# Oozie

Friday, December 31, 2021      11:46 AM

From the book on Apache Oozie

### *What is Oozie?*

Its an orchestration framework for Hadoop jobs (hive, pig, map-reduce, shell).
There are 3 kinds of Oozie jobs :
1. Jobs that are run on a standalone basis : Workflow jobs
2. Jobs that run on a scheduled basis : Coordinator jobs
3. Coordinator jobs that are bundled toether to run as a single Oozie job : Bundle job

A workflow application consists of a single file ***workflow.xml*** (XML representation of the workflow).
Associated with the workflow.xml, the XML file can obtain its properties from ***job.properties***

For all Oozie actions, we need to provide :
1. Job Tracker : Similar to resource Manager. In Hadoop 2.0, this will be YARN
2. Name node URI : This is the URI of the name node ofd Hadoop
3. Oozie workflow application path (The path that contains workflow.xml file)

How do you submit Oozie job?

We can use command line tool to monitor the progress of Oozie jobs

```
$ export OOZIE_URL=http://localhost:11000/oozie
$ oozie job -run -config target/example/job.properties
job: 0000006-130606115200591-oozie-joe-W
```

How do you monitor the progress of Oozie jobs?

```
oozie job -info 0000006-130606115200591-oozie-joe-W
```

We can also see detailed job info on oozie web interface

# *Oozie concepts*

### *Oozie application*

Oozie application is just like any other executable on Unix/Linux (Like /bin/echo)
Oozie job is like a process

### *Oozie workflow*

start and end control nodes define the start and end of a workflow. The fork and join
control nodes allow executing actions in parallel. The decision control node is like a
switch/case statement that can select a particular execution path within the workflow
using information from the job itself. Figure 2-1 represents an example workflow.

### *Oozie coordinators*

- Defined by start time, end time, frequency, and it starts when all necessary data is available
- When data is not available, workflow execution is delayed until input data is Available
- It runs periodically from start time to end time
- It is possible to configure a coordinator to wait for a max time before input shows up.. If input dosent show up before that time, we can time out the workflow.
  - When the data is delayed for more time, the coordinator job keeps a track of the missing days and processes data appropriately once data shows up in the correct input folder.



If a coordinator does not define any input data, the coordinator job is a time-based
scheduler, similar to a Unix cron job.

| Time | Input Data | Workflow | Output Data |
|---|---|---|---|
| 2013JAN02 (start) | rawlogs/2013JAN01/ | | zip_userName_interactions/2013JAN01/ username_interactions/2013JAN01/ username_ZIPs/2013JAN01/ |
| 2013JAN03 | rawlogs/2013JAN02/ | | zip_userName_interactions/2013JAN02/ username_interactions/2013JAN02/ username_ZIPs/2013JAN02/ |
| 2013JAN04 | rawlogs/2013JAN03/ | | zip_userName_interactions/2013JAN03/ username_interactions/2013JAN03/ username_ZIPs/2013JAN03/ |
| 2014JAN01 | rawlogs/2013DEC31/ | | zip_userName_interactions/2013DEC31/ username_interactions/2013DEC31/ username_ZIPs/2013DEC31/ |
| 2014JAN02 (end) | | | |

### *Oozie bundles*

- 
  - You can handle dependencies among different coordinators using a bundle job. Ex : lets say the output of a coordinator becomes the input of another coordinator.

We have three coordinator jobs: daily-logs-coordinator, weekly-aggregator-coordinator, and monthly-aggregator-coordinator. Note that we are using the same workflow application to do the reports aggregation. We are just running it using different date ranges.
A logs-processing-bundle bundle job groups these three coordinator jobs. By running the bundle job, the three coordinator jobs will run at their corresponding frequencies.
All workflow jobs and coordinator jobs are accessible and managed from a single bundle job.
This logs-processing-bundle bundle job is also known as a data pipeline job.


Oozie in general uses a lot of parameters and functions for its execution

# *Oozie architecture*



- 
  - Oozie basically executes after reading an xml file and expects all files to be present on HDFS for the workflow to function.
  - Oozie has a few main components :
    - Oozie client
      - Can be command line tool, HHTP REST API Interface, Oozie Java client API
    - Oozie Server
      - 
      - 
        e
      - Oozie server retrieves the job status from the SQL database, performs the corresponding operation and updates state of jobs to the SQL database.
        Hence, if an oozie server crashes, it can resume the jobs when the oozie server is back up, because all statuses are stored in SQL database.
      - There are 2 main entities that perform all operations :
        1. Command class :
            a) It facilitates queueing of commands that are to be executed asynchronously.
                i) There is usually a fixed size queue in Oozie. Hence even when there is too much load, only a few commands are in queue to be processed. Others are silently dropped from queue. But Oozie has a mechanism to recreate the oozie statuses and actions that were submitted.
            b) It updates database with the status of each job.
            c) *It basically manages the scheduling of Oozie jobs*
        2. ActionExecutor Class
            a) This class ACTUALLY performs the steps required for each action (Ex : Hive action, shell action, ssh action etc.)
            b) If we have to define a new Action, then we need to write an ActionExecutor class for that action, and also define the XML structure for that action.
            c) *It Manages the execution of Oozie jobs*
    - SQL Database :
      - All information about completed and in-process jobs are stored on SQL database.
      -

## _Oozie deployment_

- Oozie server : a web application that acts like a hadoop client and also a database client. It also has an optional web interface built using EXTJS.
- Oozie client :
  - Oozie CLI
  - Oozie REST API
  - Oozie Java client library (can be used with any JVM systems)


You download the Oozie application package
- Then you have a base oozie.war to setup the oozie web server
- You then inject all dependent jars (hadoop jars, database jars) and produce an enhanced .war file, which can be used to create the web server.
- You can setup the Oozie DB using a setup script

How can you upgrade support for Oozie for different databases automatically?
- Using sharelib JARS and using timestamped lib_$timestamp directories.
- oozie admin –share libupdate is the way to do it and ensures that Oozie uses the latest version.

