# Apache Big Data

projects.apache.org

**processing**

Apache Pig

HIVE

Apache Solr

Lucene

Cassandra

HBASE

OOZIE

CouchDB relax

hadoop

GORA

APACHE CAYENNE

**storage**

accumulo

Apache Zookeeper

SQOOP

WHIRR

---

| | | | |
|---|---|---|---|
| **Accumulo** - a sorted, distributed key/value store<br>http://accumulo.apache.org/ | | **Flume** - collection & import of log and event data<br>http://flume.apache.org/ | |
| **Cassandra** - column-oriented database<br>http://cassandra.apache.org/ | | **Lucene** - indexing and search<br>http://lucene.apache.org/ | |
| **Cayenne** - object-relational mapping (ORM)<br>http://cayenne.apache.org/ | | **Mahout** - library for machine learning & data mining<br>http://mahout.apache.org/ | |
| **CouchDB** - NoSQL document-oriented datastore<br>http://couchdb.apache.org/ | | **Pig** - high-level programming language for Hadoop<br>http://pig.apache.org/ | |
| **Gora** - in-memory data model & persistence<br>http://gora.apache.org/ | | **Oozie** - workflow management for Hadoop<br>http://oozie.apache.org/ | |
| **Hadoop** - a distributed computing platform:<br>• HDFS - distributed file system for Hadoop<br>• MapReduce - parallel computation on clusters<br>http://hadoop.apache.org/ | | **Solr** - Lucene-based enterprise search platform<br>http://lucene.apache.org/solr/ | |
| **HBase** - column-oriented database on top of Hadoop<br>http://hbase.apache.org/ | | **Sqoop** - imports data from RDBMS into Hadoop<br>http://sqoop.apache.org/ | |
| | | **Whirr** - cloud-agnostic deployment of clusters<br>http://whirr.apache.org/ | |
| **Hive** - data warehouse with SQL-like access<br>http://hive.apache.org/ | | **Zookeeper** - configuration & coordination<br>http://zookeeper.apache.org/ | |

# Apache Big Data

incubator.apache.org

**Apache Ambari Project**
http://incubator.apache.org/ambari

**Ambari** - deployment, configuration & monitoring of Hadoop clusters
http://incubator.apache.org/ambari/

`V` `V` `V`

BLUR

**Blur** -  search platform for searching massive amounts of data in a cloud computing environment
http://incubator.apache.org/blur/

`V` `V` `V`

**Chukwa** - log collection & analysis framework for Hadoop clusters
http://incubator.apache.org/chukwa/

`V` `V` `V`

**Crunch** -  a Java library for writing, testing & running pipelines of MapReduce jobs
http://incubator.apache.org/crunch/

`V` `V` `V`

APACHE DRILL

**Drill** - interactive analysis of large-scale dataset in seconds
http://incubator.apache.org/drill/

`V` `V` `V`

**HCatalog** - schema & data type sharing over Pig, Hive & MapReduce
http://incubator.apache.org/hcatalog/

`V` `V` `V`

**Kafka** - distributed publish-subscribe messaging system
http://incubator.apache.org/kafka/

`V` `V` `V`

**Mesos** - a cluster manager providing resource sharing & isolation across cluster applications
http://incubator.apache.org/mesos/

`V` `V` `V`

**S4** distributed stream computing platform

**S4** - platform for processing continuous unbounded streams of data
http://incubator.apache.org/s4/

`V` `V` `V`

**Tashi** - infrastructure for service providers to build applications harnessing cluster computing resources to efficiently access repositories of rich data
http://incubator.apache.org/tashi/

`V` `V` `V`

legend

`V` `V` `V`   focus on volume

`V` `V` `V`   focus on variety

`V` `V` `V`   focus on velocity