## Oozie

Complex work flow management

Relief to developers!!
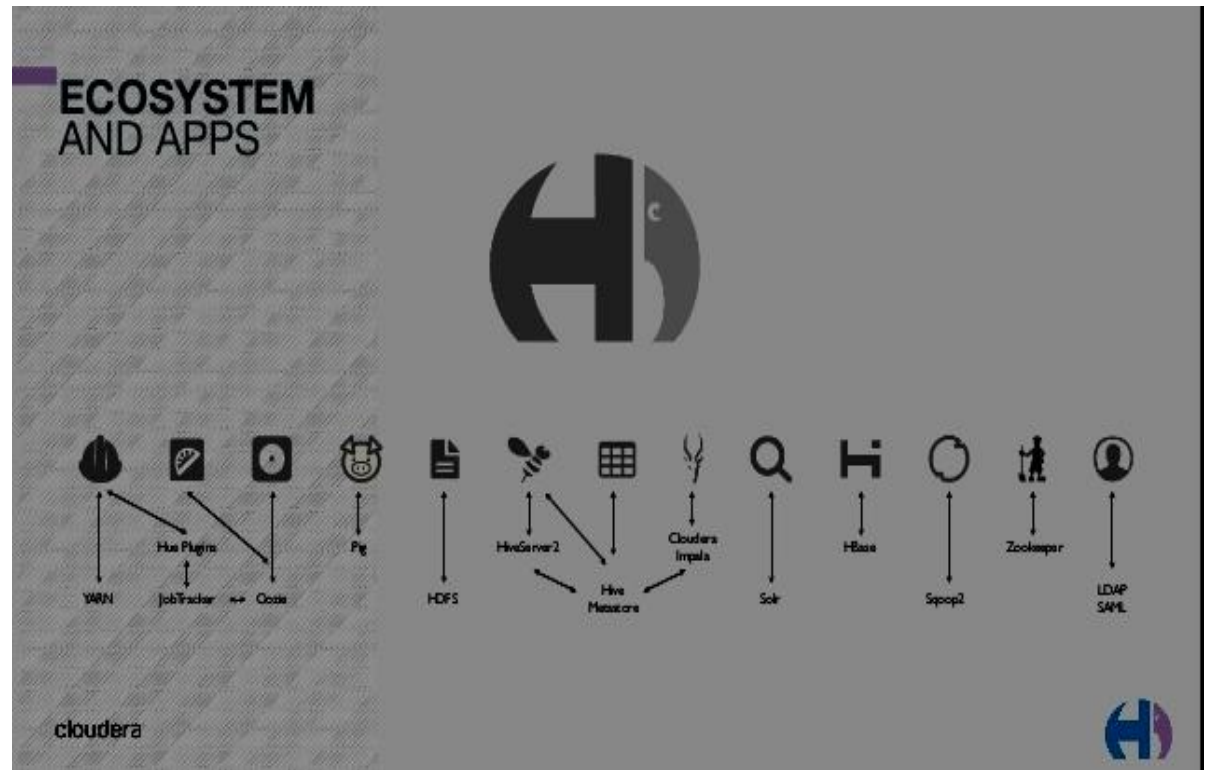


The Hadoop UI

# Agenda

- Intro to Hue:
  - What is Hue
  - When to Use it
- Oozie
  - What is Oozie
  - Life with out it
  - When to use it
- Oozie Workflow
  - Actions
  - Control flow
- Actions
  - Sqoop Action
  - Hive Action
  - HDFS Action
- Workflows:
  - Example1
  - Example2
- Scheduling Workflows – Coordinators
  - Bundles

- Cloudera designed it, initially it was a commercial tool, later made it Apache Open Source.

- Single tool, which provides multiple options to developers:

- Hive Editor

- Pig Editor

- Hive Metastore Manager

- Impala

- DB Query

- Oozie

- File Browser – HDFS

- Job Browser – Resource Manager

- One stop for developers

# When to use it

- If you are looking for the following:

- Editor to develop programs in Hive/ Pig/ Impala

- HDFS browser similar to Windows File Browser

- Track progress of :
  - Hive jobs
  - Pig Jobs
  - Spark Jobs
  - Map Reduce etc..

