

# Introduction to Pig programming



**Yahoo Search Engineering**

**陳奕瑋 ( Yiwei Chen )**



# 任務！

[殺很大、插很大\(+瑤瑤寫真性感精選54P\) @ osaki's Blog :: Xuite日誌](#)

殺不用錢～殺online瑤瑤性感變裝照+精選性感寫真童顏巨乳的娃娃音美少女瑤瑤 本名：郭書瑤  
暱稱：瑤瑤 身高：155cm 體重：42kg 三圍：33E/23/33 生日：1990/7/18 ...  
[blog.xuite.net/osaki99/blog/21865265](http://blog.xuite.net/osaki99/blog/21865265) - [頁庫存檔](#) - [類似內容](#)

[電玩美少女瑤瑤精選影音\(ヤオヤオ童顏Fカップ爆乳美少女映画videos ...](#)

2008年8月31日 ... 18歲電玩少女瑤瑤半工半讀扛家計 (內有瑤瑤男友) <http://blog.xuite.net/kaiger/daily/23136438> ... 20080913 我猜嗲嗲美少女第二段2號Kiki 3號瑤瑤 ...  
[blog.xuite.net/kaiger/daily/19128818](http://blog.xuite.net/kaiger/daily/19128818) - [頁庫存檔](#) - [類似內容](#)

[+](#) [顯示更多來自 blog.xuite.net 的結果](#)

[jays1943 分享正妹NO.24 無名瑤瑤- 樂多日誌](#)

瑤瑤也沒有哪裡得罪你們押你們為審麼這樣罵他說害女生生氣我看你們長的很醜吧不要自以為是  
喔死網友還罵人ㄟ死勒你要不要臉瑤瑤可是我的偶像你們最好是向一點 ...  
[blog.roodo.com/jays1943/archives/6850053.html](http://blog.roodo.com/jays1943/archives/6850053.html) - [頁庫存檔](#) - [類似內容](#)



# 任務！

## [瑤瑤航空 - Powered by Discuz!](#)

[瑤瑤航空 - Discuz! Board ...](#) 歡迎VIP旅客-魏如昀加入[瑤瑤航空\(2008-4-7\)](#) 歡迎VIP旅客-賴銘偉加入[瑤瑤航空\(2008-3-13\)](#) ... [瑤瑤家族](#) [瑤瑤在雅虎的第一家族](#) [瑤瑤天空部落格](#) [林佩瑤在天空的部落格](#) [林佩瑤](#) [無名網誌](#) [瑤瑤的新照片都在無名啦!](#) [無不癡齋](#) ...

[www.yaoyaofly.com](#) - [庫存頁面](#) - [更多此站結果](#)

## [瑤瑤喵小屋~ - 無名小站](#)

[瑤瑤喵小屋~ - 無名小站 Blog Album...](#) 最近好煩煩煩，我覺得我的腦容量變小了... 好多事情消化不良 好多念頭讓我無法抉擇 (More.) [goukigouki at 無名小站 at 02:39 PM post | Reply\(27\) |](#) [Trackback\(0\) | prosecute ...](#)

