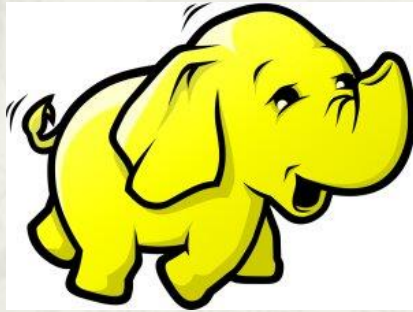# Hadoop+R應用於資料分析

# Outline

華梵大學 圖書資訊處
Office of Library and Information

➢ Hadoop 介紹
➢ R 介紹
➢ Hadoop+R 介紹
➢ Mongodb
➢ Fluentd
➢ Drill
➢ Demo

# Hadoop

# What is Hadoop?

➢ Hadoop is a software platform that lets one easily write and run applications that process vast amount of data.

  ✓ 它是軟體平台用來處理程式具巨量資料

➢ Hadoop can reliably store and process petabytes.

  ✓ 它可用來可靠地儲存和處理PB級巨量資料

➢ It distributes the data and processing across clusters of commonly available computers. These clusters can number into the thousands of nodes.

  ✓ 它將資料和處理程序分散到可以使用的電腦上，而且這些電腦的數量可以達到上千台之多

# What is Hadoop?

> By distributing the data, Hadoop can process it in parallel on the nodes where the data is located. This make it extremely rapid.

 ✓ 藉由分散資料的處理，Hadoop可以平行的運算這些資料，使得處理速度變得非常快速

> Hadoop automatically maintains multiple copies of data and automatically redeploys computing tasks based on failures.

 ✓ Hadoop可以將運算的程式和放置的資料在每一個可以運行的節點間進行複製和自動化的備份，可以避免執行中的程式或存放的資料，因為電腦的硬體或系統的上的損壞而使程式無法執行或檔案損毀

# Apache Hadoop

http://hadoop.apache.org/

# Hadoop

➢ 免費軟體

➢ 利用MapReduce作為分散式處理技術

➢ 利用HDFS作為分散式檔案系統

# MapReduce

# MapReduce

- 是一種軟體框架(Software Framework)
- 可在不同電腦組成的叢集(Clusters)上執行
- 能為巨量資料(Big Data)做分散運算處理
- 此框架的功能概念主要是映射(Map)和化簡 (Reduce)兩種
- 實作上可用JAVA、R或其他程式語言來達成

# MapReduce

➢ Map

✓ 從主節點(Master Node)輸入一組Input，此 Input是一組 key/value序對，將這組輸入切分 成好幾個小的子部分，分散到各個工作節點 (Slave Nodes)去做運算

✓ 輸入是一組 Key/Value 序對 ，輸出則為另一組 中間過程(Intermediate)的 key/value 序對
   ◆ $(K_{in}, V_{in})$ ➔ $list(K_{inter}, V_{inter})$

