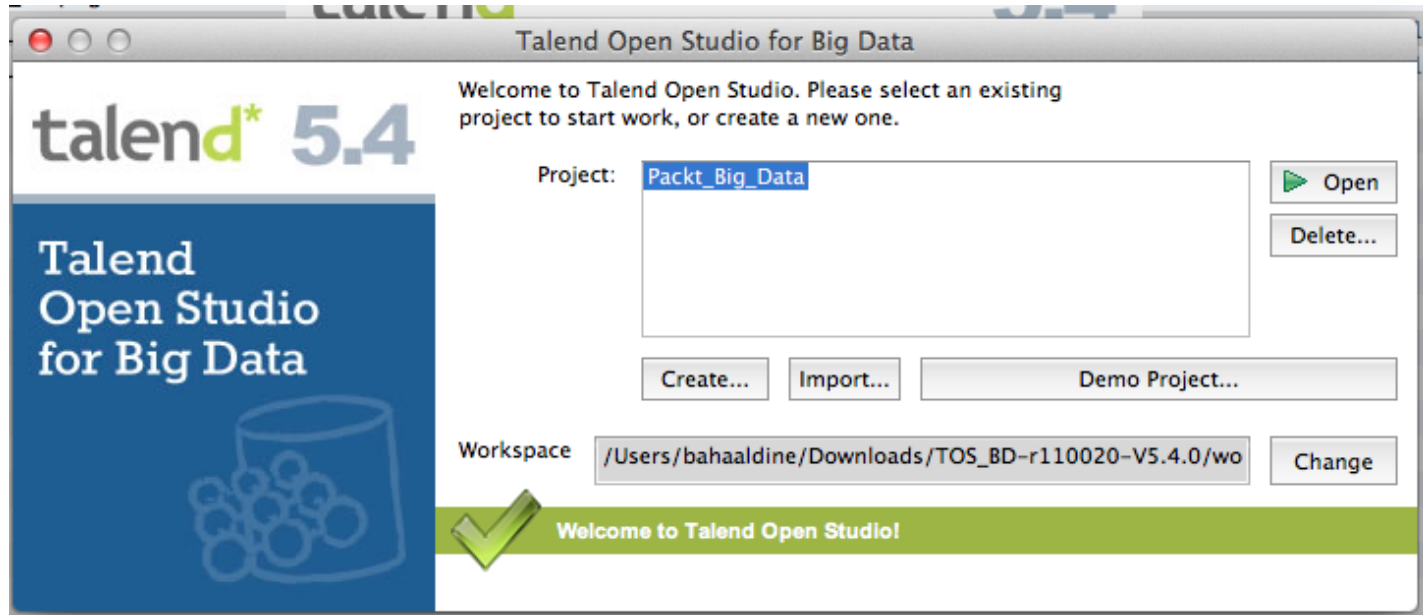# 1. Getting Started with Talend for Big Data

**New project**

Please enter the details for your new project below.

| Project Name | Packt_Big_Dat |
| Technical Name | PACKT_BIG_DAT |
| Project description | |

Cancel    Finish



Current Version    Other Releases    User Manuals

# Talend Open Studio for Big Data v5.4.0

**Download Now!**

Talend Unified Platform

Studio · Repository · Deployment · Execution · Monitoring



Lesson – 7

# Hadoop Ecosystem

**Ambari**
Provisioning, Managing and Monitoring Hadoop Clusters

**Sqoop** Data Exchange

**Flume** Log Collector

**Zookeeper** Coordination

**Oozie** Workflow

**Pig** Scripting

**Mahout** Machine Learning

**R Connectors** Statistics

**Hive** SQL Query

**Hbase** Columnar Store

**YARN Map Reduce v2**
Distributed Processing Framework

**HDFS**
Hadoop Distributed File System

**Note:** This is not an exhaustive list

http://www.facebook.com/hadoopers

## 2. Building Our First Big Data Job

Job(CH02_HDFS_WRITER )    Contexts(Job CH02_HDFS_WRITER )    Component ⊠    ⏵ Run (Job CH02_HDFS_WRITER

**tHDFSOutput_1**

| | |
|---|---|
| **Basic settings** | Property Type    Built–In ▼ |
| **Advanced settings** | Schema    Built–In ▼   Edit schema   ⋯   Sync columns |
| **Dynamic settings** | ☐ Use an existing connection |
| **View** | Version |
| **Documentation** | Distribution    Cloudera ▼ * Hadoop version   Cloudera CDH4.3+(YARN mode) ▼ * |

Version

Distribution     Cloudera ▼ * Hadoop version   Cloudera CDH4.3+(YARN mode) ▼ *

Connection

NameNode URI    "hdfs://172.16.253.202:8020/"

Authentication

☐ Use kerberos authentication
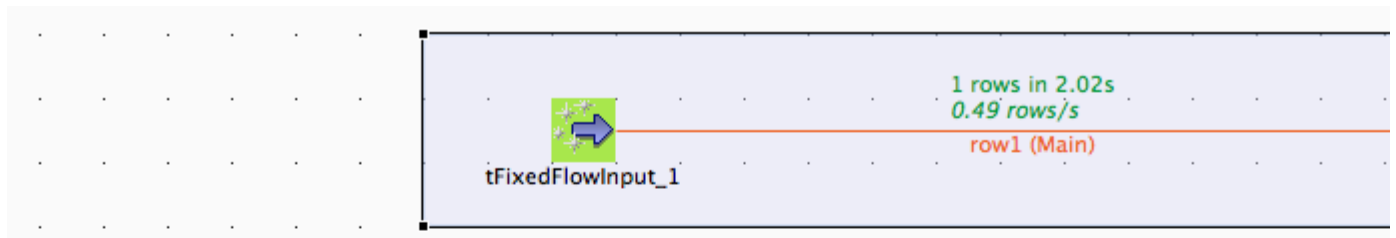
User name     context.username

File Name     "/user/bahaaldine/packt/chp01/.init"

File Type

Type     Text File ▼ *

Action     Create ▼

Row Separator     "\n"     * Field Separator   ";"

