# MongoDB Plugin & Toolchain for Hadoop

Brendan W. McAdams

10gen, Inc.

Dec. 3, 2010 @ MongoSV

# Hadoop Explained. . .

- Started in February 2006 as part of the Apache Lucene project
- Based upon Google's MapReduce and GFS Papers
- Allows distributed, scalable data processing of huge datasets
- Java based, but with support for other JVM and Non-JVM Languages
- Lots of ecosystem tools to simplify interaction such as Pig and Hive
- In use at New York Times, Last.fm, Yahoo!, Amazon, Facebook and *many more companies. . .*
- Great tools for temporary Hadoop clusters such as the Cloudera Cluster Tools, Apache Whirr and Amazon's Elastic MapReduce.

$\bigcirc$ mongoDB

# Why Integrate MongoDB?

- Language: JavaScript only; not everyone wants to write JavaScript for data processing.
  - Static Typing: JavaScript is dynamically typed which has limited value for jobs where a stronger type system is desired.
  - JVM Ecosystem: The JVM has a large number of libraries available for assisting in calculations and analysis which are available from Hadoop.
- Concurrency: Current JS Implementation is limited to one JS execution per server at a time.
- Scalability: Not a lot of ability to scale MongoDB's MapReduce except in cases of sharding.
- Integration: More development shops are integrating their data analysis into Hadoop jobs—opportunity to merge data from multiple systems including MongoDB for analysis.
- Hadoop is not a panacea—queries won't necessarily be fast(er) but it can handle larger scale and free up your MongoDB servers for long jobs.

# Why Integrate MongoDB?

- Language: JavaScript only; not everyone wants to write JavaScript for data processing.
  - Static Typing: JavaScript is dynamically typed which has limited value for jobs where a stronger type system is desired.
  - JVM Ecosystem: The JVM has a large number of libraries available for assisting in calculations and analysis which are available from Hadoop.
- Concurrency: Current JS Implementation is limited to one JS execution per server at a time.
- Scalability: Not a lot of ability to scale MongoDB's MapReduce except in cases of sharding.
- Integration: More development shops are integrating their data analysis into Hadoop jobs—opportunity to merge data from multiple systems including MongoDB for analysis.
- Hadoop is not a panacea—queries won't necessarily be fast(er) but it can handle larger scale and free up your MongoDB servers for long jobs.

# Why Integrate MongoDB?

- Language: JavaScript only; not everyone wants to write JavaScript for data processing.
  - Static Typing: JavaScript is dynamically typed which has limited value for jobs where a stronger type system is desired.
  - JVM Ecosystem: The JVM has a large number of libraries available for assisting in calculations and analysis which are available from Hadoop.
- Concurrency: Current JS Implementation is limited to one JS execution per server at a time.
- Scalability: Not a lot of ability to scale MongoDB's MapReduce except in cases of sharding.
- Integration: More development shops are integrating their data analysis into Hadoop jobs—opportunity to merge data from multiple systems including MongoDB for analysis.
- Hadoop is not a panacea—queries won't necessarily be fast(er) but it can handle larger scale and free up your MongoDB servers for long jobs.

# Why Integrate MongoDB?

- Language: JavaScript only; not everyone wants to write JavaScript for data processing.
  - Static Typing: JavaScript is dynamically typed which has limited value for jobs where a stronger type system is desired.
  - JVM Ecosystem: The JVM has a large number of libraries available for assisting in calculations and analysis which are available from Hadoop.
- Concurrency: Current JS Implementation is limited to one JS execution per server at a time.
- Scalability: Not a lot of ability to scale MongoDB's MapReduce except in cases of sharding.
- Integration: More development shops are integrating their data analysis into Hadoop jobs—opportunity to merge data from multiple systems including MongoDB for analysis.
- Hadoop is not a panacea—queries won't necessarily be fast(er) but it can handle larger scale and free up your MongoDB servers for long jobs.