# MapReduce

> Reduce

  ✓ 負責針對相同的中間過程 key 合併其所有相關聯的 Value，並產生輸出結果的 key/value 序對

  ✓ 將多對具相同 Key但不同 Value 的資料，結合為多對的 Key/Value

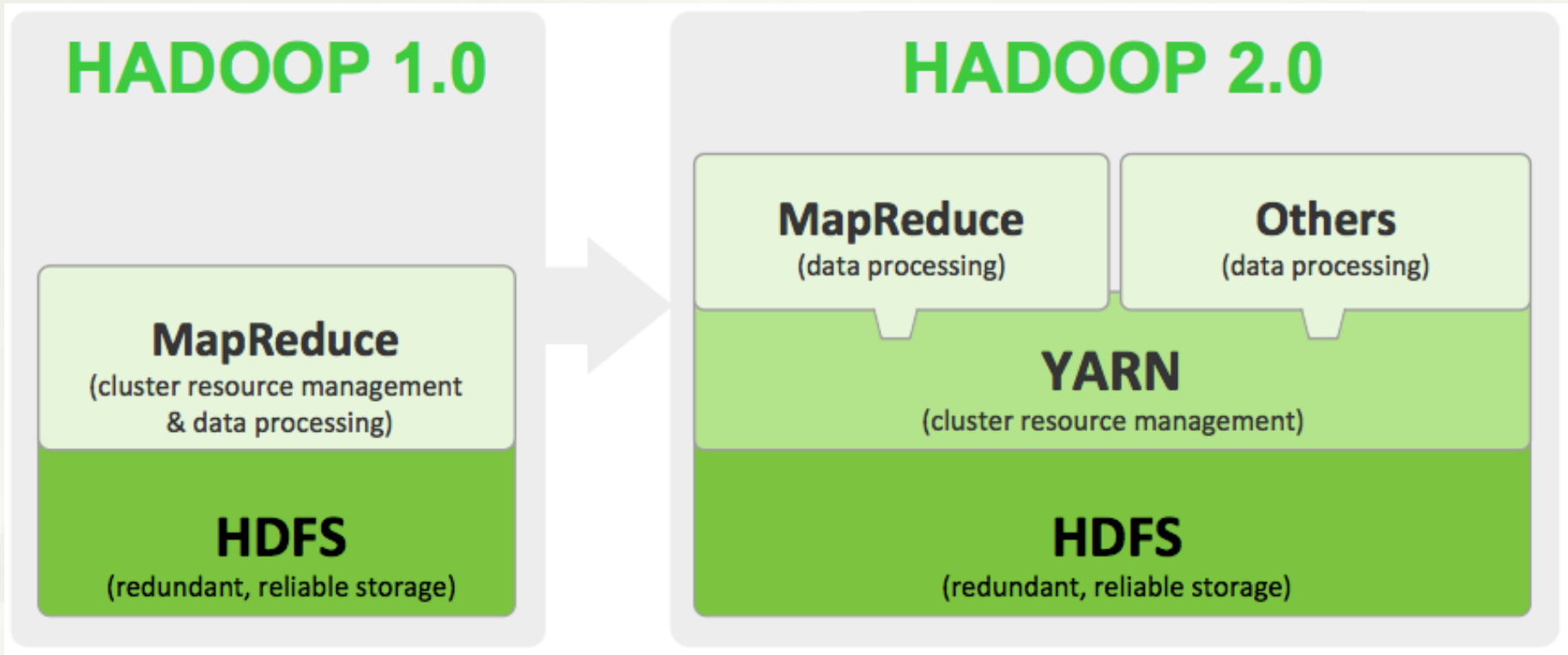    ◆ $(K_{inter}, list(V_{inter}))$ ➜ $list(K_{out}, V_{out})$

# HDFS

# Hadoop Distributed File System (HDFS)

➢ 在分散式儲存環境中，提供單一的目錄系統

➢ 資料以 Write-once-read-many 方式存取

➢ 每個檔案被分割成許多Block，每個Block複製許多複本(Replica)，並分散儲存於不同的DataNode上
  ✓ NameNode：負責維護HDFS的檔案名稱空間 (File System Namespace)
  ✓ DataNode：實際儲存檔案區塊(Blocks)的伺服器

# Hadoop 2.X



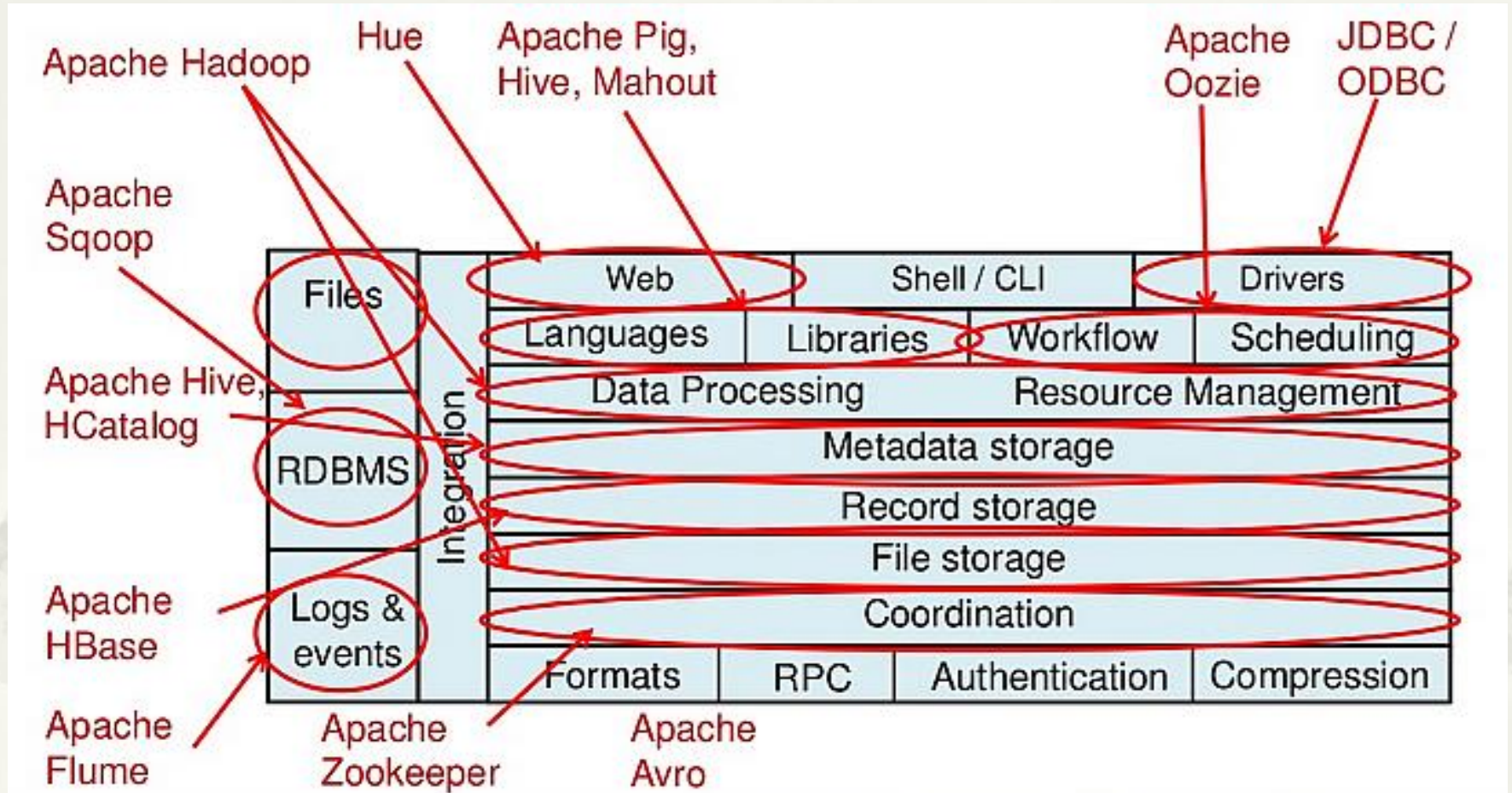http://hortonworks.com/blog/office-hours-qa-on-yarn-in-hadoop-2/

