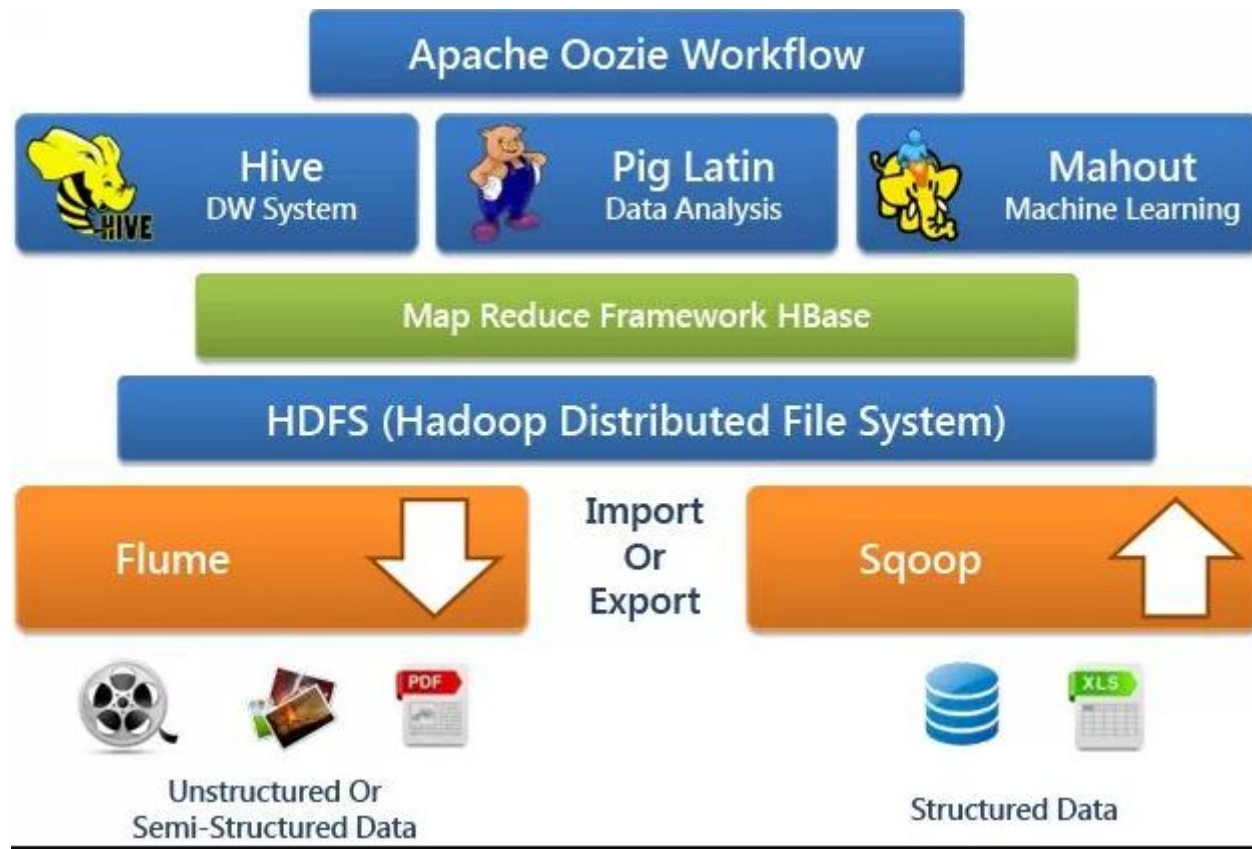