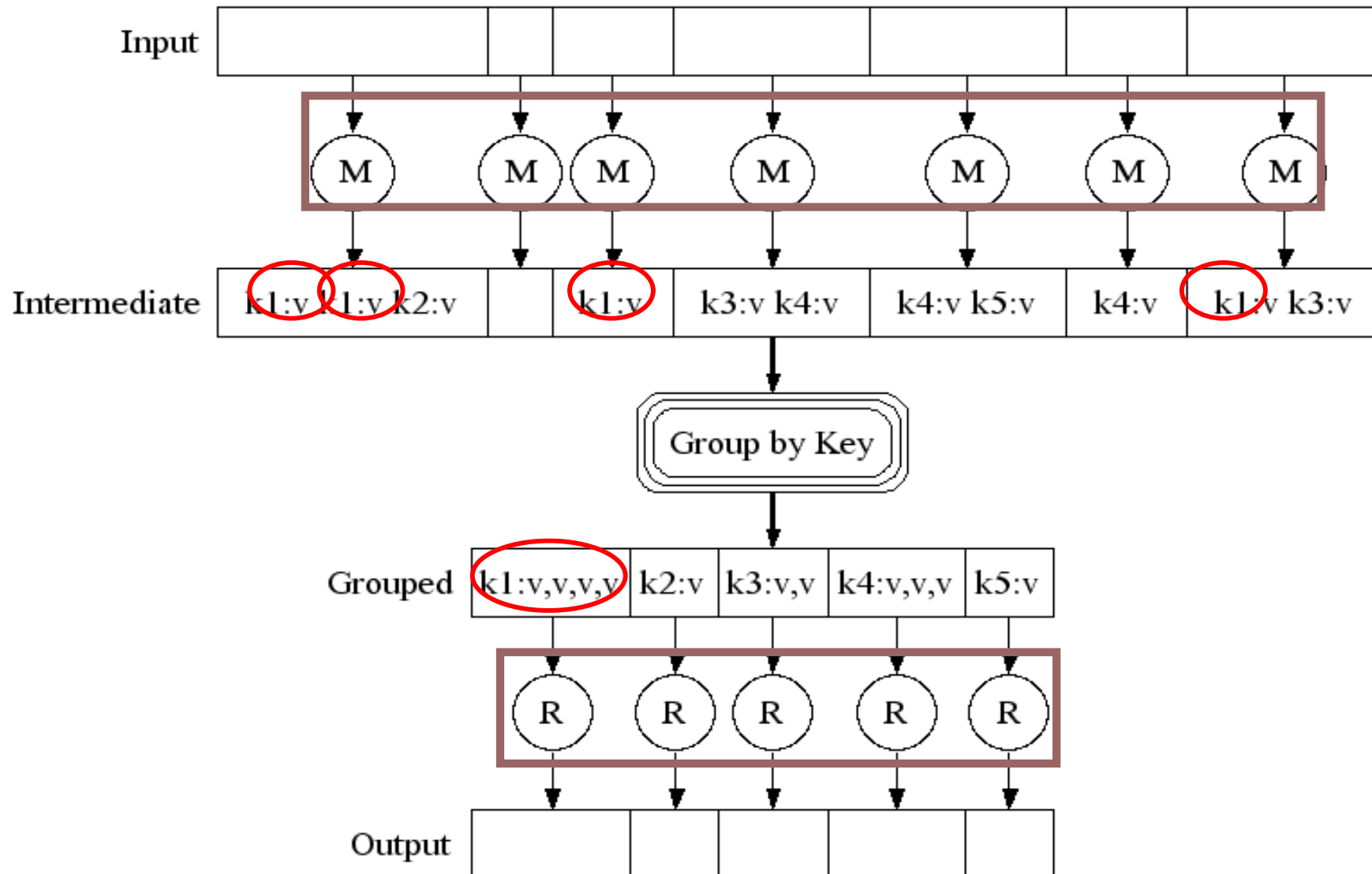
[www.wretch.cc/blog/goukigouki - 74k](#) - [庫存頁面](#) - [更多此站結果](#)



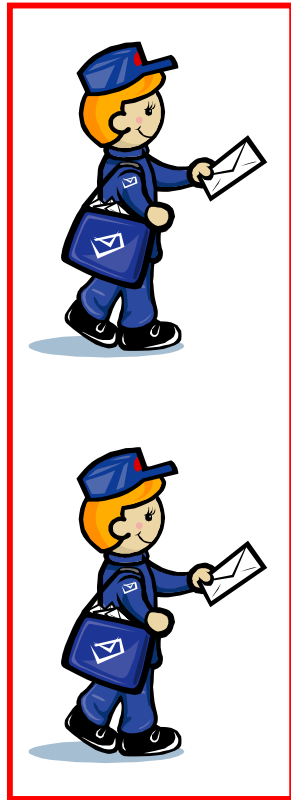
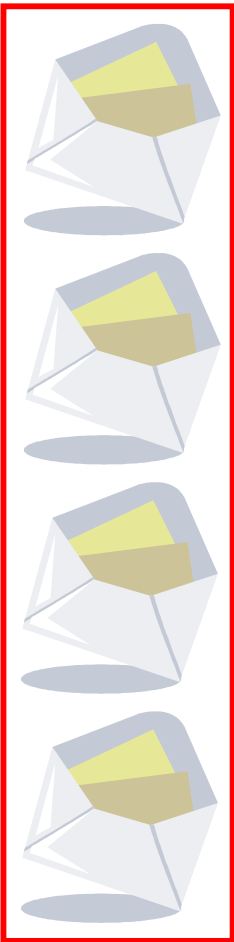
# 任務！

