big Data Batch > Simple_Conversion**

3. Rather than duplicating the job again, this time we are going to Right-click on **Blog_WordCount_NewFramework 0.1** and select **Edit Big Data Batch Properties** from the dropdown menu.

4. In the Popup window, all you need to do is choose **Spark** from the Framework dropdown. Now click **OK**.

   **Just that quickly, the MapReduce job has now been converted to run on the Spark Framework.**

   **No further configuration needs to be done, you can open the job, click on the Run Tab and execute the job on Spark!**

# Talend Big Data Sandbox
## Big Data Insights Cookbook

## Conclusion

### Simplify Big Data Integration

Talend vastly simplifies big data integration, allowing you to leverage in-house resources to use Talend's rich graphical tools that generate big data code (Spark, MapReduce, PIG, Java) for you.

Talend is based on standards such as Eclipse, Java, and SQL, and is backed by a large collaborative community.

So you can up skill existing resources instead of finding new resources.

### Built for Batch and Real-time Big Data

Talend is built for batch and real-time big data. Unlike other solutions that "map" to big data or support a few components, Talend is the first data integration platform built on Spark with over 100 Spark components.

Whether integrating batch (MapReduce, Spark), streaming (Spark), NoSQL, or in real-time, Talend provides a single tool for all your integration needs.

Talend's native Hadoop data quality solution delivers clean and consistent data at infinite scale.

### Lower Operating Costs

Talend lowers operations costs.

Talend's zero footprint solution takes the complexity out of…

✓Integration Deployment
✓Management
✓Maintenance

A usage based subscription model provides a fast return on investment without large upfront costs.