# Why Integrate MongoDB?

- Language: JavaScript only; not everyone wants to write JavaScript for data processing.
  - Static Typing: JavaScript is dynamically typed which has limited value for jobs where a stronger type system is desired.
  - JVM Ecosystem: The JVM has a large number of libraries available for assisting in calculations and analysis which are available from Hadoop.
- Concurrency: Current JS Implementation is limited to one JS execution per server at a time.
- Scalability: Not a lot of ability to scale MongoDB's MapReduce except in cases of sharding.
- Integration: More development shops are integrating their data analysis into Hadoop jobs—opportunity to merge data from multiple systems including MongoDB for analysis.
- Hadoop is not a panacea—queries won't necessarily be fast(er) but it can handle larger scale and free up your MongoDB servers for long jobs.

# Why Integrate MongoDB?

- Language: JavaScript only; not everyone wants to write JavaScript for data processing.
  - Static Typing: JavaScript is dynamically typed which has limited value for jobs where a stronger type system is desired.
  - JVM Ecosystem: The JVM has a large number of libraries available for assisting in calculations and analysis which are available from Hadoop.
- Concurrency: Current JS Implementation is limited to one JS execution per server at a time.
- Scalability: Not a lot of ability to scale MongoDB's MapReduce except in cases of sharding.
- Integration: More development shops are integrating their data analysis into Hadoop jobs—opportunity to merge data from multiple systems including MongoDB for analysis.
- Hadoop is not a panacea—queries won't necessarily be fast(er) but it can handle larger scale and free up your MongoDB servers for long jobs.

# Why Integrate MongoDB?

- Language: JavaScript only; not everyone wants to write JavaScript for data processing.
  - Static Typing: JavaScript is dynamically typed which has limited value for jobs where a stronger type system is desired.
  - JVM Ecosystem: The JVM has a large number of libraries available for assisting in calculations and analysis which are available from Hadoop.
- Concurrency: Current JS Implementation is limited to one JS execution per server at a time.
- Scalability: Not a lot of ability to scale MongoDB's MapReduce except in cases of sharding.
- Integration: More development shops are integrating their data analysis into Hadoop jobs—opportunity to merge data from multiple systems including MongoDB for analysis.
- Hadoop is not a panacea—queries won't necessarily be fast(er) but it can handle larger scale and free up your MongoDB servers for long jobs.

## Mongo + Hadoop . . .

- Integrating MongoDB and Hadoop to read & write data from/to MongoDB via Hadoop
- 10gen has been working on a plugin to integrate the two systems, written in Pure Java
- About 6 months ago I explored the idea in Scala with a project called 'Luau'
- Support for pure MapReduce as well as Pig (Currently output only - input coming soon)
- With Hadoop Streaming (soon), write your MapReduce in Python, Ruby or Perl

$\bigodot$ mongoDB

# MongoDB + MapReduce in Java I

We're going to walk through a sample application which parses US Treasury Bond historical Bid Curves since January 1990, and calculates an annual average for the 10 year Treasury.

- A sample of our dataset:

### code/sample_treasury.js

```
0 { "_id" : { "$date" : 631238400000 }, "dayOfWeek" : "TUESDAY",
    "bc3Year" : 7.9, "bc5Year" : 7.87, "bc10Year" : 7.94,
    "bc20Year" : null, "bc1Month" : null, "bc2Year" : 7.87,
    "bc3Month" : 7.83, "bc30Year" : 8, "bc1Year" : 7.81,
    "bc7Year" : 7.98, "bc6Month" : 7.89 }
1 { "_id" : { "$date" : 631324800000 }, "dayOfWeek" : "WEDNESDAY",
    "bc3Year" : 7.96, "bc5Year" : 7.92, "bc10Year" : 7.99,
    "bc20Year" : null, "bc1Month" : null, "bc2Year" : 7.94,
    "bc3Month" : 7.89, "bc30Year" : 8.039999999999999,
    "bc1Year" : 7.85, "bc7Year" : 8.039999999999999, "bc6Month"
    : 7.94 }
```