## _Chapter 4_

## _Ooze workflow actions_

Hadoop DistCp, for example, is a common tool used to pull data from S3

Workflow actions are a set of DAG actions, mentioned in workflow.xml file

### _Action- execution model in Oozie : How are actions launched in Oozie?_

- When a user submits a oozie job through Oozie CLI or any other client tools, the instruction is then passed onto the Oozie server.
- The Oozie server does not launch the action jobs (Hive, shell etc.) on its own server. Instead the Oozie server launches launcher jobs on Hadoop cluster for every action
- Each Launcher job occupies space on the Hadoop cluster
- Each launcher job in turn launches the appropriate action job on Hadoop (Hadoop job) and waits for all the hadoop jobs to complete.
- Each launcher job is a separate self-contained application that can run on any node in the hadoop cluster.
- Hence with this mechanism, Oozie server does not take the headache of running all Hadoop jobs. But the output of most hadoop jobs are redirected to the STDOUT/STDERR of Oozie server, and hence the outputs of these hadoop jobs can be viewed on Oozie web interface.



### _Action components_

Each Action has a name property and 2 other routes :
- OK <ok>
- Error <error>
Depending on if the action fails/sucedes it goes to the corresponding next action


Oozie documentation :


https://oozie.apache.org/docs/4.1.0/WorkflowFunctionalSpec.html

### _Map-reduce action_

- Every action has a defined XML structure in Oozie depending on the XST followed (XML version of Oozie followed)
- _**We can use <file> and <archive> tags to create symlinks to files, and can be used within the action code.**_
- _**Oozie automatically looks for file and folders under the lib/ sub-directory under the workflow directory. Hence you don't need to create**_

*specific symbolic links using <file> and <archive> as long as the files that you need in the actions are present inside the l ib/ sub-directory of the root workflow directory.*

- You can run other forms of map-reduce actions as well, like streaming or pipes

### Java action

- Java action can be used to run a normal java program having public static void main
- Java action can be used to run map-reduce jobs as well, because map-reduce jobs are java programs afterall
  - But this option os not optimal because :
    The launcher map task that launches the <map-reduce> action completes immediately
    and Oozie directly manages the MapReduce job. This frees up a Hadoop
    slot for a MapReduce task that would have otherwise been occupied by the
    launcher task in the case of a <java> action.

- The <capture-output> element, if present, can be used to pass the output back to the
  Oozie context.
  The oozie.action.max.output.data property defined in ooziesite.xml on the Oozie server node controls the maximum size of the output data. It is set to 2,048 by default, but users can modify it to
  suit their needs. It can be modified even on workflow.xml under the <configuration> section.

  Sometimes <capture-output> can give an error if the output size is greater than 2048 KB. In those cases, its advisable to disable it

in the job.properties file for a lot of the actions to find the required
jars and work seamlessly.

### Pig action

```
<action name=" myPigAction">
<pig>
...
<script>/mypigscript.pig</script>
<argument>-param</argument>
<argument>TempDir=${tempJobDir}</argument>
<argument>-param</argument>
<argument>INPUT=${inputDir}</argument>
<argument>-param</argument>
<argument>OUTPUT=${outputDir}/my-pig-output</argument>
</pig>
</action>
```

The code contains 2 levels of parametrization..
${tempJobDir} is a variable coming from .properties file. And the variable TempDir will be used inside the pig script as
$TempDir

We could also write it as :

The argument in the example above, -param INPUT=${inputDir},
tells Pig to replace $INPUT in the Pig script and could have also
been expressed as <param>INPUT=${inputDir}</param> in the

older style of writing Pig actions and is not recommended in newer
versions, though it is still supported.

Example of converting pig command into an oozie pig action

Command :

```
$ pig -Dmapreduce.job.queuename=research -f pig.script -param age=30
-param output=hdfs://nn.mycompany.com:8020/hdfs/user/joe/pig/output
```

Pig action :

```
<action name="myPigAction">
<pig>
<job-tracker>jt.mycompany.com:8032</job-tracker>
<name-node>hdfs://nn.mycompany.com:8020</name-node>
<prepare>
<delete path="hdfs://nn.mycompany.com:8020/hdfs/user/
joe/pig/output"/>
</prepare>
<configuration>
<property>
<name>mapred.job.queue.name</name>
<value>research</value>
</property>
</configuration>
<script>pig.script</script>
<argument>-param</argument>
<argument>age=30</argument>
```

```
<argument>-param</argument>
<argument>output=hdfs://nn.mycompany.com:8020/hdfs/user/
joe/pig/output</argument>
</pig>
<ok to="end"/>
<error to="fail"/>
</action>
```

If any jar's are included inside the pig script internally, you can include the JAR inside of the lib/ sub-directory under the subfolder of the workflow

### FS Action

- Users can run HDFS commands using Oozie using the FS action..only certain commands on HDFS can be implemented
- FS actions are launched by Oozie on the Oozie server (Wheras others are launched using launcher jobs on hadoop cluster). Hence long running FS actions can slow down oozie server and imnpact other oozie applications. This is also the reason why not all HDFS commands (e.g., copy) are supported through this action.

