# OozieProposal

## Abstract

Oozie is a server-based workflow scheduling and coordination system to manage data processing jobs for Apache Hadoop<sup>TM</sup>.

## Proposal

Oozie is an extensible, scalable and reliable system to define, manage, schedule, and execute complex Hadoop workloads via web services. More specifically, this includes:

- XML-based declarative framework to specify a job or a complex workflow of dependent jobs.
- Support different types of job such as Hadoop Map-Reduce, Pipe, Streaming, Pig, Hive and custom java applications.
- Workflow scheduling based on frequency and/or data availability.
- Monitoring capability, automatic retry and failure handing of jobs.
- Extensible and pluggable architecture to allow arbitrary grid programming paradigms.
- Authentication, authorization, and capacity-aware load throttling to allow multi-tenant software as a service.

## Background

Most data processing applications require multiple jobs to achieve their goals, with inherent dependencies among the jobs. A dependency could be sequential, where one job can only start after another job has finished. Or it could be conditional, where the execution of a job depends on the return value or status of another job. In other cases, parallel execution of multiple jobs may be permitted – or desired – to exploit the massive pool of compute nodes provided by Hadoop.

These job dependencies are often expressed as a Directed Acyclic Graph, also called a workflow. A node in the workflow is typically a job (a computation on the grid) or another type of action such as an eMail notification. Computations can be expressed in map/reduce, Pig, Hive or any other programming paradigm available on the grid. Edges of the graph represent transitions from one node to the next, as the execution of a workflow proceeds.

Describing a workflow in a declarative way has the advantage of decoupling job dependencies and execution control from application logic. Furthermore, the workflow is modularized into jobs that can be reused within the same workflow or across different workflows. Execution of the workflow is then driven by a runtime system without understanding the application logic of the jobs. This runtime system specializes in reliable and predictable execution: It can retry actions that have failed or invoke a cleanup action after termination of the workflow; it can monitor progress, success, or failure of a workflow, and send appropriate alerts to an administrator. The application developer is relieved from implementing these generic procedures.

Furthermore, some applications or workflows need to run in periodic intervals or when dependent data is available. For example, a workflow could be executed every day as soon as output data from the previous 24 instances of another, hourly workflow is available. The workflow coordinator provides such scheduling features, along with prioritization, load balancing and throttling to optimize utilization of resources in the cluster. This makes it easier to maintain, control, and coordinate complex data applications.

Nearly three years ago, a team of Yahoo! developers addressed these critical requirements for Hadoop-based data processing systems by developing a new workflow management and scheduling system called Oozie. While it was initially developed as a Yahoo!-internal project, it was designed and implemented with the intention of open-sourcing. Oozie was released as a GitHub project in early 2010. Oozie is used in production within Yahoo and since it has been open-sourced it has been gaining adoption with external developers

## Rationale

Commonly, applications that run on Hadoop require multiple Hadoop jobs in order to obtain the desired results. Furthermore, these Hadoop jobs are commonly a combination of Java map-reduce jobs, Streaming map-reduce jobs, Pipes map-reduce jobs, Pig jobs, Hive jobs, HDFS operations, Java programs and shell scripts.

Because of this, developers find themselves writing ad-hoc glue programs to combine these Hadoop jobs. These ad-hoc programs are difficult to schedule, manage, monitor and recover.

Workflow management and scheduling is an essential feature for large-scale data processing applications. Such applications could write the customized solution that would require separate development, operational, and maintenance overhead. Since it is a prevalent use-case for data processing, the application developer would surely prefer a generalized solution with little or no such overhead. Oozie addresses the challenge by providing an execution framework to flexibly specify the job dependency, data dependency, and time dependency. In addition, Oozie provides a multi-tenant-based centralized service and the opportunity to optimize load and utilization while respecting SLAs.

Oozie is built on Apache Hadoop<sup>TM</sup> to schedule jobs related to various Apache projects such as Hadoop, Pig, and Hive. As an Apache Open source project, Oozie is expected to attract the larger and more diversified community that currently uses such Apache sponsored projects. Additionally, users of the Hadoop ecosystem can influence Oozie's roadmap, and contribute to it. Likewise, Oozie, as part of the Apache Hadoop ^TM^ecosystem, will be a great benefit to the current Hadoop/Pig/Hive/HBase/HCatalog community.

## Current Status

### Meritocracy

Oozie currently is a github-based open sourced project where developers from multiple companies are contributing to the project. Our intent with this incubator proposal is to further extend this diverse developer community around Oozie following the Apache meritocracy model. We plan to continue to provide adequate support to new developers and to quickly recruit those who make solid contributions to committer status. In addition, Oozie will expect, accept, and work to attract contributions from amateurs as well.

## Community

While an efficient workflow management and scheduling system is critical for large companies with huge data processing in multi-tenant clusters, it is equally necessary for any non-trivial deployment. Different companies are currently using Oozie as a workflow scheduler for Hadoop-based data processing. At Yahoo! it is being used extensively in production clusters to process thousand of jobs. Like the Oozie user community, the Oozie developer community is also very strong. Developers from Yahoo! provided the initial code base, and they are still the most active contributors. In late 2010, developers from Cloudera also started contributing, and currently other companies (e.g., IBM) are beginning to participate.