1 rows in 2.02s
*0.49 rows/s*
row1 (Main)
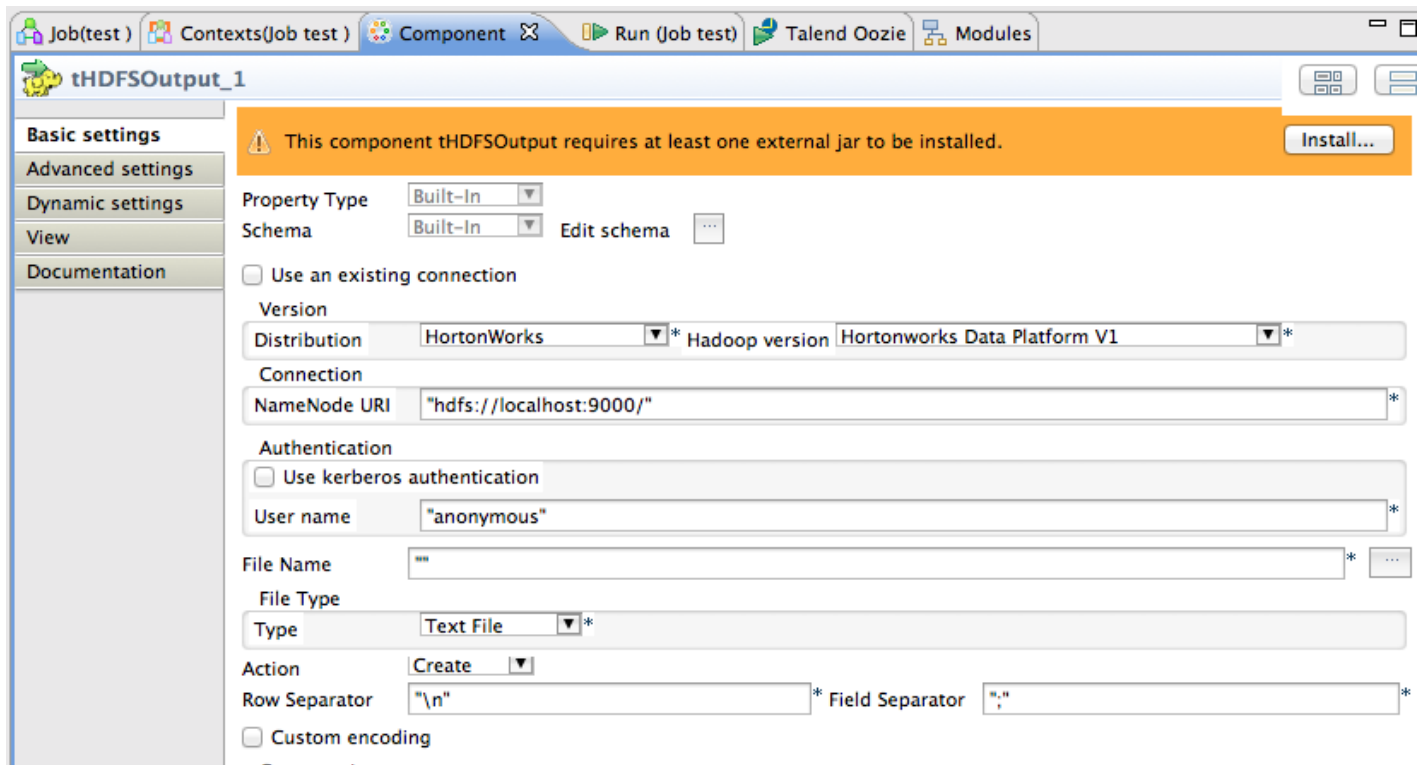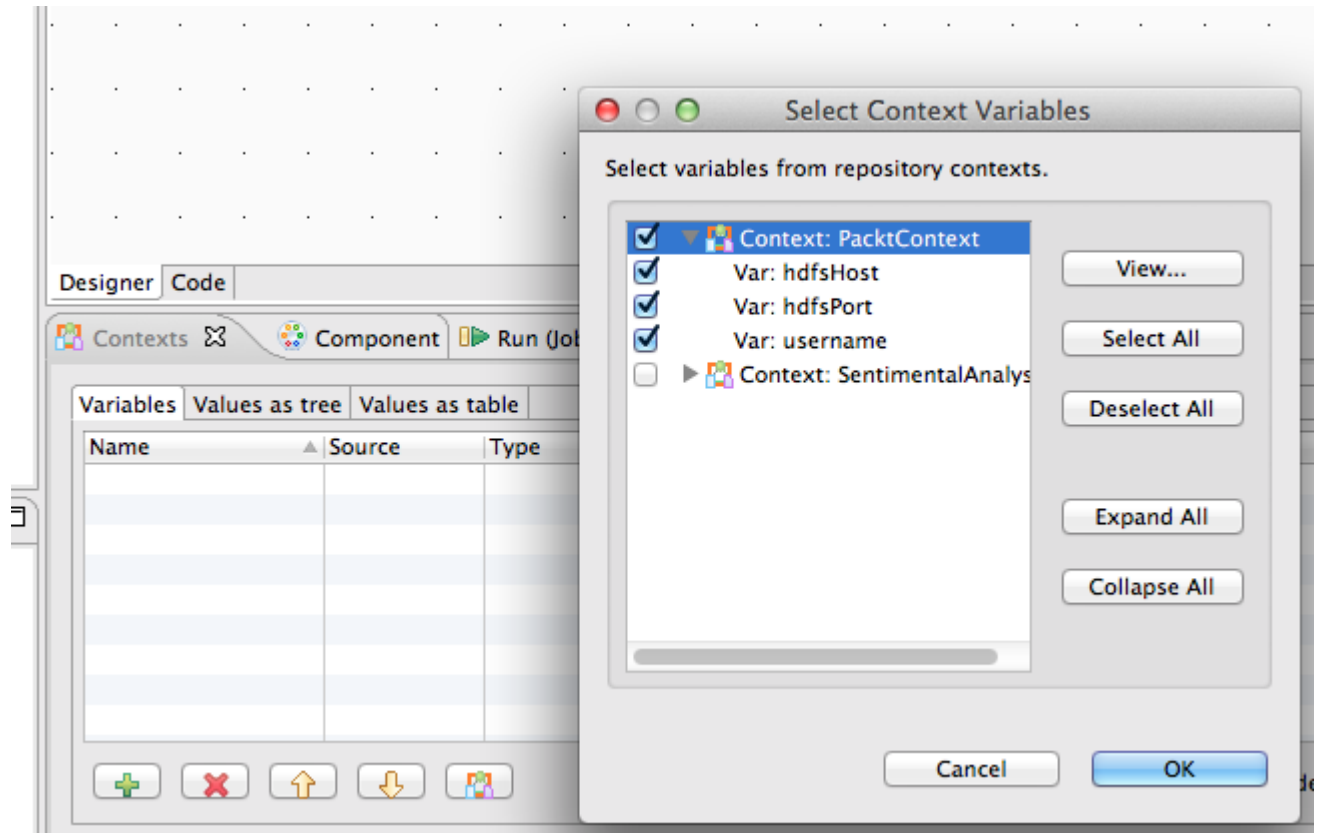
tFixedFlowInput_1