- Hive Metastore Manager

- Better Access to Databases and Tables on Hive

- Download results of Hive queries
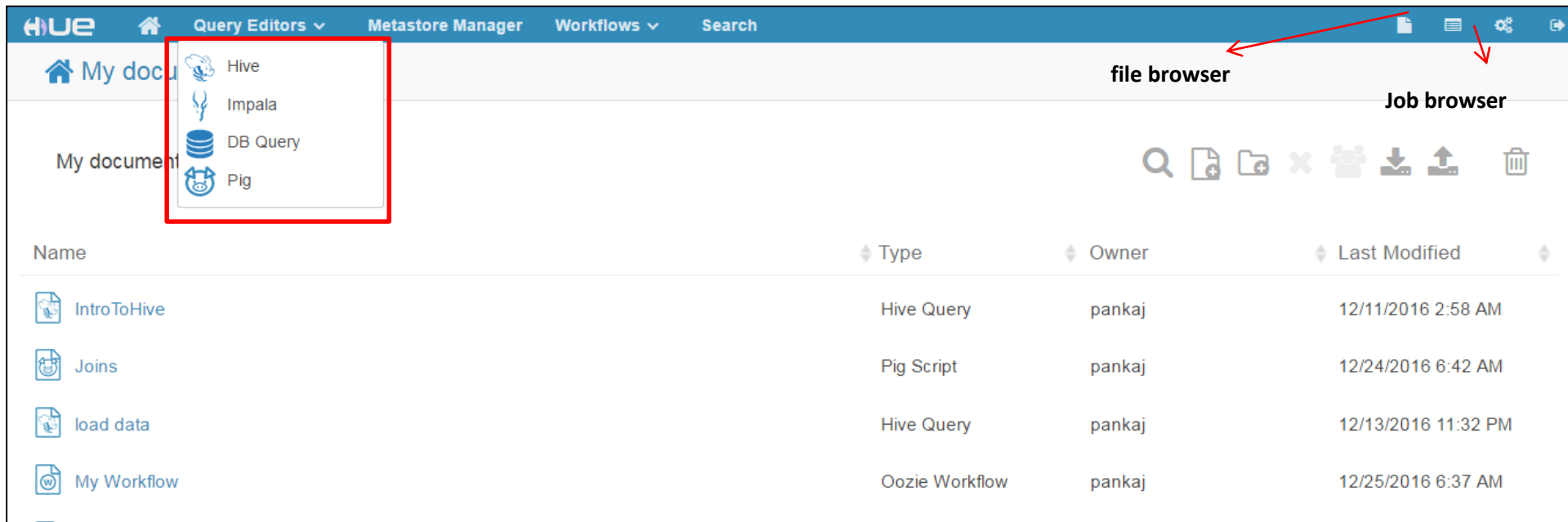
- Visualize results of Hive Queries

# Query Editors, Job and File Browsers

# Oozie Workflow

# Oozie

- Oozie as a tool is known for designing workflows, schedule and track them.

- In many industries, there are technology workflows, where output from one tool is consumed by another tool.

- Consider a example where you want to do the following:

- Clean data in HDFS

- Sqoop Import

- Hive / Pig scripts

- Sqoop export

- Drop a mail up on completion

- Answer to this example is Oozie. If you carefully observe, there is direction for data flow, and these flows are called as DAG( Directed Acyclic Graphs)

- The flow of the graph can be controlled using some control nodes, which helps in setting the start and end of the flow and some decision making nodes based on some intermediate predicate based values

## Apache Oozie Workflow Scheduler for Hadoop

### Overview

Oozie is a workflow scheduler system to manage Apache Hadoop jobs.

Oozie Workflow jobs are Directed Acyclical Graphs (DAGs) of actions.

Oozie Coordinator jobs are recurrent Oozie Workflow jobs triggered by time (frequency) and data availability.

Oozie is integrated with the rest of the Hadoop stack supporting several types of Hadoop jobs out of the box (such as Java map-reduce, Streaming map-reduce, Pig, Hive, Sqoop and Distcp) as well as system specific jobs (such as Java programs and shell scripts).

Oozie is a scalable, reliable and extensible system.

- Imagine, you have to do the following manually:
  - Connect hadoop components
  - If something goes wrong, drop a mail to stake holders
  - Up on completion, send a detailed report
  - Launch parallel jobs daily, weekly, monthly..

- All the mentioned points are hard to manage, schedule, coordinate, track the progress.
- Oozie can do this for you!

Apache Oozie Market Share in Big Data

We use the best scanning and sleuthing tech in the world to track the install bases of over 3,000 technology products, including Big Data (e.g. databases). In the Big Data category, Apache Oozie has a market share of about 3.3%. Other major products in this category include:

Market Share:

## 3.3%

2,262 Companies

Market-Share for Apache Oozie

Other top Products

- 23,621 companies using Informatica
- 17,295 companies using Apache Hadoop
- 7,448 companies using Teradata
- 6,271 companies using Apache Hbase
- 4,080 companies using Cloudera
- View all other top products

Reference: https://idatalabs.com/tech/products/apache-oozie

# Other tools similar to Oozie for Hadoop/ Big Data

**Commercial tool, reliable fast and easy to use**

**Not a good tool for Big Data, too costly!**

**Good Competitor to talend**

**Azkaban, open source workflow manager, similar to Oozie**

# OOZIE WORKFLOW

- Oozie workflows contain control flow nodes and action nodes.