- 你怎麼知道我們放的網頁比較好？
- 你怎麼知道第一筆結果應該要多熱門？

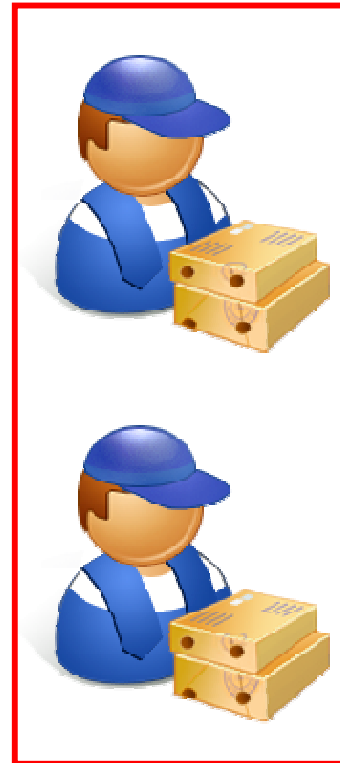
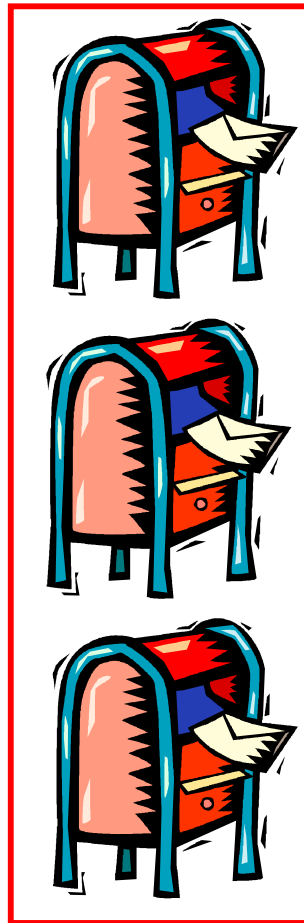
# Hadoop Programming – Map/Reduce



# Map / Reduce



**mappers**



**reducers**



100:  
37



220:  
28

# Map-Reduce

- 全新想法
- 須分別撰寫 mappers & reducers
- 會有超級無敵霹靂多的 mapper/reducer 要維護！

# We usually do ...

- 大部份時候：
  - filtering, projecting
  - grouping, aggregation, joining
- 今天有多少人搜尋「美國生」



# Pig (Latin)

- Procedural dataflow language (Pig Latin) for Map-Reduce