Code

exts  Component  Run (Job CH01_HDFS_WRITER)  ⊠  Oozie(CH01_HDFS_WRITER)  Modules
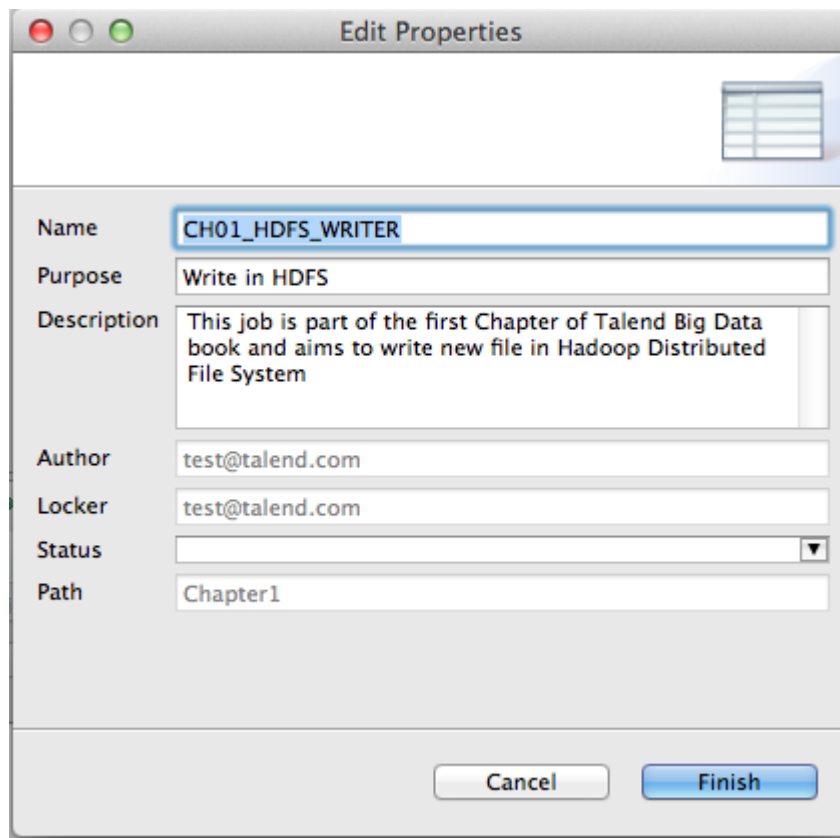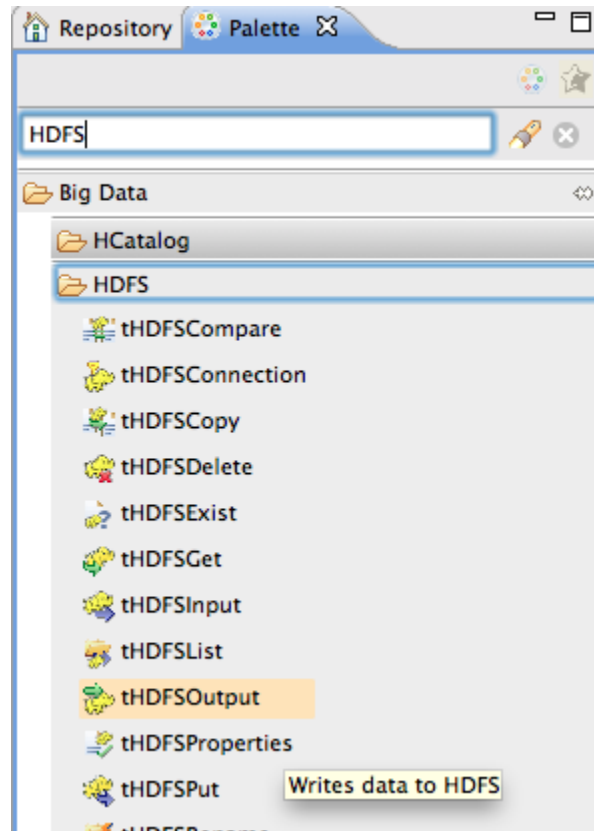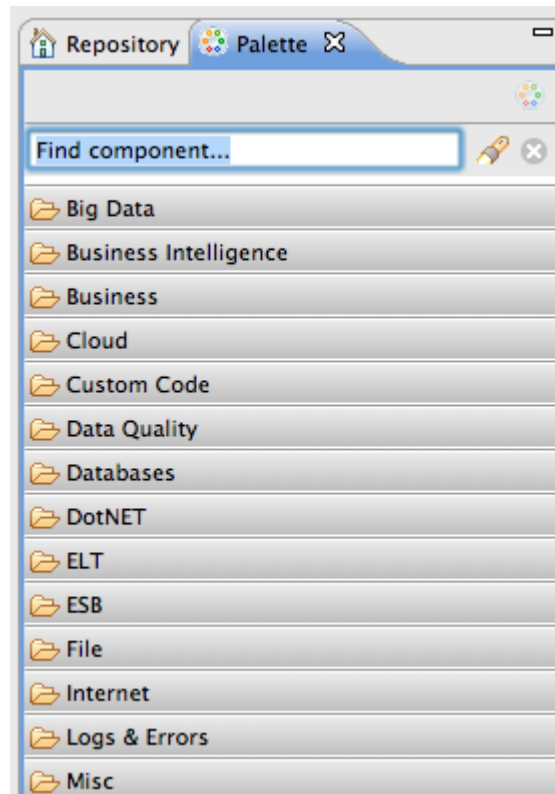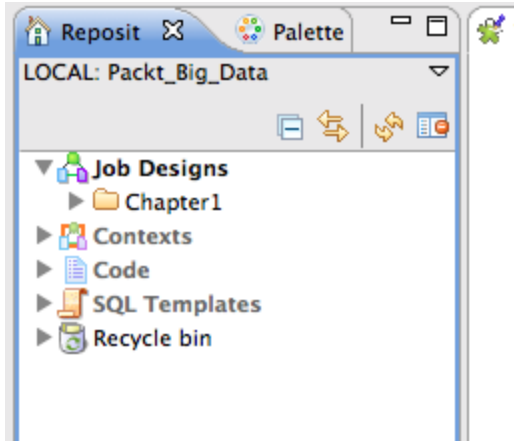
01_HDFS_WRITER