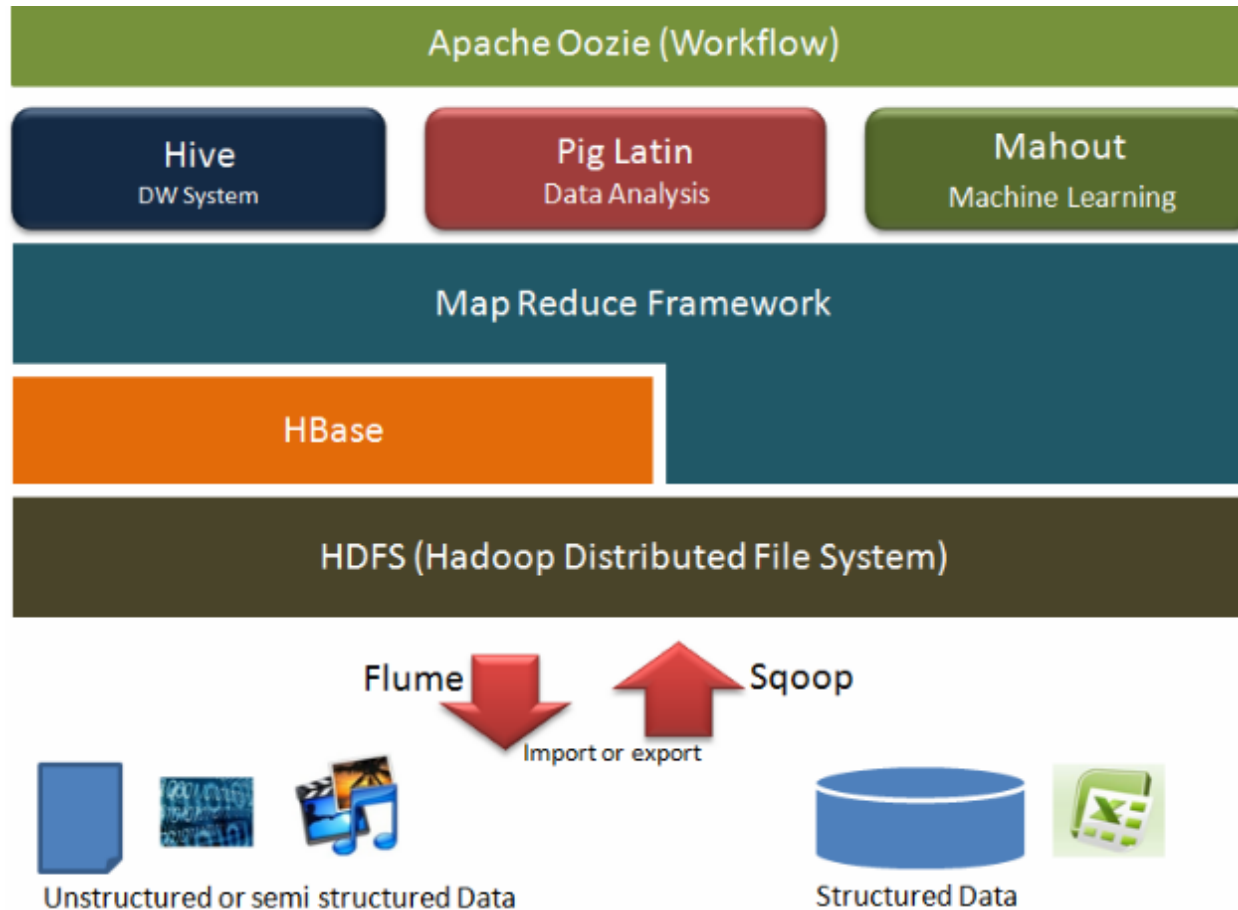


