

## Azkaban Hadoop – A workflow scheduler for Hadoop

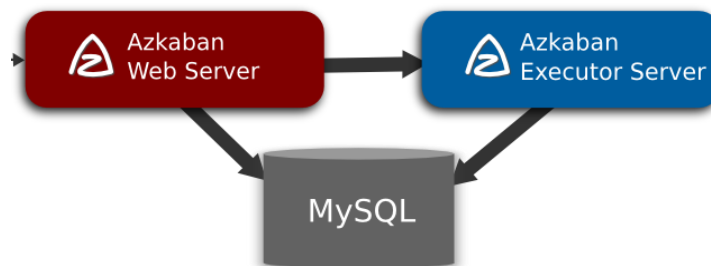
Azkaban is simple batch scheduler for constructing and running Hadoop jobs or other offline processes.

A workflow scheduler allows you to string together a group of processes to run in an order that respects the dependencies between the jobs.

Batch jobs need to be scheduled to run periodically. They also typically have intricate dependency chains—for example, dependencies on various data extraction processes or previous steps. Larger processes might have 50 or 60 steps, of which some might run in parallel and others must wait for the output of earlier steps. Azkaban is a workflow scheduler that allows the independent pieces to be declaratively assembled into a single workflow, and for that workflow to be scheduled to run periodically.

### **Azkaban consists of 3 key components:**

- Relational Database (MySQL)
- AzkabanWebServer
- AzkabanExecutorServer



### **Relational Database (MySQL)**

Azkaban uses MySQL to store much of its state. Both the AzkabanWebServer and the AzkabanExecutorServer access the DB.

### **How does AzkabanWebServer use the DB?**

The web server uses the db for the following reasons:

- **Project Management** - The projects, the permissions on the projects as well as the uploaded files.
- **Executing Flow State** - Keep track of executing flows and which Executor is running them.
- **Previous Flow/Jobs** - Search through previous executions of jobs and flows as well as access their log files.
- **Scheduler** - Keeps the state of the scheduled jobs.
- **SLA** - Keeps all the sla rules

## How does the AzkabanExecutorServer use the DB?

The executor server uses the db for the following reasons:

- **Access the project** - Retrieves project files from the db.
- **Executing Flows/Jobs** - Retrieves and updates data for flows and that are executing
- **Logs** - Stores the output logs for jobs and flows into the db.
- **Interflow dependency** - If a flow is running on a different executor, it will take state from the DB.

## AzkabanWebServer

The AzkabanWebServer is the main manager to all of Azkaban. It handles project management, authentication, scheduler, and monitoring of executions. It also serves as the web user interface.

Using Azkaban is easy. Azkaban uses \*.job key-value property files to define individual tasks in a work flow, and the \_dependencies\_ property to define the dependency chain of the jobs. These job files and associated code can be archived into a \*.zip and uploaded through the web server through the Azkaban UI or through curl.

## AzkabanExecutorServer

Previous versions of Azkaban had both the AzkabanWebServer and the AzkabanExecutorServer features in a single server. The Executor has since been separated into its own server.

## Getting Started

In version 3.0 there are three modes: the stand alone "solo-server" mode, the heavier weight two server mode and distributed multiple-executor mode. The following describes the differences between the two modes.

In **solo server mode**, the DB is embedded H2 and both web server and executor server run in the same process. This should be useful if one just wants to try things out. It can also be used on small scale use cases.

The **two server mode** is for more serious production environment. Its DB should be backed by MySQL instances with master-slave set up. The web server and executor server should run in different processes so that upgrading and maintenance shouldn't affect users.

The **multiple executor mode** is for most serious production environment. Its DB should be backed by MySQL instances with master-slave set up. The web server and executor servers should ideally run in different hosts so that upgrading and maintenance shouldn't affect users. This multiple host setup brings in robust and scalable aspect to Azkaban.