mongoDB

# MongoDB + MapReduce in Java II

```
2  { "_id" : { "$date" : 631411200000 }, "dayOfWeek" : "THURSDAY",
       "bc3Year" : 7.93, "bc5Year" : 7.91, "bc10Year" : 7.98,
       "bc20Year" : null, "bc1Month" : null, "bc2Year" : 7.92,
       "bc3Month" : 7.84, "bc30Year" : 8.039999999999999,
       "bc1Year" : 7.82, "bc7Year" : 8.02, "bc6Month" : 7.9 }
3  { "_id" : { "$date" : 631497600000 }, "dayOfWeek" : "FRIDAY",
       "bc3Year" : 7.94, "bc5Year" : 7.92, "bc10Year" : 7.99,
       "bc20Year" : null, "bc1Month" : null, "bc2Year" : 7.9,
       "bc3Month" : 7.79, "bc30Year" : 8.06, "bc1Year" : 7.79,
       "bc7Year" : 8.029999999999999, "bc6Month" : 7.85 }
```

mongoDB

# MongoDB + MapReduce in Java III

- The MongoDB JavaScript version of our demo:

code/mongo_treasury_mr.js

```
0  function m() {
1      emit( this._id.getYear(), { count: 1, sum: this.bc10Year })
2  }
3
4  function r( year, values ) {
5      var n = { count: 0,  sum: 0 }
6      for ( var i = 0; i < values.length; i++ ){
7          n.sum += values[i].sum;
8          n.count += values[i].count;
9      }
10
11     return n;
12 }
13
14 function f( year, value ){
15     value.avg = value.sum / value.count;
16     return value.avg;
17 }
```

mongoDB

```
18
19 result = db.yield_historical.in.mapReduce( m , r, { finalize: f
      });
```

mongoDB

# MongoDB + MapReduce in Java V

- The same work can be done in Hadoop. . .

$\bigcirc$ mongoDB

footer_navigationB.W. McAdams (10gen, Inc.)     MongoDB Plugin & Toolchain for Hadoop     Dec. 3, 2010 @ MongoSV     9 / 40

The Configuration is specified in XML

## code/mongo–treasury_yield.xml

```
0  <?xml version="1.0"?>
1  <configuration>
2    <property>
3      <!-- run the job verbosely ? -->
4      <name>mongo.job.verbose</name>
5      <value>true</value>
6    </property>
7    <property>
8      <!-- Run the job in the foreground and wait for response, or
         background it? -->
9      <name>mongo.job.background</name>
10     <value>false</value>
11   </property>
12   <property>
13     <!-- If you are reading from mongo, the URI -->
14     <name>mongo.input.uri</name>
15     <value>mongodb://localhost/demo.yield_historical.in</value>
```

mongoDB

```
16    </property>
17    <property>
18      <!-- If you are writing to mongo, the URI -->
19      <name>mongo.output.uri</name>
20      <value>mongodb://localhost/demo.yield_historical.out</value>
21    </property>
22    <property>
23      <!-- The query, in JSON, to execute [OPTIONAL] -->
24      <name>mongo.input.query</name>
25      <!--<value>{"x": {"$regex": "^eliot", "$options": ""}}</value>-->
26      <value></value>
27    </property>
28    <property>
29      <!-- The fields, in JSON, to read [OPTIONAL] -->
30      <name>mongo.input.fields</name>
31      <value></value>
32    </property>
33    <property>
34      <!-- A JSON sort specification for read [OPTIONAL] -->
35      <name>mongo.input.sort</name>
36      <value></value>
```