Pig, Making Hadoop Easy

# Hadoop Ecosystem

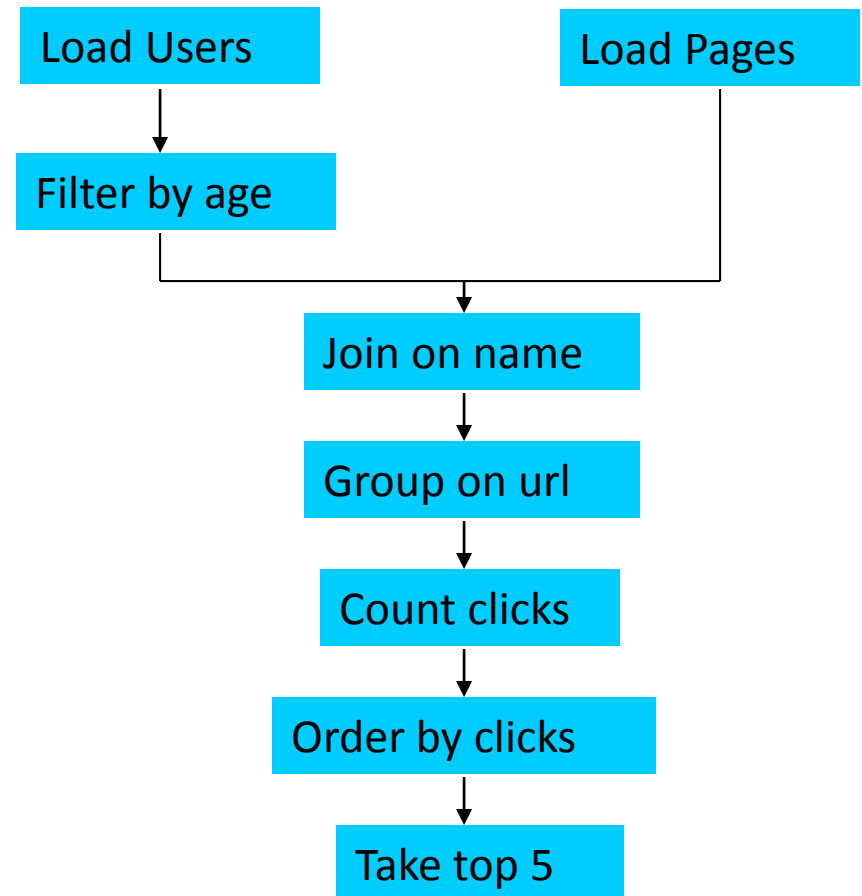


# Hadoop Ecosystem



# Motivation By Example

Suppose you have user data in one file, website data in another, and you need to find the top 5 most visited pages by users aged 18 - 25.



