



# Hadoop Workflow System

**Arun C. Murthy**

CCDI, Yahoo!

[acm@yahoo-inc.com](mailto:acm@yahoo-inc.com)

YAHOO!



# Hadoop Overview

---

- Open Source Apache Project
  - <http://hadoop.apache.org/core>
- Hadoop Core includes:
  - Hadoop Distributed File System - distributes data
  - Map/Reduce – parallel processing framework
- Written in Java
- Runs on
  - Linux, Mac OS/X, Windows, and Solaris



# Workflow for Hadoop

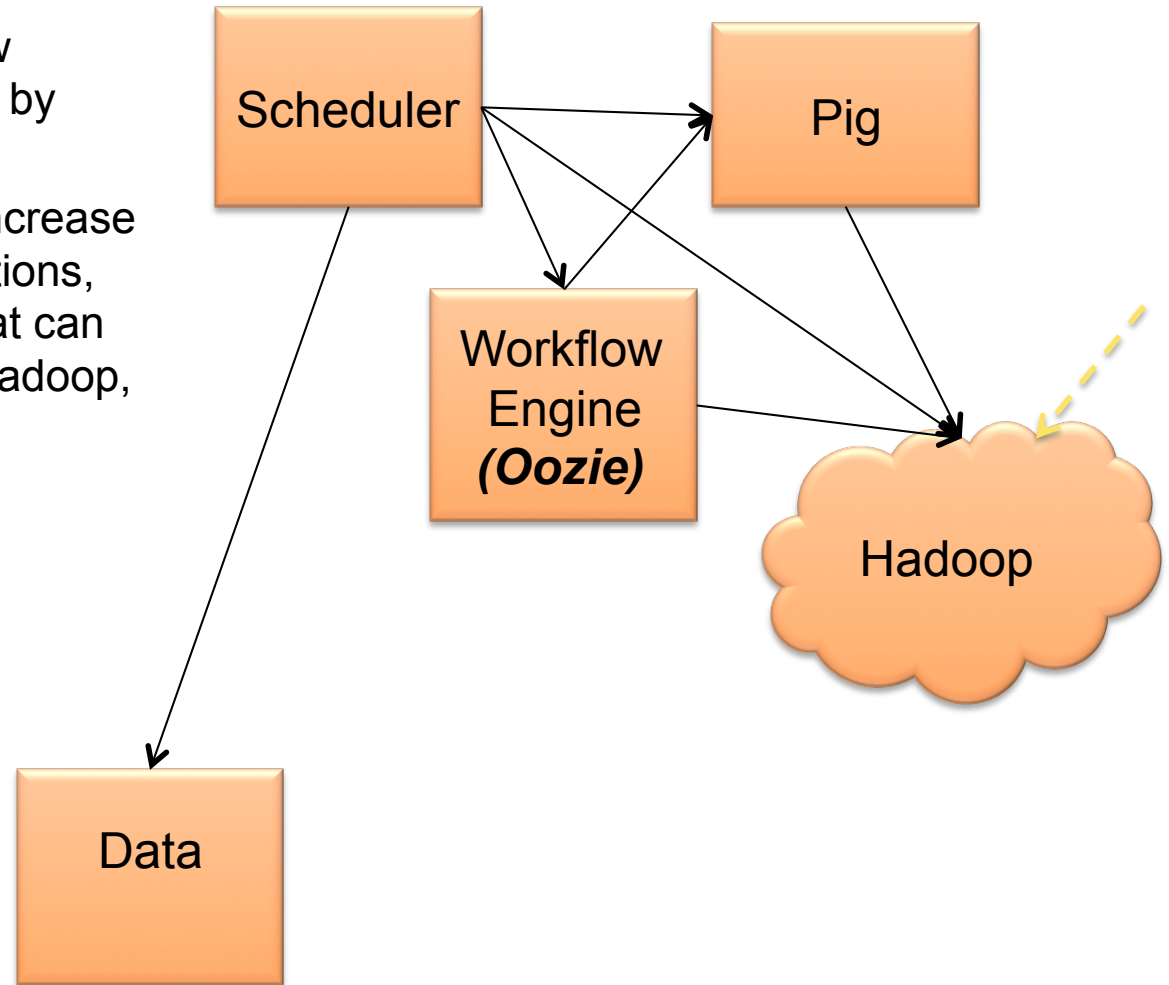
---

- Tens of thousands of Hadoop Map-Reduce, Pig jobs are being run in production today at Yahoo! and other companies
- Adhoc submissions, crontabs etc.
- Applications have complicated dependency graphs
- Need an automated system for managing complex workflows
- The next step in evolution of Hadoop:
- Work in Progress:
  - <http://issues.apache.org/jira/browse/HADOOP-5303>

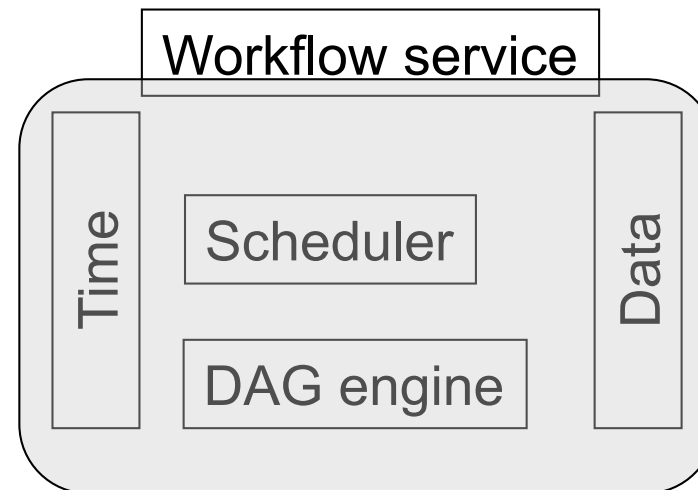


# The Bigger Picture

- The primary role of the workflow engine is to coordinate jobs run by external job-execution engines
- To improve operability, and to increase our chances for future optimizations, we will *limit* the kinds of jobs that can be run from workflows to Pig, Hadoop, and perhaps a few others



- The *scheduler* starts workflow jobs
  - Event-based triggers: *time*, *data* availability, other workflows, external events
- The *DAG engine* runs those jobs
- The *state* includes the following
  - Workflow definitions for execution
  - State of workflow jobs being executed
- User view
  - Externally, the workflow service exports one, simple workflow-service API





YAHOO!