# YARN

> **Y**et **A**nother **R**esource **N**egotiator

> YARN is a more general purpose framework of which classic MapReduce is one application.

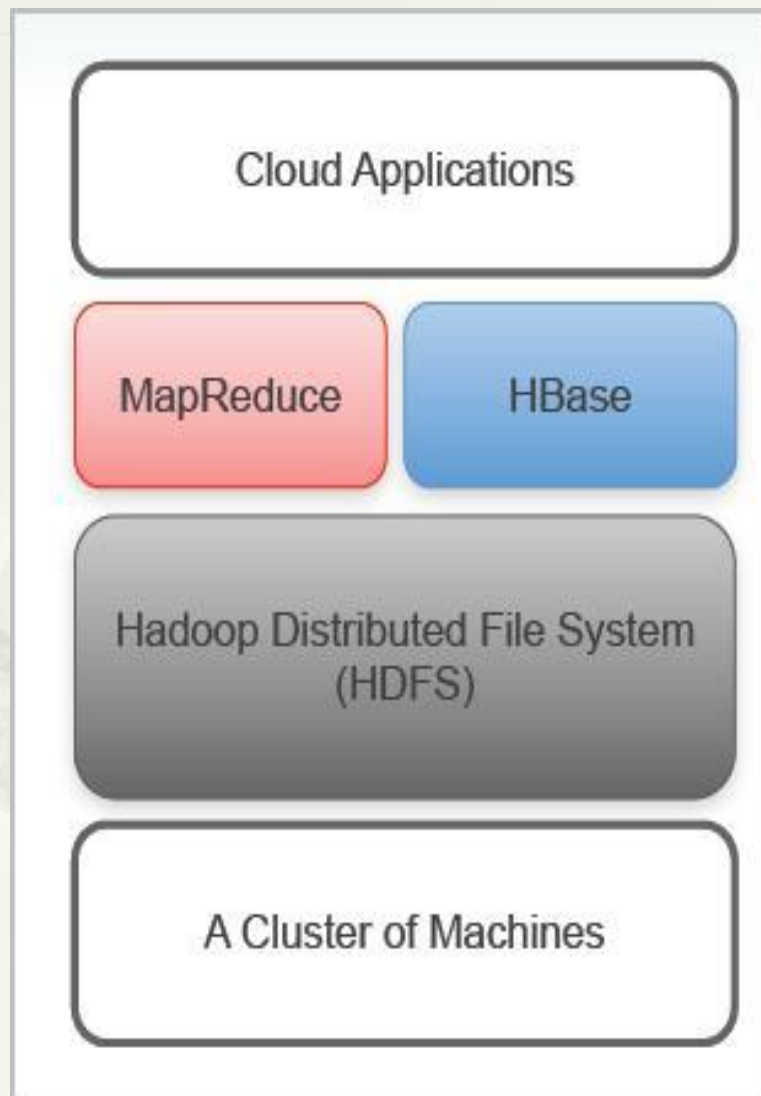  ✓ YARN 是更通用的軟體框架，而 MapReduce 只是其中的一個應用

# Hadoop Ecosystem

# HBase

Cloud Applications

MapReduce

HBase

Hadoop Distributed File System
(HDFS)

A Cluster of Machines

HBase:

➢是HDFS上的資料庫。

➢沒有正規化與Join的觀念

➢利用Family Columns將相
似的欄位群聚在一起，用
於強化效率。

17

# Hive

- ➢ Developed at Facebook

- ➢ "Relational database" built on Hadoop
  - ✓ Maintains list of table schemas
  - ✓ SQL-like query language (HiveQL)
  - ✓ Can call Hadoop Streaming scripts from HiveQL

# Sqoop

➢ 是一個用來將 Hadoop 和關聯式資料庫中的資料相互轉移的工具，可以將一個關聯式資料庫（MySQL,SQL Server等）中的資料導入到 Hadoop 的 HDFS 中，也可以將 HDFS 的資料導入到關聯式資料庫中

# Fluentd

➤ 是一個日誌收集系統，它的特點是可通過簡單的配置，將日誌收集到不同的地方
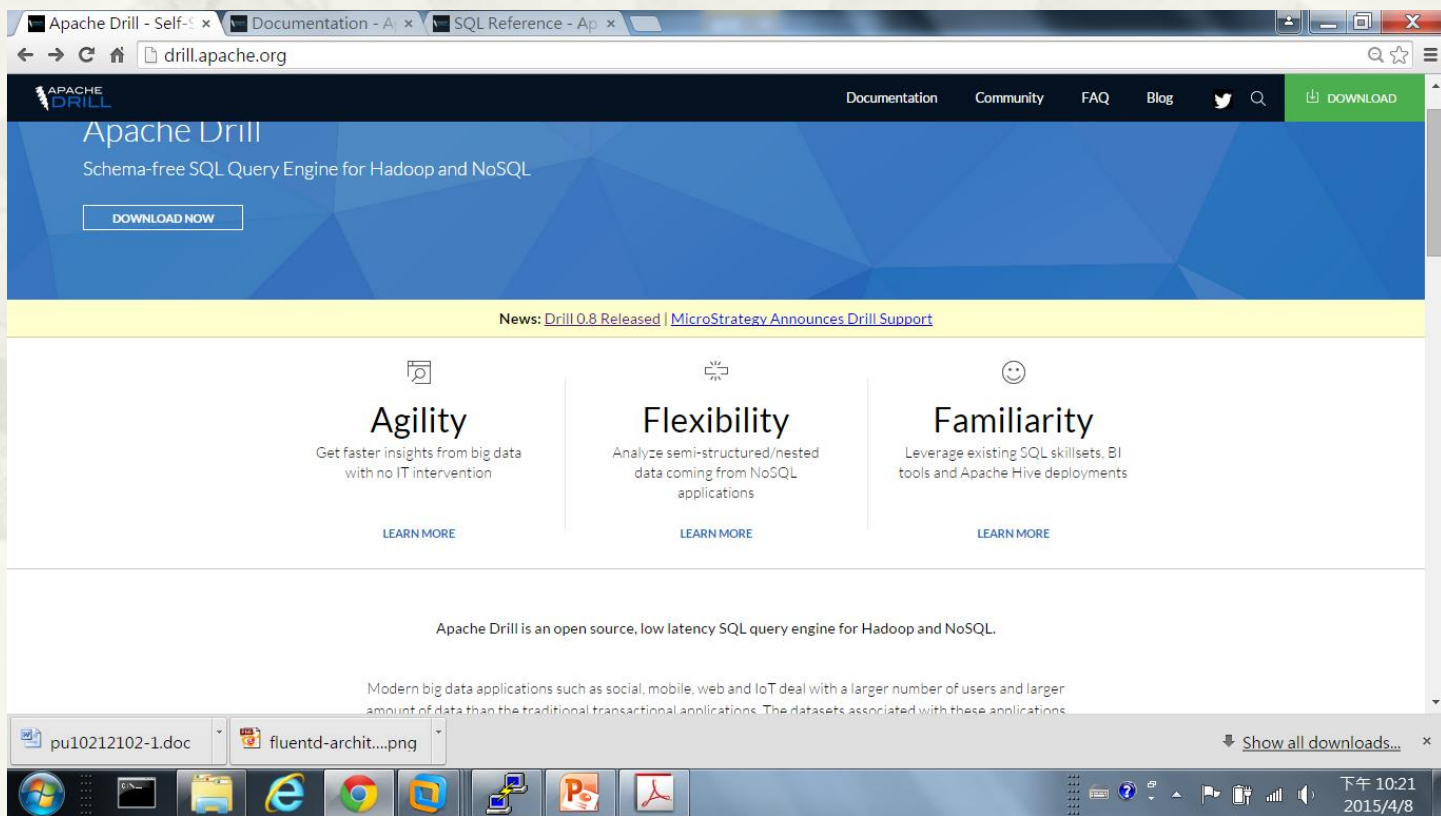
➤ 使用 Fluentd + MongoDB 構建即時日誌收集系統

# MongoDB

> 隨著資訊爆炸時代的來臨，RDBMS 的效能與彈性遭遇相當大的瓶頸，因此我們開始需要一個可以處理巨量資料(Big Data)的儲存方式，這時 NoSQL 已經悄悄誕生

