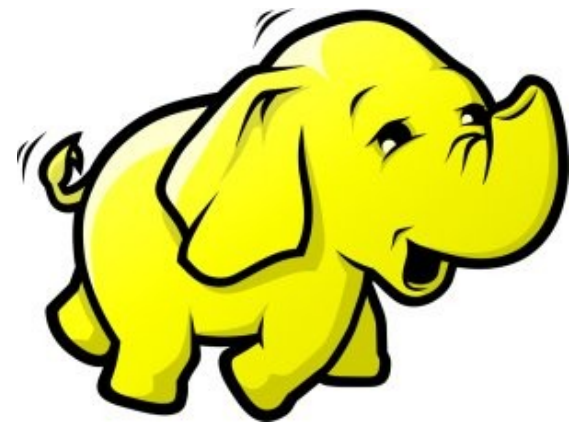




# 淺談海量資料的趨勢、挑戰與因應對策

Big Data : the Trends, Challenges and Solutions

**Jazz Wang**  
**Yao-Tsung Wang**  
**jazz@nchc.org.tw**



# Agenda 演講大綱

**What is Big Data ?** 何謂海量資料

**Why should we care?** 為何需要關切

**When to deploy it ?** 何時導入技術

**How to handle it ?** 三大因應策略

**Who is key player ?** 誰是成功關鍵

# WHAT



## What is Big Data ?

## 何謂海量資料

趨勢

Trends

定義

Definitions

挑戰：管理維度

The Six Dimensions

# Trends .... It's all about **Buzzwords** ..... 「趨勢」亦或「流行語」？ Web 3.0, Cloud Computing, Social Network, Big Data, ....

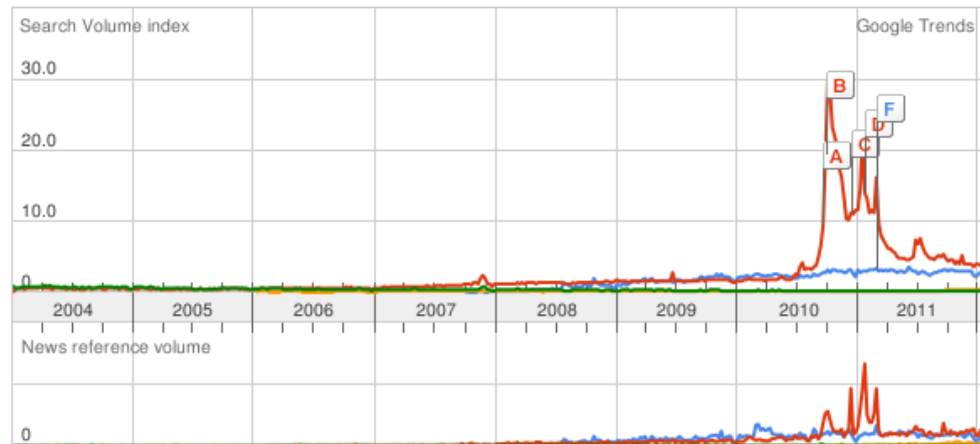
Google Trends

Tip: Use commas to compare multiple search terms.

Searches [Websites](#)

- Scale is based on the average worldwide traffic of **cloud computing** in all years. [Learn more](#)
- An improvement to our geographical assignment was applied retroactively from 1/1/2011. [Learn more](#)

cloud computing 1.00 social network 2.40 big data 0.20  
semantic web 0.40



語意網 ( Semantic Web ) 從 2001 年開始制定標準後，逐漸下滑。而同義詞 Web 3.0 也呈現相似趨勢。海量資料 ( Big Data ) 與其關鍵技術 Hadoop，則仍在上揚中。

A  
B

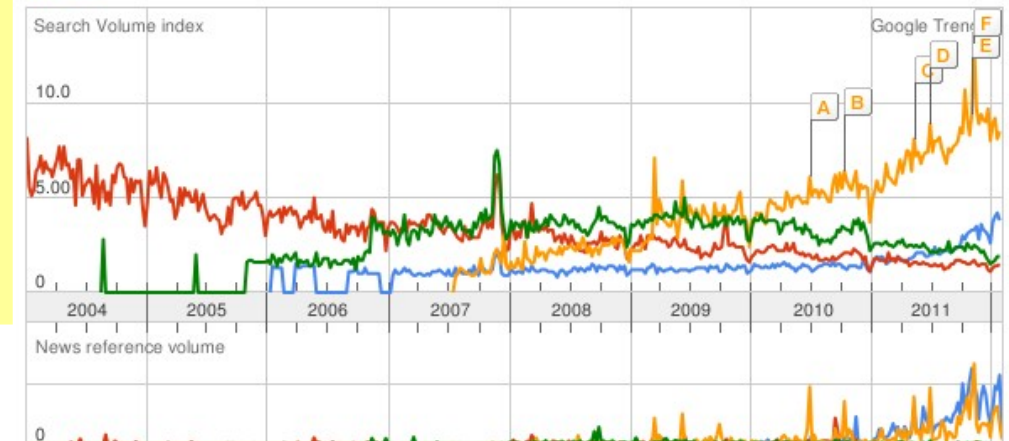
Google Trends

Tip: Use commas to compare multiple search terms.

Searches [Websites](#)

- Scale is based on the average worldwide traffic of **big data** in all years. [Learn more](#)
- An improvement to our geographical assignment was applied retroactively from 1/1/2011. [Learn more](#)

big data 1.00 semantic web 3.30 hadoop 2.50  
web 3.0 2.40



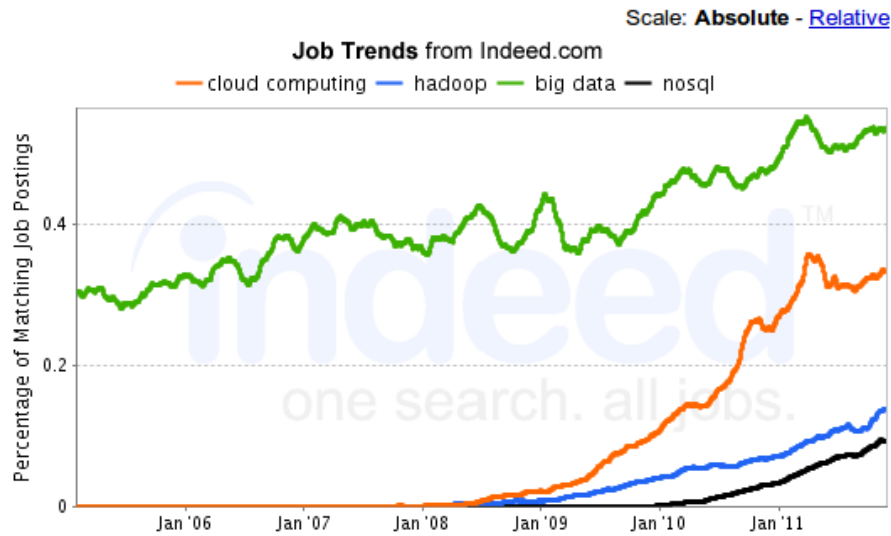
A  
B  
C  
D  
E  
F

整體而言，雲端運算 ( Cloud Computing ) 與社交網路 ( Social Network ) 呈現上揚。且社交網路比雲端運算還引人注目。



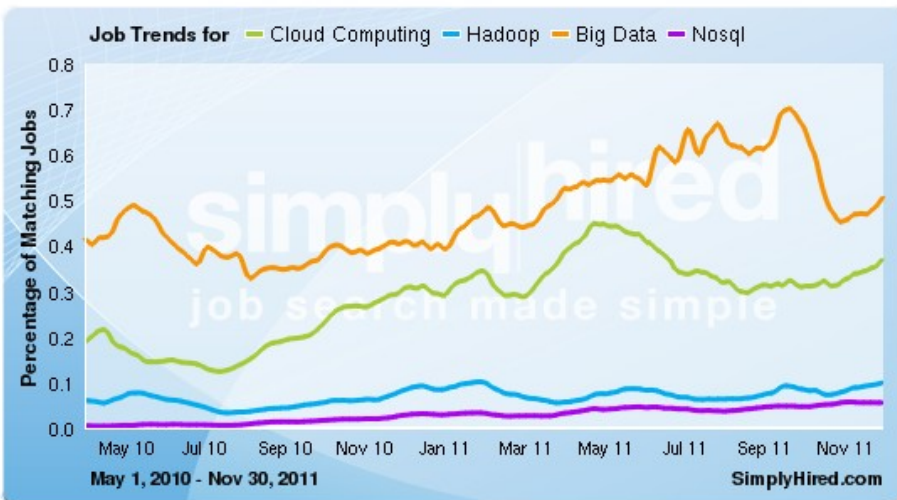
# Trends of Market Needs 市場需求趨勢

cloud computing, hadoop, big data, nosql Job Trends



Indeed.com searches millions of jobs from thousands of job sites. This job trends graph shows the percentage of jobs we find that contain your search terms.

Find [Cloud Computing jobs](#), [Hadoop jobs](#), [Big Data jobs](#), [Nosql jobs](#)



美國軟體就業市場分析，根據 indeed 與 simply hired 兩間公司的趨勢觀察，都得到一樣的結果：

Big Data > Cloud Computing > Hadoop > NoSQL

To

CIO technologies	Ranking of technologies CIOs selected as one of their top 3 priorities in 2012			
Ranking	2012	2011	2010	2009
Analytics and business intelligence	1	5	5	1
Mobile technologies	2	3	6	12
Cloud computing (SaaS, IaaS, PaaS)	3	1	2	16
Collaboration technologies (workflow)	4	8	11	5
Virtualization	5	2	1	3
Legacy modernization	6	7	15	4
IT management	7	4	10	*
Customer relationship management	8	18	*	*
ERP applications	9	13	14	2
Security	10	12	9	8
Social media/Web 2.0	11	10	3	15