Execution

Run    Kill    Clear

```
Starting job CH01_HDFS_WRITER at 11:21 20/11/2013.


[statistics] connecting to socket on port 3341
[statistics] connected
[WARN ]: org.apache.hadoop.conf.Configuration - fs.default.name is deprecated. Instead, use
fs.defaultFS
2013-11-20 11:21:37.309 java[30590:f07] Unable to load realm info from SCDynamicStore
[WARN ]: org.apache.hadoop.util.NativeCodeLoader - Unable to load native-hadoop library for
your platform... using builtin-java classes where applicable
[statistics] disconnected

Job CH01_HDFS_WRITER ended at 11:21 20/11/2013. [exit code=0]
```
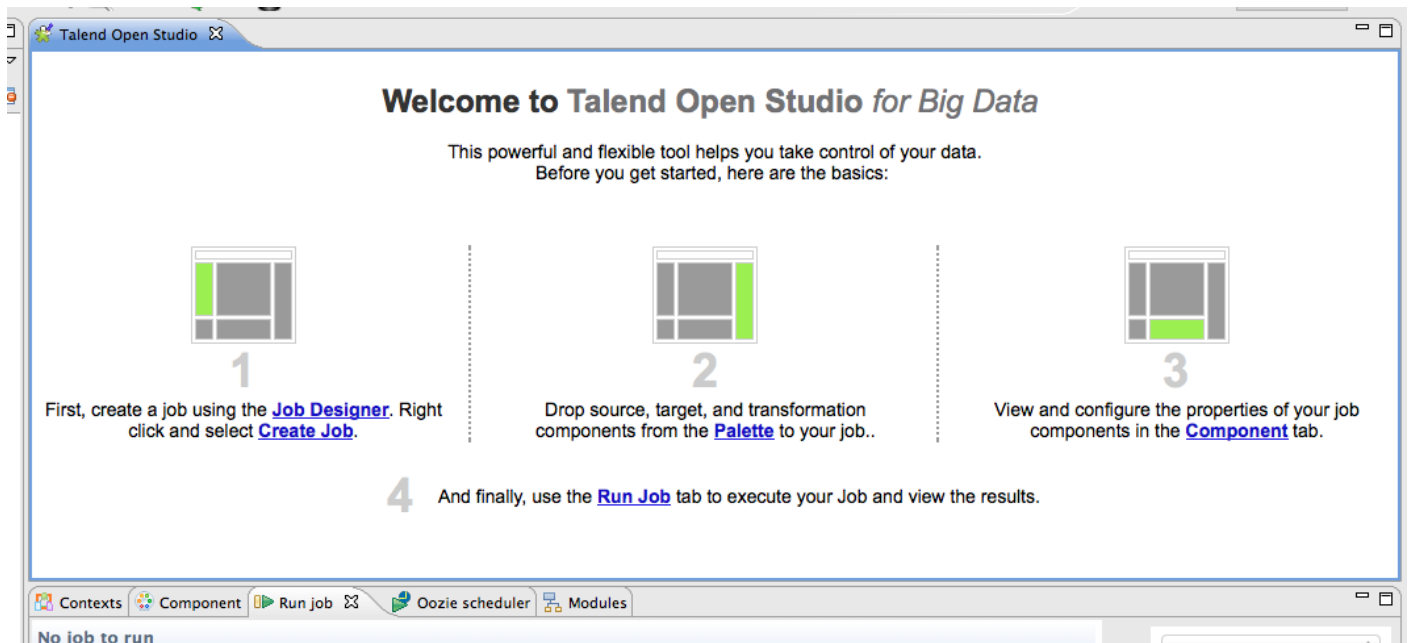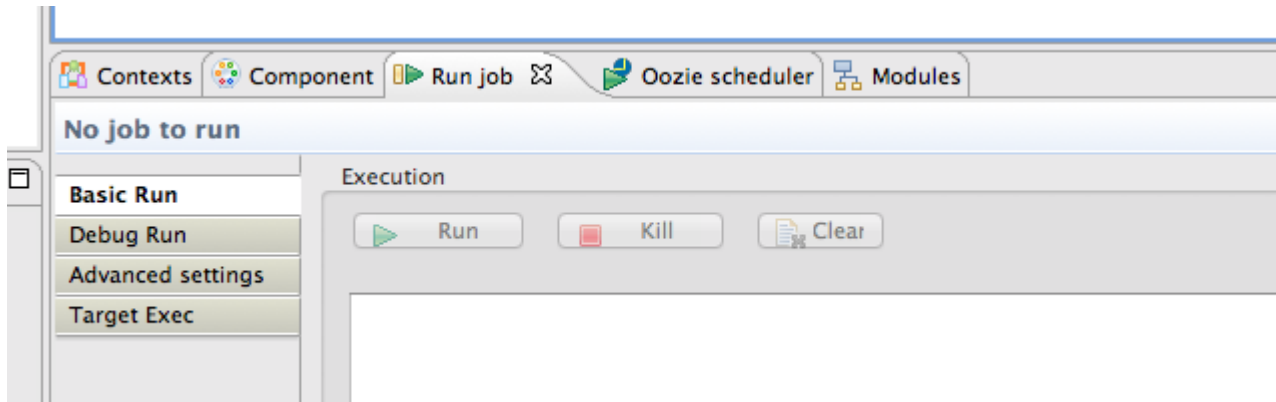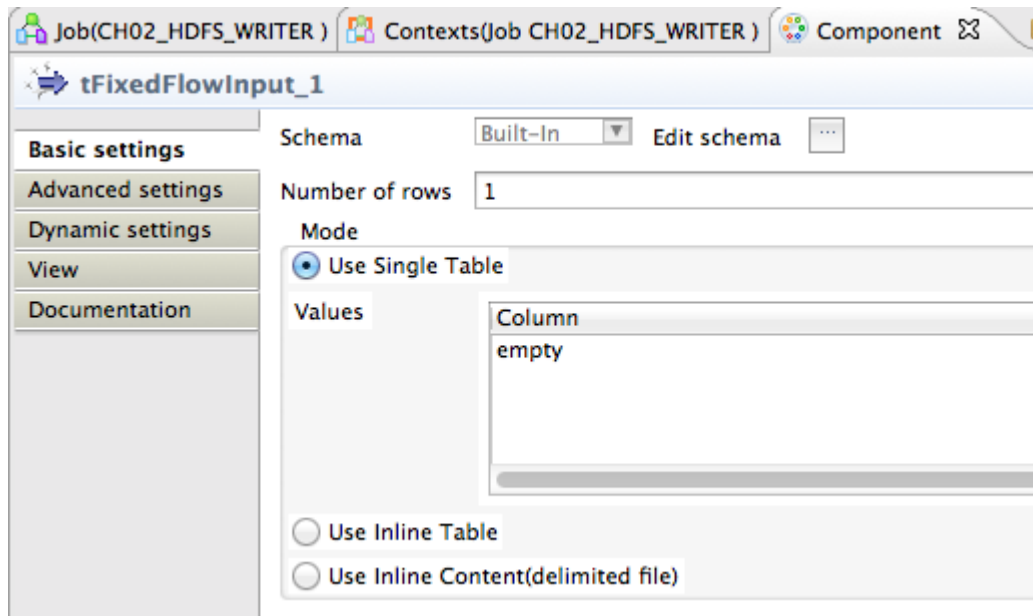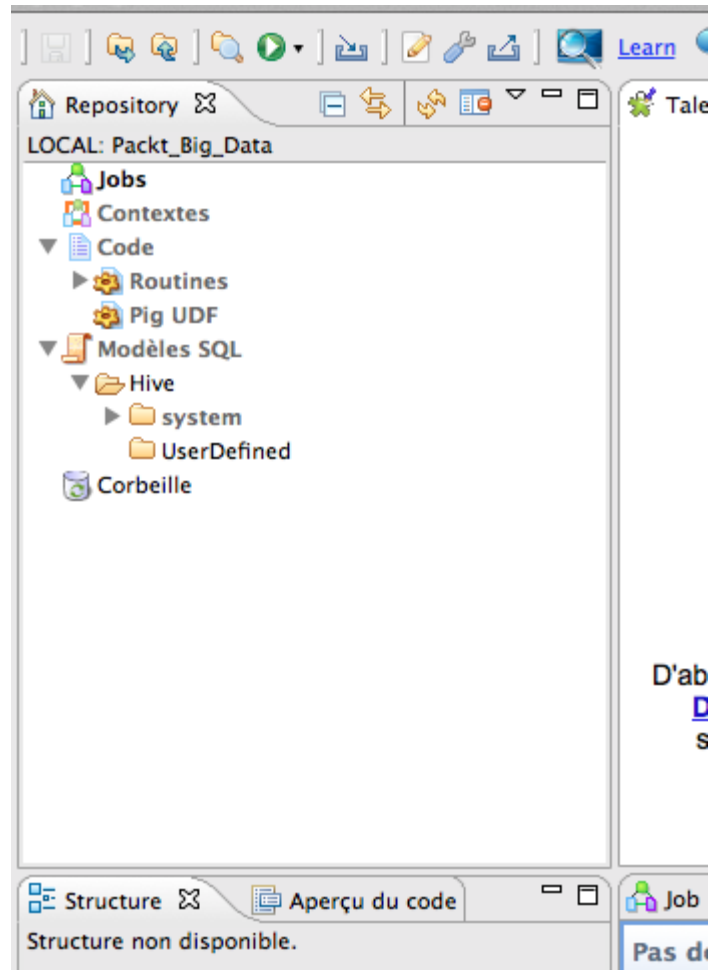
un
Run
ed settings
Exec