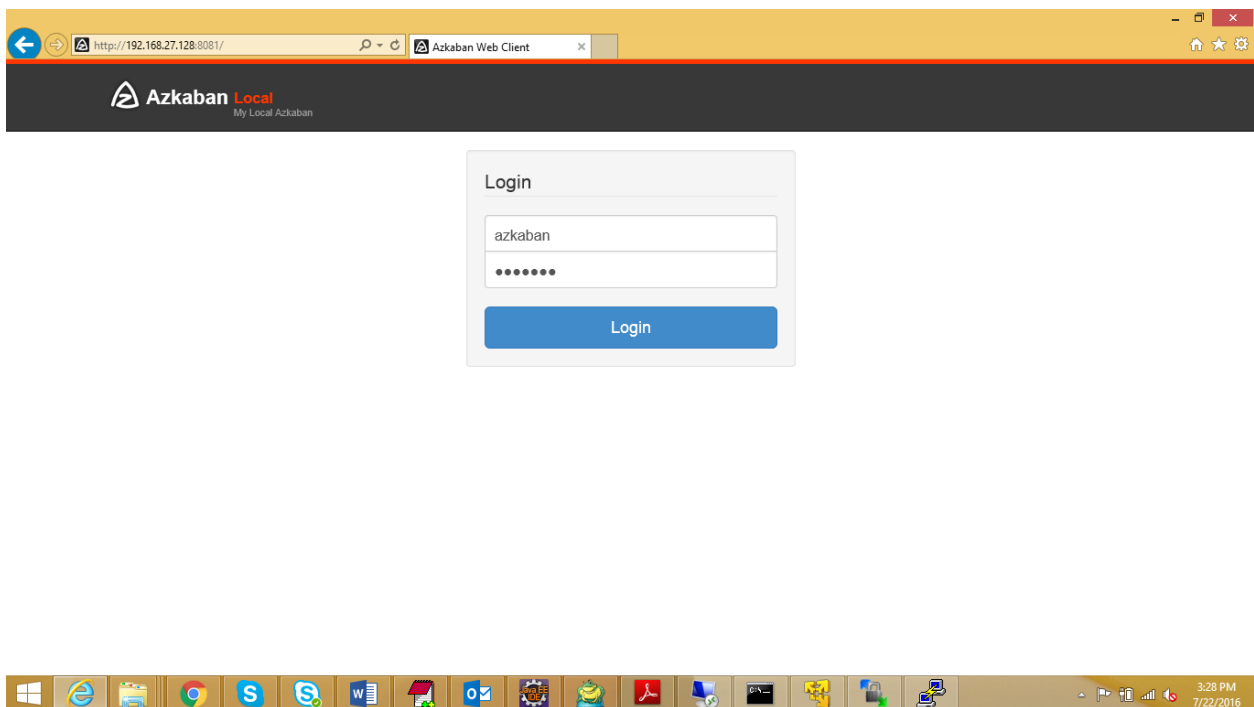
## **Azkaban Features**

- Time based dependency scheduling of workflows for hadoop jobs and compatible with any version of hadoop.
- Provides easy to use web UI and rich set of virtualizations for displaying interactive graph in browser.
- Provides support for e-mail alert on failure and successes of workflow and also supports SLA(Service Level Agreement) alerting via e-mails.
- We can retry running the failed jobs again via browser itself.
- Tracks user actions and supports authentication.
- Provides separate project work spaces for each project for easier future references.