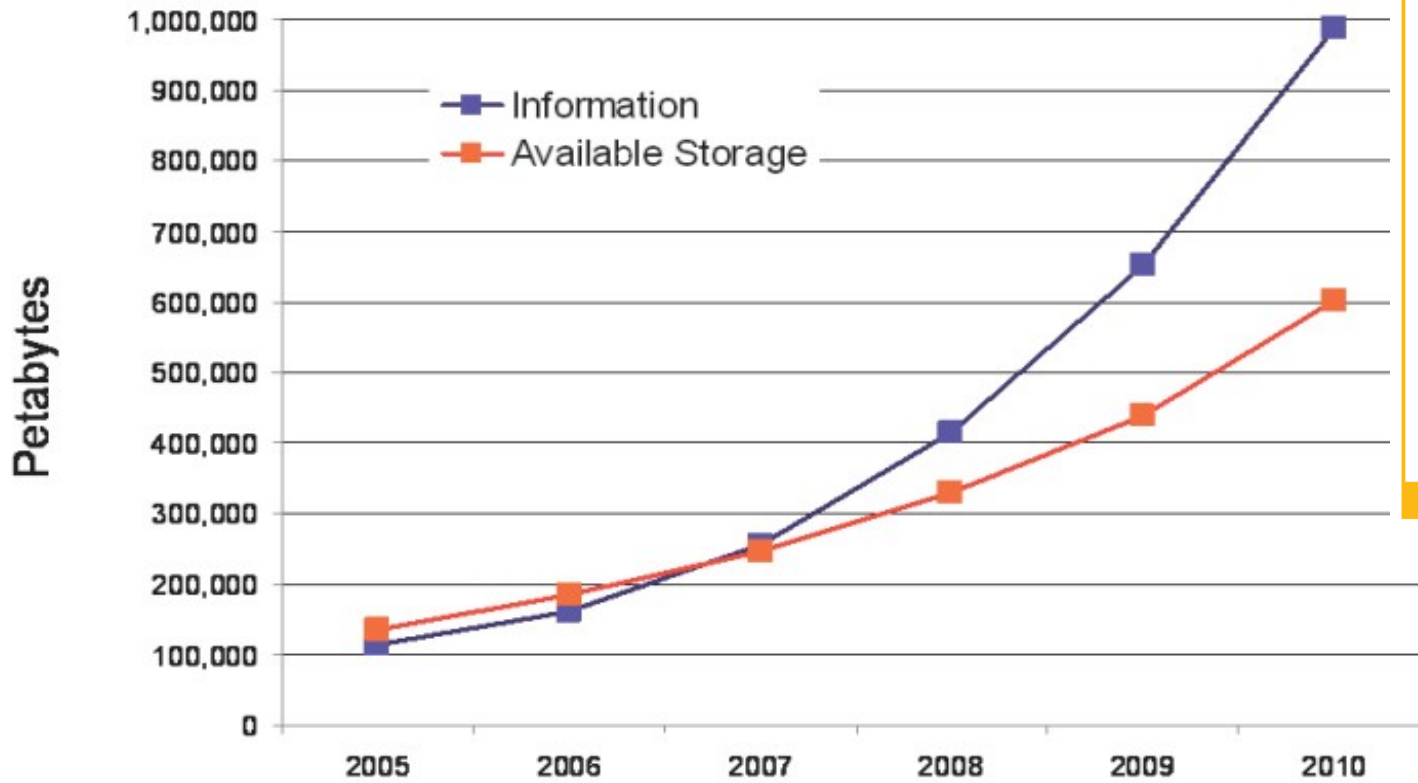
Gartner CIO Agenda 2012 前三名：  
 [1] Business Intelligence (Big Data)  
 [2] Mobile technology  
 [3] Cloud Computing

# How BIG? 讓我們先來認識一下容量單位

Bit (b)	1 or 0
Byte (B)	8 bits
Kilobyte (KB)	1,000 bytes
Megabyte (MB)	1,000 KB
Gigabyte (GB)	1,000 MB
Terabyte (TB)	1,000, GB
Petabyte (PB)	1,000 TB
Exabyte (EB)	1,000 PB
Zettabyte (ZB)	1,000 EB

# Data Explosion!! 始於 2007 的「資料大爆炸」時代

## Information Versus Available Storage

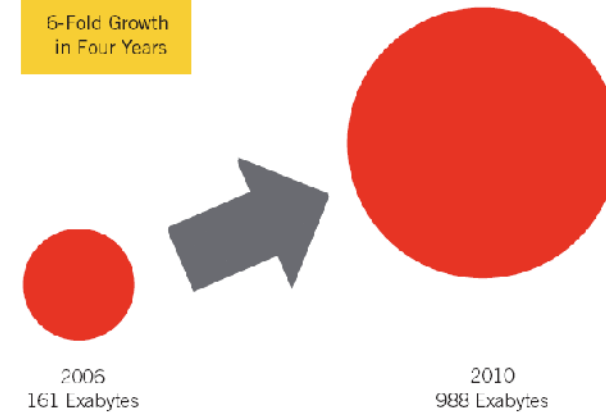


Source: IDC, 2007

Figure 1

### Information Created, Captured and Replicated

6-Fold Growth  
in Four Years



Source: IDC, 2007

2007 年，IDC 預估  
2010 年會成長**六倍**！  
(相較 2006 年)

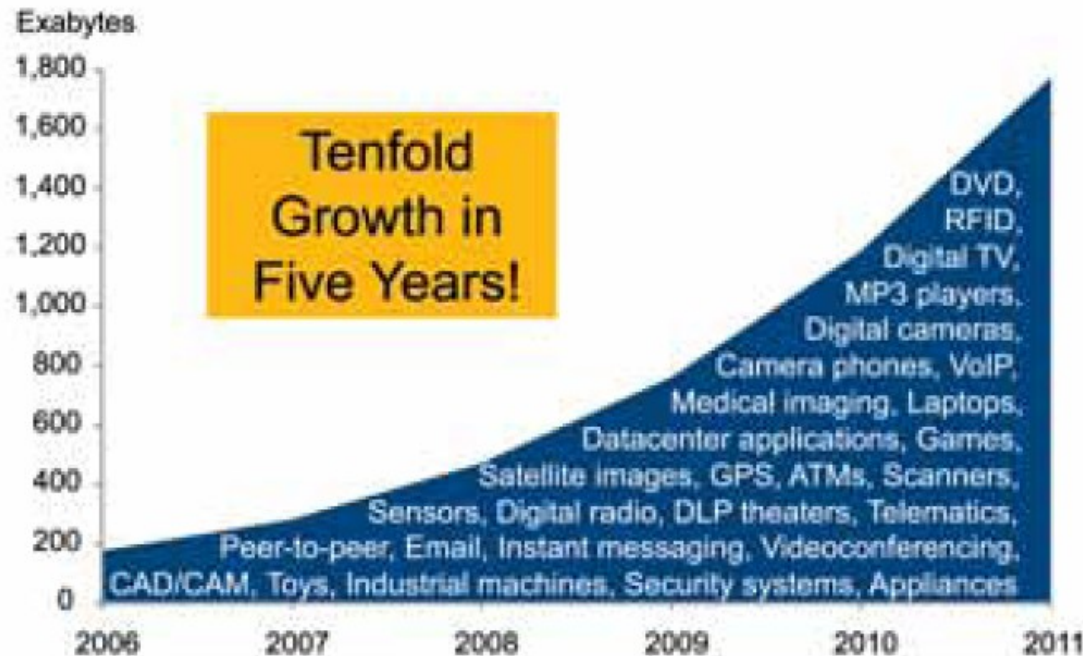
2006 161 EB  
2010 988 EB (預測)

出處：The Expanding Digital Universe,  
A Forecast of Worldwide Information Growth Through 2010,  
March 2007, An IDC White Paper - sponsored by EMC  
<http://www.emc.com/collateral/analyst-reports/expanding-digital-idc-white-paper.pdf>

# Data Explosion!! 始於 2007 的「資料大爆炸」時代

Figure 1

Digital Information Created, Captured, Replicated Worldwide



Source: IDC, 2008

2009 年，IDC 預估  
2011 年會成長**十倍**！  
(相較 2006 年)

2006	161	EB
2007	281	EB
2010	988	EB (預測)
2011	1773	EB (預測)

出處：[The Diverse and Exploding Digital Universe](#),  
An Updated Forecast of Worldwide Information Growth Through 2011  
[March 2008](#), An IDC White Paper - [sponsored by EMC](#)  
<http://www.emc.com/collateral/analyst-reports/diverse-exploding-digital-universe.pdf>



# Data expanded 2x each year !! 每年約略兩倍



追蹤歷年的 IDC 數據：

2006 161 EB

2007 281 EB

2008 487 EB

2009 800 EB (0.8 ZB)

2010 988 EB (預測)

2010 1200 EB (1.2 ZB)

2011 1773 EB (預測)

2011 1800 EB (1.8 ZB)

景氣差而成長趨緩？  
或受新技術抑制？

出處：[Extracting Value from Chaos](#),  
June 2011, An IDC White Paper - sponsored by EMC

<http://www.emc.com/collateral/about/news/idc-emc-digital-universe-2011-infographic.pdf>

# What is Big Data?! 何謂『海量資料』？

海量資料泛指資料大小已無法用一般軟體擷取、管理與處理；  
單一資料集大小介於數十 TB 至數 PB 的資料。

'Big Data' = few dozen TeraBytes to PetaBytes in single data set.

## Definition

[edit]

Big data is a term applied to data sets whose size is beyond the ability of commonly used software tools to capture, manage, and process the data within a tolerable elapsed time. Big data sizes are a constantly moving target currently ranging from a few dozen terabytes to many petabytes of data in a single data set.