$\bigcirc$ mongoDB

```
37    </property>
38    <property>
39      <!-- The number of documents to limit to for read [OPTIONAL] -->
40      <name>mongo.input.limit</name>
41      <value>0</value> <!-- 0 == no limit -->
42    </property>
43    <property>
44      <!-- The number of documents to skip in read [OPTIONAL] -->
45      <!-- TODO - Are we running limit() or skip() first? -->
46      <name>mongo.input.skip</name>
47      <value>0</value> <!-- 0 == no skip -->
48    </property>
49    <property>
50      <!-- Class for the mapper -->
51      <name>mongo.job.mapper</name>
52      <value>com.mongodb.hadoop.examples.TreasuryYieldMapper</value>
53    </property>
54    <property>
55      <!-- Reducer class -->
56      <name>mongo.job.reducer</name>
57      <value>com.mongodb.hadoop.examples.TreasuryYieldReducer</value>
```

mongoDB

```
58    </property>
59    <property>
60      <!-- InputFormat Class -->
61      <name>mongo.job.input.format</name>
62      <value>com.mongodb.hadoop.MongoInputFormat</value>
63    </property>
64    <property>
65      <!-- OutputFormat Class -->
66      <name>mongo.job.output.format</name>
67      <value>com.mongodb.hadoop.MongoOutputFormat</value>
68    </property>
69    <property>
70      <!-- Output key class for the output format -->
71      <name>mongo.job.output.key</name>
72      <value>org.apache.hadoop.io.IntWritable</value>
73    </property>
74    <property>
75      <!-- Output value class for the output format -->
76      <name>mongo.job.output.value</name>
77      <value>org.apache.hadoop.io.DoubleWritable</value>
78    </property>
```

```
79    <property>
80      <!-- Output key class for the mapper [optional] -->
81      <name>mongo.job.mapper.output.key</name>
82      <value></value>
83      <value>org.apache.hadoop.io.IntWritable</value>
84    </property>
85    <property>
86      <!-- Output value class for the mapper [optional] -->
87      <name>mongo.job.mapper.output.value</name>
88      <value>org.apache.hadoop.io.DoubleWritable</value>
89    </property>
90    <property>
91      <!-- Class for the combiner [optional] -->
92      <name>mongo.job.combiner</name>
93      <value>com.mongodb.hadoop.examples.TreasuryYieldReducer</value>
94    </property>
95    <property>
96      <!-- Partitioner class [optional] -->
97      <name>mongo.job.partitioner</name>
98      <value></value>
99    </property>
```

mongoDB

```
100    <property>
101      <!-- Sort Comparator class [optional] -->
102      <name>mongo.job.sort_comparator</name>
103      <value></value>
104    </property>
105
106  </configuration>
```

The Job is loaded & configured via a Java Tool:

**code/TreasuryYieldXMLConfig.java**

```java
// TreasuryYieldXMLConfig.java
/*
 * Copyright 2010 10gen Inc.
 *
 * Licensed under the Apache License, Version 2.0 (the "License");
 * ...
 */
package com.mongodb.hadoop.examples;

import java.io.*;
import java.util.*;

import org.apache.hadoop.conf.*;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
import org.apache.hadoop.util.*;
import org.bson.*;
```

```
17
18  import com.mongodb.hadoop.util.*;
19
20  public class TreasuryYieldXMLConfig extends MongoTool {
21
22      static {
23          // Load the XML config defined in hadoop-local.xml
24          Configuration.addDefaultResource(
            "src/examples/hadoop-local.xml" );
25          Configuration.addDefaultResource(
            "src/examples/mongo-defaults.xml" );
26      }
27
28      public static void main( String[] args ) throws Exception{
29          final int exitCode = ToolRunner.run( new
            TreasuryYieldXMLConfig(), args );
30          System.exit( exitCode );
31      }
32  }
```

mongoDB

# MongoDB Hadoop Job for Treasury Data I
Configuration

MongoTool is a helper utility for doing XML configs with queries, etc:

### code/MongoTool.java

```java
 0  // MongoTool.java
 1  /*
 2   * Copyright 2010 10gen Inc.
 3   *
 4   * Licensed under the Apache License, Version 2.0 (the "License");
 5   * ...
 6   */
 7
 8  package com.mongodb.hadoop.util;
 9
10  import java.util.Map.Entry;
11
12  import org.apache.commons.logging.*;
13  import org.apache.hadoop.conf.*;
14  import org.apache.hadoop.mapreduce.*;
15  import org.apache.hadoop.util.*;
16
```