tFixedFlowInput_1                tHDFSOutput_1

## Select Context Variables

Select variables from repository contexts.

- ☑ ▼ Context: PacktContext
  - ☑ Var: hdfsHost
  - ☑ Var: hdfsPort
  - ☑ Var: username
- ☐ ▶ Context: SentimentalAnalys

[View...]  [Select All]  [Deselect All]  [Expand All]  [Collapse All]

[Cancel]  [OK]

---

**Designer** | **Code**

🔲 Contexts ✕ | 🔲 Component | ▶ Run (Job

**Variables** | Values as tree | Values as table

| Name ▲ | Source | Type |
|--------|--------|------|

---

🔲 Job(test ) | 🔲 Contexts(Job test ) | 🔲 Component ✕ | ▶ Run (Job test) | Talend Oozie | Modules

### tHDFSOutput_1

- Basic settings
- Advanced settings
- Dynamic settings
- View
- Documentation

⚠ This component tHDFSOutput requires at least one external jar to be installed.   [Install...]

Property Type    Built-In ▼
Schema           Built-In ▼    Edit schema  [...]

☐ Use an existing connection

**Version**
Distribution     HortonWorks ▼ *  Hadoop version  Hortonworks Data Platform V1 ▼ *

**Connection**
NameNode URI     "hdfs://localhost:9000/"    *

**Authentication**
☐ Use kerberos authentication
User name        "anonymous"    *

File Name        ""    *  [...]

**File Type**
Type             Text File ▼ *
Action           Create ▼
Row Separator    "\n"    *  Field Separator  ";"    *
☐ Custom encoding

## Repository | Palette ☒

HDFS

📁 **Big Data**                                              ◇

    📁 HCatalog

    📂 HDFS

        🔧 tHDFSCompare

        🐸 tHDFSConnection

        🔧 tHDFSCopy

        🐸 tHDFSDelete

        🐸 tHDFSExist

        🐸 tHDFSGet

        🐸 tHDFSInput

        🐸 tHDFSList

        🐸 tHDFSOutput

        📄 tHDFSProperties

        🐸 tHDFSPut      Writes data to HDFS

        tHDFSRename

---

## Edit Properties

| | |
|---|---|
| Name | CH01_HDFS_WRITER |
| Purpose | Write in HDFS |
| Description | This job is part of the first Chapter of Talend Big Data book and aims to write new file in Hadoop Distributed File System |
| Author | test@talend.com |
| Locker | test@talend.com |
| Status | ▼ |
| Path | Chapter1 |

Cancel    Finish

| Contexts | Component | Run job ⊠ | Oozie scheduler | Modules |

**No job to run**

| Basic Run |
| Debug Run |
| Advanced settings |
| Target Exec |

Execution

▶ Run    ■ Kill    Clear



Talend Open Studio ⊠

# Welcome to Talend Open Studio *for Big Data*

This powerful and flexible tool helps you take control of your data.
Before you get started, here are the basics:

**1**

First, create a job using the **Job Designer**. Right click and select **Create Job**.