# In Map Reduce

```
import java.io.IOException;
import java.util.ArrayList;
import java.util.Iterator;
import java.util.List;

import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.io.Writable;
import org.apache.hadoop.io.WritableComparable;
import org.apache.hadoop.mapred.FileInputFormat;
import org.apache.hadoop.mapred.FileOutputFormat;
import org.apache.hadoop.mapred.JobConf;
import org.apache.hadoop.mapred.KeyValueTextInputFormat;
import org.apache.hadoop.mapred.Mapper;
import org.apache.hadoop.mapred.MapReduceBase;
import org.apache.hadoop.mapred.OutputCollector;
import org.apache.hadoop.mapred.RecordReader;
import org.apache.hadoop.mapred.Reducer;
import org.apache.hadoop.mapred.Reporter;
import org.apache.hadoop.mapred.SequenceFileInputFormat;
import org.apache.hadoop.mapred.SequenceFileOutputFormat;
import org.apache.hadoop.mapred.jobcontrol.Job;
import org.apache.hadoop.mapred.jobcontrol.JobC;
import org.apache.hadoop.mapred.lib.IdentityMapper;

public class MRExample {
    public static class LoadPages extends MapReduceBase
    implements Mapper<LongWritable, Text, Text, Text> {

        public void map(LongWritable k, Text val,
            OutputCollector<Text, Text> oc,
            Reporter reporter) throws IOException {
            // Pull the key out
            String line = val.toString();
            int firstComma = line.indexOf(',');
            String key = line.substring(0, firstComma);
            String value = line.substring(firstComma + 1);
            Text outKey = new Text(key);
            // Prepend an index to the value so we know which file
            // it came from.
            Text outVal = new Text(value);
            oc.collect(outKey, outVal);
        }

        public static class LoadAndFilterUsers extends MapReduceBase
        implements Mapper<LongWritable, Text, Text, Text> {

            public void map(LongWritable k, Text val,
                OutputCollector<Text, Text> oc,
                Reporter reporter) throws IOException {
                // Pull the key out
                String line = val.toString();
                int firstComma = line.indexOf(',');
                String value = line.substring(firstComma + 1);
                int age = Integer.parseInt(value);
                if (age < 18 || age > 25) return;
                String key = line.substring(0, firstComma);
                Text outKey = new Text(key);
                // Prepend an index to the value so we know which file
                // it came from.
                Text outVal = new Text("2" + value);
                oc.collect(outKey, outVal);
            }

            public static class Join extends MapReduceBase
            implements Reducer<Text, Text, Text, Text> {

                public void reduce(Text key,
                    Iterator<Text> iter,
                    OutputCollector<Text, Text> oc,
                    Reporter reporter) throws IOException {
                    // For each value, figure out which file it's from and
                    // accordingly.
                    List<String> first = new ArrayList<String>();
                    List<String> second = new ArrayList<String>();

                    while (iter.hasNext()) {
                        Text t = iter.next();
                        String value = t.toString();
                        if (value.charAt(0) == '1')
                            first.add(value.substring(1));
                        else second.add(value.substring(1));
                    }

                    reporter.setStatus("OK");
                }

                // Do the cross product and collect the values
                for (String s1 : first) {
                    for (String s2 : second) {
                        String outval = key + "," + s1 + "," + s2;
                        oc.collect(null, new Text(outval));
                        reporter.setStatus("OK");
                    }
                }
            }

            public static class LoadJoined extends MapReduceBase
            implements Mapper<Text, Text, Text, LongWritable> {

                public void map(
                    Text k,
                    Text val,
                    OutputCollector<Text, LongWritable> oc,
                    Reporter reporter) throws IOException {
                    // Find the url
                    String line = val.toString();
                    int firstComma = line.indexOf(',');
                    int secondComma = line.indexOf(',', firstComma + 1);
                    String key = line.substring(firstComma, secondComma);
                    // Drop the rest of the record, I don't need it anymore,
                    // just pass a 1 for the combiner/reducer to sum instead.
                    Text outKey = new Text(key);
                    oc.collect(outKey, new LongWritable(1L));
                }

                public static class ReduceUrls extends MapReduceBase
                implements Reducer<Text, LongWritable, WritableComparable<Text>,
                    Writable> {

                    public void reduce(
                        Text key,
                        Iterator<LongWritable> iter,
                        OutputCollector<WritableComparable, Writable> oc,
                        Reporter reporter) throws IOException {
                        // Add up all the values we see
                        long sum = 0;
                        while (iter.hasNext()) {
                            sum += iter.next().get();
                            reporter.setStatus("OK");
                        }

                        oc.collect(key, new LongWritable(sum));
                    }

                    public static class LoadClicks extends MapReduceBase
                    implements Mapper<WritableComparable, Writable, LongWritable,
                        Text> {

                        public void map(
                            WritableComparable key,
                            Writable val,
                            OutputCollector<LongWritable, Text> oc,
                            Reporter reporter) throws IOException {
                            oc.collect((LongWritable)val, (Text)key);
                        }

                        public static class LimitClicks extends MapReduceBase
                        implements Reducer<LongWritable, Text, LongWritable, Text> {

                            int count = 0;
                            public void reduce(
                                LongWritable key,