$\bullet$ mongoDB

```
17  /**
18   * Tool for simplifying the setup and usage of Mongo Hadoop jobs
19   * using the Tool / Configured interfaces for use w/ a ToolRunner
20   * Primarily useful in cases of XML Config files.
21   *
22   * @author Brendan W. McAdams <brendan@10gen.com>
23   */
24  public class MongoTool extends Configured implements Tool {
25      private static final Log log = LogFactory.getLog( MongoTool.class
        );
26
27      public int run( String[] args ) throws Exception{
28          /**
29           * ToolRunner will configure/process/setup the config
30           * so we need to grab the classlevel one
31           * This will be inited with any loaded xml files or -D
        prop=value params
32           */
33          final Configuration conf = getConf();
34
```

$\bigcirc$ mongoDB

```
35      log.info( "Created a conf: '" + conf + "' on {" +
    this.getClass() + "} as job named '" + _jobName + "'" );

36
37      for ( final Entry<String, String> entry : conf ) {
38          log.trace( String.format( "%s=%s\n", entry.getKey(),
    entry.getValue() ) );
39      }

40
41      final Job job = new Job( conf , _jobName );
42      /**
43       * Any arguments specified with -D <property>=<value>
44       * on the CLI will be picked up and set here
45       * They override any XML level values
46       * Note that -D<space> is important - no space will
47       * not work as it get spicked up by Java itself
48       */
49      // TODO - Do we need to set job name somehow more
    specifically?
50      // This may or may not be correct/sane
51      job.setJarByClass( this.getClass() );
52      final Class mapper = MongoConfigUtil.getMapper( conf );
```

```
53
54        log.info( "Mapper Class: " + mapper );
55        job.setMapperClass( mapper );
56        job.setCombinerClass( MongoConfigUtil.getCombiner( conf ) );
57        job.setReducerClass( MongoConfigUtil.getReducer( conf ) );
58
59        job.setOutputFormatClass( MongoConfigUtil.getOutputFormat(
     conf ) );
60        job.setOutputKeyClass( MongoConfigUtil.getOutputKey( conf )
     );
61        job.setOutputValueClass( MongoConfigUtil.getOutputValue( conf
     ) );
62
63        job.setInputFormatClass( MongoConfigUtil.getInputFormat( conf
     ) );
64
65        /**
66         * Determines if the job will run verbosely e.g. print debug
     output
67         * Only works with foreground jobs
68         */
```

mongoDB

```
69      final boolean verbose = MongoConfigUtil.isJobVerbose( conf );
70      /**
71       * Run job in foreground aka wait for completion or
    background?
72       */
73      final boolean background = MongoConfigUtil.isJobBackground(
    conf );
74      try {
75          if ( background ) {
76              log.info( "Setting up and running MapReduce job in
    background." );
77              job.submit();
78              return 0;
79          }
80          else {
81              log.info( "Setting up and running MapReduce job in
    foreground, will wait for results.  {Verbose? " + verbose + "}"
    );
82              return job.waitForCompletion( true ) ? 0 : 1;
83          }
84      }
```

```
85          catch ( final Exception e ) {
86              log.error( "Exception while executing job... ", e );
87              return 1;
88          }
89      }
90
91      /**
92       * Main will be a necessary method to run the job - suggested
         implementation
93       * template:
94       * public static void main(String[] args) throws Exception {
95       * int exitCode = ToolRunner.run(new <YourClass>(), args);
96       * System.exit(exitCode);
97       * }
98       *
99       */
100
101     /**
102      * SET ME
103      * Defines the name of the job on the cluster.
104      * Left non-final to allow tweaking with serial #s, etc
```

```
105        **/
106        String _jobName = "<unnamed MongoTool job>";
107   }
```