  - 很像 SQL
    - group, join, filter, sort ...
  - 人人都會 SQL

# Pig Script Example

- Top sites visited by users aged 18 to 25

```
Users = LOAD 'users.in' AS (name, age);
Fltrd = FILTER Users by age >= 18 and age <= 25;

Pages = LOAD 'pages.in' 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;

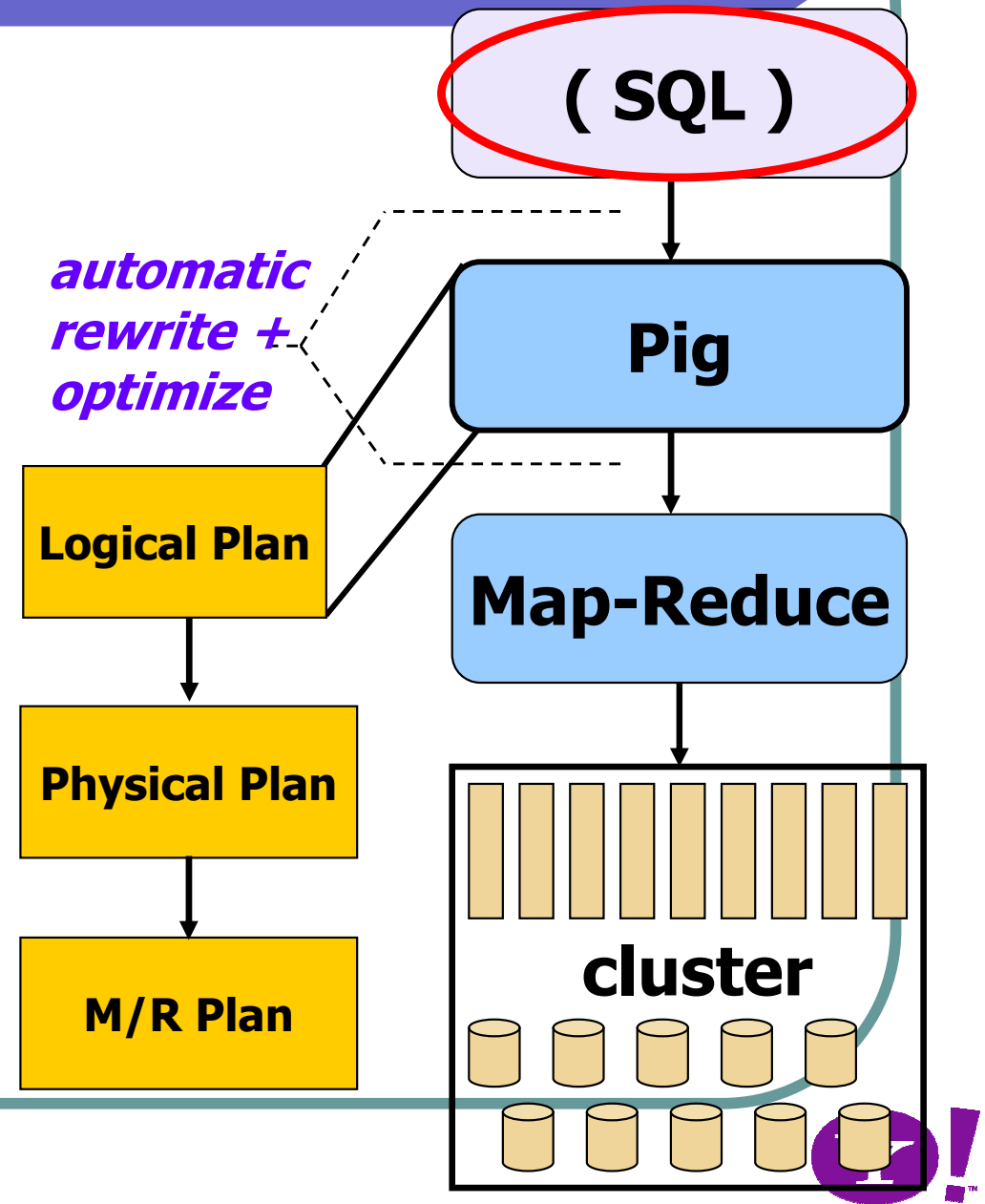
Srttd  = ORDER Smmd BY clicks;
Top100 = LIMIT Srttd 100;

STORE Top100 INTO 'top100sites.out';
```