In a 2001 research report<sup>[14]</sup> and related conference presentations, then META Group (now Gartner) analyst, Doug Laney, defined data growth challenges (and opportunities) as being three-dimensional, i.e. increasing volume (amount of data), velocity (speed of data in/out), and variety (range of data types, sources). Gartner continues to use this model for describing big data.<sup>[15]</sup>

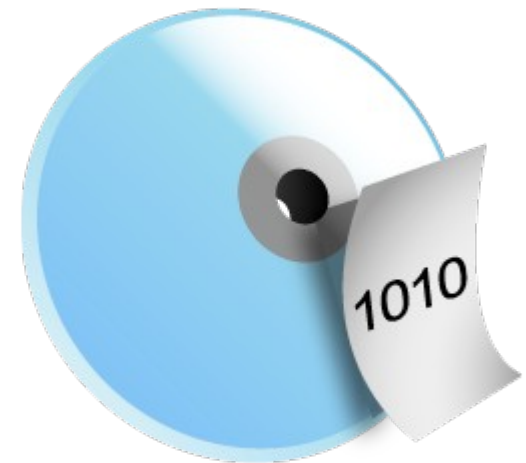
出處：[http://en.wikipedia.org/wiki/Big\\_data](http://en.wikipedia.org/wiki/Big_data)



多個檔案，容量 100TB



一個資料庫，容量 100TB

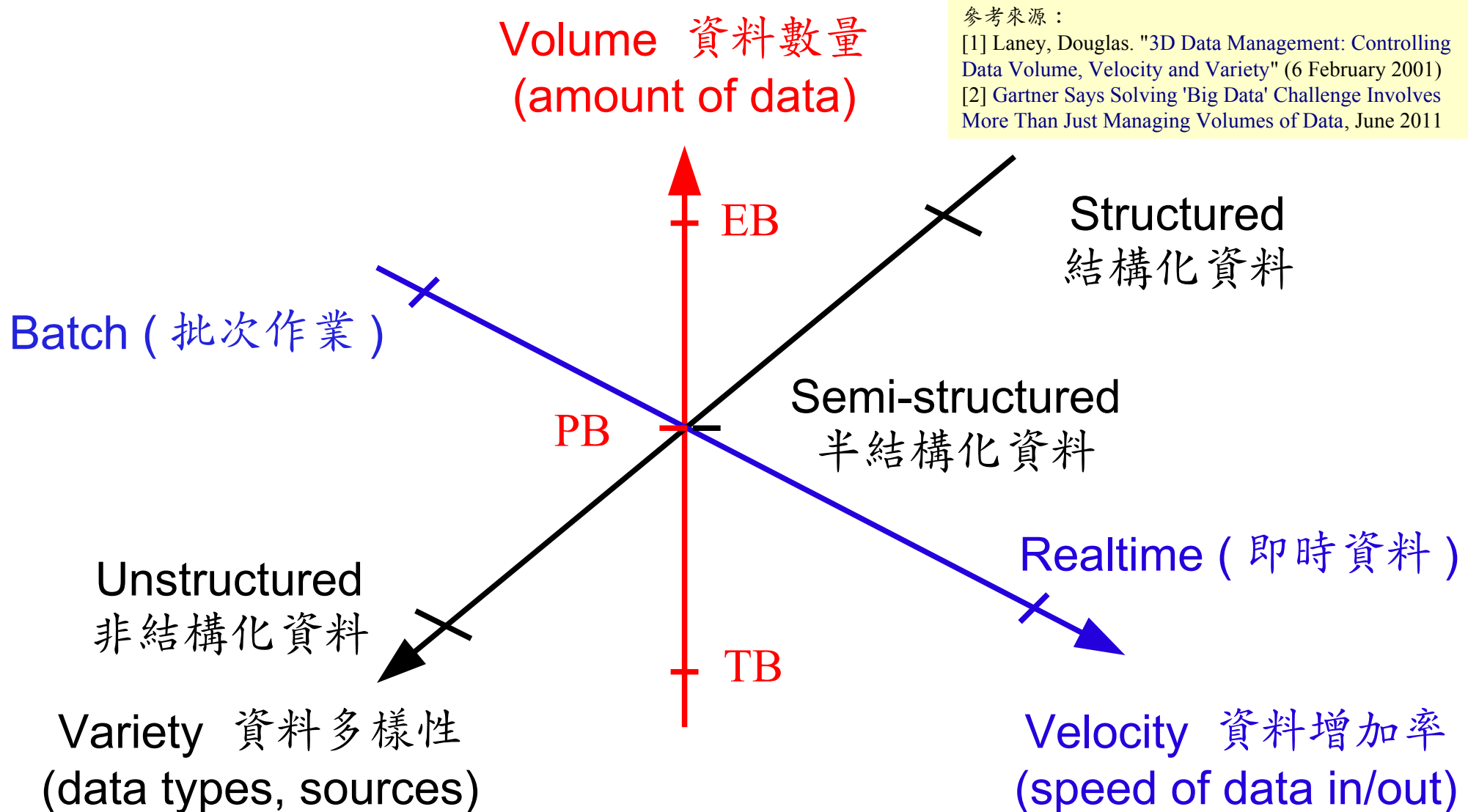


一個檔案，容量 100TB



# Gartner Big Data Model? 海量資料的模型?

海量資料的挑戰在於如何管理「數量」、「增加率」與「多樣性」



參考來源:

[1] Laney, Douglas. "3D Data Management: Controlling Data Volume, Velocity and Variety" (6 February 2001)

[2] Gartner Says Solving 'Big Data' Challenge Involves More Than Just Managing Volumes of Data, June 2011

# Six Dimensions of Big Data? 六個維度?



# 12D of Information Management? 12 個維度?



Source: Gartner (March 2011), 'Big Data' Is Only the Beginning of Extreme Information Management, 7 April 2011, <http://www.gartner.com/id=1622715>

# Agenda 演講大綱

What is Big Data ?

何謂海量資料

Why should we care?

為何需要關切

資料 Data

知識 Knowledge

智慧 Wisdom

WHY



花精靈-小葵

# Why we call it “SMART” !!

智慧打哪兒來？！

**Smart Phone**

智慧手機

**Smart Car**

智慧車輛

**Smart Grid**

智慧電網

**SMART**

哪裡長  
智慧了？

**Smart City**

智慧城市

**Smart Home**

智慧家庭

**Smart Meter**

智慧電錶

資料

Data

知識

Knowledge

智慧

Wisdom



# Can Machine understand You? 讓機器更懂你?

iPhone

Features Built-in Apps



Siri. Beta

Your wish is  
its command.

Siri on iPhone 4S lets you use your voice to send messages, schedule meetings, place phone calls, and more. Ask Siri to do things just by talking the way you talk. Siri understands what you say, knows what you mean, and even talks back. Siri is so easy to use and does so much, you'll keep finding more and more ways to use it.



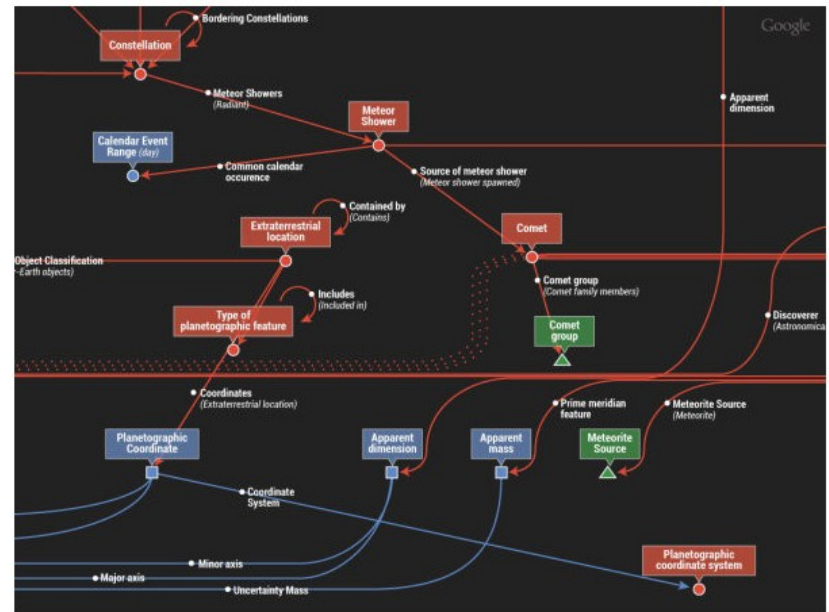
## Google將發展「人工智慧」 永久改變搜尋引擎

2012年02月15日 00:11

點評: 超級阿斯拉, 衝啊! (阿斯拉: 好的, 華人!)

記者黃郁棋 / 綜合報導

每個人都在猜, 下一波網路革命是什麼? 每個人都在猜, 未來的世界會如何運作? Google的資深副總Amit Singhai透露了一點訊息。「Google正努力從『單字』層面進展到『意義』層面, 未來搜尋引擎提供的不只是關鍵字搜尋, 搜尋引擎甚至會『明白』你到底要什麼。」

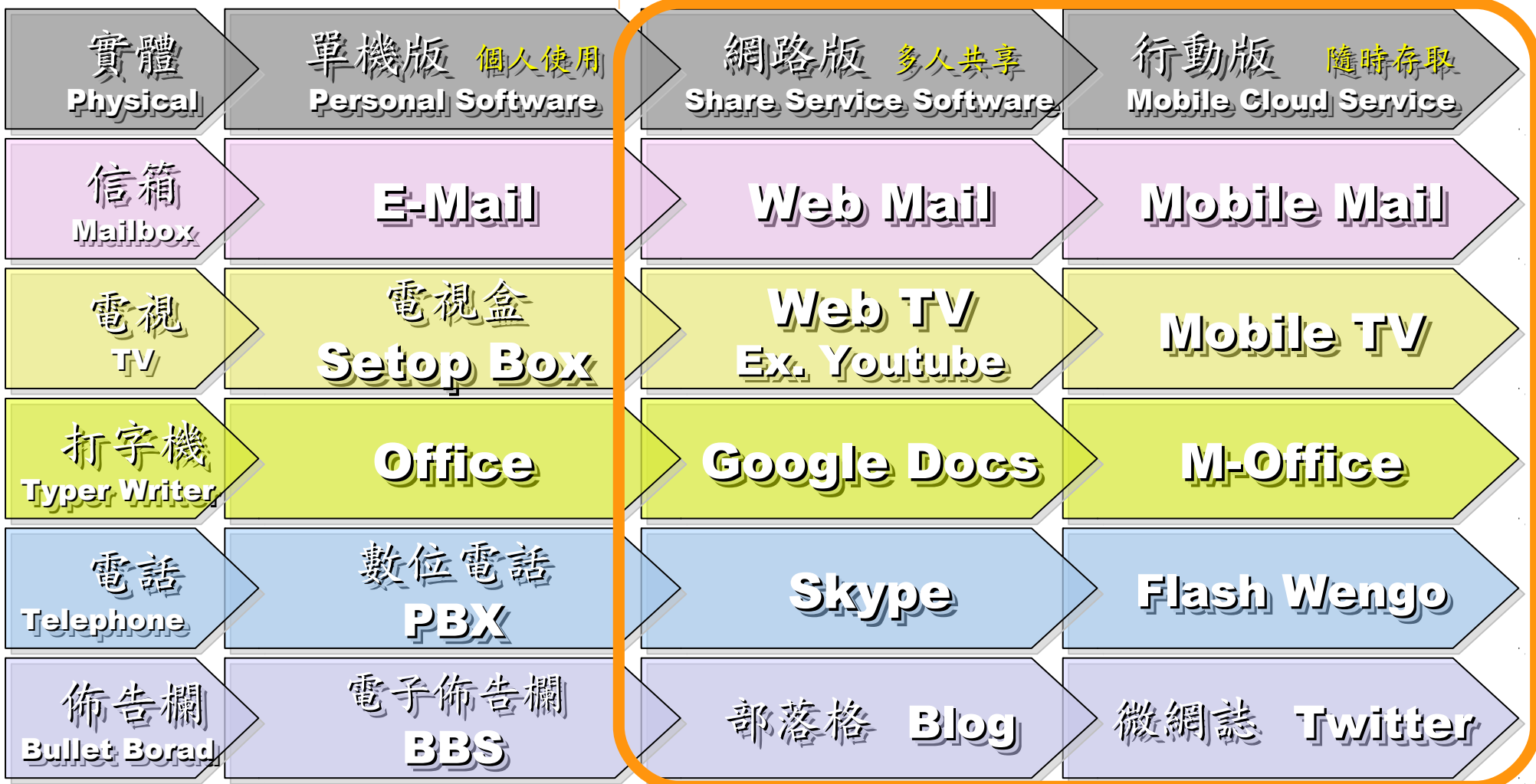


▲ Google未來將會朝「人工智慧」前進。(圖 / 取自mashable.com)