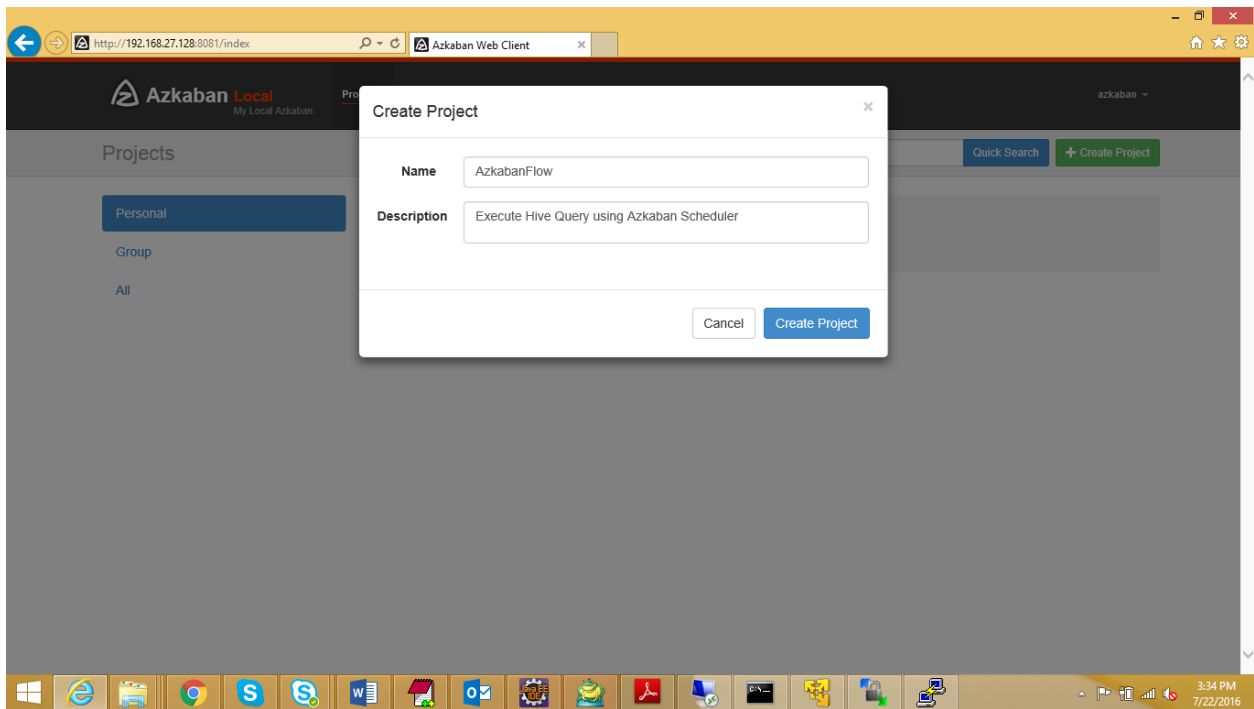
## **Creating Workflows in Azkaban**

Hit Url: <http://192.168.27.128:8081/manager?project=test>

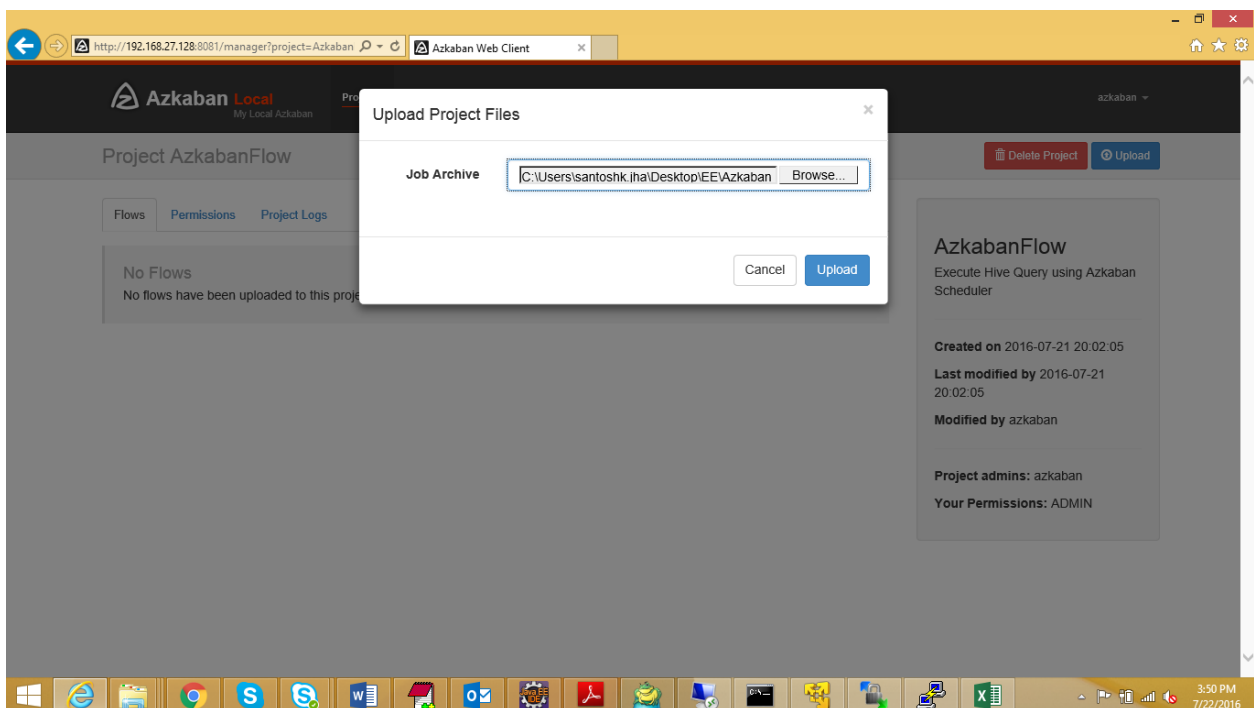