Example of an FS Action

```
<action name="myFSAction">
<fs>
<delete path='hdfs://foo:8020/usr/joe/temp-data'/>
<mkdir path='myDir/${wf:id()}'/>
<move source='${jobInput}' target='myDir/${wf:id()}/input'/>
<chmod path='${jobOutput}' permissions='-rwxrw-rw-'
dir-files='true'/>
</fs>
</action>
```

Examle 2 :

```
$ hadoop fs -rm -r /hdfs/user/joe/logs
$ hdfs dfs -mkdir /hdfs/user/joe/logs
$ hdfs dfs -chmod -R 755 /hdfs/user/joe/
```

Hadoop action :

```
<action name="myFSAction">
<fs>
<delete path='/hdfs/user/joe/logs'/>
<mkdir path='/hdfs/user/joe/logs'/>
<chmod path=' /hdfs/user/joe/' permissions='755' dir-files='true'>
<recursive/>
</chmod>
</fs>
</action>
```

### Sub-Action wf

- We can run a sub-workflow as a part of a parent workflow using the sub-workflow action
- From the parent's perspective, the child wf is an action and the next action in the parent wf is triggered only when the child wf is completed in its entirety
- You can propagate configuration (job.properties) from parent to the child wf

```
<action name="mySubWorkflow">
<sub-workflow>
<app-path>hdfs://nn.mycompany.com:8020/hdfs/user/joe/
sub_workflow</app-path>
<propagate-configuration/>
</sub-workflow>
<ok to="success"/>
<error to="fail"/>
</action>
```

### Hive action

- Hive action is similar to a pig action.. You can use <param> for hive parameters.. Instead you can also use <arguement> to pass in arguments

Example :

```
<action name="myHiveAction">
<hive>
<job-tracker>jt.mycompany.com:8032</job-tracker>
<name-node>hdfs://nn.mycompany.com:8020</name-node>
<job-xml>hive-config.xml</job-xml>
<script>hive.hql</script>
<argument>-hivevar</argument>
<argument>age=30</argument>
</hive>
<ok to="success"/>
```

```
<error to="fail"/>
</action>


<action name="myHiveAction">
<hive>
<job-tracker>jt.mycompany.com:8032</job-tracker>
<name-node>hdfs://nn.mycompany.com:8020</name-node>
<job-xml>hive-config.xml</job-xml>
<script>hive.hql</script>
<param>age=30</param>
</hive>
<ok to="success"/>
<error to="fail"/>
</action>
```

### Shell action

- If shell action calls shell_file.sh, we need to pass shell_file.sh as a file in the <file> parameter.
- If the shell script calls other files, and if the other files are present under the lib/ sub directory of the workflow root directory, they can be referenced inside of the shell script with relative paths.
- While Oozie does run the shell command on a Hadoop node, it runs it via the launcher job. It does not invoke another MapReduce job to accomplish this task.
- <exec> should execute a binary.. Like /usr/bin/python

For ex :

```
<action name="Impala_job" cred="hcat,hs2-creds">
        <shell xmlns="uri:oozie:shell-action:0.1">
            <job-tracker>${jobTracker}</job-tracker>
            <name-node>${nameNode}</name-node>
            <exec>${invalidate_metadata_shell_script}</exec>
            <argument>${mapdt}</argument>
            <argument>${invalidate_metadata_impala_script}</argument>
            <argument>${keytab_path}</argument>
            <file>${hiveSite}#hive-site.xml</file>
            <file>${hiveConfig}#hive-config.xml</file>
            <!--file>${keytab_path}#sunda29.keytab</file-->
            <file>${invalidate_metadata_shell_script}#invalidate_metadata_impala.sh</file>
            <!--file>${invalidate_metadata_impala_script}#invalidate_metadata_impala_script.sql</file-->
            <!--capture-output/-->
        </shell>
        <ok to="EMAIL_SUCCESS"/>
        <error to="KILL"/>
    </action>
```

Example 2 :

```
$ export TZ=PST
$ python test.py 07/21/2014
```

```
<action name="Impala_job" cred="hcat,hs2-creds">
        <shell xmlns="uri:oozie:shell-action:0.1">
            <job-tracker>${jobTracker}</job-tracker>
            <name-node>${nameNode}</name-node>
            <exec>/usr/bin/python</exec>
              <argument>test.py</argument>
            <argument>07/21/2014</argument>
            <env-var>TZ=PST</env-var>
             <file>test.py#test.py</file>

             <!--capture-output/-->
        </shell>
        <ok to="EMAIL_SUCCESS"/>
        <error to="KILL"/>
    </action>
```

                                                                    -directory of the workflow directory.

### Distcp action

- Distcp supports hadoop distributed copy tool , that can be used to copy data from one cluster to another, from s3 to hadoop etc.