# MongoDB Hadoop Job for Treasury Data I
## Mapper

The Map function:

### code/TreasuryYieldMapper.java

```
 0 // TreasuryYieldMapper.java
 1 /*
 2  * Copyright 2010 10gen Inc.
 3  *
 4  * Licensed under the Apache License, Version 2.0 (the "License");
 5  * ...
 6  */
 7 package com.mongodb.hadoop.examples;
 8
 9 import java.io.*;
10 import java.util.*;
11
12 import org.apache.commons.logging.*;
13
14 import org.apache.hadoop.conf.*;
15 import org.apache.hadoop.io.*;
16
```

mongoDB

```
17  import org.apache.hadoop.mapreduce.*;
18  import org.apache.hadoop.util.*;
19  import org.bson.*;
20
21  import com.mongodb.hadoop.util.*;
22
23  public class TreasuryYieldMapper extends Mapper<java.util.Date,
       BSONObject, IntWritable, DoubleWritable> {
24
25      private static final Log log = LogFactory.getLog(
        TreasuryYieldMapper.class );
26
27      public void map( java.util.Date key , BSONObject value , Context
        context ) throws IOException, InterruptedException{
28
29          int year = key.getYear() + 1900;
30          double bid10Year = ((Number) value.get( "bc10Year"
        )).doubleValue();
31
32          context.write( new IntWritable( year ), new DoubleWritable(
        bid10Year ) );
```

```
33
34        }
35
36 }
```

The Reduce function:

## code/TreasuryYieldReducer.java

```
0  // TreasuryYieldMapper.java
1  /*
2   * Copyright 2010 10gen Inc.
3   *
4   * Licensed under the Apache License, Version 2.0 (the "License");
5   * ...
6   */
7  package com.mongodb.hadoop.examples;
8
9  import java.io.*;
10 import java.util.*;
11
12 import org.apache.commons.logging.*;
13
14 import org.apache.hadoop.conf.*;
15 import org.apache.hadoop.io.*;
16
```

```
17  import org.apache.hadoop.mapreduce.*;
18  import org.apache.hadoop.util.*;
19  import org.bson.*;
20
21  import com.mongodb.hadoop.util.*;
22
23  public class TreasuryYieldReducer extends Reducer<IntWritable,
        DoubleWritable, IntWritable, DoubleWritable> {
24
25      private static final Log log = LogFactory.getLog(
        TreasuryYieldReducer.class );
26
27      public void reduce( IntWritable key , Iterable<DoubleWritable>
        values , Context context ) throws IOException,
        InterruptedException{
28
29          int count = 0;
30          double sum = 0;
31          for ( final DoubleWritable value : values ) {
32              log.debug( "Key: " + key + " Value: " + value );
33              sum += value.get();
```

mongoDB

```
34            count++;
35        }
36
37        double avg = sum / count;
38
39        log.info( "Average 10 Year Treasury for " + key.get() + " was
    " + avg );
40
41        context.write( key, new DoubleWritable( avg ) );
42    }
43
44 }
```

- Make sure the Java Driver is in the 'lib' directory on every cluster member, e.g. `sudo cp lib/mongo-java-driver-2.3.jar /usr/lib/hadoop-0.20/lib/` (Bounce hadoop after adding libraries)
- Build the Hadoop Jar (`ant resolve; ant jar`)
- Run the Job!