- Control flow node define the beginning and the end of the workflow (start, end and fail nodes)

- They also provide a way to control the workflow execution path (decision, fork and join)

- Action nodes can trigger the execution of a computation/processing task

- These tasks include – map-reduce, HDFS commands, Pig, SSH, HTTP, eMail and Oozie sub-workflow

# Oozie Architecture



Reference: yahoo

# Oozie on Hue



- Under workflow, you see two options:
  - Dashboards for monitoring jobs.
  - Editors for designing, scheduling and bundling jobs.
- Click on workflows to design workflows.
- Click on coordinators to schedule workflows.
- Bundle to batch a set of coordinator applications.

- **Start node**
  It indicates the first workflow node, through which the workflow job will start. This is the starting point for the workflow job. The "to" attribute points to the node where the job starts

- **End node**
  It is the end of the workflow job. When a workflow job reaches its end, it has completed successfully. Even if some other workflow jobs are in running state, when an end node is reached, these are killed forcefully and the program still exits successfully

# Control Flow Nodes

- **Kill node**

  Kill node allows a workflow job to kill itself. All the running actions of the workflow job would also be killed and a message as mentioned in the tags will be entered in the log file

# Control Flow nodes

- **Fork and join control nodes**
  Fork node is used to split a path into multiple concurrent nodes. It allows tasks to be run in parallel. Join nodes then waits for every concurrent execution paths to reach to it. Fork and Join should be used in pairs

- **Decision node**
  Decision node allows a workflow job to make a selection between its execution paths based on a list of predicate. The predicates are evaluated on order of appearance until one of them comes true. In case all return false, default transition is taken. These predicate can contain logics such as the size of file being greater than some threshold, or the file being completely loaded or the exit status of an action node

# Oozie Actions

- Action nodes are generally used for performing computation tasks

- No Computation performed by Action nodes takes within oozie, all compution are performed remotely

- All computations performed by action node are asynchronous, but for most of the computations the workflow jobs waits for the action task to complete by polling and callbacks

- Ok attribute provides the path to follow on successful completion

- Error attribute provides the path to follow in case of error

# Hive Action



Provide the path of Hive Script located on HDFS.

Optional config files can be loaded.

# Hive Server2 Action



Provide the path of Hiveserver 2 Script located on HDFS.

Optional config files can be loaded.

# Pig Action



Provide the path of Pig Script located on HDFS.

# Spark Action



Spark Master IP Address
Mode of execution
Provide the path of Jar Files to be executed
Provide the main class path

# Java Action



Java scripts in Oozie are used to set parameters real time.

# Sqoop Action



Provide sqoop command

# Map reduce Action



Provide the path of runnable Jar File

# Sub Workflow, Fork and Joins control nodes



Like Informatica, in Oozie, one workflow can be triggered from another workflow.

Here choose the workflows which are already designed.

# Shell Action



Provide the path of shell scripts.

The shell action runs a Shell command.

# SSH action



Remote shell commands execution.

The ssh action starts a shell command on a remote machine as a remote secure shell in background. The workflow job will wait until the remote shell command completes before continuing to the next action.

The shell command must be present in the remote machine and it must be available for execution via the command path.

Ref:https://oozie.apache.org/

# HDFS file action



HDFS file level actions like delete, create a directory, touch file, move a file or directory can be execute.

Select one option, and provide the path.

# Email Action



Email Action, to update the status of the job up on failures, or completion.

Provide the list of email addresses, subject and message to be sent up on completion or failure of an activity.

# Hadoop Streaming Action



Hadoop streaming programs action to run map reduce programs written in other than Java language.

Example: Python Map Reduce programs.