# Evolution of Software / Service

## 軟體演化勢必走向『智能化』



# The wisdom of Clouds (Crowds)

雲端序曲：雲端的智慧始終來自於群眾的智慧

2006年8月9日

Google 執行長施密特 ( Eric Schmidt ) 於SES'06會議中首次使用「雲端運算 ( Cloud Computing ) 」來形容無所不在的網路服務

2006年8月24日

Amazon 以 Elastic Compute Cloud 命名其虛擬運算資源服務



# Data is the source of Wisdom !!

## 用雲掌握資料，加以分析，形成智能給端用



# 雲

資料中心  
提供服務

## 雲端設計新思維：端的智能來自於雲的服務

## Devices share the wisdom of Cloud

# 端



各類裝置  
存取服務

# Agenda 演講大綱

What is Big Data ? 何謂海量資料

Why should we care? 為何需要關切

**When to deploy it ?** 何時導入技術

基礎建設 IaaS

分析平台 PaaS

智慧服務 SaaS

## WHEN



花精靈-小蠻



# National Definition of Cloud Computing

## 美國國家標準局 NIST 給雲端運算所下的定義

### 5 Characteristics

五大基礎特徵

### 4 Deployment Models

四個佈署模型

### 3 Service Models

三個服務模式

#### 1. On-demand self-service.

隨需自助服務

#### 2. Broad network access

隨時隨地用任何網路裝置存取

#### 3. Resource pooling

多人共享資源池

#### 4. Rapid elasticity

快速重新佈署靈活度

#### 5. Measured Service

可被監控與量測的服務

# 4 Deployment Models of Cloud Computing

## 雲端運算的四種佈署模型

Public Cloud

公用雲端



Target Market

is **S.M.B.**

主要客戶為  
中小企業

**Dynamic Resource Provisioning  
between public and private cloud**

私有雲端動態根據計算需求  
調用公用雲端的資源

Hybrid  
Cloud

以大型企業  
為主要客戶

**Enterprise is  
key market**

Community Cloud

社群雲端

Academia 學術為主



私有雲端

Private Cloud



# 3 Service Models of Cloud Computing

## 雲端運算的三種服務模式 (市場區隔)

### IaaS

Infrastructure as a Service

架構即服務

### PaaS

Platform as a Service

平台即服務

### SaaS

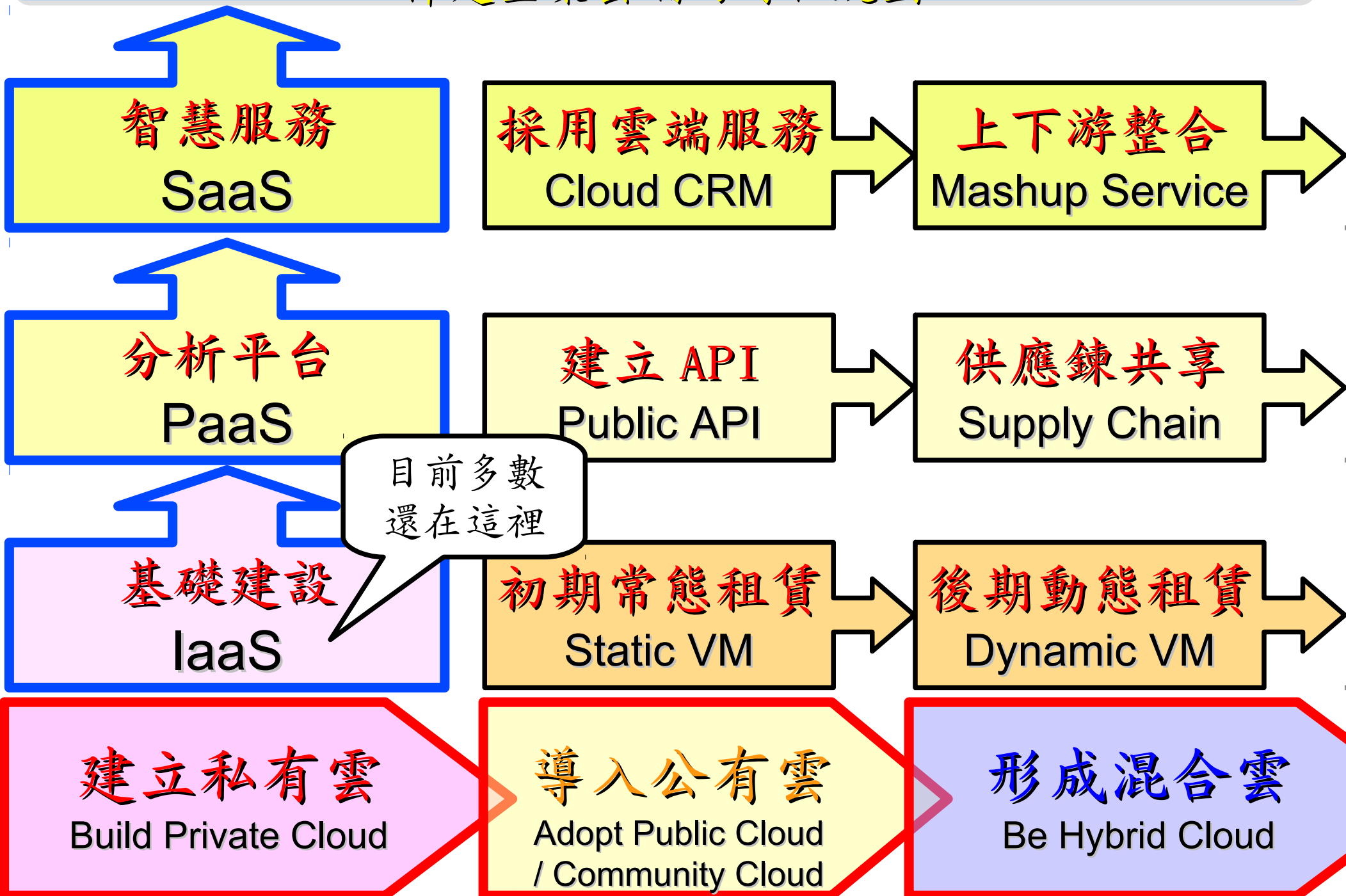
Software as a Service

軟體即服務



# Roadmap to build Your Enterprise Cloud !!

佈建企業雲端的時程規劃



# Agenda 演講大綱

What is Big Data ? 何謂海量資料

Why should we care? 為何需要關切

When to deploy it ? 何時導入技術

**How to handle it ?** 三大因應策略

儲存虛擬化 Dedup.

資料安全 Security

智慧服務 SaaS

## HOW



花精靈-麗兒

# Three Solutions !! 三種服務模式 vs. 三類因應對策

## SaaS

Software as a Service

軟體即服務

## Web 2.0

網頁服務

(A) 提供 API 介面

(B) 分散式資料庫

## PaaS

Platform as a Service

平台即服務

## Data Analysis

資料分析

(A) 資料整合

(B) 資料探勘

## IaaS

Infrastructure as a Service

架構即服務

## Virtualization

虛擬化技術

(A) 儲存虛擬化

(B) 備援與加密

# Agenda 演講大綱

**What is Big Data ?** 何謂海量資料

**Why should we care?** 為何需要關切

**When to deploy it ?** 何時導入技術

**How to handle it ?** 三大因應策略

**Who is key player ?** 誰是成功關鍵



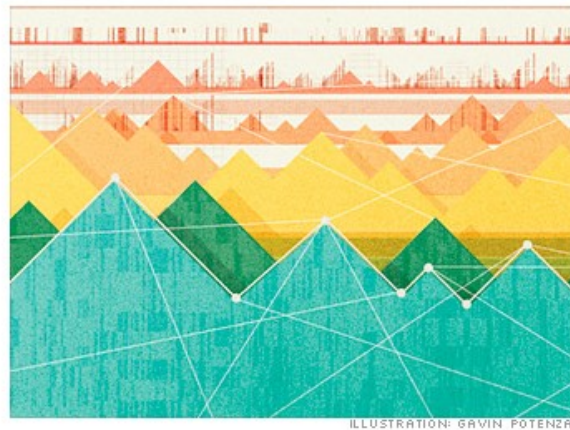


# Data Scientist !! 資料科學家 !!

## Data scientist: The hot new gig in tech

By Michal Lev-Ram, writer September 6, 2011: 5:00 AM ET

Companies that want to make sense of all their bits and bytes are hiring so-called data scientists - if they can find any.



FORTUNE -- The unemployment rate in the U.S. continues to be abysmal (9.1% in July), but the tech world has spawned a new kind of highly skilled, nerdy-cool job that companies are scrambling to fill: data scientist.

會「統計」的人照過來！

財星雜誌 (FORTUNE) 等均報導今年最熱門的職缺是「資料科學家」！

## What is data science?

Data science can be broken down into four essential parts.

### Mining data



Collecting and formatting the information

### Statistics



Information analysis

### Interpret



Representation or visualization in the form of presentations, infographics, graphs or charts

### Leverage



Implications of the data, application of the data, interaction using the data and predictions formed from studying it





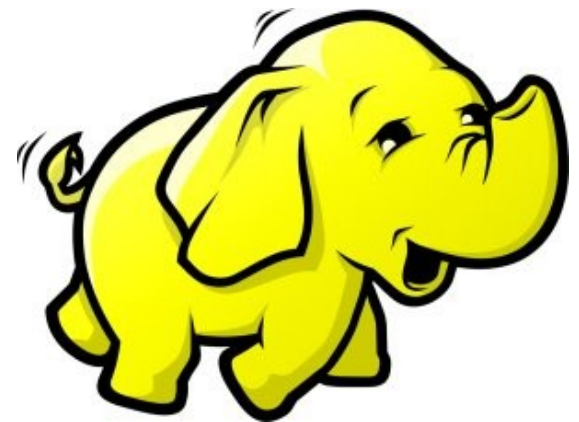
# 處理海量資料的資訊架構與關鍵技術

Technologies to build IT Stack for Big Data

**Jazz Wang**

**Yao-Tsung Wang**

**[jazz@nchc.org.tw](mailto:jazz@nchc.org.tw)**



# Hot Jobs in Big Data

## 從海量資料的熱門工作談起

**Data Mining**

**資料探勘**

**Data Visualization**

**資料視覺化**

**Data Analysis**

**資料分析**

**Data Manipulation**

**資料操控**

**Data Discovery**

**資料鑑識**

How to Get a Hot Job in Big Data, Dan Tynan, InfoWorld, March 19, 2012  
出處：<http://www.cio.com/article/print/702388>

# Applications of Data Mining

## 資料探勘的應用 - 搜尋引擎

搜尋結果

### 檔案搜尋

網址(D) 搜尋結果


搜尋小幫手

您想要搜尋什麼?