We currently use JIRA for issue tracking, github for code hosting and Yahoo! Groups for developer and user communications.

### Core Developers

Oozie is currently being designed and developed by four engineers from Yahoo! – Mohammad Islam, Angelo Huang, Mayank Bansal, and Andreas Neumann. In addition, many outside contributors are actively contributing in design and development. Among them, Alejandro Abdelnur from Cloudera and Chao Wang from IBM are very important contributors. All of these core developers have deep expertise in Hadoop and the Hadoop Ecosystem in general.

### Alignment

The ASF is a natural host for Oozie given that it is already the home of Hadoop, Pig, Hive, and other emerging cloud software projects. Oozie was designed to support Hadoop from the beginning in order to solve data processing challenges in Hadoop clusters. Oozie complements the existing Apache cloud computing projects by providing a flexible framework for managing complex data processing tasks.

# Known Risks

## Orphaned Products

The core developers plan to work full time on the project. There is very little risk of Oozie getting orphaned since large companies like Yahoo! are extensively using it on their production Hadoop clusters. For example, there are nearly 400 Yahoo! internal Oozie users and thousands of jobs are processed hourly through Oozie in production. In addition, there are nearly 400 active users (including Yahoo! internal and external) in the email community where nearly 15 emails are exchanged per day. Furthermore, there were more than 1500 downloads of the Oozie binary in the last eight months from the github site and a large number of downloads were conducted by other companies such as Cloudera. Oozie has three major releases and more than 15 patch releases in the last couple of years which further demonstrates Oozie as a very active project. We plan to extend and diversify this community further through Apache.

## Inexperience with Open Source

The core developers are all active users and followers of open source. They are already committers and contributors to the Oozie Github project. In addition, they are very familiar with Apache principals and philosophy for community driven software development.

## Homogeneous Developers

The core developers are from Yahoo! as well as from several other corporations, including Cloudera and IBM.

## Reliance on Salaried Developers

Currently, the developers are paid to do work on Oozie. Companies like Yahoo! and Cloudera are invested in Oozie as the solution to the workflow management and scheduling problem in Hadoop clusters, and that is not likely to change. In addition, since workflow management is very important for most hadoop based data processing, non-salaried developers and researchers from various institutes are expected to contribute to the project.

## Relationships with Other Apache Products

Oozie is based on Apache Hadoop to manage jobs created by different Apache projects such as Hadoop, Pig, and Hive. Users of these products are extensively using Oozie as their workflow scheduler.

## An Excessive Fascination with the Apache Brand

We deeply respect the reputation of Apache and have had great success with other Apache projects such as Pig and HCatalog. We are motivated to expand and increase the adoption and development of Oozie following Apache's established open source model. We have also given reasons in the Rationale and Alignment sections.

# Documentation

Information about Oozie can be found at http://yahoo.github.com/oozie/. The following links provide more information about Oozie in open source:

- Codebase at GitHub: https://github.com/yahoo/oozie.
- JIRA : http://oozie-jira.hadoop.developer.yahoo.net
- Continuous Integration (CI) build: http://oozie-ci.hadoop.developer.yahoo.net/
- Yahoo user community: http://tech.groups.yahoo.com/group/Oozie-users/

# Initial Source

Oozie has been under development since 2009 by a team of engineers at Yahoo!. It is currently hosted on GitHub under an Apache license at https://github.com/yahoo/oozie.

# External Dependencies

The required external dependencies are all Apache License or compatible licenses. Following the components with non-Apache licenses are enumerated:

- HSQLDB License: HSQLDB
- JDOM license: JDOM
- BSD: Serp
- CCDL v1: jaxb-api, ejb, JAF

NOTE: With the exception of HSQLDB and JDOM that are directly used by Oozie, the other listed components are transitive dependencies of other Apache components used by Oozie.

# Cryptography

Oozie supports the Kerberos authentication mechanism to access secured Hadoop services.

# Required Resources

## Mailing Lists

- oozie-private for private PMC discussions (with moderated subscriptions)
- oozie-dev
- oozie-commits
- oozie-user

## Subversion Directory

https://svn.apache.org/repos/asf/incubator/oozie

## Issue Tracking

JIRA Oozie (OOZIE)

## Other Resources

The existing code already has unit tests, so we would like a Hudson instance to run them whenever a new patch is submitted. This can be added after project creation.

# Initial Committers

- Mohammad K Islam (mislam77 at yahoo dot com)
- Angelo K Huang (angelohuang at gmail dot com)
- Mayank Bansal (mabansal at gmail dot com)
- Andreas Neumann (neunand at gmail dot com)
- Alejandro Abdelnur (tucu00 at gmail dot com)
- Chao Wang (brookwc at gmail dot com)

# Affiliations

- Mohammad K Islam (Yahoo!)
- Angelo Huang (Yahoo!)
- Mayank Bansal (Yahoo!)
- Andreas Neumann (Yahoo!)
- Alejandro Abdelnur (Cloudera)
- Chao Wang (IBM)

# Sponsors

## Champion

Alan Gates

## Nominated Mentors

- Owen O'Malley (Incubator PMC member)
- Alan Gates (Incubator PMC member)
- Christopher Douglas(Incubator PMC member)
- Devaraj Das (Hadoop PMC member)

## Sponsoring Entity

We are requesting the Incubator to sponsor this project.