                                Iterator<Text> iter,
                                OutputCollector<LongWritable, Text> oc,
                                Reporter reporter) throws IOException {
                                    // Only output the first 100 records
                                    while (count < 100 && iter.hasNext()) {
                                        oc.collect(key, iter.next());
                                        count++;
                                    }
                                }

                                public static void main(String[] args) throws IOException {
                                    JobConf lp = new JobConf(MRExample.class);
                                    lp.setJobName("Load Pages");
                                    lp.setInputFormat(TextInputFormat.class);

                                    lp.setOutputKeyClass(Text.class);
                                    lp.setMapperClass(LoadPages.class);
                                    FileInputFormat.addInputPath(lp, new
                                        Path("/user/gates/pages"));
                                    FileOutputFormat.setOutputPath(lp, new
                                        Path("/user/gates/tmp/indexed_pages"));
                                    lp.setNumReduceTasks(0);
                                    Job loadPages = new Job(lp);

                                    JobConf lfu = new JobConf(MRExample.class);
                                    lfu.setJobName("Load and Filter Users");
                                    lfu.setInputFormat(TextInputFormat.class);
                                    lfu.setOutputKeyClass(Text.class);
                                    lfu.setOutputValueClass(Text.class);
                                    lfu.setMapperClass(LoadAndFilterUsers.class);
                                    FileInputFormat.addInputPath(lfu, new
                                        Path("/user/gates/users"));
                                    FileOutputFormat.setOutputPath(lfu,
                                        new Path("/user/gates/tmp/filtered_users"));
                                    lfu.setNumReduceTasks(0);
                                    Job loadUsers = new Job(lfu);

                                    JobConf join = new JobConf(MRExample.class);
                                    join.setJobName("Join Users and Pages");
                                    join.setInputFormat(KeyValueTextInputFormat.class);
                                    join.setMapperClass(Join.class);
                                    join.setOutputValueClass(Text.class);
                                    join.setMapperClass(IdentityMapper.class);
                                    join.setReducerClass(Join.class);
                                    FileInputFormat.addInputPath(join, new
                                        Path("/user/gates/tmp/indexed_pages"));
                                    FileInputFormat.addInputPath(join, new
                                        Path("/user/gates/tmp/filtered_users"));
                                    FileOutputFormat.setOutputPath(join, new
                                        Path("/user/gates/tmp/joined"));
                                    join.setNumReduceTasks(50);
                                    Job joinJob = new Job(join);
                                    joinJob.addDependingJob(loadPages);
                                    joinJob.addDependingJob(loadUsers);

                                    JobConf group = new JobConf(MRExample.class);
                                    group.setJobName("Group URLs");
                                    group.setInputFormat(KeyValueTextInputFormat.class);
                                    group.setMapperClass(LoadJoined.class);
                                    group.setOutputKeyClass(Text.class);
                                    group.setOutputValueClass(LongWritable.class);
                                    group.setReducerClass(ReduceUrls.class);
                                    group.setCombinerClass(ReduceUrls.class);
                                    FileInputFormat.addInputPath(group, new
                                        Path("/user/gates/tmp/joined"));
                                    FileOutputFormat.setOutputPath(group, new
                                        Path("/user/gates/tmp/grouped"));
                                    group.setNumReduceTasks(50);
                                    Job groupJob = new Job(group);
                                    groupJob.addDependingJob(joinJob);

                                    JobConf top100 = new JobConf(MRExample.class);
                                    top100.setJobName("Top 100 sites");
                                    top100.setInputFormat(SequenceFileInputFormat.class);
                                    top100.setMapperClass(LimitClicks.class);
                                    top100.setOutputValueClass(Text.class);
                                    top100.setOutputFormat(SequenceFileOutputFormat.class);
                                    top100.setMapperClass(LoadClicks.class);
                                    top100.setCombinerClass(LimitClicks.class);
                                    top100.setReducerClass(LimitClicks.class);
                                    FileInputFormat.addInputPath(top100, new
                                        Path("/user/gates/tmp/grouped"));
                                    FileOutputFormat.setOutputPath(top100, new
                                        Path("/user/gates/top100/sitesforusers18to25"));
                                    top100.setNumReduceTasks(1);
                                    Job limit = new Job(top100);
                                    limit.addDependingJob(groupJob);

                                    JobControl jc = new JobControl("Find sites for use:
                                        18 to 25");
                                    jc.addJob(loadPages);
                                    jc.addJob(loadUsers);
                                    jc.addJob(joinJob);
                                    jc.addJob(groupJob);
                                    jc.addJob(limit);
                                    jc.run();
                                }
                            }
                        }
                    }
                }
            }
        }
    }
}
```