- 圖片、音樂、或視訊(P)
- 文件(文字處理、試算表, 等等)(O)
- 所有檔案和資料夾(L)
- 電腦或人員(C)
- 說明和支援中心裡的資訊(I)

您也可能想要...

- 搜尋網際網路(S)
- 變更喜好(G)



0 個物件

Gmail Calendar Documents Photos Sites Web More -

Search

All Mail

From

To

Subject

Has the words

Doesn't have

Has attachment

Date within 1 day of

Examples: f

Create

### 信件搜尋

發的交談

larwin.nchc.org.tw 於 2011年12月02日 (週五) 10時53分46秒 的交談

日 (週五)

- (10時53分48秒) Shunfa 楊順發
- (10時53分51秒) Jazz Yao-Tsung
- (10時54分08秒) Shunfa 楊順發
- (10時54分42秒) Jazz Yao-Tsung
- (10時54分49秒) Jazz Yao-Tsung
- (10時54分51秒) Jazz Yao-Tsung
- (10時55分02秒) Shunfa 楊順發
- (10時55分04秒) Shunfa 楊順發
- (10時55分39秒) Jazz Yao-Tsung

3 KiB

尋找(F)

關閉(C)

### 即時通訊搜尋

IEEE Xplore DIGITAL LIBRARY

BROWSE

- Journals & Magazines
- Conference Proceedings
- Standards
- Books
- Educational Courses

SIGN IN

Search 3,076,887 documents

SEARCH

Advanced Search | Preferences | Search Tips

### 資料庫搜尋

「網頁搜尋」

設Yahoo!奇摩為首頁 資訊展PK線上搶先

# YAHOO! 奇摩

網頁 | 知識+ | 圖片 | 影片 | 部落格 | 字典 | 新聞 | 購物 BETA

網頁搜尋

熱門: 第一美腿 12歲父親 嫩模女神 幼稚病 51區 花心星座 解夢 知識: 傷口癢竟是 電鍋料理

2011 資訊月 ONLINE 3G特展搶先看!!



# Applications of Data Visualization

## 資料視覺化的應用 - Infographics

### Data Scientist Study



The explosion in digital data, bandwidth and processing power — combined with new tools for analyzing the data — has sparked massive interest in the emerging field of data science. Organizations of all sizes are turning to people who are capable of translating this trove of data — created by mobile sensors, social media, surveillance, medical imaging, smart grids and the like — into predictive insights that lead to business value. Despite the growing opportunity, demand for data scientists has outpaced supply of talent, and will for the next five years. Who are data science practitioners, what skills do they need, and why are they so different?

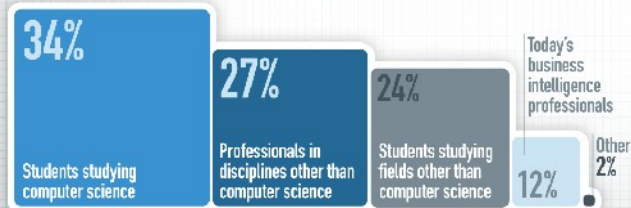
Over 2/3 believe demand for talent will outpace the supply of data scientists

OVER THE NEXT FIVE YEARS, DEMAND FOR DATA SCIENTISTS WILL:



Only 12% see today's BI professional as the best source for new data scientists

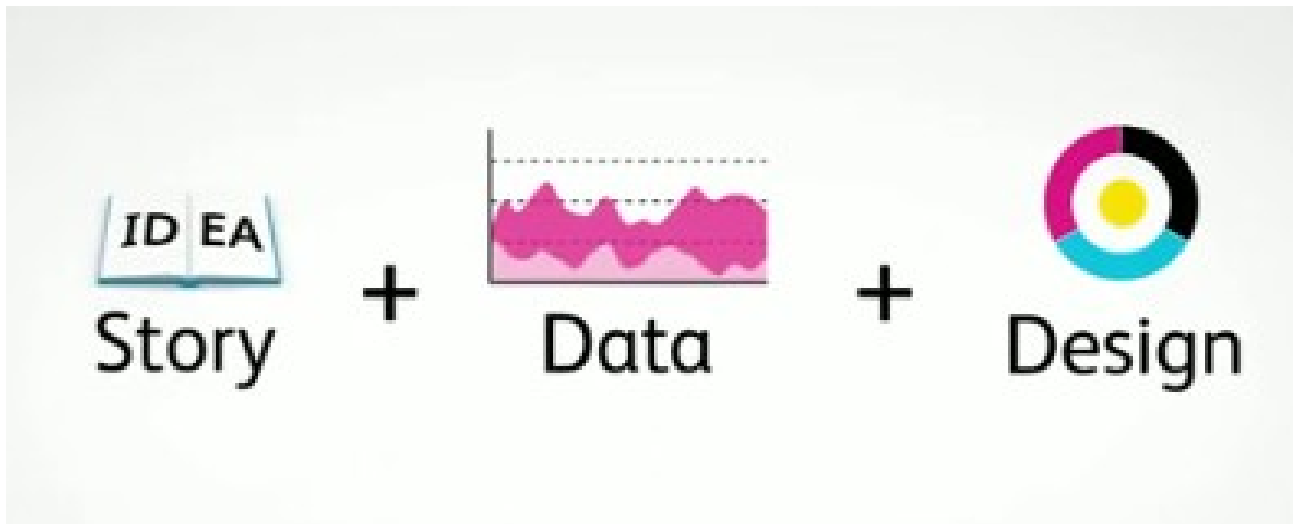
THE BEST SOURCE OF NEW DATA SCIENCE TALENT IS:



DUET TO THE ROUNDING, SOME PERCENTAGES MAY NOT ADD UP TO 100

Lack of training and resources are the biggest obstacle to data science in organizations

THE BIGGEST OBSTACLE TO DATA SCIENCE ADOPTION IN OUR ORGANIZATION IS:



參考來源：未來「夯」職業：資料科學家  
淺談超吸睛的資訊圖表

<http://www.bnext.com.tw/print/article/id/21740>  
<http://www.inside.com.tw/2011/04/13/infographics>

# Applications of Data Analysis

## 資料分析的應用 - 商業智慧 (BI)





# Applications of Data Discovery

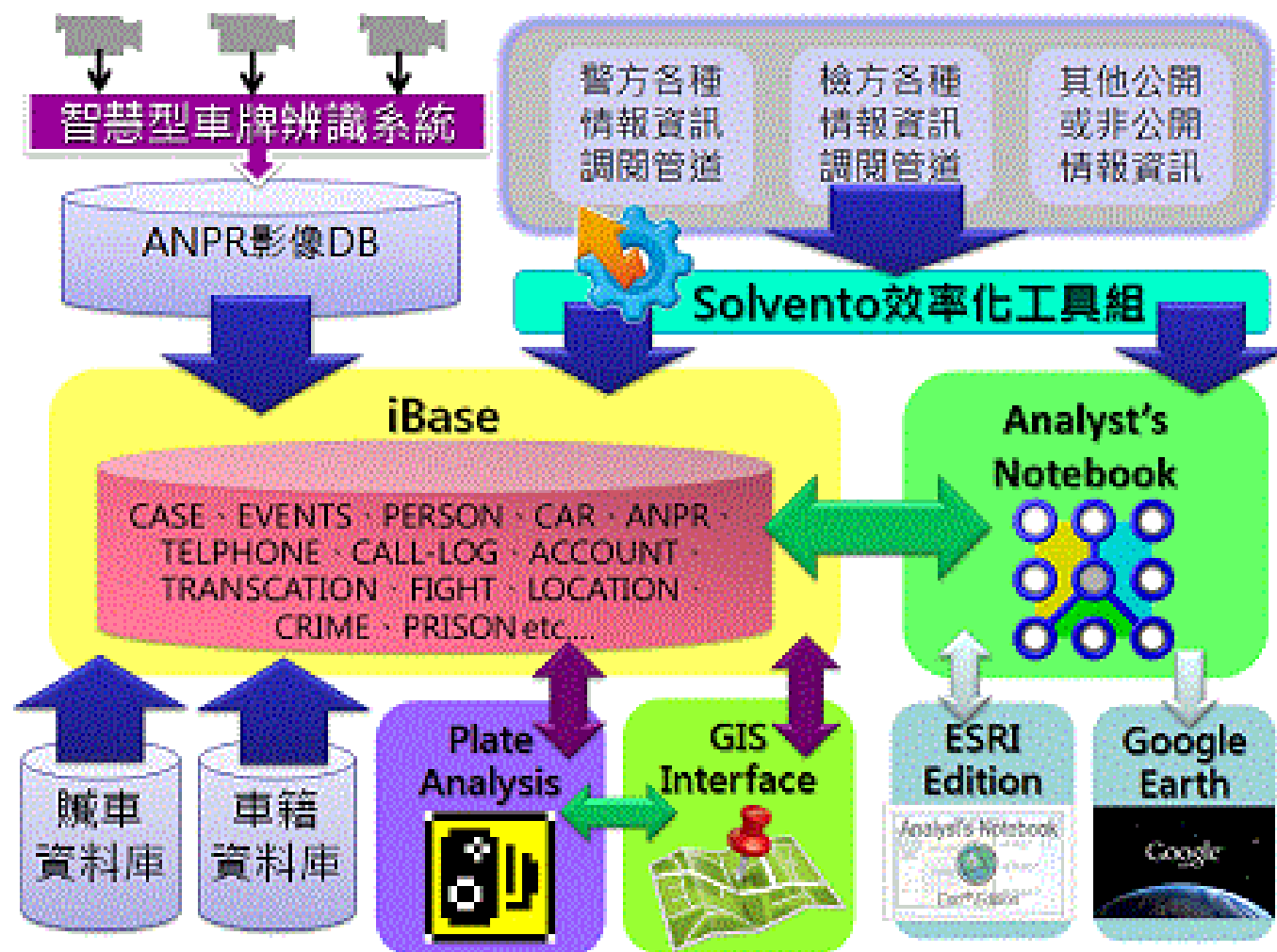
## 數位鑑識 - 資訊與法律的結合

### 電腦鑑識 & 會計鑑識

[http://www.solventosoft.com/upload/ANPR\\_02s.gif](http://www.solventosoft.com/upload/ANPR_02s.gif)