$\bigcirc$ mongoDB

# A Modified Pig Tutorial with MongoDB Output I

- We took the stock Pig Tutorial and modified it to save to MongoDB
- Input file is an anonymized Search Engine Log:

### code/excite_log_snippet.txt

| | | | |
|---|---|---|---|
| 0 | 2A9EABFB35F5B954 | 970916105432 | +md foods +proteins |
| 1 | BED75271605EBD0C | 970916001949 | yahoo chat |
| 2 | BED75271605EBD0C | 970916001954 | yahoo chat |
| 3 | BED75271605EBD0C | 970916003523 | yahoo chat |
| 4 | BED75271605EBD0C | 970916011322 | yahoo search |
| 5 | BED75271605EBD0C | 970916011404 | yahoo chat |
| 6 | BED75271605EBD0C | 970916011422 | yahoo chat |
| 7 | BED75271605EBD0C | 970916012756 | yahoo caht |
| 8 | BED75271605EBD0C | 970916012816 | yahoo chat |
| 9 | BED75271605EBD0C | 970916023603 | yahoo chat |
| 10 | BED75271605EBD0C | 970916025458 | yahoo caht |
| 11 | BED75271605EBD0C | 970916025516 | yahoo chat |
| 12 | BED75271605EBD0C | 970916030348 | yahoo chat |
| 13 | BED75271605EBD0C | 970916034807 | yahoo chat |
| 14 | BED75271605EBD0C | 970916040755 | yahoo chat |
| 15 | BED75271605EBD0C | 970916090700 | hawaii chat universe |

mongoDB

# A Modified Pig Tutorial with MongoDB Output II

```
16  BED75271605EBD0C    970916094445    yahoo chat
17  BED75271605EBD0C    970916191427    yahoo chat
18  BED75271605EBD0C    970916201045    yahoo chat
19  BED75271605EBD0C    970916201050    yahoo chat
20  BED75271605EBD0C    970916201927    yahoo chat
21  824F413FA37520BF    970916184809    garter belts
22  824F413FA37520BF    970916184818    garter belts
23  824F413FA37520BF    970916184939    lingerie
24  824F413FA37520BF    970916185051    spiderman
```

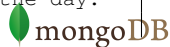- The Script is only lightly modified from stock pig:

### code/pigtutorial.pig

```
0  /*
1   * Licensed to the Apache Software Foundation (ASF) under one
2   * or more contributor license agreements.  See the NOTICE file
3   * distributed with this work for additional information
4   * regarding copyright ownership.  The ASF licenses this file
5   * to you under the Apache License, Version 2.0 (the
6   * "License"); you may not use this file except in compliance
```

mongoDB

```
 7  * with the License.  You may obtain a copy of the License at
 8  *
 9  *       http://www.apache.org/licenses/LICENSE-2.0
10  *
11  * Unless required by applicable law or agreed to in writing,
      software
12  * distributed under the License is distributed on an "AS IS"
      BASIS,
13  * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express
      or implied.
14  * See the License for the specific language governing
      permissions and
15  * limitations under the License.
16  */
17
18  -- Query Phrase Popularity (local mode)
19
20  -- This script processes a search query log file from the Excite
         search engine and finds search phrases that occur with
         particular high frequency during certain times of the day.
21
```

mongoDB

```
22 -- Register the tutorial JAR file so that the included UDFs can
      be called in the script.
23
24 -- Based on the Pig tutorial ,modified for Mongo support tests
25 REGISTER examples/pigtutorial/lib/pigtutorial.jar;
26 REGISTER mongo-hadoop.jar;
27 REGISTER lib/mongo-java-driver-2.3.jar;
28
29 -- Use the PigStorage function to load the excite log file into
      the raw bag as an array of records.
30 -- Input: (user,time,query)
31 raw = LOAD 'excite-small.log' USING PigStorage('\t') AS (user,
      time, query);
32
33 -- Call the NonURLDetector UDF to remove records if the query
      field is empty or a URL.
34 clean1 = FILTER raw BY
      org.apache.pig.tutorial.NonURLDetector(query);
35
36 -- Call the ToLower UDF to change the query field to lowercase.
37 clean2 = FOREACH clean1 GENERATE user, time,
      org.apache.pig.tutorial.ToLower(query) as query;
```