Example of conversion between distcp command to oozie command

$ /opt/hadoop/bin/hadoop distcp -m 100 s3n://my-logfiles/2014-04-15/*
/hdfs/user/joe/logs/2014-04-15/

```
<action name="myDistcpAction">
<distcp xmlns="uri:oozie:distcp-action:0.1">
<job-tracker>jt.mycompany.com:8032</job-tracker>
<name-node>hdfs://nn.mycompany.com:8020</name-node>
<prepare>
```

```
<delete path="hdfs://nn.mycompany.com:8020/hdfs/user/joe/
logs/2014-04-15/"/>
</prepare>
<arg>-Dfs.s3n.awsAccessKeyId=XXXX</arg>
<arg>-Dfs.s3n.awsSecretAccessKey=YYYY</arg>
<arg>-m</arg>
<arg>100</arg>
<arg>s3n://my-logfiles/2014-04-15/*</arg>
<arg>/hdfs/user/joe/logs/2014-04-15/</arg>
</distcp>
<ok to="success"/>
```

### Email action

- This is one of the few actions that gets launched on Oozie server

### Ssh action

- Ssh action can be used to run shell commands on a remote host (which is usually not the hadoop cluster)
- However, the oozie.action.ssh.allow.user.at.host should be set to true in oozie-site.xml for this to be enabled.
-                                         o run shell commands or some custom scripts
  
  Hadoop cluster.
- Also, the <shell> action runs through an Oozie launcher while the <ssh> action is initiated from the Oozie server.

```
<action name=" mySSHAction ">
<ssh>
<host>foo@bar.com<host>
<command>uploaddata</command>
<args>jdbc:derby://bar.com:1527/myDB</args>
<args>hdfs://foobar.com:8020/usr/joe/myData</args>
</ssh>
</action>
```

### Sqoop action

- Sqoop is used to import/export data to/from a relational database
- You can either have the entire command under <command> or have the sqoop arguments present under <arg> elements

```
$ /opt/sqoop-1.4.5/bin/sqoop import --connect jdbc:mysql://mysqlhost.mycompany
.com/MY_DB --table test_table -username mytestsqoop -password password
--target-dir /hdfs/joe/sqoop/output-data -m 1
```

```
<action name="sqoop-import">
<sqoop xmlns="uri:oozie:sqoop-action:0.2">
<job-tracker>jt.mycompany.com:8032$lt;/job-tracker>
<name-node>hdfs://nn.mycompany.com:8020$lt;/name-node>
<prepare>
<delete path=" hdfs://nn.mycompany.com:8020/hdfs/joe/sqoop/
output-
```

```
<arg>SELECT count(*) FROM test_table</arg>
</sqoop>
<ok to="end"/>
<error to="fail"/>
</action>
```

*__Synchronous vs Asynchronous actions__*

page 40. These are called asynchronous actions because they are launched via a
launcher as Hadoop jobs. But the filesystem action, email action, SSH action, and
sub-workflow action are executed by the Oozie server itself and are called synchronous
actions.

# *__Chapter 5__*

# *__Workflow applications__*

## __Outline of a basic workflow__

Each wf is comprised of :
1. Global configurations
2. Control nodes
3. Action nodes

## *__Control nodes__*

<start>, <end>

### *__<fork> and <join>__*

- These nodes are used for parallel executions
- There can be unlimited nesting of fork-join pairs
- Fork should always end with a corresponding join node



### *__<decision> node__*

A <decision> node behavior is best described as an if-then-else-if-then-else...sequence, where the first predicate that resolves to true will determine the execution path. Unlike a
<fork> node where all execution paths are followed, only one execution path will be followed in a <decision> node

```
<decision name="decision">
<switch>
      <case to="mapReduce">${jobType eq "mapReduce"}</case>
      <case to="hive">${jobType eq "hive"}</case>
      <case to="pig">${jobType eq "pig"}</case>
      <default to="mapReduce"/>
</switch>
</decision>
```

### *__<kill> node__*

- <kill> node will make sure that the oozie workflow is killed.
- Any downstream actions will be killed, however the currently running hadoop job will be allowed to complete
- Using a <kill> node in a workflow is similar to doing a System.exit(1) (any nonzero

exit code) in Java

### *Job XML*

<job-xml> tag can be used to send configuration parameters to be passed to the actions
We can also specify multiple conf files in job-xml
<map-reduce>
<job-tracker>${jobTracker}</job-tracker>
<name-node>${nameNode}</name-node>
<job-xml>/conf_A_job.xml</job-xml>
<job-xml>/conf_B_job.xml</job-xml>
<job-xml>/conf_C_job.xml</job-xml>

In this case,
What happens if the same configuration property is defined in
multiple <job-xml> files? Settings from the later files in the list of
files override the earlier ones. In this example, if the property
mapred.queue.name is defined in all three job XML files, the value
in conf_C_job.xml will take precedence over the value in the first
two files.