# In Pig Latin

```
Users = load 'users' as (name, age);
Fltrd = filter Users by
    age >= 18 and age <= 25;
Pages = load 'pages' as (user, url);
Jnd = join Fltrd by name, Pages by user;
Grpd = group Jnd by url;
Smmd = foreach Grpd generate group,
    COUNT(Jnd) as clicks;
Srted = order Smmd by clicks desc;
Top5 = limit Srted 5;
store Top5 into 'top5sites';
```

# Pig and Pig Latin

# What is Pig?

- Apache Pig is a Hadoop platform for creating MapReduce jobs. Pig uses a high-level, SQL-like programming language named Pig Latin.

The benefits of Pig include:

- Run a MapReduce job with a few simple lines of code.
- Process structured data with a schema, or Pig can process unstructured data without a schema. (Pigs eat anything!)
- Pig Latin uses a familiar SQL-like syntax.
- Pig scripts read and write data from HDFS.
- Pig Latin is a data flow language, a logical solution for many MapReduce algorithms.



# Pig Latin

- Pig Latin is a high-level data flow scripting language.

Pig Latin scripts can be executed one of three ways:

- **Pig script:** write a Pig Latin program in a text file and execute it using the pig executable.
- **Grunt shell:** enter Pig statements manually one-at-a-time from a CLI tool known as the Grunt interactive shell.
- **Embedded in Java:** use the PigServer class to execute a Pig query from within Java code.

# The Grunt Shell

- Grunt is an interactive shell that enables users to enter Pig Latin statements and also interact with HDFS.
- To enter the Grunt shell, run the pig executable in the PIG\_HOME\bin folder:

```
[root@sandbox ~]# hadoop fs -put movies.txt /user/hadoop/
[root@sandbox ~]# pig
2014-07-13 11:45:17,017 [main] INFO org.apache.pig.Main - Apache Pig version 0.12.1.2.1.1.0-385 (rexported) compiled Apr 16 2014, 15:59:00
2014-07-13 11:45:17,019 [main] INFO org.apache.pig.Main - Logging error messages to: /root/pig_1405277117014.log
2014-07-13 11:45:17,056 [main] INFO org.apache.pig.impl.util.Utils - Default bootstrap file /root/.pigbootstrap not found
2014-07-13 11:45:17,852 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2014-07-13 11:45:17,853 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2014-07-13 11:45:17,853 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: hdfs://sandbox.hortonworks.com:8020
2014-07-13 11:45:19,234 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
grunt> Movies = Load '/user/hadoop/movies.txt' USING PigStorage(',') as (id,name,year,rating,duration);
grunt> Movies1 = LOAD '/user/hadoop/movies.txt' USING PigStorage(',') as (id:int,name:chararray,year:int,rating:float, duration:int);
grunt> _
```

# Pig Latin Types

Category	Type	Description
Numeric	int	32-bit signed integer
	long	64-bit signed integer
	float	32-bit floating-point number
	double	64-bit floating-point number
Text	chararray	Character array
Binary	bytearray	Byte array
Complex	tuple	Sequence of fields of any type
	bag	An unordered collection of tuples, possibly with duplicates
	map	A set of key-value pairs. Keys must be character arrays; values may be any type

# Who uses Pig for What?

- 70% of production jobs at Yahoo
- Also used by Twitter, LinkedIn, Ebay, AOL, ...
- Used to
  - Process web logs
  - Build user behavior models
  - Process images
  - Build maps of the web
  - Do research on raw data sets