Login credential : azkaban/Azkaban



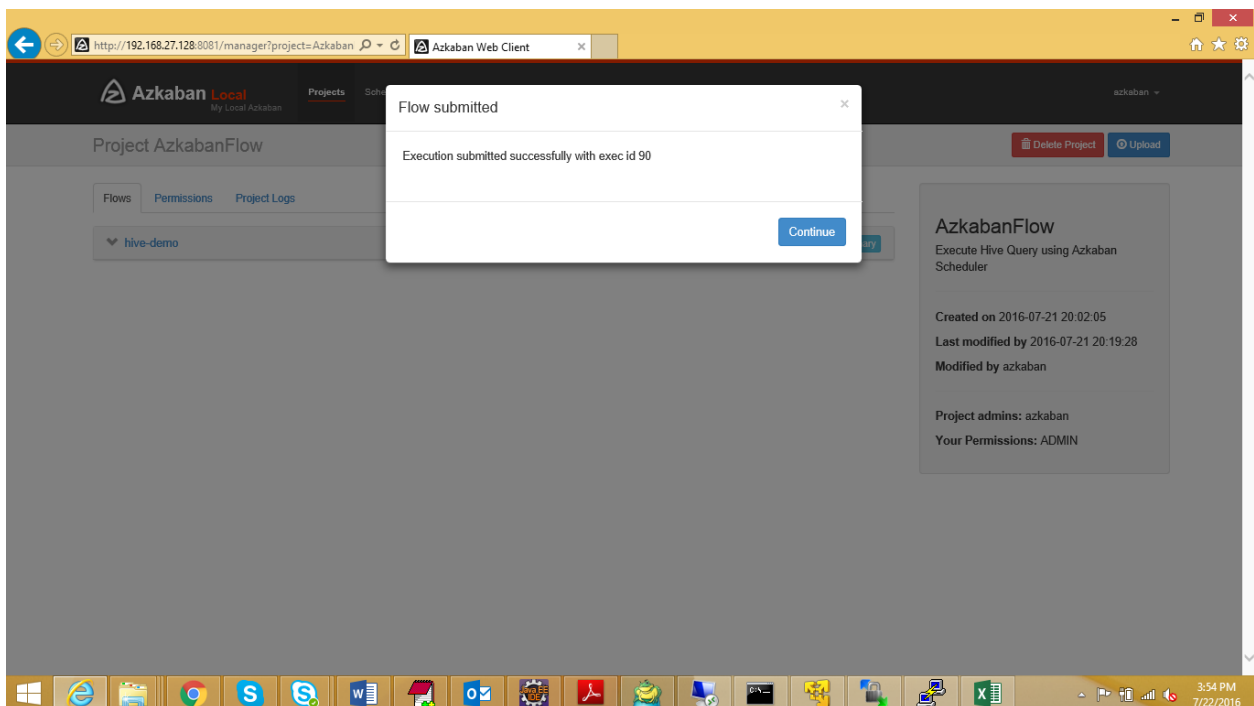
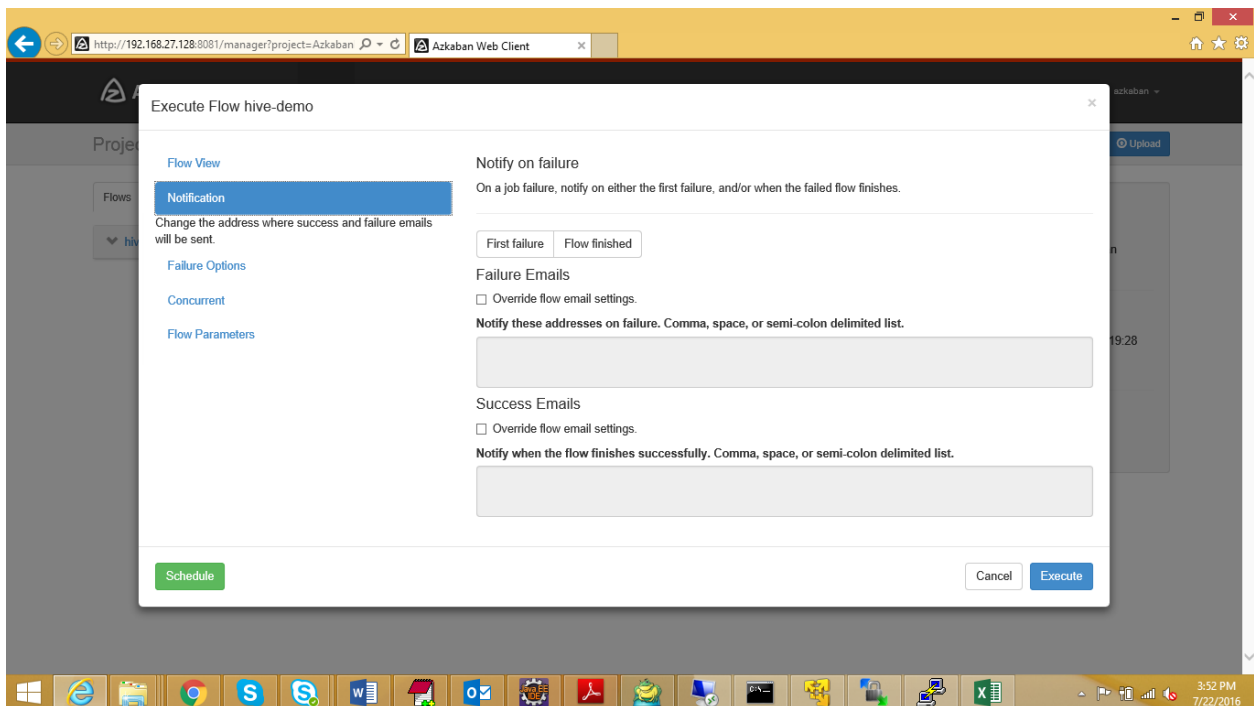
Create Project :