**2**

Drop source, target, and transformation components from the **Palette** to your job..

**3**

View and configure the properties of your job components in the **Component** tab.

**4**

And finally, use the **Run Job** tab to execute your Job and view the results.

| Contexts | Component | Run job ⊠ | Oozie scheduler | Modules |

**No job to run**

# 3. Formatting Data



| Column | Db Column | Key | DB Type |
|---|---|---|---|
| day_of_week | day_of_week | | STRING |
| mont | mont | | STRING |
| day_of_month | day_of_month | | STRING |
| time | time | | STRING |
| zone | zone | | STRING |
| year | year | | STRING |
| username | username | | STRING |
| content | content | | STRING |

```
hive> desc tweets
    > ;
OK
day_of_week     string
mont    string
day_of_month    string
time    string
zone    string
year    string
content string
hours   int
Time taken: 0.351 seconds
```

**Schema of parse tweets**

parse tweets

| Column | Key | Type | Nullab | Date Pattern ( | Length | Precisior | Defaul | Commer |
|--------|-----|------|--------|----------------|--------|-----------|--------|--------|
| day_of_week | ☐ | String | ☑ | | 3 | 0 | | |
| month | ☐ | String | ☑ | | 4 | 0 | | |
| day_of_month | ☐ | String | ☑ | | 3 | 0 | | |
| time | ☐ | String | ☑ | | 9 | 0 | | |
| zone | ☐ | String | ☑ | | 4 | 0 | | |
| year | ☐ | String | ☑ | | 5 | 0 | | |
| content | ☐ | String | ☑ | | 157 | 0 | | |



**parse tweets(tFileInputPositional_1)**

| | |
|---|---|
| **Basic settings** | |
| Advanced settings | |
| Dynamic settings | |
| View | |
| Documentation | |

Property Type [Built-In ▼]  💾

☐ Use existing dynamic

File name/Stream [context.tweetFile]

Row Separator ["\n"]

☐ Use byte length as the cardinality

☐ Customize  Pattern ["3,4,3,9,5,5,*"]

☐ Skip empty rows  ☐ Uncompress as zip file

☐ Die on error

Header [0]  Footer

Schema [Built-In ▼]  Edit schema [...]

## 4. Processing Tweets with Apache Hive



connection to hive  →OnSubjobOk→  drop table sentiments  →OnSubjobOk→  create the table sentiments  →OnSubjobOk→  build the sentiments

build the sentiments(tHiveRow_3)

build the sentiments(tHiveRow_

| Paramètres simples | ☑ Utiliser une connexion existante   Liste des composants   tHiveConnection_1 – connection to hive ▼ |
|---|---|
| Advanced settings | Schéma   Built-In ▼   Editer le schéma  …   Nom de la table  "" |
| Paramètres dynamiques | Type de requête   Built-In ▼   Guess Query |
| View | Requête   "insert into table sentiments select hash_tags.hash_tags_label, hash_tags.time, emoticons.emoticons_label from hash_tags INNER JOIN emoticons ON hash_tags.hash_tags_id == emoticons.emoticons_id" |
| Documentation | |

connection to hive → OnSubjobOk → add UDF jar into the hive runtime classpath → OnSubjobOk → register the UDF → OnSubjobOk → drop table emoticons → OnSubjobOk → create the table emoticons → OnSubjobOk → feed the emoticons

feed the emoticons(tHiveRow_5)

| Paramètres simples | ☑ Utiliser une connexion existante   Liste des composants   tHiveConnection_1 – connection to hive ▼ |
|---|---|
| Advanced settings | Schéma   Built-In ▼   Editer le schéma  …   Nom de la table  "" |
| Paramètres dynamiques | Type de requête   Built-In ▼   Guess Query |
| View | Requête   "insert into table emoticons |
| Documentation | select concat(tweets.day_of_week, tweets.day_of_month, tweets.time, tweets.mont) as emoticons_id, tweets.day_of_week, tweets.day_of_month, substr(tweets.time,1,2), tweets.mont, emoticons_label from tweets LATERAL VIEW explode( extract_pattern(tweets.content, '((?:\\;|:|=|8|\\\\^|X|:\\'|<|>)(?:-|O|c|\\\\.|_|\\\\^) ?(?:D|P|<|>|\\\\)\\\\)\\\\(\\\\(\\\\|\\\\|\\\\;|= |X|O|\\\\*|S|\\\\|\\\\|{|\\\\}|\\\\([\\\\]|\\\\(|\\\\)))' ) ) emoticonTable as emoticons_label" |

connection to hive → OnSubjobOk → add UDF jar into the hive runtime classpath → OnSubjobOk → register the UDF → OnSubjobOk → drop table hash_tags → OnSubjobOk → create the table hash_tags → OnSubjobOk → feed the hash tags

Designer | Code

add UDF jar into the hive runtime classpath(tHiveRow_1)

| Basic settings | ☑ Use an existing connection   Component List   tHiveC |
|---|---|
| Advanced settings | Schema   Built-In ▼   Edit schema  …   Tab |
| Dynamic settings | Query Type   Built-In ▼   Guess Query |
| View | Query   "add jar "+context.custom_udf_jar |
| Documentation | |

"Orange Ball in the Sky- fullmoon at dawn| #photography #Fullmoon http://t.co/pZylCf1XtB"

**Custom UDF
Extract_pattern**

**Hive UDF
explode**

[ #photography, #Fullmoon ]

#photography
#Fullmoon

## 5. Aggregating Data with Apache Pig

## Top Twitters Timeline



SkyF1GP (17)  candyharreh (4)
MaybeKerr (6)  dyelahott (6)  uiharu2351 (10)  fas_almscm (6)  SportsAB (7)
rllyzaynrlly (7)  raquela_1D (7)  GaFraBo_ (21)  dbfla2505 (4)  MacguiverFilms (6)
Robert_Testone (5)  Tasiaps67 (7)  hipsta_chapel (3)  PPEMR (28)  ThompsongkgJazm (5)
ORWaitTime (11)  SpeedTracs (16)  our_cars (53)  SwisherGangEnt (5)  SEVENCLOVERS (21)  WandaDeosa (