Inline Configuration

ns and job-xml

Parametrization

Most EL functions are present in
https://oozie.apache.org/docs/4.0.1/WorkflowFunctionalSpec.html#a4.2_Expression_Language_Functions

- This can be achieved using :
  - User defined variables : ${jobTracker}
  - EL constants and system defined variables
    - KB, MB, GB, TB, PB : predefined long integers
    - Oozie supports variables like ${YEAR}, ${MONTH} and ${DAY} : Used commonly in coordinators
  - Hadoop Counters :
    - Can be used in all hadoop nodes (map reduce, hive, pig) to make descisions
    - ${hadoop:counters("myMRNode")["FileSystemCounters"]["FILE_BYTES_READ"]}
    - ${hadoop:counters("myPigNode")["RECORD_WRITTEN"]}
    - There are system variables representing some of the common Hadoop counters: RECORDS, MAP_IN, MAP_OUT, REDUCE_IN, and REDUCE_OUT. RECORDS is
      -and-out counters for mappers and reducers

The preferred syntax for Oozie variables is ${VAR}. But this only
works for variable names that follow Java naming conventions.

the {wf:conf('VAR')} function. If a variable name has spaces or

consistent with the syntax for EL functions while most users like

between variables and functions.

### *Job.properties*

- We can either have job.properties file or a job.xml file

Example 5-8. Sample job.properties file
nameNode= hdfs://localhost:8020
jobTracker=localhost:8032
queueName=research
oozie.use.system.libpath=true
oozie.wf.application.path=${nameNode}/user/joe/oozie/mrJob/firstWorkflow.xml

Example 5-9. Sample job.xml file
<configuration>
<property>
<name>nameNode</name>
<value>hdfs://localhost:8020</value>
</property>
<property>
<name>jobTracker</name>
<value>localhost:8032</value>
</property>
<property>

```xml
<name>queueName</name>
<value>research</value>
</property>
<property>
<name>oozie.use.system.libpath</name>
<value>true</value>
</property>
<property>
<name>oozie.wf.application.path</name>
<value>${nameNode}/user/joe/oozie/mrJob/firstWorkflow.xml</value>
</property>
</configuration>
```

## *Command line option*

- You can use the -D parameter to override the property file

oozie job -oozie http://localhost:4080/oozie/ -DqueueName=research -config job.properties –run


Config-default.xml

config-default.xml file can be used to keep all default values for all workflows
- We can have separate job.properties for each of the other workflows
- This config-default.xml file is most probably overridden by other inline configurations and -D option

## *Parameters section inside workflow.xml*

```xml
<workflow-app name="parametersWF" xmlns="uri:oozie:workflow:0.5">
<parameters>
        <property>
        <name>queueName</name>
        <value>production</value>
        </property>
</parameters>
```

NOTE :

Please do not use config-default.xml and parameters together. Use one or the other

## *Lifecycle of a workflow :*

When we give -run command on oozie, the wf will start executing right away. When we give -submit command the workflow is in the PREP stage, and will execute when we provide -start option


1. PREP
2. RUNNING
3. SUCCEEDED
4. KILLED : When there is an error in the wf
5. FAILED : When there is an Oozie error
6. SUSPENDED : When the workflow is suspended using the suspend command



# *Chapter 6*

# *Oozie coordinator*

They are scheduled execution of workflows based on trggers.. The triggers include :
1. Data trigger
2. Time trigger

Time trigger

- Basically represents a CRON job
- 3 basic components :
  - Start time
  - End time
  - Frequency

Coordinator application : coordinator.xml
Coordinator job : job launched by Oozie because of the corodinator
Oozie executes coordinator jobs whereas users write coordinator applications.

### Coordinator action :

If a coordinator job has start date as Jan 1 2021 and end date as Dec 31 2021, with a frequency of daily, then there are 365 coordinator actions.
Each coordinator action checks for data availability before submitting the workflow.

### Nominal time of a coordinator action

If a coordinator job has start date as Jan 1 2021 and end date as Dec 31 2021, with a frequency of daily, then there are 365 coordinator actions.
Lets say each coordinator action hhas to start at 8AM in the morning.
Then each of the times,
Jan 1 , 8AM
Jan 2, 8AM
Jan 3 , 8AM
Until Dec 31, 8AM is called NOMINAL TIME for each of the coordinator actions.

Nominal time of each coordinator action remains the same irrespective of the actual execution time of the coordinator action.

Upon successful submission, Oozie returns a unique coordinator job ID. Each coordinator ID has a -C at the end. At the start time for this job, Oozie initiates the creation of the coordinator action. Oozie also assigns an ID for each new action. coordinator action IDs are generated by concatenating the coordinator job ID, the @ sign, and a sequentially incrementing action number. For example, if the coordinator job ID is 0000003-140329120933279-ooziejoe- C, the first two action IDs will be  0000003-140329120933279- oozie-joe-C@1 and 0000003-140329120933279-oozie-joe-C@2.

### Coordinator action lifecycle

This is the lifecycle of every coordinator action.



- At first the action is in WAITING state, waiting for data dependency. If the data dosent appear before the timeout period, then the action reaches TIMEDOUT state.
- But if the action finds the data required, it moves to READY state, and checks the concurrency requirements set in the coordinator job (The number of coordinator actions that can run at the same time). If the concurrency requirements is met, then it its moved to SUBMITTED, and only then the workflow starts executing.