# Pig script → Map/Reduce

- 不需懂底下 Map-Reduce 運作
- Pig 幫忙翻譯



# Why Pig?

- 容易學
- 開發快
- 一目瞭然

# Why Pig?

```
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.mapred.FileInputFormat;
import org.apache.hadoop.mapred.Mapper;
import org.apache.hadoop.mapred.MapperContext;
import org.apache.hadoop.mapred.MapperRunner;
import org.apache.hadoop.mapred.RecordReader;
import org.apache.hadoop.mapred.Reporter;
import org.apache.hadoop.mapred.Reducer;
import org.apache.hadoop.mapred.ReducerContext;
import org.apache.hadoop.mapred.ReducerRunner;
import org.apache.hadoop.mapred.SequenceFileInputFormat;
import org.apache.hadoop.mapred.TextInputFormat;
import org.apache.hadoop.mapred.JobControl;
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 {
            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("1:" + line);
            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 into a LongWritable(sum));
            String line = val.toString();
            int firstComma = line.indexOf(',');
            String value = line.substring(firstComma);
            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();
            }
        }
    }

    public static class LoadJoined extends MapReduceBase
        implements Mapper<Text, Text, Text, LongWritable> {
        public void map(
            Text key,
            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);
            String key = line.substring(firstComma, secondComma);
            String value = line.substring(secondComma);
            Text outKey = new Text(key);
            oc.collect(outKey, new LongWritable(1L));
        }
    }

    public static class ReduceUrls extends MapReduceBase
        implements Reducer<Text, LongWritable, WritableComparable,
            WritableComparable> {
        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");
            }
        }
    }

    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(key, val.toString());
            }
        }

    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) {
        JobConf jp = new JobConf(MRExample.class);
        jp.setJobName("Load Pages");
        jp.setInputFormat(TextInputFormat.class);
        jp.setOutputKeyClass(Text.class);
        jp.setOutputValueClass(Text.class);
        FileInputFormat.addInputPath(jp, Path("user/gates/pages"));
        FileOutputFormat.setOutputPath(jp, Path("user/gates/tmp/"));
        jp.setNumReduceTasks(1);
        Job loadPages = new Job(jp);

        JobConf lfu = new JobConf(MRExample.class);
        lfu.setJobName("Load and Filter");
        lfu.setInputFormat(TextInputFormat.class);
        lfu.setOutputKeyClass(Text.class);
        lfu.setOutputValueClass(Text.class);
        lfu.setMapperClass(LoadAndFilterUsers.class);
        FileInputFormat.addInputPath(lfu, Path("user/gates/users"));
        FileOutputFormat.setOutputPath(lfu, Path("user/gates/tmp/"));
        lfu.setNumReduceTasks(1);
        Job loadUsers = new Job(lfu);

        JobConf join = new JobConf(MRExample.class);
        join.setJobName("Join Users and Pages");
        join.setInputFormat(KeyValueTextInputFormat.class);
        join.setOutputKeyClass(Text.class);
        join.setOutputValueClass(Text.class);
        join.setMapperClass(IdentityMapper.class);
        join.setReducerClass(Join.class);
        FileInputFormat.addInputPath(join, Path("user/gates/tmp/indexd_pages"));
        FileInputFormat.addInputPath(join, Path("user/gates/tmp/indexd_urls"));
        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.setOutputKeyClass(Text.class);
        group.setOutputValueClass(LongWritable.class);
        group.setOutputFormat(SequenceFileOutputFormat.class);
        group.setMapperClass(LoadJoined.class);
        group.setCombinerClass(Join.class);
        group.setReducerClass(ReduceUrls.class);
        FileInputFormat.addInputPath(group, Path("user/gates/tmp/joined"));
        FileOutputFormat.setOutputPath(group, 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.setOutputKeyClass(LongWritable.class);
        top100.setOutputValueClass(Text.class);
        top100.setOutputFormat(SequenceFileOutputFormat.class);
        top100.setMapperClass(LoadClicks.class);
        top100.setCombinerClass(LimitClicks.class);
        top100.setReducerClass(LimitClicks.class);
        FileInputFormat.addInputPath(top100, Path("user/gates/tmp/grouped"));
        FileOutputFormat.setOutputPath(top100, Path("user/gates/top100sitesforusers1"));
        top100.setNumReduceTasks(1);
        Job limit = new Job(top100);
        limit.addDependingJob(groupJob);

        JobControl jc = new JobControl();
        jc.addJob(loadPages);
        jc.addJob(loadUsers);
        jc.addJob(joinJob);
        jc.addJob(groupJob);
        jc.addJob(limit);
    }
}
```