00:00  04:00  08:00  12:00  16:00  20:00



▼ Jobs
  ▶ Chapter2
  ▶ Chapter3
  ▶ Chapter4
  ▼ Chapter5
      CH05_01_PIG_TOP_TWITTERS
      CH05_02_PIG_TOP_HASH_TAGS
      CH05_03_PIG_TOP_EMOTICONS
      CH05_04_PIG_TOP_SENTIMENTS



load formatted sentiments stores in json — row1 (Pig) → group and count the sentiments — row3 (Pig) → sort the top sentiments — row4 (Pig) → store the top sentiments



### Schema of filter on the hash tag and the time

filter on the hash tag and the time (Output)

| ng | Preci | Defa | Com | | Column | Key | Type | ✓ |
|---|---|---|---|---|---|---|---|---|
| | 0 | | | | hash_tags_label | ☐ | String | |
| | 0 | | | | time | ☐ | String | |
| | 0 | | | | | | | |
| | 0 | | | | | | | |
| | 0 | | | | | | | |

load the formatted hash tags    row1 (Pig)    filter on the hash tag and the time    row2 (Pig)    group and count the hash tags    row3 (Pig)    get top hash tags    row4 (Pig)    store the top hash tags



load the formatted tweets    row1 (Pig)    filter on the twitters username and the time    row2 (Pig)    Group and count the tweets    row3 (Pig)    get the top twitters    row6 (Pig)    sort twitters    row4 (Pig)    store twitters



05_01_PIG_TOP_TWITTERS )   Component   Run (Job CH05_01_PIG_TOP_TWITTERS)

rs(tPigStoreResult_1)

| | |
|---|---|
| Property Type | Built-In |
| Schema | Built-In   Edit schema   ···   Sync columns |
| Result Folder URI | "/user/"+context.username+"/packt/chp05/twitters" |
| ☑ Remove result directory if exists | |
| Store function | PigStorage ▼ * |
| Field separator | ";" |



PigSort_1)

| | |
|---|---|
| Schema | Built-In ▼   Edit schema   ···   Sync columns |

Sort key

| Column | Order |
|---|---|
| username_count | DESC |

| Schema | Built-In ▼ | Edit schema | ⋯ | Sync columns |

**Group by**

| Column |
|---|
| username |
| time |

[➕] [✖] [⬆] [⬇] [📄] [📋]

**Operations**

| Additional Output Column | Function | Input Column |
|---|---|---|
| username_count | max | username_count |

[➕] [✖] [⬆] [⬇] [📄] [📋] [➕]

---

🗒 **Group and count the tweets(tPigAggregate_1)**

| **Basic settings** |
|---|
| Advanced settings |
| Dynamic settings |
| View |
| Documentation |

| Schema | Built-In ▼ | Edit schema | ⋯ | Sync columns |

**Group by**

| Column |
|---|
| username |
| time |

[➕] [✖] [⬆] [⬇] [📄] [📋]

**Operations**

| Additional Output Column | Function | Input Column |
|---|---|---|
| username_count | count | username |

[➕] [✖] [⬆] [⬇] [📄] [📋] [➕]

## Schema of filter on the twitters username and the time

**load the formatted tweets (Input – Pig)**

| Column | Key | Type | Nullab | Date Patte | Length | Precis | Defa | Comm |
|---|---|---|---|---|---|---|---|---|
| day_of_week | ☐ | String | ☑ | | 5 | 0 | | |
| mont | ☐ | String | ☑ | | 5 | 0 | | |
| day_of_mont | ☐ | String | ☑ | | 2 | 0 | | |
| time | ☐ | String | ☑ | | 10 | 0 | | |
| zone | ☐ | String | ☑ | | 5 | 0 | | |
| year | ☐ | String | ☑ | | 5 | 0 | | |
| username | ☐ | String | ☑ | | 30 | 0 | | |
| content | ☐ | String | ☑ | | 150 | 0 | | |

**filter on the twitters username and the time (Output)**

| Column | Key | Type | Nullab | Date Patte | Length | Precis | Defa | Comm |
|---|---|---|---|---|---|---|---|---|
| username | ☐ | String | ☑ | | 30 | 0 | | |
| time | ☐ | String | ☑ | | 10 | 0 | | |

OK    Cancel

▼ 📄 Code
   ▶ 🔧 Routines
   ▼ 🔧 Pig UDF
      📄 MyJSONLoad
      📄 MyJSONStore
▼ 📄 SQL Templates
   ▶ 📁 Hive
▶ 🗑 Recycle bin

row1 (Pig)

filter on the twitters username and the time

_TOP_TWITTERS )    Contexts(Job CH05_01_PIG_TOP_TWITTERS )    Com

twitters username and the time(tPigFilterColumns_1)

Schema    Built-In ▼    Edit schema  ⋯  Sync columns

Operations

| Additional Output Column | Function | Input Column |
|---|---|---|
| username_count | count ▼ | username |
| | count | |
| | avg | |
| | sum | |
| | max | |
| | min | |

# 6. Using the SQL Database

export the top twitters from hdfs to mysql

twitters from hdfs to mysql(tSqoopExport_1)

**Version**

| Distribution | Cloudera ▼ * Hadoop version | Cloudera CDH4.X (MR 1 mode) ▼ * |

**Configuration**

| NameNode URI | "hdfs://"+context.hdfs_host+":"+context.hdfs_port |
| JobTracker Host | context.JT_host+":"+context.JT_port |

**Authentication**

☐ Use kerberos authentication

| Hadoop user name | "" |

| Connection | "jdbc:mysql://"+context.mysql_host+":"+context.mysql_port+"/ | Table Name | context.mysql_twitters_table |
| Export Dir | context.sqoop_twitters_export_dir |
| Username | context.mysql_user |
| Password | context.mysql_password |

☐ Specify Number of Mappers

☐ Print Log

---

**tSqoopImport_1**

| Basic settings |
| Advanced settings |
| Dynamic settings |
| View |
| Documentation |

**Mode**
- ⦿ Use Commandline
- ○ Use Java API

| Connection | "jdbc:mysql://"+context.mysql_host+":"+context.mysql_port+"/"+context.mysql_database |
| Username | context.mysql_user |
| Password | context.mysql_password |
| Table Name | context.mysql_twitters_table |

☐ Append

| File Format | textfile ▼ |

☐ Compress

☑ Print Log  ☐ Verbose

load the mysql jdbc driver

export the top twitters from hdfs to mysql



Contexts(Job CH06_01_SQOOP_EXPORT_TWI    Component ⊠    Run (Job CH06_01_SQOOP_EXPORT_TWITTE

load the mysql jdbc driver(tLibraryLoad_1)

Basic settings

Advanced settings

Dynamic settings

View

Documentation

Library        mysql-connector-java-5.1.22-bin.jar        ▼ *    ...