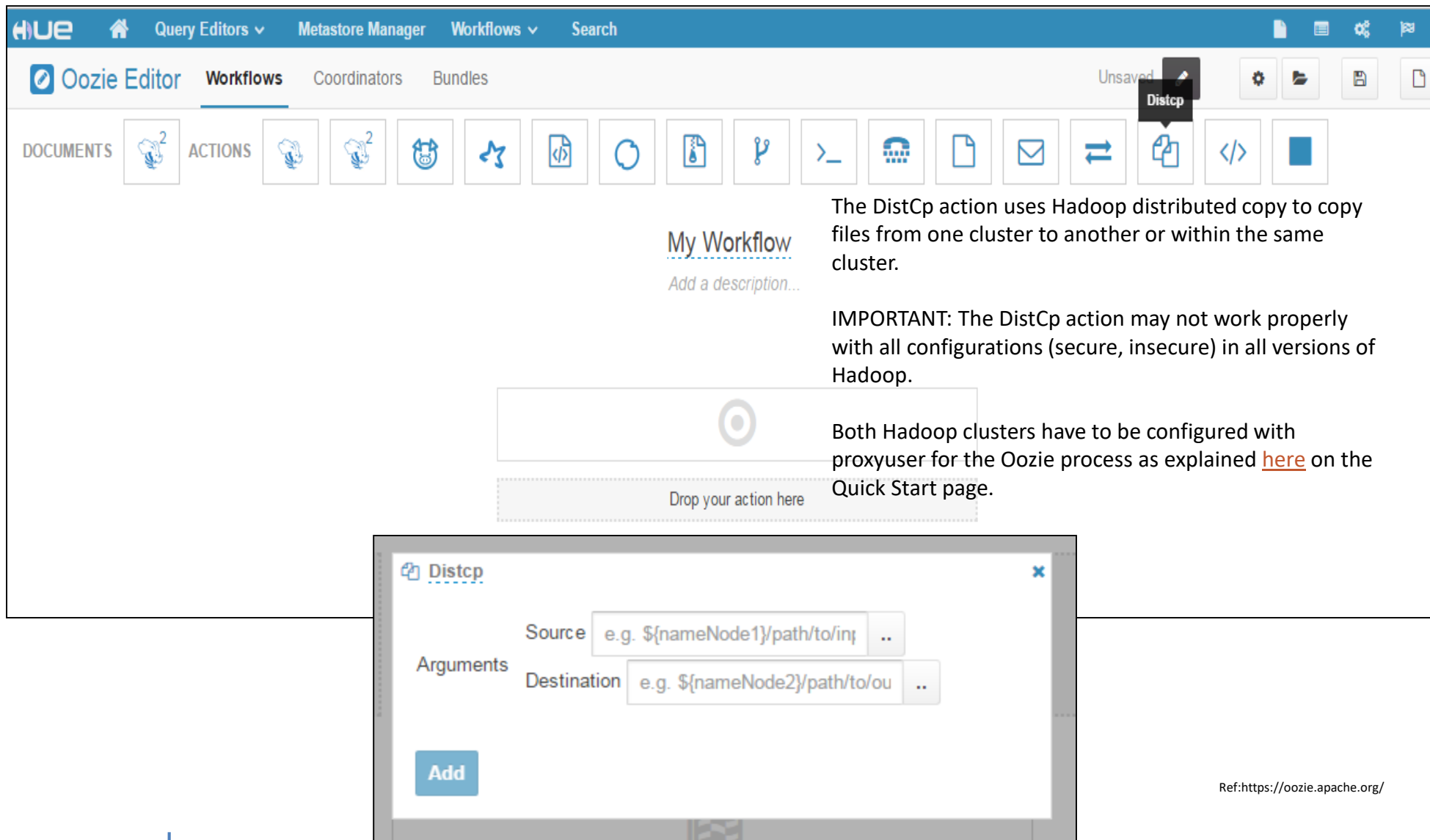Developer should provide Mapper and Reducer python files.

# Distcp Action



The DistCp action uses Hadoop distributed copy to copy files from one cluster to another or within the same cluster.

IMPORTANT: The DistCp action may not work properly with all configurations (secure, insecure) in all versions of Hadoop.