**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;**  
**Smmr = FOREACH Grpd GENERATE group, COUNT(jnd) AS clicks;**  
**Srtd = ORDER Smmr BY clicks;**  
**Top100 = LIMIT Srtd 100;**  
**STORE Top100 INTO 'top100sites';**



# Why (NOT) Pig?

- 不是史上究極霹靂大無敵武器
  - Focus: aggregation, filter, join,...
- 另一種做分散運算工作的方式

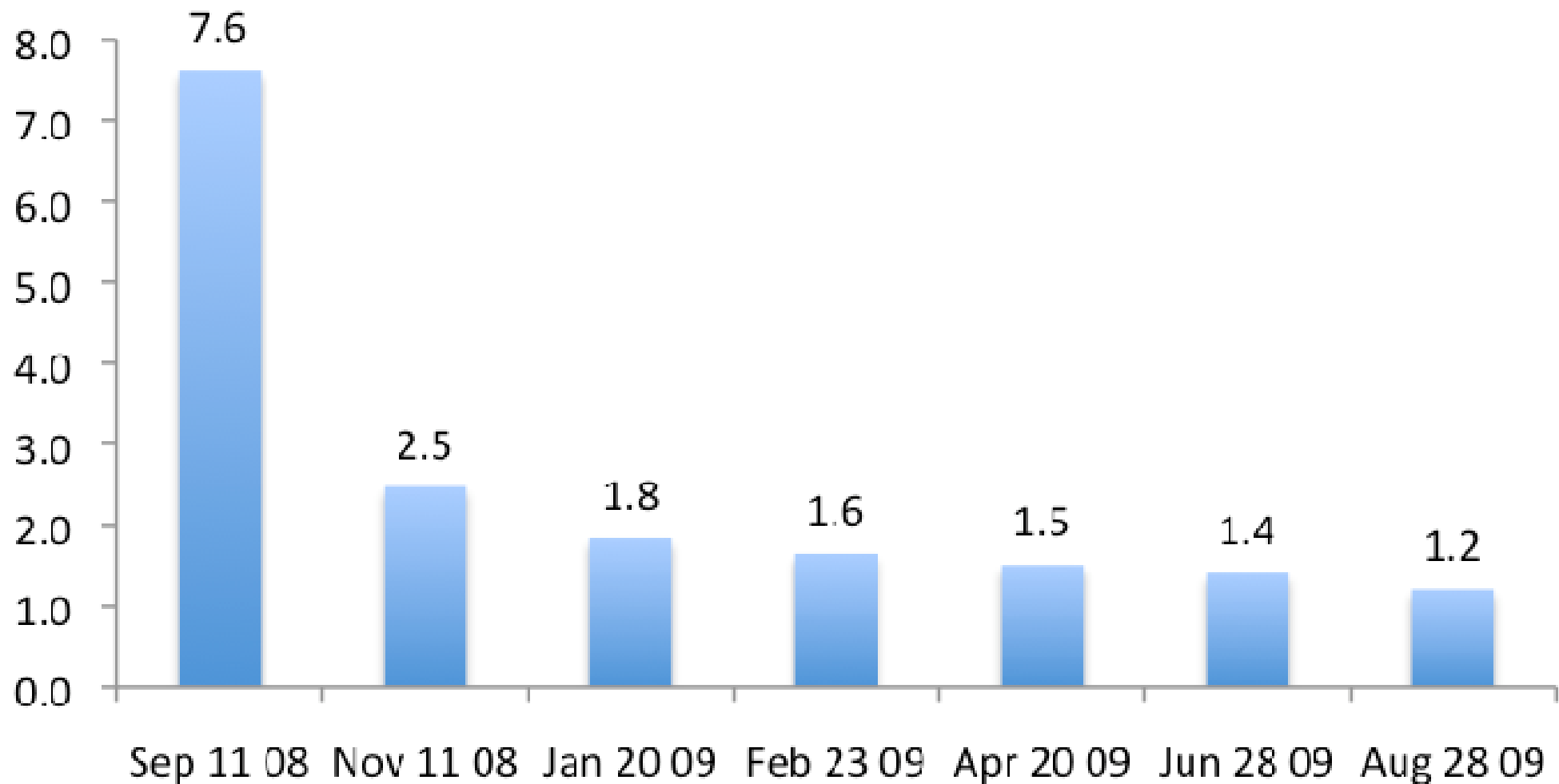


# Sweet spot between SQL – M/R

	SQL	Pig	Map-Reduce
<i>Programming style</i>	Large blocks of declarative constraints	→	“Plug together pipes”
<i>Built-in data manipulations</i>	Group-by, Sort, Join, Filter, Aggregate, Top-k, etc...	←	Group-by, Sort
<i>Execution model</i>	Fancy; trust the query optimizer	→	Simple, transparent
<i>Opportunities for automatic optimization</i>	Many	←	Few (logic buried in map() and reduce())
<i>Data Schema</i>	Must be known at table creation	→	Not required, may be defined at runtime



# Pig Performance vs Map-Reduce







# Execution and Syntax

# Pig Example

- Show users aged 18-25

```
Users = LOAD 'users.txt'  
        USING PigStorage(',') AS (name, age);  
Fltrd = FILTER Users  
        BY age >= 18 AND age <= 25;  
Names = FOREACH Fltrd GENERATE name;  
  
STORE Names INTO 'names.out';
```



# How to execute

- Local:

- `pig -x local foo.pig`

- Hadoop (HDFS):

- `pig foo.pig`

- `pig -Dmapred.job.queue.name=xxx foo.pig`

- `hadoop queue -showacIs`



# How to execute

- Interactive pig shell
  - `$ pig`
  - `grunt> _`

# Load Data

```
Users = LOAD 'users.txt'  
        USING PigStorage(',') AS (name, age);
```

- LOAD ... AS ...
- PigStorage(',') to specify separator

```
John,18  
Mary,20  
Bob,30
```