Repository    Palette ⊠

sqoop

Big Data

Sqoop

tSqoopExport

tSqoopImport

tSqoopImportAllTables

tSqoopMerge

# 7. Big Data Architecture and Integration Patterns



```
9 /user/bahaaldine/data/00
9 /user/bahaaldine/data/01
9 /user/bahaaldine/data/02
9 /user/bahaaldine/data/03
9 /user/bahaaldine/data/04
9 /user/bahaaldine/data/05
9 /user/bahaaldine/data/06
9 /user/bahaaldine/data/07
9 /user/bahaaldine/data/08
9 /user/bahaaldine/data/09
9 /user/bahaaldine/data/10
9 /user/bahaaldine/data/11
9 /user/bahaaldine/data/12
9 /user/bahaaldine/data/13
9 /user/bahaaldine/data/14
9 /user/bahaaldine/data/15
9 /user/bahaaldine/data/16
9 /user/bahaaldine/data/17
9 /user/bahaaldine/data/18
9 /user/bahaaldine/data/19
9 /user/bahaaldine/data/20
9 /user/bahaaldine/data/21
9 /user/bahaaldine/data/22
9 /user/bahaaldine/data/23
```

TwitterConfig — Twitter Stream — route1 — 1696968 rows – 34736.82s / 48.85 rows/s — Add Carriage Return — route5 — 1696968 rows – 34736.83s / 48.85 rows/s — FILE — Tweets Hourly File



1  Streaming Tweets

**Talend ESB**

2  Push to HDFS

**Talend Big Data**

3  Hive / Pig Processing

**cloudera**
Ask Bigger Questions

4



use default database — OnComponentOk — drop the table tweets — OnComponentOk — create the table tweets — OnComponentOk — for each hour

Iterate

create the partitions

## 8. Appendix-A

### Configuration

| Category | Property | Value |
|---|---|---|
| **Default** | Bind DataNode to Wildcard Address | ☑ Reset to the default value: false ↰ Override Instances |
| ▸ Service-Wide | | |
| ▸ Balancer (Default) | Use DataNode Hostname | ☑ |
| ▾ DataNode (Default) | dfs.datanode.use.datanode.hostname | Reset to the default value: false ↰ |
|     Resource Management | | |
|     Performance | | |
|     **Ports and Addresses** | DataNode Protocol Port | 50020 |
|     Security | dfs.datanode.ipc.address | default value |

**cloudera manager**   Home   **Services** ▾   Hosts   Activities ▾   Diagnose ▾   Audits   Charts ▾   Administration

Services » Service hdfs1 »

### 🖴 hdfs1

🏠 Status    ☰ Instances    ⊙ Commands    🔧 Configuration ▾    👁 Audits    🖼 Charts Library    ⤴ NameNod...

### Configuration

🔍 dfs.permission ✖

✔ 1 validation check. ❯

| Category | Property | Value |
|---|---|---|
| Service-Wide / Security | Superuser Group<br>dfs.permissions.supergroup,<br>dfs.permissions.superusergroup | supergroup<br>default value |
| Service-Wide | Check HDFS Permissions<br>dfs.permissions | ☐<br>Reset to the default value: true ↰ |

## Cluster 1 - CDH4

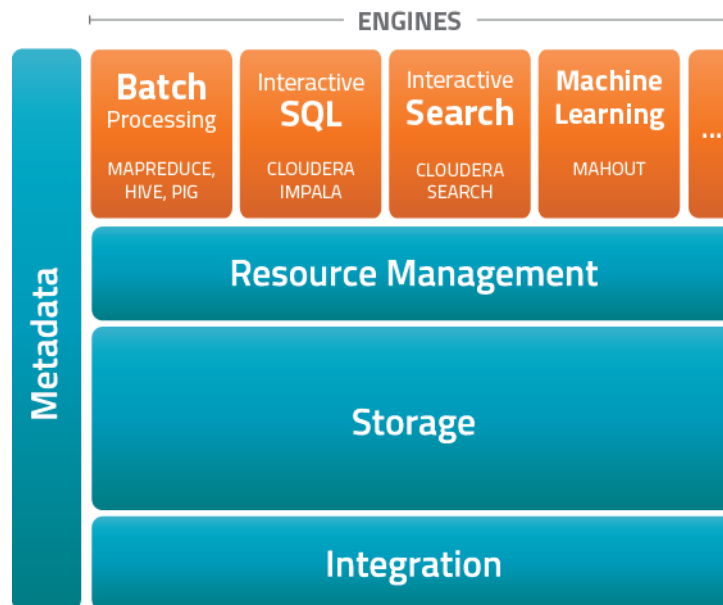| Name | ▲ | Status |
|------|---|--------|
| 🔶 flume1 ▾ | | ◉ Stopped |
| **H** hbase1 ▾ | | ◉ Stopped |
| 📦 hdfs1 ▾ | | ◯ Good Health |
| 🐝 hive1 ▾ | | ...h |
| **H)** hue1 ▾ | | ...h |
| ⚡ impala | | |
| 📘 ks_ind | | |
| ⠿ mapreduce1 ▾ | | ◯ Good Health |

🏠 Status
☰ Instances
⊙ Commands
🔧 Configuration
👁 Audits
🖼 Charts Library

# Login

Username:

Password:

☐ Remember me on this computer.

Login

ENGINES

| | **Batch** Processing | Interactive **SQL** | Interactive **Search** | **Machine Learning** | ... |
|---|---|---|---|---|---|
| **Metadata** | MAPREDUCE, HIVE, PIG | CLOUDERA IMPALA | CLOUDERA SEARCH | MAHOUT | |

Resource Management

Storage

Integration

**cloudera**

**Cloudera's Distribution Including Apache Hadoop (CDH)**

## Cloudera QuickStart VM

**The Cloudera Experience on a Single Machine.**

This VM contains a sample of Cloudera's Platform for Big Data. Although the true power of Hadoop comes when it can be distributed across hundreds, even thousands of nodes, this VM makes it easy for you to learn without having to set up a full cluster.

**Components in the VM:**

Cloudera Standard

- CDH
- Cloudera Manager (limited features)

**Downloads & Instructions** >

**Hadoop Tutorial**