Both Hadoop clusters have to be configured with proxyuser for the Oozie process as explained here on the Quick Start page.
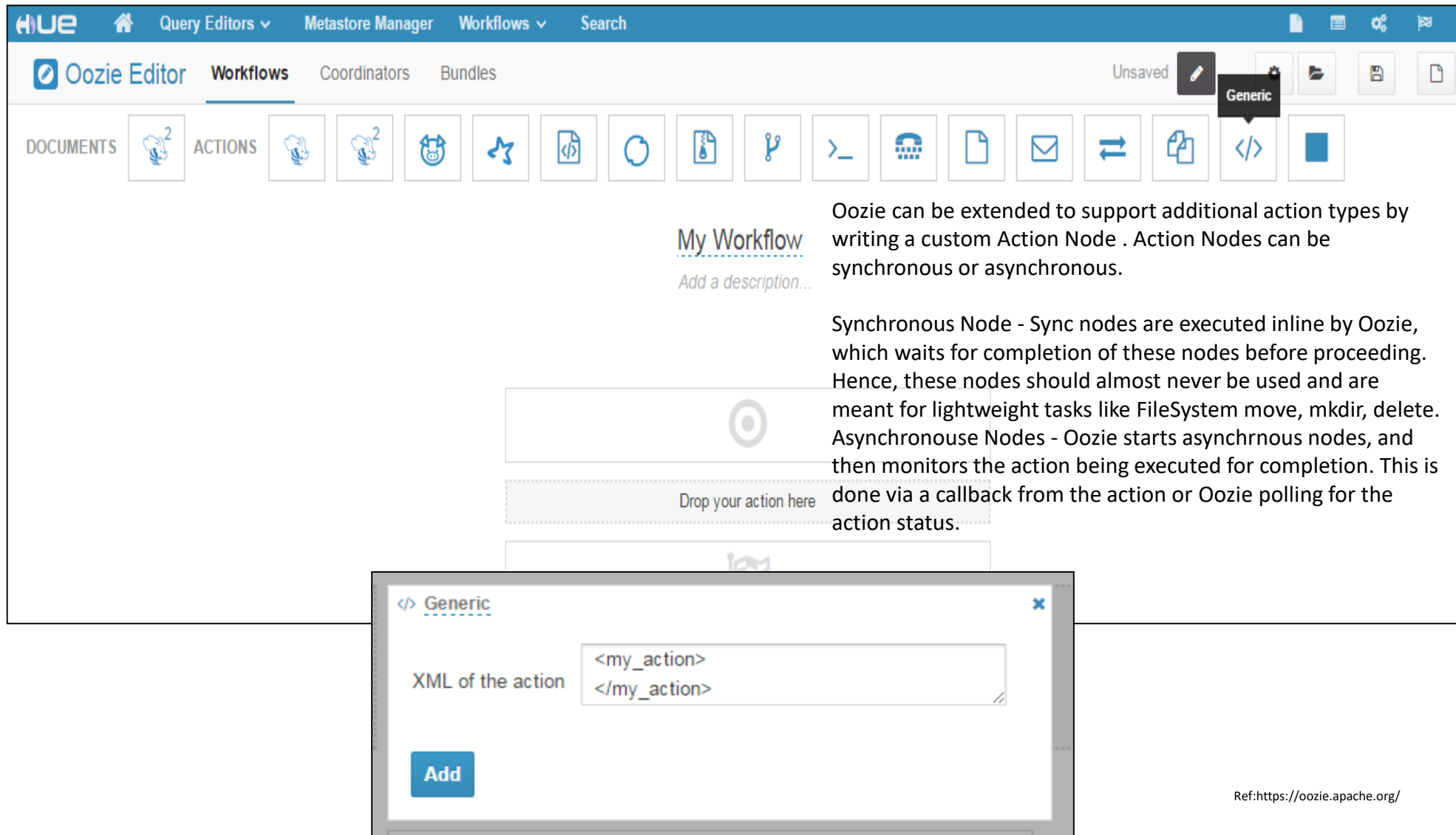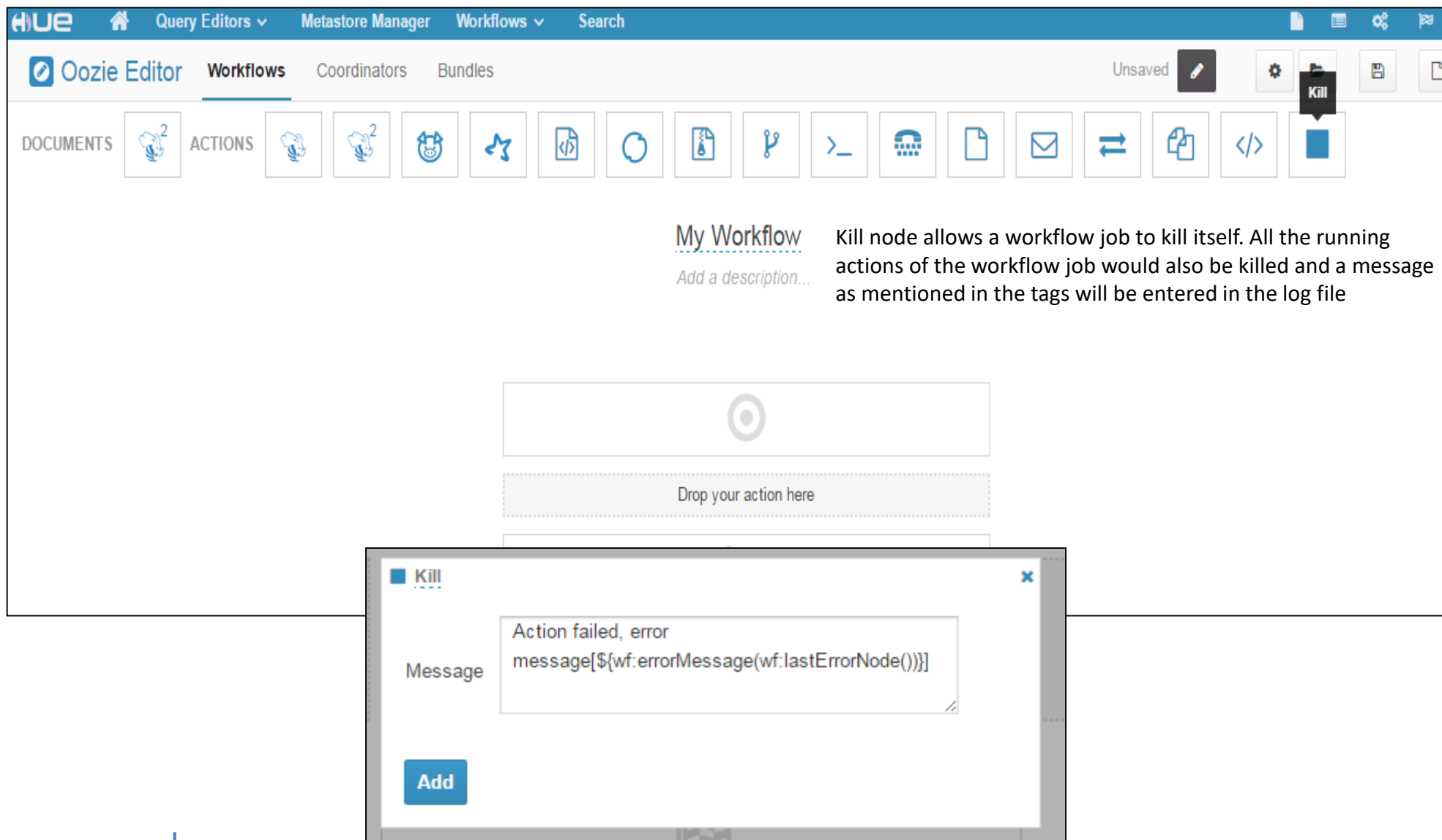
Ref:https://oozie.apache.org/