> Mongodb是由10gen團隊所開發的一套NoSQL程式，它是文件式的資料庫系統，也就是說你的每一個Mongo資料庫都會以「檔案」的形式存放在資料夾中，如果要刪除資料庫，就直接把檔案刪掉就可以了

# Drill

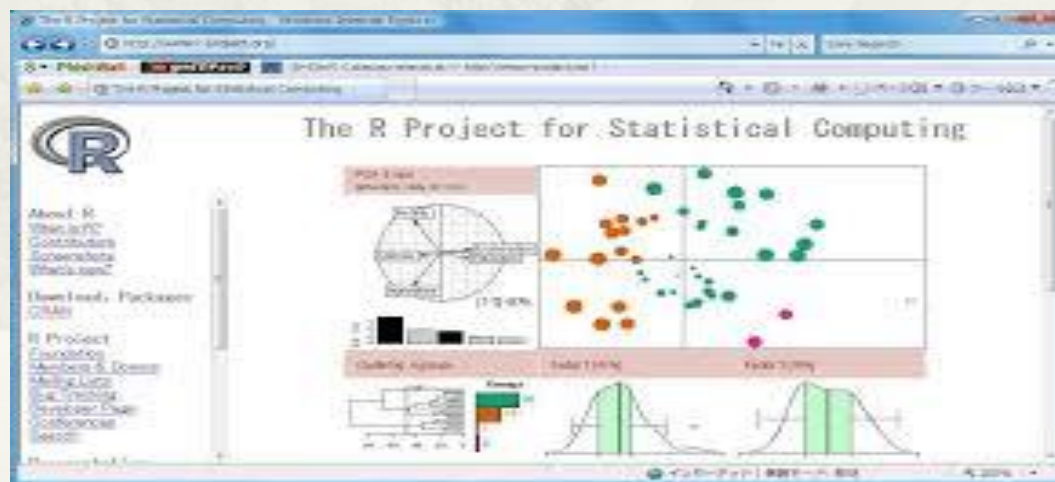➢ Apache Drill 將有助於Hadoop用戶實現以
  SQL指令更快查詢巨量資料

# R 語言

# R 簡介

> Ross Ihaka 與 Robert Gentleman（1966）所開發出來之相似於 AT & T 貝爾實驗室所開發之 S 語言

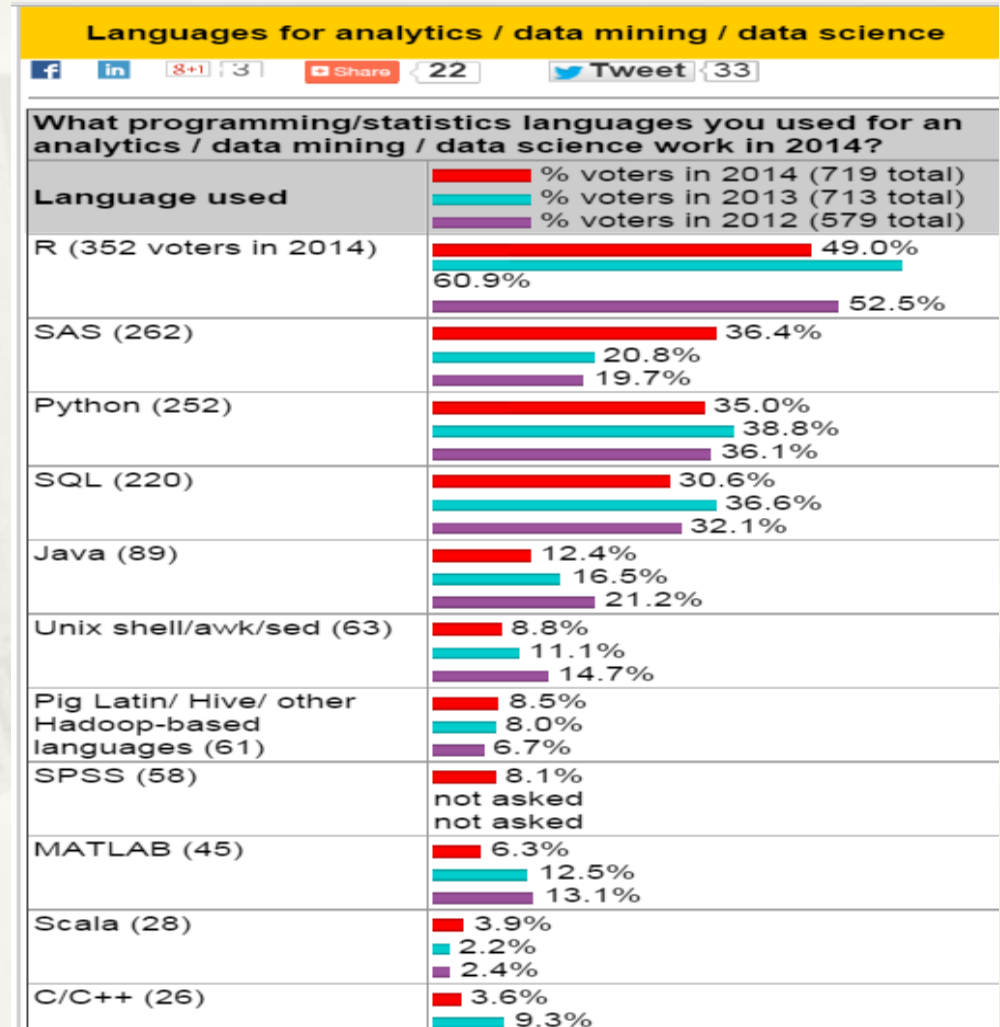> R 有 Windows、Unix、Linux 及 Apple MacOS 等不同作業系統的版本