name	age
John	18
Mary	20
Bob	30

# Filter

```
Fltrd = FILTER Users  
      BY age >= 18 AND age <= 25;
```

- **FILTER ... BY ...**
  - constraints can be composite

name	age
John	18
Mary	20
Bob	30

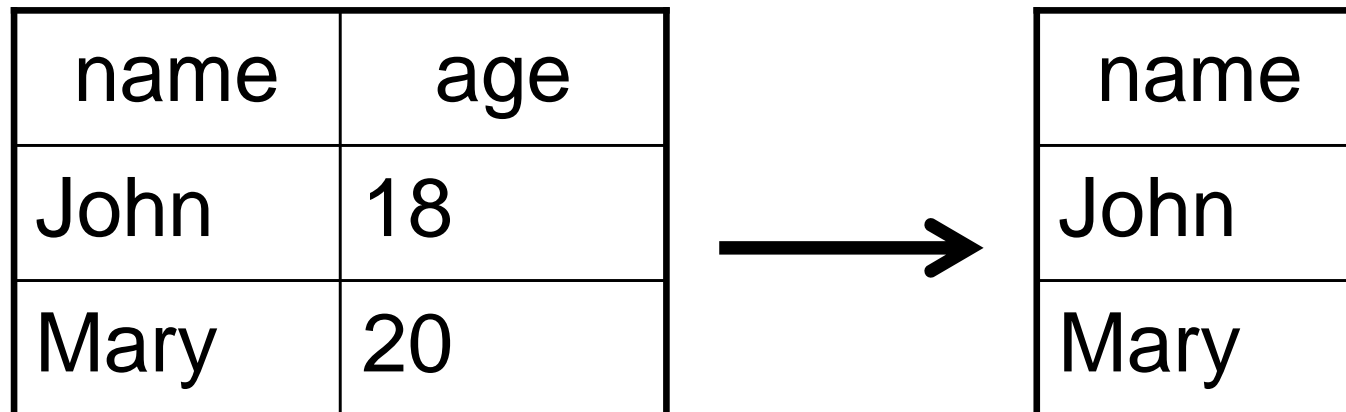


name	age
John	18
Mary	20

# Generate / Project

```
Names = FOREACH Fltrd GENERATE name;
```

- FOREACH ... GENERATE



# Store Data

```
STORE Names INTO 'names.out';
```

- **STORE ... INTO ...**
  - PigStorage(',') to specify separator if multiple fields



# Command - JOIN

```
Users = LOAD 'users' AS (name, age);  
Pages = LOAD 'pages' AS (user, url);  
Jnd   = JOIN Users BY name, Pages BY user;
```

name	age
John	18
Mary	20
Bob	30

user	url
John	yaho
Mary	goog
Bob	bing



name	age	user	url
John	18	John	yaho
Mary	20	Mary	goog
Bob	30	Bob	bing

# Command - GROUP

```
Grpd = GROUP Jnd by url;  
describe Grpd;
```

name	age	url
John	18	yhoo
Mary	20	goog
Dee	25	yhoo
Kim	40	bing
Bob	30	bing



yhoo	(John, 18, yhoo) (Dee, 25, yhoo)
goog	(Mary, 20, goog)
bing	(Kim, 40, bing) (Bob, 30, bing)

# Other Commands

- `PARALLEL` – controls `#reducer`
- `ORDER` – sort by a field
- `COUNT` – eval: count `#elements`
- `COGROUP` – structured JOIN
- More at  
[http://hadoop.apache.org/pig/docs/r0.5.0/piglatin\\_reference.html](http://hadoop.apache.org/pig/docs/r0.5.0/piglatin_reference.html)





# Features

# Parameter Substitution

```
%default TYPE 'view'
```

```
%declare ID '18987'
```

```
A = load '/data/$DATE/$ID/$TYPE'
```

- `$ pig a.pig`
- `$ pig -param DATE=20091009 a.pig`
- `$ pig -param DATE=20091009 -param  
TYPE=click a.pig`



# RegEx Comparison

- `itsyou = FILTER` urls by  
(`$0 MATCHES 'http://.*\\.yahoo\\.com.*'`)
- **MATCHES** matches 'whole' string
  - `'aaaa' MATCHES 'aaa.*'` is true
  - `'bbaaaa' MATCHES 'aaa.*'` is false
- **pattern syntax:** `java.util.regex.Pattern`



# User-defined Function (UDF)

(John,171)  
(Mary,165)  
(Bob,183)



(**JOHN**,171)  
(**MARY**,165)  
(**BOB**,183)

# UDF – user function part

```
package myudf;
import java.io.IOException;
import org.apache.pig.EvalFunc;
import org.apache.pig.data.Tuple;

public class UPPER extends EvalFunc<String>
{
    public String exec(Tuple in) throws IOException {
        if (in == null || in.size() == 0) return null;
        String str = (String)in.get(0);
        return str.toUpperCase();
    }
}
```





# UDF

- <http://hadoop.apache.org/pig/docs/r0.3.0/udf.html>
- <http://hadoop.apache.org/pig/javadoc/docs/api/>
- **PiggyBank**
  - Pig users UDF repo
  - <http://wiki.apache.org/pig/PiggyBank>



# Embedded in Java

```
/* create a pig server in the main class*/
{
    PigServer pigserver = new PigServer(args[0]);
    runMyQuery(pigServer, "/user/viraj/mydata.txt")
}

/* submit in function runMyQuery */

runMyQuery(PigServer pigServer, String inputFile) throws
IOException {
    pigServer.registerQuery("A = load '" + inputFile +
"' as (f1,f2,f3);");
    pigServer.registerQuery("B = group A by f1;");
    pigServer.registerQuery("C = foreach B generate
flatten(group);");
    pigServer.store("C", "/user/viraj/myoutput");
}
```



# References

- **FAQ**
  - <http://wiki.apache.org/pig/FAQ>
- **Documentation**
  - <http://hadoop.apache.org/pig/docs/r0.5.0/>
- **Talks & papers**
  - <http://wiki.apache.org/pig/PigTalksPapers>
  - <http://www.cloudera.com/hadoop-training-pig-introduction>



Questions?