# Generic action



Oozie can be extended to support additional action types by writing a custom Action Node . Action Nodes can be synchronous or asynchronous.

Synchronous Node - Sync nodes are executed inline by Oozie, which waits for completion of these nodes before proceeding. Hence, these nodes should almost never be used and are meant for lightweight tasks like FileSystem move, mkdir, delete. Asynchronouse Nodes - Oozie starts asynchrnous nodes, and then monitors the action being executed for completion. This is done via a callback from the action or Oozie polling for the action status.

Ref:https://oozie.apache.org/

# Kill Node



Kill node allows a workflow job to kill itself. All the running actions of the workflow job would also be killed and a message as mentioned in the tags will be entered in the log file

The Oozie **Coordinator** system allows the user to define and execute recurrent and interdependent workflow jobs (data application pipelines).

Choose a workflow, and you will find the following options:

# Bundle

Bundle is a higher-level oozie abstraction that will batch a set of coordinator applications. The user will be able to start/stop/suspend/resume/rerun in the bundle level resulting a better and easy operational control.

More specififcally, the oozie Bundle system allows the user to define and execute a bunch of coordinator applications often called a data pipeline.

My Bundle **Name of bundle**

Add a description...

Which schedules to bundle?

➕ Add a coordinator **Add the coordinators to bundle**

sample ↗                                                                                    ✖

| start_date | 2017-01-12T08:54 | — |
| end_date | 2017-01-19T08:54 | — |

**Select the time period for which this bundle should be active.**

➕ Add a parameter

In this lab, you will learn the follwing:
1. Creating sqoop action
2. Creating HDFS action
3. Creating Hive action
4. Creating Pig Action
5. Creating Workflows
6. Integrating sub-workflows
7. Coordinating workflows

1. Sqoop data to HDFS
2. Pig consumes this data, and generates top10Cust. Sqoop exports the results back to MySQL.
3. Hive creates external tables, and generates chain_stats. Sqoop exports the results back to MySQL.

# Create required tables in MySQL

**Login:**

```
mysql -h 54.149.41.179 -u username -p --local-infile

use database;

CREATE TABLE transactions(id varchar(20),chain varchar(20), dept varchar(20),
category varchar(20), company varchar(20), brand varchar(20), date1 varchar(10),
productsize int, productmeasure varchar(10), purchasequantity int, purchaseamount FLOAT);

LOAD DATA LOCAL INFILE '/home/training/Desktop/transactions.csv'
INTO TABLE transactions FIELDS TERMINATED BY ',' ENCLOSED BY '"'
LINES TERMINATED BY '\r\n';

-- This table will be imported to HDFS


-- Results will be exported back to mysql

create table chain_stats(chain varchar(3), deptcin int,categorycin int,
companycin int,brandcin int,totalspent float);

create table chainTop10Cust(chain varchar(3), id varchar(10),totalSales float);
```

Login to Hive, and create database oozie_username. Replace your name under username.

In oozieScripts folder, change the DBName in hive script to oozie_username.

Up on making relavent changes:

Copy oozieScripts folder to home folder using winscp.

Push the folder to hadoop:

hadoop fs -put oozieScripts

```
pankaj@ip-172-31-4-182:~/oozie$ hadoop fs -put oozieScripts
pankaj@ip-172-31-4-182:~/oozie$ hadoop fs -ls oozieScripts
Found 2 items
-rw-r--r--   1 pankaj labusers       1627 2017-01-13 07:08 oozieScripts/reporting.hql
-rw-r--r--   1 pankaj labusers        948 2017-01-13 07:08 oozieScripts/storeTop10cust.pig
```

# Sqoop Import workflow



**Name of the workflow**

**Drag and drop HDFS fs action.**
Delete the directory, if exists before sqoop starts.

**Drag and drop sqoop action.**

**Add the below command:**
import --connect jdbc:mysql://54.149.41.179/pankaj --username username--password pwd--table transactions --split-by id --target-dir TransactionsData

Add username and pwd.