mongoDB

```
38
39 -- Because the log file only contains queries for a single day,
        we are only interested in the hour.
40 -- The excite query log timestamp format is YYMMDDHHMMSS.
41 -- Call the ExtractHour UDF to extract the hour (HH) from the
        time field.
42 houred = FOREACH clean2 GENERATE user,
        org.apache.pig.tutorial.ExtractHour(time) as hour, query;
43
44 -- Call the NGramGenerator UDF to compose the n-grams of the
        query.
45 ngramed1 = FOREACH houred GENERATE user, hour,
        flatten(org.apache.pig.tutorial.NGramGenerator(query)) as
        ngram;
46
47 -- Use the DISTINCT command to get the unique n-grams for all
        records.
48 ngramed2 = DISTINCT ngramed1;
49
50 -- Use the GROUP command to group records by n-gram and hour.
51 hour_frequency1 = GROUP ngramed2 BY (ngram, hour);
52
```

mongoDB

```
53  -- Use the COUNT function to get the count (occurrences) of each
        n-gram.
54  hour_frequency2 = FOREACH hour_frequency1 GENERATE flatten($0),
        COUNT($1) as count;
55
56  -- Use the GROUP command to group records by n-gram only.
57  -- Each group now corresponds to a distinct n-gram and has the
        count for each hour.
58  uniq_frequency1 = GROUP hour_frequency2 BY group::ngram;
59
60  -- For each group, identify the hour in which this n-gram is
        used with a particularly high frequency.
61  -- Call the ScoreGenerator UDF to calculate a "popularity" score
        for the n-gram.
62  uniq_frequency2 = FOREACH uniq_frequency1 GENERATE flatten($0),
        flatten(org.apache.pig.tutorial.ScoreGenerator($1));
63
64  -- Use the FOREACH-GENERATE command to assign names to the
        fields.
65  uniq_frequency3 = FOREACH uniq_frequency2 GENERATE $1 as hour,
        $0 as ngram, $2 as score, $3 as count, $4 as mean;
66
```

**mongoDB**

```
67 -- Use the FILTER command to move all records with a score less
       than or equal to 2.0.
68 filtered_uniq_frequency = FILTER uniq_frequency3 BY score > 2.0;
69
70 -- Use the ORDER command to sort the remaining records by hour
       and score.
71 ordered_uniq_frequency = ORDER filtered_uniq_frequency BY hour,
       score;
72
73 -- Use the PigStorage function to store the results.
74 -- Output: (hour, n-gram, score, count,
       average_counts_among_all_hours)
75 STORE ordered_uniq_frequency INTO
       'mongodb://localhost/demo.pig.output' USING
       com.mongodb.hadoop.pig.MongoStorage;
```

mongoDB

# Mongo Hadoop Plugin

- Starting today, you can download and use the Hadoop Plugin:
  https://github.com/mongodb/mongo-hadoop
- Still an Alpha—*not* production ready.
- Continuing Development. . .
    - Working on support for Hadoop Streaming, which allows Python, Ruby, and Perl to be used instead of Java (Waiting to migrate to .21 for Binary support)
    - Pig Input support.
    - Support for Sharding concurrency (Once one of you asks for it <hint-hint>)
    - Exploring other complimentary toolchains like Cascading.
    - Tell us how you are using Hadoop or how you want to integration MongoDB + Hadoop. . . So we can focus on the areas you, the users, want and need.

$\bigcirc$ mongoDB

## Questions?

- Twitter: **@rit** | mongodb: **@mongodb** | 10gen: **@10gen**
- email: **brendan@10gen.com**
- mongo-hadoop ... Available in a Pre-Production Alpha:
  https://github.com/mongodb/mongo-hadoop
- Report Bugs, request features:
  https://github.com/mongodb/mongo-hadoop/issues
- Pressing Questions?
  - IRC - freenode.net **#mongodb**
  - MongoDB Users List -
    http://groups.google.com/group/mongodb-user
- 10gen is *hiring!* We need smart engineers in both NY and Bay Area: http://10gen.com/jobs
- Up Next: 15 Minute Break followed by Mathias Stearn on "MongoDB Internals: The Storage Engine"

$\bullet$ mongoDB