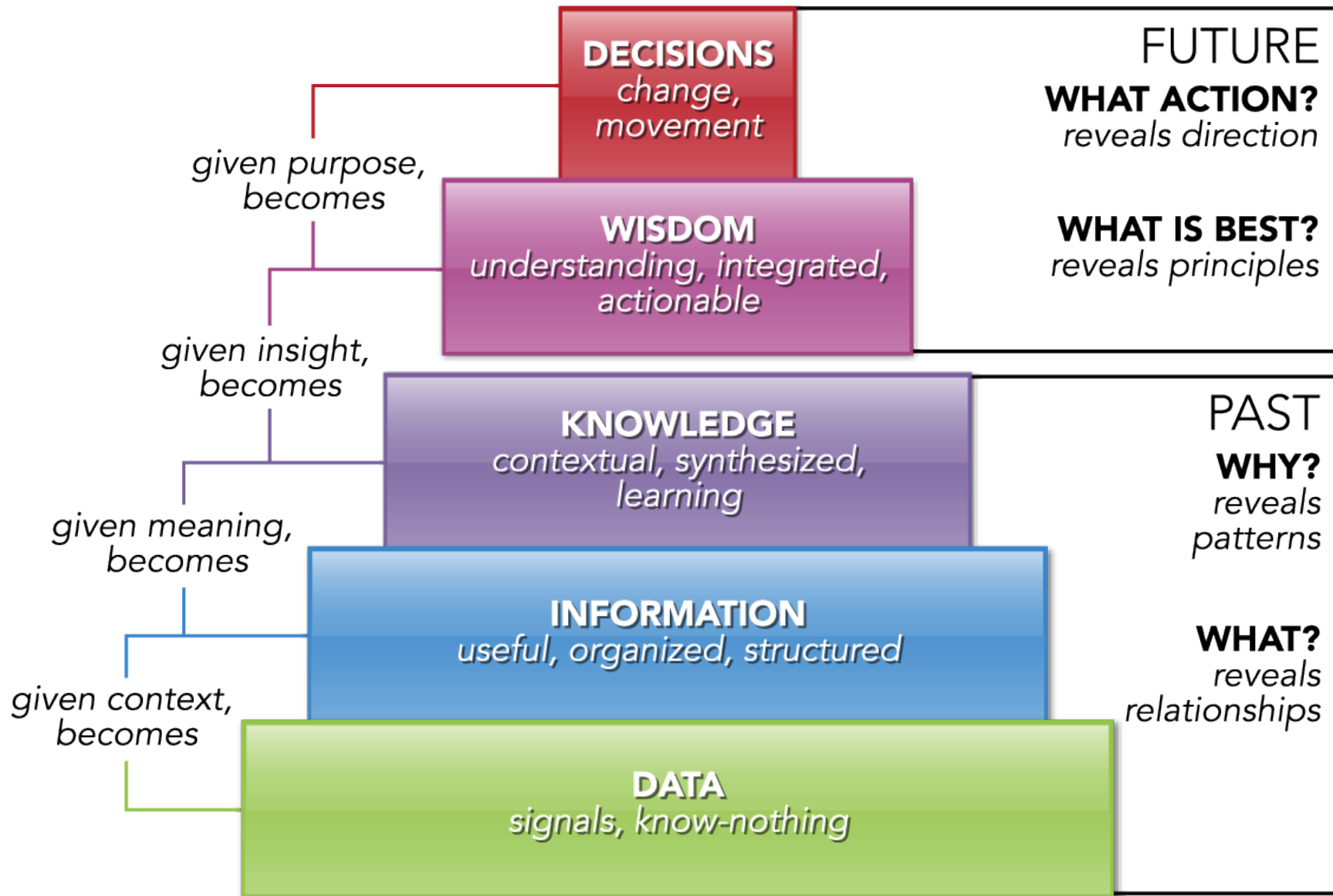


<http://blog.udn.com/kf0630/6018593>



# Data, Information, Knowledge, Wisdom

## 知識管理模型：資料、資訊、知識與智慧



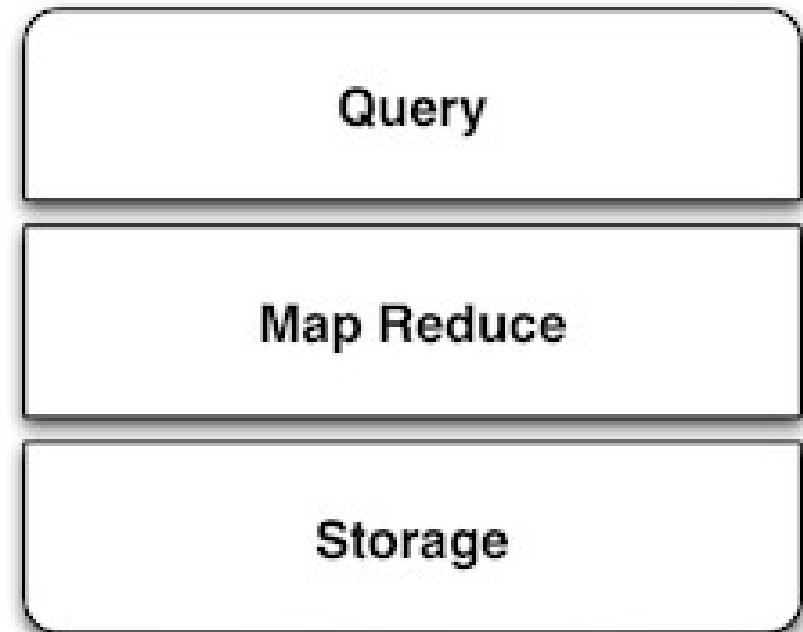
# The SMAQ stack for big data

## 海量資料處理的資訊架構

做網頁相關的人可能聽過 LAMP



未來處理海量資料的人必需知道  
SMAQ (Storage, MapReduce and Query)

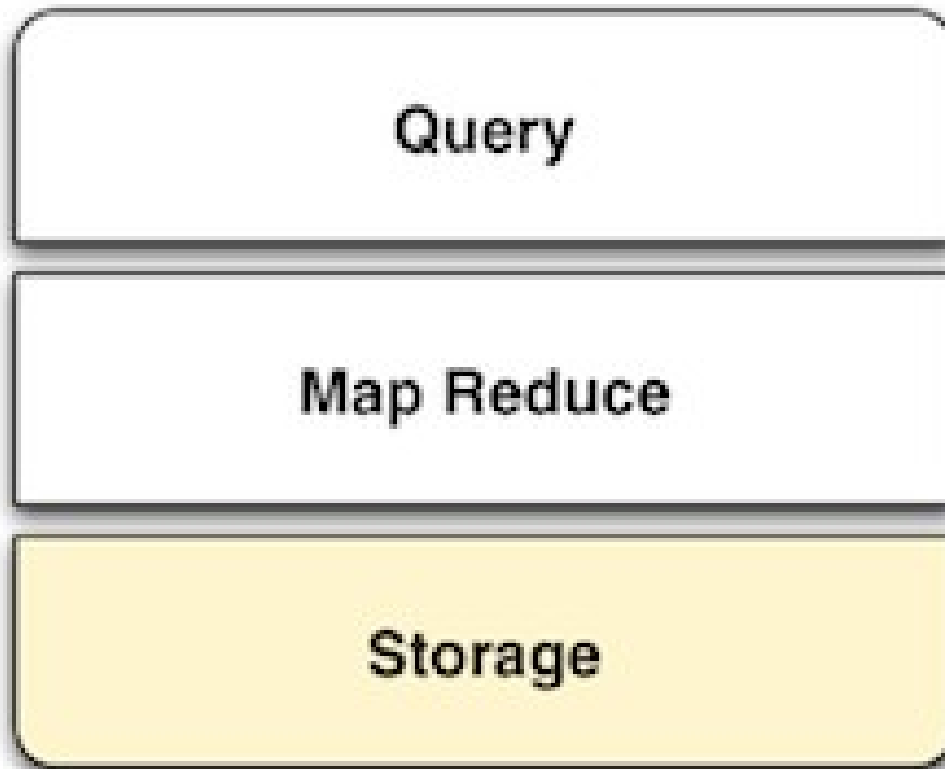


參考來源：The SMAQ stack for big data，Edd Dumbill，22 September 2010，  
<http://radar.oreilly.com/2010/09/the-smaq-stack-for-big-data.html>

圖片來源：<http://smashingweb.ge6.org/wp-content/uploads/2011/10/apache-php-mysql-ubuntu.png> 37