# Sqoop Import workflow - Submit

Save and submit.

Submit          Save

# Sqoop Import workflow – Track progress

# Data Check up on completion



You should see data imported to HDFS

# Pig <--> Sqoop <--> MySQL – top 10 customers in each chain - report

**DOCUMENTS** | **ACTIONS**

pig MIS Reporitng

⊙

📄 HDFS Fs      ⊙

Delete
/user/pankaj/chainTop10Cust/

🐷 Pig Script      ⊙

storeTop10cust.pig

○ Sqoop 1      ⊙

export --connect jdbc:mysql://54.149.41.179/pankaj --username pankaj -..

🏁

## Name of the workflow

**Clean the directories.**
Delete /user/username/chainTop10Cust

*Replace username with your username*

**Path of pig script:**

Browse and select the pig file in HDFs.

**Sqoop Export o MYSQL**

Add the following command to export results to MySQL:

export --connect jdbc:mysql://54.149.41.179/pankaj --
username username--password pwd --table chainTop10Cust --
export-dir chainTop10Cust --input-fields-terminated-by "\t"

Add username and pwd
After completion, check the table in MySQL

Save and submit.

Submit    Save

# Pig <--> Sqoop <--> MySQL– top 10 customers in each chain – Track progress

# Pig <--> Sqoop <--> MySQL, pig output check

# Pig <--> Sqoop <--> MySQL,  Sqoop export check

**Login to MySQL** :

mysql -h 54.149.41.179 -u username-p

```
mysql> select * from chainTop10Cust;
+-------+----------+------------+
| chain | id       | totalSales |
+-------+----------+------------+
| 4     | 96841999 |    9552.77 |
| 4     | 13744500 |    9551.77 |
| 58    | 42937475 |    10223.8 |
| 88    | 73507112 |    12875.6 |
| 88    | 73850140 |    10351.8 |
| 88    | 18854215 |    9139.25 |
| 88    | 49806426 |    8786.69 |
| 88    | 88852308 |    8037.27 |
| 88    | 85358490 |    7767.95 |
| 88    | 50791864 |    7753.57 |
| 88    | 77989055 |    7313.03 |
| 88    | 66650733 |    7135.31 |
| 88    | 70462513 |    7001.42 |
| 95    | 83868868 |    15302.2 |
| 95    | 49522674 |    11603.2 |
| 95    | 83938442 |    11153.1 |
| 95    | 61933479 |    10753.9 |
```
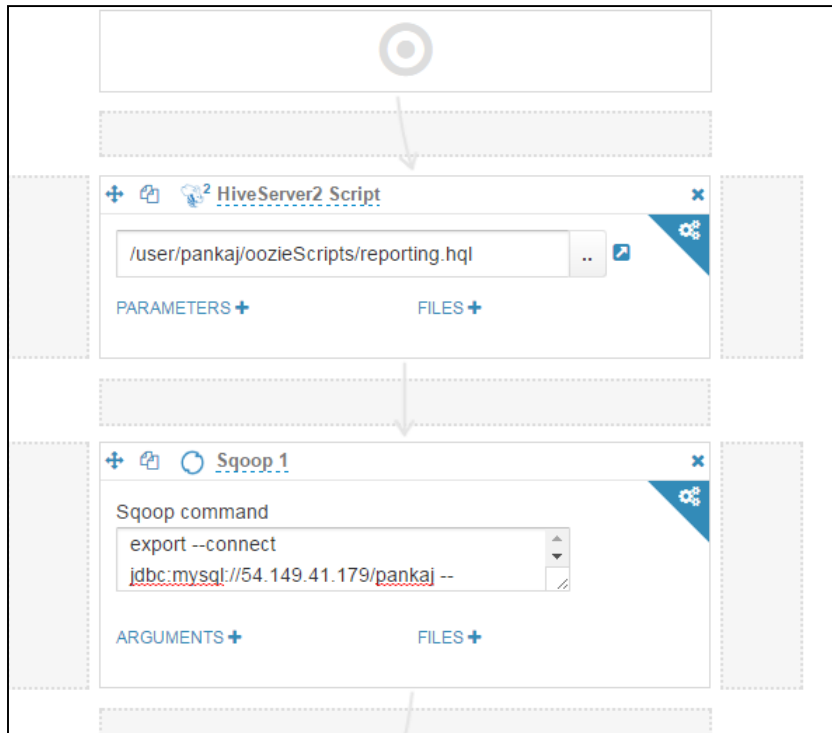
# Hive<--> Sqoop <--> MySQL – Chain Report

**Hive MIS Reporting**

*Add a description...*

**Name of the workflow**

**Path of hive script:**

Browse and select the hive file in HDFs.

**Sqoop Export o MYSQL**

Add the following command to export results to MySQL:

export --connect jdbc:mysql://54.149.41.179/pankaj --username username -password pwd --table chain_stats --export-dir /user/hive/warehouse/oozie_test.db/chain_stats/ --input-fields-terminated-by "\001"

Add username and pwd
After completion, check the table in MySQL

# Hive<--> Sqoop <--> MySQL Track Progress

Login to hive, or use Hue.