Parametrization of coordinator

Time related EL functions
Frequency related EL functions

Day and month based frequency :

coord:days(N)
coord:endOfDays(N)
coord:months(N)
coord:endOfMonths()

## Coordinator execution controls

We provide these controls under &lt;controls&gt; section inside coordinator.xml

```
<coordinator-app name="my_second" start="${startTime}" end="${endTime}"
frequency="${coord:days(1)}" timezone="UTC"
xmlns="uri:oozie:coordinator:0.4">
<controls>
        <timeout>${my_timeout}</timeout>
        <concurrency>${my_concurrency}</concurrency>
        <execution>${execution_order}</execution>
        <throttle>${materialization_throttle}</throttle>
</controls>
```

1. Throttle :
    a. The maximum number of coordinator actions that can be in WAITING state (i.e. number of coordinator actions that are waiting for data)
    b. System default = 12, you can also mention a user-set upper limit or a system set upper limit
2. timeout :
    a. Time in minutes that coordinator actions can be in waiting state before timing out
    b. Default is 7 days, and max default timeout limit is 60 days
3. execution order :
    a. FIFO : executes the earliest coordinator action first
    b. LIFO : executes the last coordinator action first
    c. LAST_ONLY : Executes only the last one and discards the rest
4. concurrency :
    a. The number of coordinatr actions that can be in RUNNING state at the same time

# Chapter 7

# Data trigger coordinators

## Dataset :
- Logical entity used to represent a set of data produced by an application
- You can define dataset with :
    ○ Directory-based
    ○ Metadata-based
        ▪ H-Catalog metadata based data dependency

### Defining a dataset

5 attributes to define a dataset :
1. name
2. initial-instance : defines the first valid instance in the dataset.
3. frequency : frequency of succesive data instances
4. uri-template : template of the dataset, with EL functions (like {YEAR}- {MONTH}- {DAY} and so on)
5. done-flag : the file the indicates that the data is complete. Ex : _SUCCESS file or a _trigger file.
    a. The done-                                                                                    -size file called _SUCCESS at the end of
       processing to indicate data completeness. If done-flag exists, but the value is specified as empty, Oozie just checks for the existence of the directory and uses that as a
       signal for completion

Example :

```
<dataset name="ds_input1" frequency="${coord:hours(6)}"
        initial-instance="2014-12-29T02:00Z">
        <uri-template>
                ${baseDataDir}/revenue_feed/${YEAR}-${MONTH}-${DAY}-${HOUR}
        </uri-template>
        <done-flag>_trigger</done-flag>
</dataset>
```

- There could be multiple datasets defined in a coordinator.
    ○ For this, we can use &lt;datasets&gt; section inside oozie coordinator.xml file
    ○ We could also define all datasets in a separate XML file and include it in the &lt;include&gt; tag in &lt;datasets&gt; tag inside coordinator.xml

```
<datasets>
<include>hdfs://localhost:8020/user/joe/shares/common_datasets.xml</include>
        <dataset name="ds_input1" frequency="${coord:hours(6)}"
                initial-instance="2014-12-29T02:00Z">
                <uri-template>
                        ${baseDataDir}/revenue_feed/${YEAR}-${MONTH}-${DAY}-${HOUR}
                </uri-template>
                <done-flag>_trigger</done-flag>
        </dataset>
```

```
</datasets>
```

The example common_datasets.xml could be as follows:
```
<datasets>
    <dataset name="ds_input2" frequency="${coord:hours(6)}"
    initial-instance="2014-12-29T02:00Z">
    <uri-template>
    ${baseDataDir}/revenue_feed/${YEAR}-${MONTH}-${DAY}-${HOUR}
    </uri-template>
    <done-flag>_trigger</done-flag>
    </dataset>
</datasets>
```

NOTE:
1. There are 2 independent timelines,
    a. The coordinator timeline
    b. Dataset timelines..


## Input events

- Whereas datasets declare data items of interest, <input-events> describe the actual instance(s) of dependent dataset for this coordinator. More specifically, a workflow will not start until all the data instances defined in the input-events are available.
- There is only one <input-events> section, but there can be multiple <data-in> sections inside it.
- <data-in > definition requires 3 attributes :
    ○ name
    ○ dataset
    ○ instance definition
        ▪ There are 2 ways to specify instances :
            Using an individual <instance > tag
            Using <start-instance> and <end-instance>
- In summary, the input-events allows a user to define the list of required datasets and the corresponding data instances


In this case, the coordinator action will be triggered only when it finds the success file in 4 previous instances of the dat aset.

```
<input-events>
    <data-in name="event_input1" dataset="ds_input1">
    <start-instance>${coord:current(-4)}</start-instance>
    <end-instance>${coord:current(-1)}</end-instance>
    </data-in>
</input-events>
```

### Output events

<output-events> specifies the data instance produced by a coordinator action. Its very similar to input -events.

```
<output-events>
<data-out name="event_output1" dataset="daily-feed">
<instance>${coord:current(0)}</instance>
</data-out>
</output-events>
```

Example :

```
<coordinator-app name="my_first_rollup_job" start="2014-01-01T02:00Z"
end="2014-12-31T02:00Z" frequency="${coord:days(1)}"
xmlns="uri:oozie:coordinator:0.4">
<datasets>
    <dataset name="ds_input1" frequency="${coord:hours(6)}"
    initial-instance="2014-12-29T02:00Z">
    <uri-template>
    ${baseDataDir}/revenue_feed/${YEAR}-${MONTH}-${DAY}-${HOUR}
    </uri-template>
    <done-flag>_trigger</done-flag>
    </dataset>
</datasets>
<input-events>
    <data-in name="event_input1" dataset="ds_input1">
            <start-instance>${coord:current(-4)}</start-instance>
            <end-instance>${coord:current(-1)}</end-instance>
    </data-in>
</input-events>
<action>
    <workflow>
            <app-path>${appBaseDir}/basic-cron</app-path>
            <property>
                    <name>nameNode</name>
                    <value>hdfs://localhost:8020</value>
            </property>
```

```xml
<property>
    <name>jobTracker</name>
    <value>localhost:8032</value>
</property>
```

```xml
<property>
    <name>jobTracker</name>
    <value>localhost:8032</value>
</property>
```