# The SMAQ stack for big data

## 海量資料處理的資訊架構



用來儲存分散、沒有關聯  
的非結構化資料

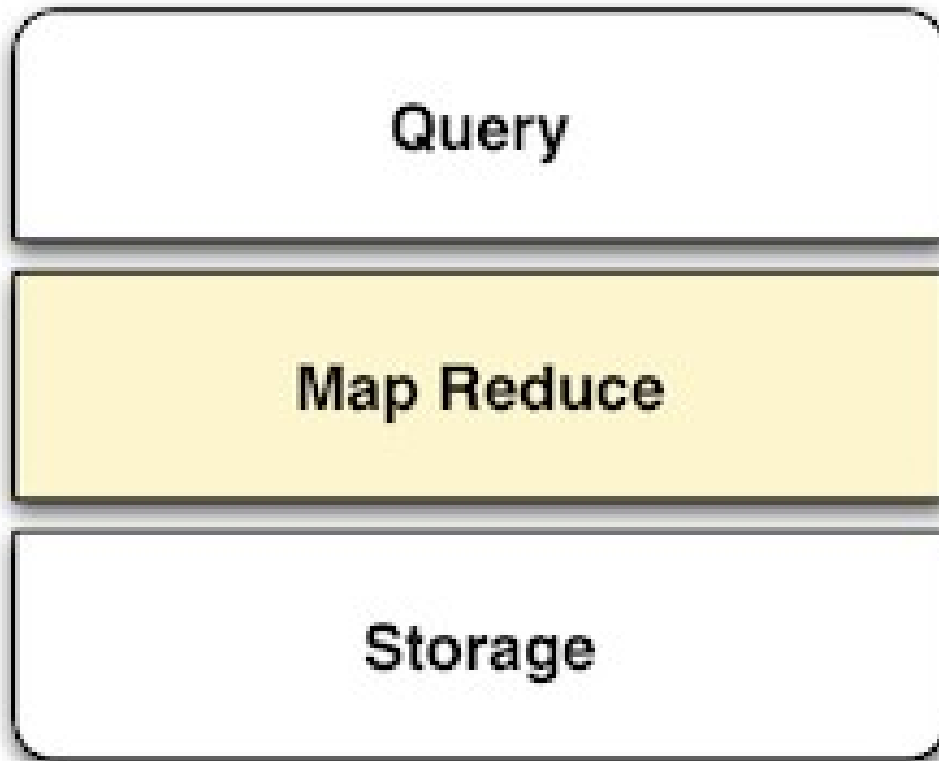
### Key features

- Distributed
- Non-relational or unstructured

# The SMAQ stack for big data

## 海量資料處理的資訊架構

運用批次處理的方式，將  
運算工作平均分散到許多  
的伺服器做運算。



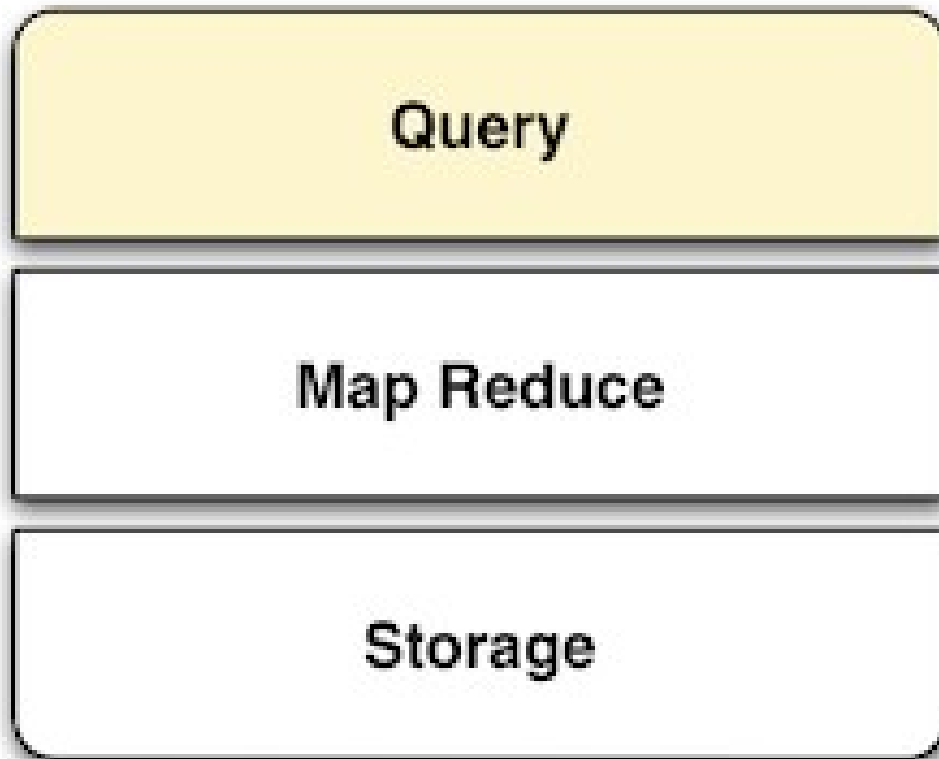
### Key features

- Distributes computation over many servers
- Batch processing model



# The SMAQ stack for big data

## 海量資料處理的資訊架構



### Key features

- Efficient way of defining computation
- Platform for user friendly analytical systems

將算完的結構化資料儲存到可供查詢的資料庫系統

# Three Core Technologies of Google ....

## Google 的三大關鍵技術 .....

- Google 在一些會議分享他們的三大關鍵技術
- Google shared their design of web-search engine
  - SOSP 2003 :
    - “The Google File System”
    - <http://labs.google.com/papers/gfs.html>
  - OSDI 2004 :
    - “MapReduce : Simplified Data Processing on Large Cluster”
    - <http://labs.google.com/papers/mapreduce.html>
  - OSDI 2006 :
    - “Bigtable: A Distributed Storage System for Structured Data”
    - <http://labs.google.com/papers/bigtable-osdi06.pdf>



# Open Source Mapping of Google Core Technologies

## Google 三大關鍵技術對應的自由軟體

### BigTable

A huge key-value datastore

HBase, Hypertable  
Cassandra, ....

### MapReduce

To parallel process data

Hadoop MapReduce API  
Sphere MapReduce API, ...

### Google File System

To store petabytes of data

Hadoop Distributed File System (HDFS)  
Sector Distributed File System

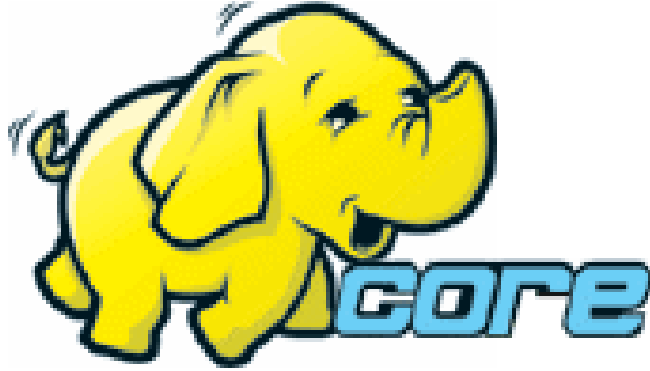
更多不同語言的 MapReduce API 實作：

<http://trac.nchc.org.tw/grid/intertrac/wiki%3Ajazz/09-04-14%23MapReduce>

其他值得觀察的分散式檔案系統：

- IBM GPFS - <http://www-03.ibm.com/systems/software/gpfs/>
- Lustre - <http://www.lustre.org/>
- Ceph - <http://ceph.newdream.net/>

# Hadoop

- <http://hadoop.apache.org>
  - Hadoop 是 Apache Top Level 開發專案
  - **Hadoop is Apache Top Level Project**
  - 目前主要由 Yahoo! 資助、開發與運用
  - **Major sponsor is Yahoo!**
  - 創始者是 Doug Cutting，參考 Google Filesystem
  - **Developed by Doug Cutting, Reference from Google Filesystem**
  - 以 Java 開發，提供 HDFS 與 MapReduce API。
  - **Written by Java, it provides HDFS and MapReduce API**
  - 2006 年使用在 Yahoo 內部服務中
  - **Used in Yahoo since year 2006**
  - 已佈署於上千個節點。
  - **It had been deploy to 4000+ nodes in Yahoo**
  - 處理 Petabyte 等級資料量。
  - **Design to process dataset in Petabyte**
- 
- Facebook、Last.fm  
、Joost are also  
powered by Hadoop**

# Sector / Sphere

- <http://sector.sourceforge.net/>
- 由美國資料探勘中心研發的自由軟體專案。
- **Developed by National Center for Data Mining, USA**
- 採用 C/C++ 語言撰寫，因此效能較 Hadoop 更好。
- **Written by C/C++, so performance is better than Hadoop**
- 提供「類似」Google File System 與 MapReduce 的機制
- **Provide file system similar to Google File System and MapReduce API**
- 基於UDT高效率網路協定來加速資料傳輸效率
- **Based on UDT which enhance the network performance**
- Open Cloud Testbed有提供測試環境，並開發Ma1Stone效能評比軟體
- **Open Cloud Consortium provide Open Cloud Testbed and develop Ma1Stone toolkit for benchmark**