Check in the oozie database created.

```
hive> select * from oozie_test.chain_stats;
OK
14      82      700     2026    2486    365785.2999999863
15      82      749     2715    3244    1077677.4699999483
17      81      706     1960    2317    394737.3999999857
18      82      740     2393    2931    693190.2099999533
2       65      240     227     285     3814.63999999999
20      82      721     1833    2193    291686.07999999437
205     80      632     1115    1563    106421.01999999971
3       82      682     1703    2044    208104.0299999987
4       82      737     2394    2884    767159.359999946
58      76      375     322     449     10223.759999999973
88      82      714     1921    2269    324290.9899999864
95      82      717     2160    2546    441778.4099999827
Time taken: 1.928 seconds, Fetched: 12 row(s)
```
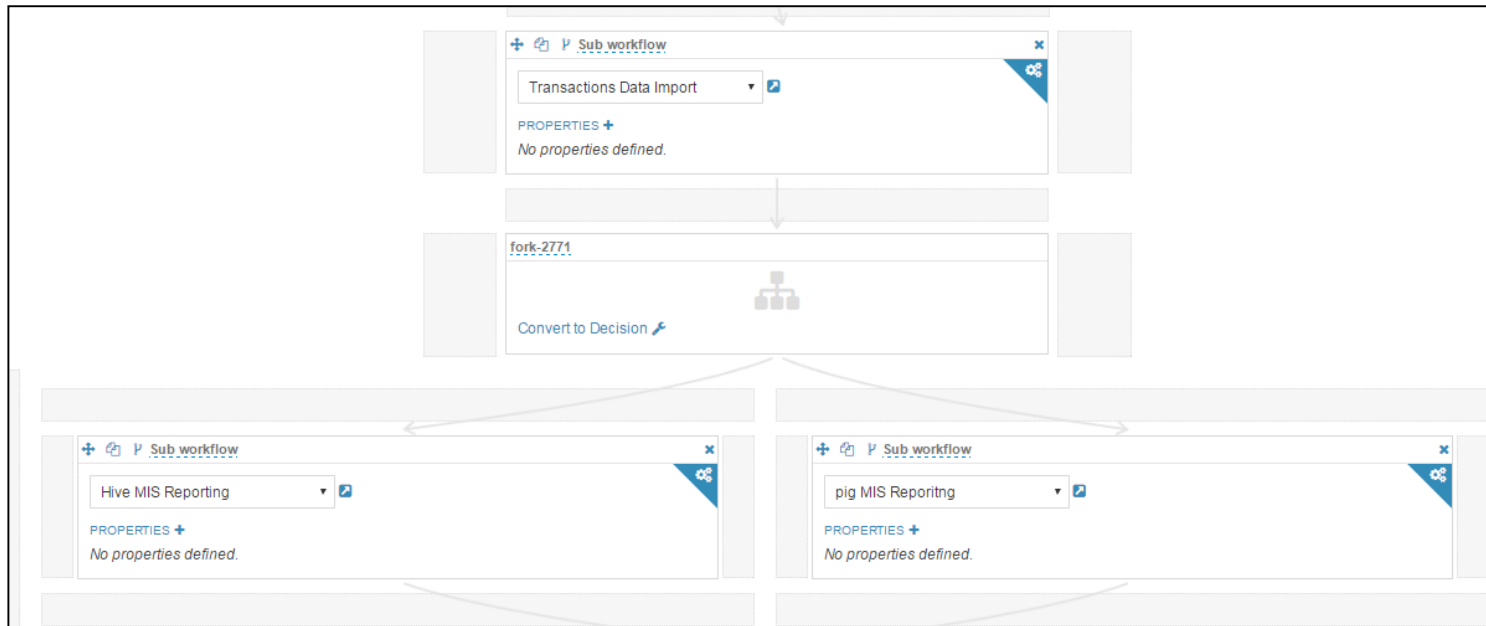
# Combining all three workflows in to single workflow.

Create a new workflow.
Use sub workflow action.
Step1: Sqoop import work flow.
Step2: Hive and Pig reporting workflows in parallel.



Save it as reporting workflow.
Run and see what happens!

# Submit

# Progress

Add the previous workflow and schedule it.



Create a new one

Select this workflow



Name of the coordinator

**ReportGenCoordinator**

*Add a description...*

## Which workflow to schedule?

ReporintWorkflow 🔗

## How often?

Every day at 0 : 0

⇄ Options

## Workflow Parameters

➕ Add parameter

Save

# Conclusion

You learnt the following:

1. Creating sqoop action
2. Creating HDFS action
3. Creating Hive action
4. Creating Pig Action
5. Creating Workflows
6. Integrating sub-workflows
7. Coordinating workflows