## Table 7-6. current() versus latest() comparison

| Topics | current(n) | latest(n) |
|---|---|---|
| Data checking starts from | Action nominal time | Action actual time OR the present wall clock time |
| Fixed versus Variable | Fixed. Returns the same timestamp for the same action irrespective of when it checks. | Variable. Returns different timestamps based on when the check happens. |
| Gaps in data availability | Disregards gaps in data availability. Always returns the same instance(s) of data for a given action and does not skip any data whether it exists or not. | Accounts for the gaps in data availability. Skips missing data instances. Only considers the available instances. |
| Range of 'n' | Any integer | Only '0' OR negative integer. |

Parameters passing to the workflow

1. dataIn(eventName) :
   Ex : ${coord:dataIn('event_input1')} gives
   hdfs://localhost:8020/user/joe/revenue_feed/2014-12-31-02,
   hdfs://localhost:8020/user/joe/revenue_feed/2014-12-31-08,
   hdfs://localhost:8020/user/joe/revenue_feed/2014-12-31-14,
   hdfs://localhost:8020/user/joe/revenue_feed/2014-12-31-20

   A comma concatenated string

2. dataOut(eventName)
3. nominalTime()
4. actualTime()
5. dateOffset(baseTimeStamp, skipInstance, timeUnit)
6. formatTime(timeStamp, formatString)

Examples ;

```
<workflow>
 <app-path>${wfAppPath_dpe}</app-path>

<property>
        <name>nominalTime</name>
        <value>${coord:nominalTime()}</value>
    </property>
    <property>
        <name>ENV_NAME_L</name>
        <value>${ENV_NAME_L}</value>
    </property>
    <property>
        <name>LAST_MONTH</name>
        <value>${coord:formatTime(coord:dateOffset(coord:nominalTime(), -1, 'MONTH'), 'yyyy-MM')}</value>
    </property>
    <property>
        <name>STATEMENT_INCLUDE_HISTORY</name>
        <value>N</value>
    </property>
    <property>
        <name>STATEMENT_END_DATE</name>
        <value>${coord:formatTime(coord:dateOffset(coord:nominalTime(), -2, 'DAY'), 'yyyy-MM-dd')}</value>
    </property>

</workflow>
```

## A complete coordinator

```
<coordinator-app name="my_rollup_job" start="2014-01-01T02:00Z"
end="2014-12-
xmlns="uri:oozie:coordinator:0.4">
<datasets>
    <dataset name="ds_input1" frequency="${coord:hours(6)}"
    initial-instance="2013-12-29T02:00Z">
        <uri-template>
        hdfs://localhost:8020/user/joe/revenue_feed/${YEAR}-${MONTH}-${DAY}-
        ${HOUR}
        </uri-template>
        <done-flag>_trigger</done-flag>
```

```xml
        </dataset>
        <dataset name="daily-feed" frequency="${coord:days(1)}"
        initial-instance="2013-12-29T02:00Z">
                <uri-template>
                hdfs://localhost:8020/user/joe/revenue_daily_feed/${YEAR}-${MONTH}-
                ${DAY}
                </uri-template>
        </dataset>
</datasets>
<input-events>
        <data-in name="event_input1" dataset="ds_input1">
                <start-instance>${coord:current(-4)}</start-instance>
                <end-instance>${coord:current(-1)}</end-instance>
        </data-in>
</input-events>
<output-events>
        <data-out name="event_output1" dataset="daily-feed">
                <instance>${coord:current(0)}</instance>
        </data-out>
</output-events>
<action>
        <workflow>
        <app-path>${myWFHomeInHDFS}/app</app-path>
                <property>
                        <name>myInputDirs</name>
                        <value>${coord:dataIn('event_input1')}</value>
                </property>
                <property>
                        <name>myOutputDirs</name>
                        <value>${coord:dataOut('event_output1')}</value>
                </property>
                <property>
                        <name>myNominalTime</name>
                        <value>${coord:nominalTime()}</value>
                </property>
                <property>
                        <name>myActualTime</name>
                        <value>${coord:actualTime()}</value>
                </property>
                <property>
                        <name>myPreviousInstance</name>
                        <value>${coord:dateOffset(coord:nominalTime(), -1, 'DAY')}</value>
                </property>
                <property>
                        <name>myFutureInstance</name>
                        <value>${coord:dateOffset(coord:nominalTime(), 1, 'DAY')}</value>
                </property>
                <property>
                        <name>nameNode</name>
                        <value>hdfs://localhost:8020</value>
                </property>
                <property>
                        <name>jobTracker</name>
                        <value>localhost:8032</value>
                </property>
        </workflow>
</action>
</coordinator-app>
```

# Chapter 8

# Oozie bundles

- Its mainly for operaational flexibility that can deal with stopping/resuming all coordinators related to a data pipeline all at once.
- You have to have a bundle application specification (bundle.xml), which will contain a bunch of coordinator applications

**Oozie Bundle**

Coordinator A — Workflow A — Actions: hive, pig, mr, email

Coordinator N — Workflow N — Actions: mr, hive, hive, java, shell

# Chapter 9

## Advanced Concepts

### Managing action jars

- In Oozie, the system-provided JARs are known as sharelib
- In order to use sharelib that is defined at system level use : oozie.use.system.libpath=true in workflow.properties
- In general, you can override the sharelib of any action at three levels: action, job, and system. For defining at the action level, set the property ozie.action.sharelib.for.
  #action Type# in the configuration section of the action in workflow.xml. For job level, you can define it in job.properties file as
  a key-value pair. For a system-level change, the admin can define the property in oozie-site.xml.