> 免費軟體，其網站位於 http://www.r-project.org

# R 語言

➤ 內建許多函式(Function)及約5000多個免費套件

➤ R 是直譯式語言（Interpreted Language）
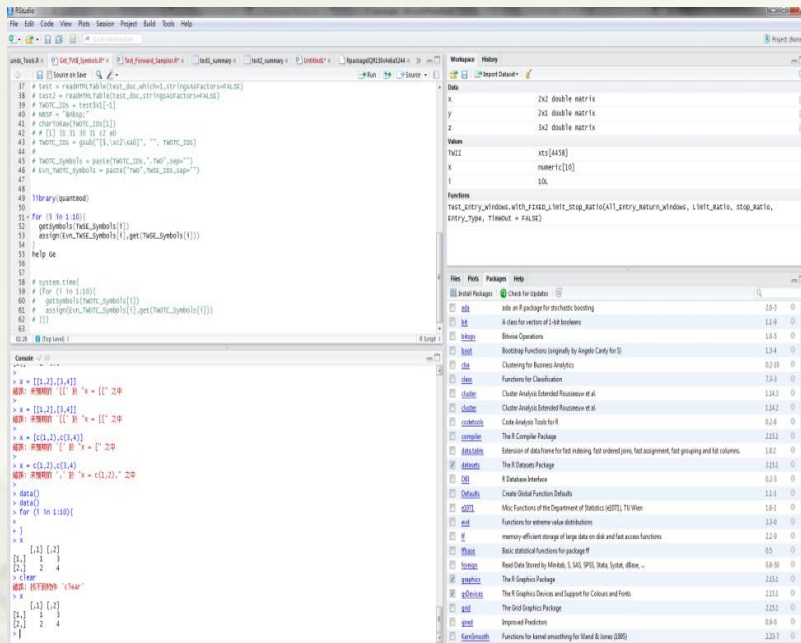  ✓ 一行行執行，可直接看到執行結果