**Sector-Sphere**

National Center for Data Mining  
University of Illinois at Chicago



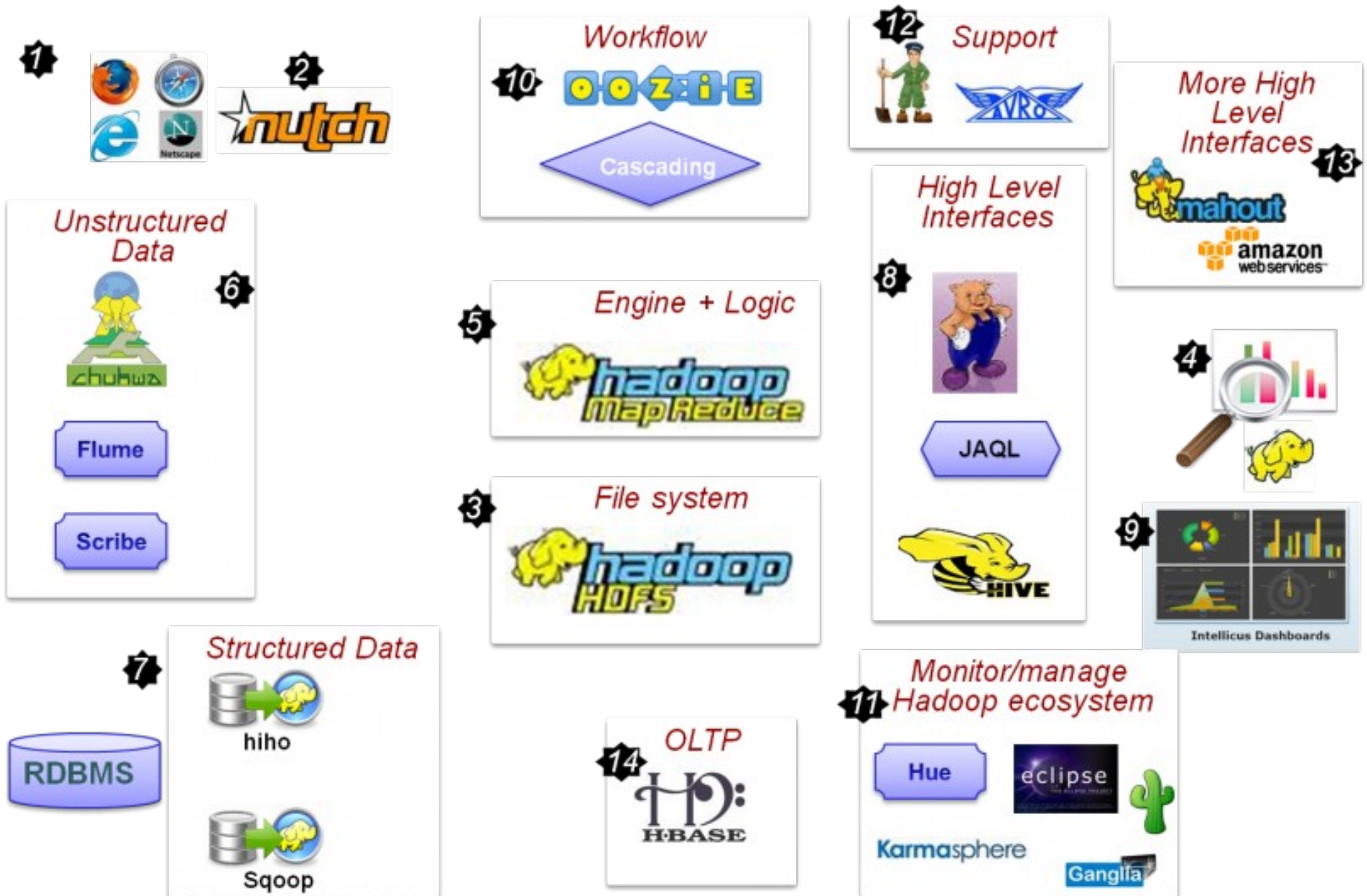
Open Data Group

<http://www.opendatagroup.com/>



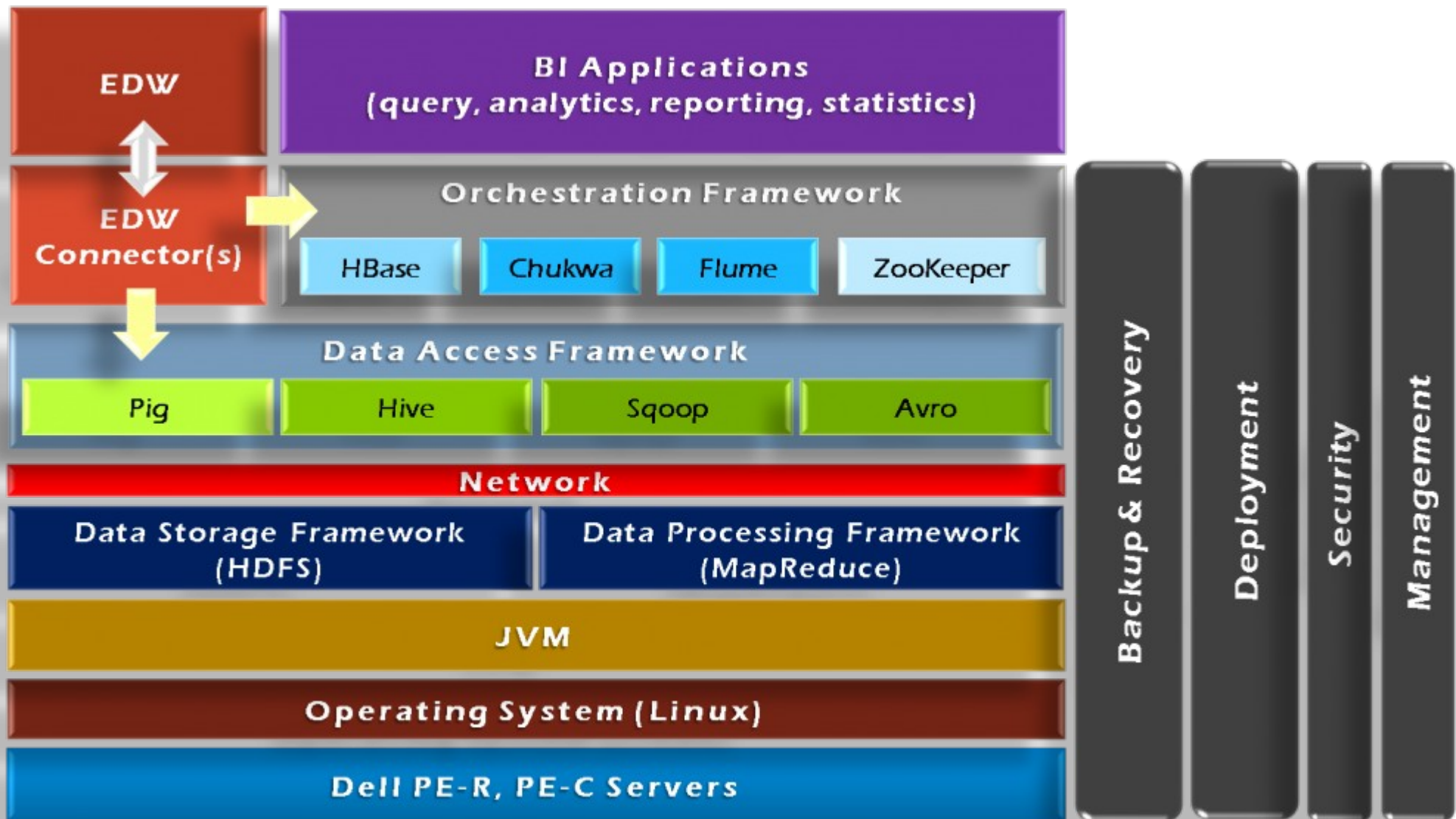
# Why we choice Hadoop? Good Ecosystem!

豐富的生態系建構出處理海量資料的工具庫



# BI and EDW build on Hadoop Ecosystem

運用 Hadoop 生態系搭建資料倉儲與商業智慧分析



# Build your own search engine, too

您也能用 **Hadoop** 搭建自己的搜尋引擎

Web UI ( Crawlzilla Website + Search Engine)

JSP + Servlet + JavaBean

Nutch

Lucene

Crawlzilla System Management

Tomcat

Hadoop

PC1

PC2

PC3



# Microsoft love Hadoop, too

## 微軟幫 Azure 還有 SQL Server 都接上 Hadoop



SQL Server | All Microsoft Sites | United States | Change | Search Microsoft | bing | Web

Microsoft SQL Server

Contact Us > | Facebook | Twitter | YouTube

About SQL Server | Solutions & Technologies | Editions | Get SQL Server | Learning Center | Partners

### Business Intelligence

Share this page

#### Big Data Analytics

#### Big Data Solution

Unlock business insights from all your structured and unstructured data, including large volumes of data not previously activated, with Microsoft's Big Data solution. Microsoft's end-to-end roadmap for Big Data embraces Apache Hadoop™ by distributing enterprise class Hadoop based solutions on both Windows Server and Windows Azure. Our solution is also integrated into the Microsoft BI tools such as SQL Server Analysis Services, Reporting Services and even PowerPivot and Excel. This enables you to do BI on all your data, including those in Hadoop.

#### Key Benefits

- Broader access of Hadoop to end users, IT professionals and Developers, through easy installation and configuration and simplified programming with JavaScript.
- Enterprise ready Hadoop distribution with greater security, performance, ease of management and options for Hybrid IT usage.

參考來源：Big Data Solution | Microsoft SQL Server 2008 R2

<http://www.microsoft.com/sqlserver/en/us/solutions-technologies/business-intelligence/big-data-solution.aspx>

# Oracle love Hadoop, too

## Oracle 也接上 Hadoop



CNET > News > Software, Interrupted

## Cloudera teams up to connect Oracle and Hadoop

Cloudera and Quest software are partnering to provide connectivity between Oracle and Hadoop.



by [Dave Rosenberg](#) | June 21, 2010 5:30 AM PDT



This week [Cloudera](#), a provider of software and services for the Apache Hadoop project, is set to announce a partnership with [Quest Software](#) to develop, support, and distribute an Oracle connector for Hadoop.





# Hinet Application of Big Data

## 中華電信已經在做的海量資料應用

Business  
Next 數位時代

### 中華電信：分析駭客行為，拓展對外新服務

撰文者：趙郁竹

發表日期：2012-03-06



[214期雜誌精選]

全球最大的中華電信提供行動電話、市話、寬頻固網、MOD……，各種業務服務，加起來的用戶數就有3000萬，比全台灣人口還多，光是單月帳務數量就高達100億筆資料。除了電信、寬頻服務，還有日益增加的數位服務、行動增值服務，從服務內容到客戶端，累積出的資料相當驚人。

「資料量越來越大，日常分析工作需要很多時間，但新的運算技術有效解決了這個問題，」中華電信資訊處處長陳明仕說。2010年開始，因為中華電信本身的資料運算需求，採用分散式運算架構Hadoop技術，打造出大資料運算平台，不但解決了自身的資料問題，還能對外提供資料運算應用。

以MOD為例，一天有幾千萬筆資料，如何找出使用者在什麼時段做了什麼事？廣告效益又如何？「用傳統的方法，需要400分鐘才能分析完；用Hadoop大資料平台，13分鐘就能解決，節省非常多時間，」他說。

#### 追蹤再拆解

大資料運算技術除了節省時間，還能防止駭客入侵。「駭客的攻擊行為都有模式可循，」陳明仕解釋，就像球賽一樣，了解進攻模式就能防守。用戶的資料保護是第一要務，因此透過行為模式分析，能有效保護企業資訊安全，也保障客戶的個資安全。

參考來源：中華電信：分析駭客行為，拓展對外新服務，發表日期：2012-03-06

<http://www.bnext.com.tw/print/article/id/22333>

# Hinet Application of Big Data

## 中華電信已經在做的海量資料應用

IT ithome.com.tw

### 中華電信用Hadoop技術分析通話明細

READ LATER

面對資料快速成長以及非結構性資料的增加，中華電信資訊處第四科科長楊秀一表示，中華電信近來利用Hadoop雲端運算技術自行開發了一個專門用來分析非結構化資料的巨量資料（Big Data）運算平臺，嘗試在資料進到資料倉儲系統之前，先進行資料的分析與處理以減少資料倉儲的資料量。

近年來行動語音市場趨於飽和，為了掌握用戶特性進行客製化行銷，一份資料要進行分析，就會被多次複製，因此即使用戶增加趨緩，但中華電信擁有的資料量仍快速暴增。

中華電信用來分析的資料模型最早於10多年前已有雛形，但當初主要用於行動語音分析。一直到2009年，他們完整導入Teradata的電信業邏輯資料模型cLDM 9.0版，整合更多電信服務的用戶資料。楊秀一表示，當初導入該模型的目的主要是為了整合行動語音、固網、數據的資料，進行以人為中心的分析模式。在導入之前，中華電信的資料模型是以設備為中心，因為不同設備的記錄資料儲存在不同的資料庫，無法進行整合性的分析。

參考來源：中華電信用 Hadoop 技術分析通話明細，發表日期：2011-06-12  
<http://www.ithome.com.tw/itadm/article.php?c=68023>