### Supporting User Jars
1. Using lib/ folder of workflow
2. Using the variable oozie.libpath in job.properties

### JAR precedence :
1. Application lib directory : Those defined under lib/ directory of the wf receives highest priority
2. User-level shared library : those jars defined under the path mentioned in oozie.libpath in job.properties receives 2nd priotity
3. System level shared library : those jars defined using sharelib receives lowest priority

### Oozie CRON specification
- Uses Quartz scheduler
- -based frequency definition is a string of five fields separated by white space. These five fields denote various time components such as Minute, Hour, Dayof-Month, Month, and Day-of-Week, in that order.

### Table 9-1. Oozie's cron syntax

| Field name | Allowed values | Allowed special chars |
|---|---|---|
| Minute | 0-59 | Commas (,), dashes (-), asterisks, and slashes (/) |
| Hour | 0-23 | Commas (,), dashes (-), asterisks, and slashes (/) |
| Day of Month | 1-31 | Commas (,), dashes (-), asterisks, question marks (?), slashes (/), and the letters "L" and "W" |
| Month | 1-12 or JAN-DEC | Commas (,), dashes (-), asterisks, and slashes (/) |
| Day of Week | 1-7 or SUN-SAT | Commas (,), dashes (-), asterisks, question marks (?), slashes (/), and the letter "L" |

# Chapter 10

## Developer topics

### Developing Custom EL functions

- EL functions run on Oozie server, hence its important to make efficient, fast running EL functions.
- It is actually recommended to use shell or java action to implement custom EL functions.
- The Oozie admin should approve of any new EL functions
  - The corresponding JAR has to be injected into the oozie server, and the server must be restarted.

- The custom EL function should be a normal JAVA class with functions.

Supporting custom action types

- Addition of new synchronous action types is highly discouraged
- So the recommended way of writing a heavy-duty action is to use the asynchronous model
- The steps to create a new action are :
    - Extend the ActionExecutor class and implement 5 methods
    - Writing the XML schema
- Package the class into a JAR and inject into oozie_setup.war file
- Then restart Oozie server

Overriding an Asynchronous Action Type


Chapter 11

## *Oozie CLI*

$ oozie validate my_workflow.xml
$ oozie validate my_coordinator.xml

$ oozie job -config ./job.properties –submit

$ oozie job –suspend 0000006-130606115200591-oozie-joe-W
$ oozie job –resume 0000006-130606115200591-oozie-joe-W
$ oozie job –kill 0000006-130606115200591-oozie-joe-W

$ oozie job -info 0000084-141219003455004-oozie-joe-C -len 10

$ oozie job -dryrun -config wf_job.properties

$ oozie job -

For example, the following command shows how to rerun a set of coordinator actions based on date. It also removes the old files and recalculates the data dependencies. This command reruns the actions with the nominal time between 2014-10-20T05:00Z to 2014-10-25T20:00Z and individual actions with nominal time 2014-10-28T01:00Z and 2014-10-30T22:00Z:

```
$ oozie job -rerun 0000673-120823182447665-oozie-hado-C-refresh
-date 2014-10-20T05:00Z::2014-10-25T20:00Z, 2014-10-28T01:00Z,
2014-10-30T22:00Z
```

The next command demonstrates how to rerun coordinator actions using action


the action number 4 and 7 through 10:

```
$ oozie job -rerun 0000673-120823182447665-oozie-hado-C-nocleanup
-action 4,7-10
```

***Bundle reprocessing***


```
$ oozie job -rerun 0000094-141219003455004-oozie-joe-B-coordinator test-coord
Coordinators [test-coord] of bundle 0000094-141219003455004-oozie-joe-B
are scheduled to rerun on date ranges [null].
```

```
$ oozie job -rerun 0000094-141219003455004-oozie-joe-B-coordinator test-coord
-date 2014-12-28T01:28Z
Coordinators [test-coord] of bundle 0000094-141219003455004-oozie-joe-B
are scheduled to rerun on date ranges [2014-12-28T01:28Z].
```

```
$ oozie job -rerun 0000094-141219003455004-oozie-joe-B-coordinator test-coord
-date 2014-12-28T01:28Z::2015-01-06T00:30Z
Coordinators [test-coord] of bundle 0000094-141219003455004-oozie-joe-B
are scheduled to rerun on date ranges [2014-12-28T01:28Z::2015-01-06T00:30Z].
```

```
$ oozie job -rerun 0000094-141219003455004-oozie-joe-B-date 2014-12-28T01:28Z
All coordinators of bundle 0000094-141219003455004-oozie-joe-B are scheduled
to rerun on the date ranges [2014-12-28T01:28Z].
```