➤ R 是物件導向語言（Object Oriented Language）

# 常用的資料分析語言

最近最受歡迎的資料分析語言 R

# 整合式開發環境 IDE



R Studio

http://www.rstudio.com/

R

http://www.r-project.org/
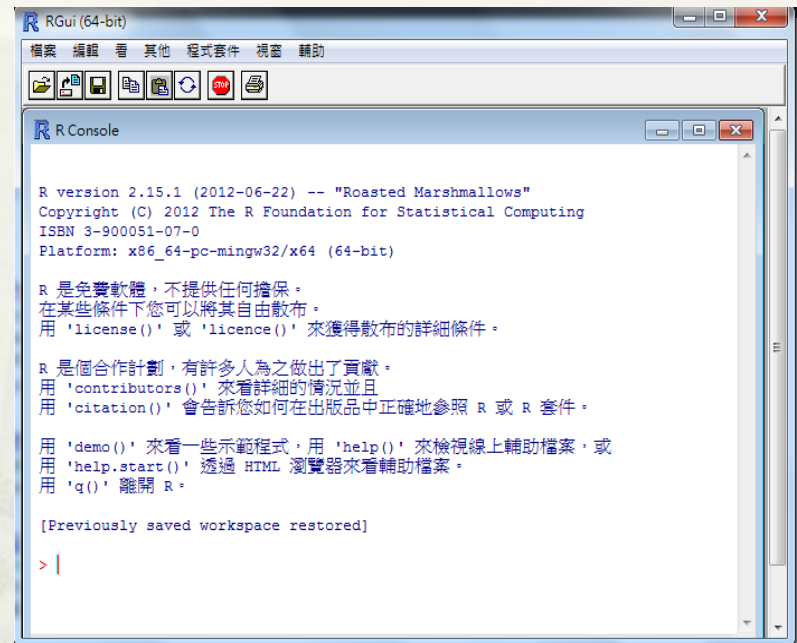
# R 應用領域

- Big Data
- 統計分析
- 資料探勘
- 機器學習
- 推薦系統
- 文字探勘
  …

# Hadoop+R

# Hadoop+R

- ➢ 擴大 R 處理資料能力
  - ✓ R 將資料全部讀進 Memory（無法讀入巨量資料）
  - ✓ Hadoop 讓 R 可以進行分散式運算

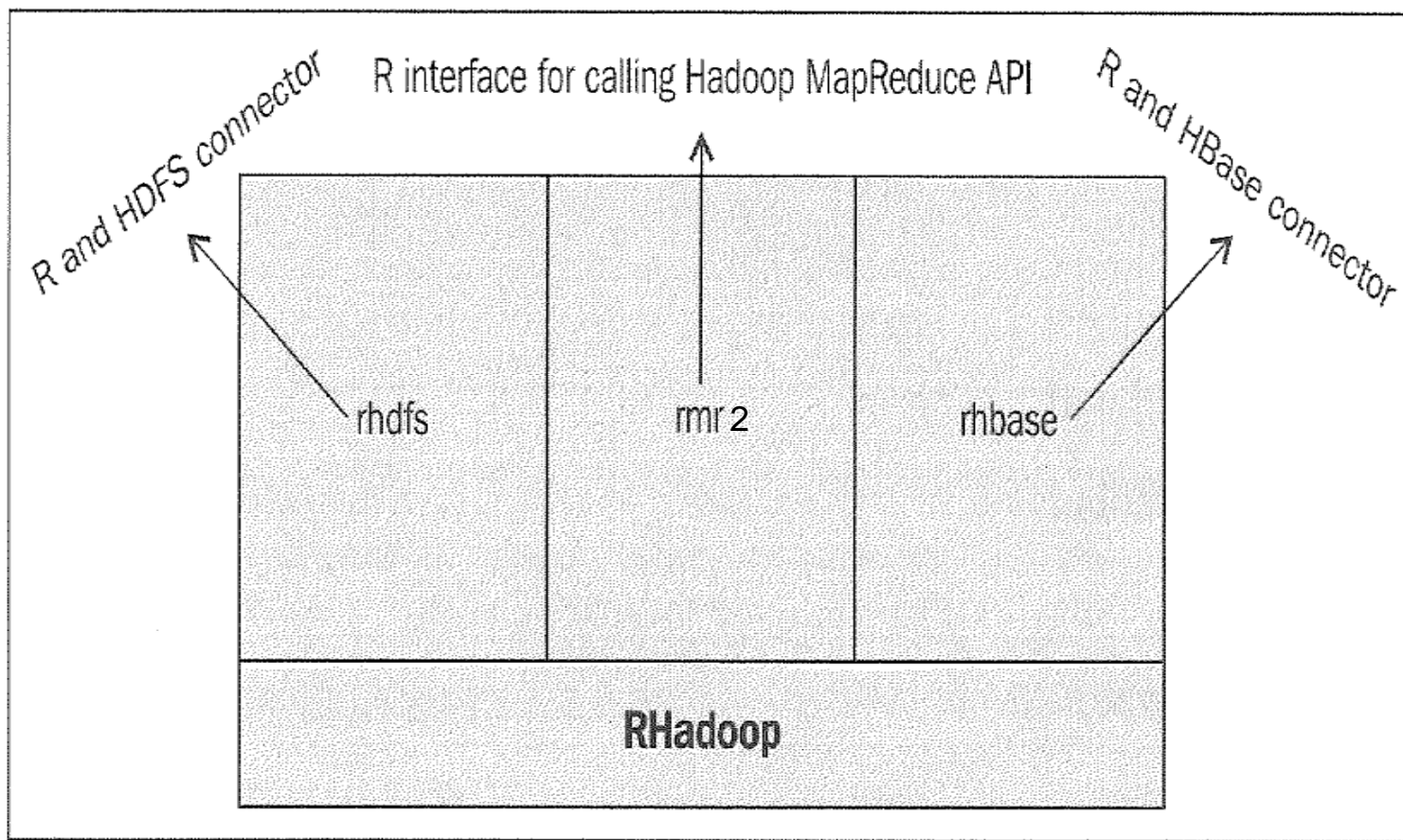- ➢ 使用 R 語言就可輕易使用Hadoop功能

# RHadoop 套件

➢ 由 Revolution Analytics 開發

➢ 針對 MapReduce、HDFS 及 HBase發展三個免費套件
- ✓ rmr2
- ✓ rhdfs
- ✓ rhbase

# RHadoop 套件功能

- ➢ rhdfs
  - ✓ 讓使用者可以透過 R 存取 HDFS

- ➢ rmr2
  - ✓ 可以讓使用者發展並呼叫 MapReduce 工作
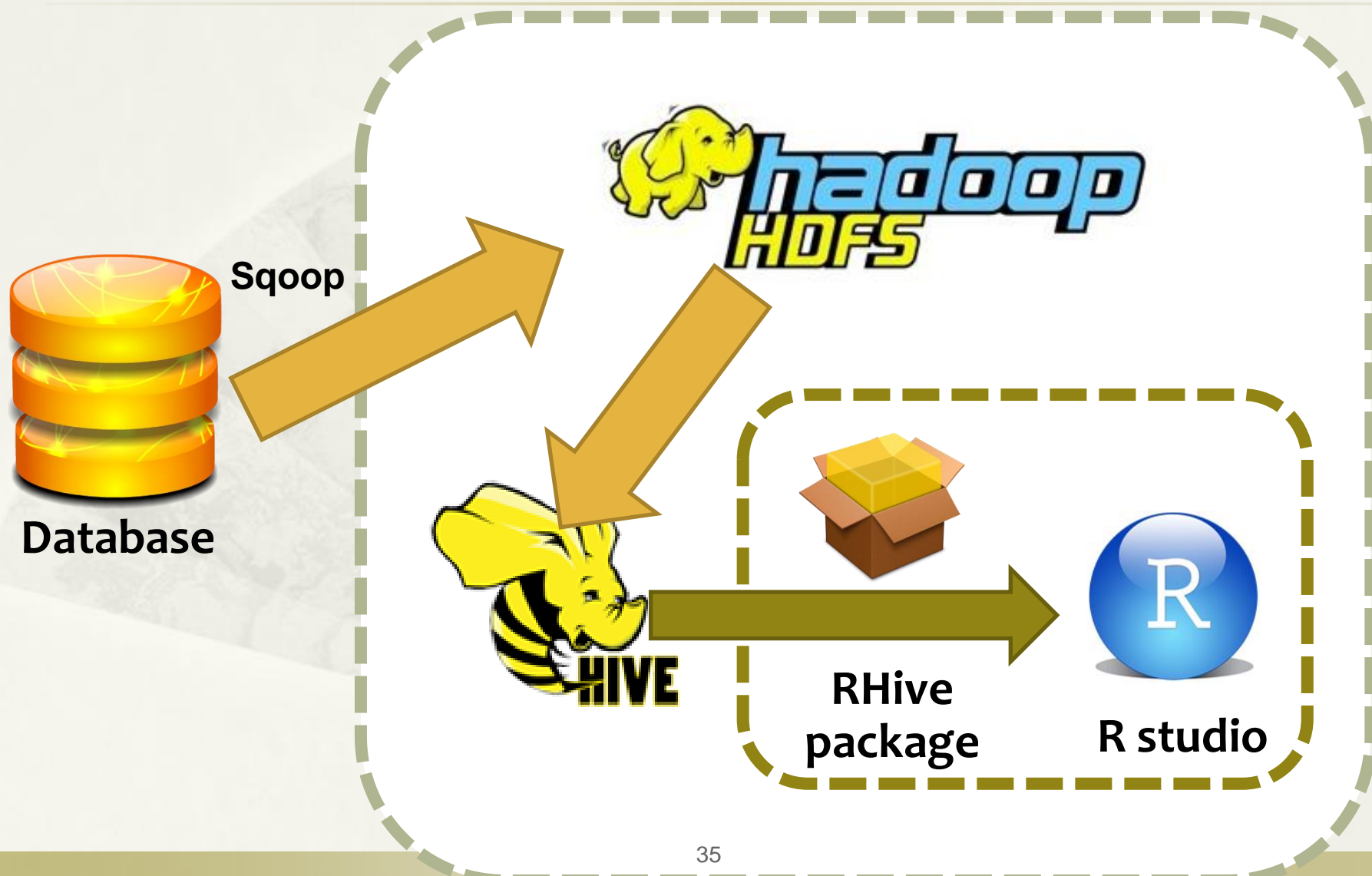
- ➢ rhbase
  - ✓ 可以操作 HBase 資料

# RHadoop 套件架構



RHadoop Ecosystem

# D e m o

# 執行流程



**Sqoop**

**Database**

**hadoop HDFS**

**HIVE**

**RHive package**

**R studio**

# Sqoop 指令

sqoop import --connect
"jdbc:sqlserver://$SQL_SRV:$PORT;database=mitopac;username=$U
;password=$P" --hive-import -m 1 --table reader --warehouse-dir
$DEST_DIR --map-column-hive reader11=String --hive-overwrite

# Demo

# Log 資料收集及分析

## Fluentd+Mongodb(192.168.244.131)

```
1. # login root
2. # mongo httpd
3. > db.accesslog.count()
4. > db.accesslog.find()
5. quit()
```

## Drill (192.168.244.131:8047, select  * from mongo.httpd. accesslog  limit 10)

1. bin/sqlline -u jdbc:drill:zk=local (./start)
2. use mongo.httpd;
3. show tables;
4. select * from accesslog; select host from accesslog;
5. !quit