Backup slides

# Parameter Substitution

```
$ pig -param myparam=val foo.pig
```

```
B = filter A by ($0 eq '$myparam')
```

- `pig -dryrun` produces processed script

```
B = filter A by ($0 eq 'val')
```



# Parameter Substitution

- Params in file instead of command line
- `$ pig -param_file myparams.txt a.pig`

```
#myparams.txt  
DATE=20081009  
TYPE=clicks
```



# UDF – build user function

- `javac`
  - `-cp $PIG_HOME/lib/pig.jar`
  - `-sourcepath src`
  - `-d classes`
  - `src/myudf/UPPER.java`
- `jar cf myudf.jar -C classes myudf/UPPER.class`





# UDF – pig latin part

- **register** myudf.jar;
- B =  
foreach A generate  
    **myudf.UPPER(name)**, height;

# SQL vs. Pig Latin

<u>SQL</u>	<u>Pig</u>	<u>Example</u>
<b>From table</b>	<b>Load file(s)</b>	<b>SQL:</b> from X; <b>Pig:</b> A = load 'mydata' using PigStorage('\t') as (col1, col2, col3);
<b>Select</b>	<b>Foreach ... generate</b>	<b>SQL:</b> select col1 + col2, col3 ... <b>Pig:</b> B = foreach A generate col1 + col2, col3;
<b>Where</b>	<b>Filter</b>	<b>SQL:</b> select col1 + col2, col3 from X where col2>2; <b>Pig:</b> C = filter B by col2 > '2';

(adapted from Viraj's slide)



# SQL vs. Pig Latin

<u>SQL</u>	<u>Pig</u>	<u>Example</u>
<b>Group by</b>	<b>Group + foreach ... generate</b>	<b>SQL:</b> select col1, col2, sum(col3) from X group by col1, col2; <b>Pig:</b> D = group A by (col1, col2); E = foreach D generate flatten(group), SUM(A.col3);
<b>Having</b>	<b>Filter</b>	<b>SQL:</b> select col1, sum(col2) from X group by col1 having sum(col2) > 5; <b>Pig:</b> F = filter E by \$1 > '5';
<b>Order By</b>	<b>Order ... By</b>	<b>SQL:</b> select col1, sum(col2) from X group by col1 order by col1; <b>Pig:</b> H = ORDER E by \$0;

(adapted from Viraj's slide)



# SQL vs. Pig Latin

<u>SQL</u>	<u>Pig</u>	<u>Example</u>
<b>Distinct</b>	<b>Distinct</b>	<b>SQL:</b> select distinct col1 from X; <b>Pig:</b> I = foreach A generate col1; J = distinct I;
<b>Distinct Agg</b>	<b>Distinct in foreach</b>	<b>SQL:</b> select col1, count (distinct col2) from X group by col1; <b>Pig:</b> K = foreach D { L = distinct A.col2; generate flatten(group), SUM(L); }

(adapted from Viraj's slide)



# SQL vs. Pig Latin

<u>SQL</u>	<u>Pig</u>	<u>Example</u>
<b>Join</b>	<b>Cogroup + flatten</b>  <b>(also shortcut: JOIN)</b>	<b>SQL:</b> select A.col1, B.col3 from A join B using (col1);  <b>Pig:</b> A = load 'data1' using PigStorage('\t') as (col1, col2); B = load 'data2' using PigStorage('\t') as (col1, col3); C = cogroup A by col1 <b>inner</b> , B by col1 <b>inner</b> ; D = foreach C generate flatten(A), flatten(B); E = foreach D generate A.col1, B.col3;

(adapted from Viraj's slide)



# Debug Tips

- Use small data and `pig -x local`
- LIMIT
  - `A = LOAD 'data' AS (a1,a2,a3)`
  - `B = LIMIT A 3;`
- `DUMP` , `DESCRIBE`

# FAQ

- <http://wiki.apache.org/pig/FAQ>
  - can assign #reducer
  - support regex
  - can use allocated HOD cluster

# pig.vim

- [http://www.vim.org/scripts/script.php?script\\_id=2186](http://www.vim.org/scripts/script.php?script_id=2186)

```
A = load 'data.txt' as (f1,f2,f3);  
dump A;  
B = foreach A generate f1,f3;  
dump B;  
store B into 'output.txt' using PigStorage('-');
```