Upload Job file including hive query in zip format :



## Execute or Schedule Job :



A Windows taskbar with a light blue background. It contains several application icons: Internet Explorer, a folder icon, Google Chrome, Skype, a green checkmark icon, Microsoft Word, a calendar icon, Outlook, a puzzle icon, a red Adobe Reader icon, a yellow cartoon character, a blue folder icon, a black icon with a white 'X', a yellow puzzle icon, a blue folder icon, and a blue folder icon. On the right side, there is a system tray with icons for network, volume, and battery, along with the time and date: 3:58 PM, 7/22/2016.

The image is a screenshot of a web browser displaying the Azkaban Web Client. The browser's address bar shows the URL 'http://192.168.27.128:8081/executor?execid=908job=...'. The page header includes the Azkaban logo, navigation tabs for 'Projects', 'Scheduling', 'Executing' (which is active), and 'History'. A user profile dropdown for 'azkaban' is visible in the top right. The main content area is titled 'Job Execution hive-demo' and includes a 'Job Properties' button. Below the title, there are links for 'Project AzkabanFlow', 'Flow hive-demo', 'Execution 90', and 'Job hive-demo'. The 'Job Logs' tab is selected, showing a list of log entries. Each entry includes a timestamp, job name, and log message. The logs show the job starting, building the executor, setting up secure proxy info, and then executing commands like 'java -Dhive.querylog.location...' and 'hadoop jar ...'. The logs also include deprecation warnings for various Hadoop configuration properties. The bottom of the image shows a Windows taskbar with various application icons and a system clock indicating 3:58 PM on 7/22/2016.