



財團法人國家實驗研究院

國家高速網路與計算中心
NATIONAL CENTER FOR HIGH-PERFORMANCE COMPUTING

雲端運算簡介


王耀聰 陳威宇

Jazz@nchc.org.tw

waue@nchc.org.tw

2008. 04 . 27-28

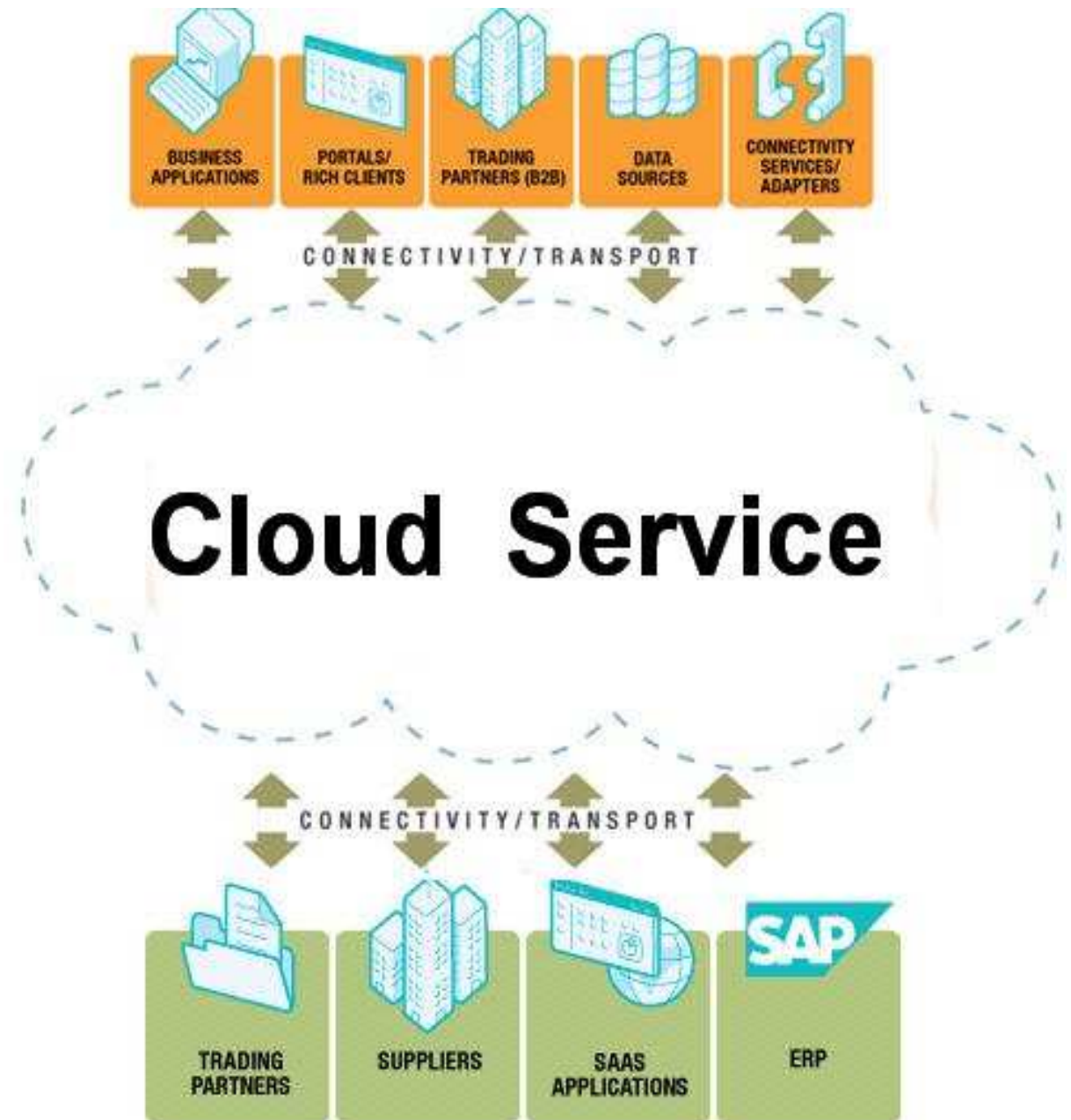
國家高速網路與計算中心(NCHC)



雲端運算??

雲端服務

- Web Email
- 線上掃毒
- YouTube
- 線上文件
- 部落格
- ...



雲端運算特色

虛擬化

超大規模

高可靠度

使用者付費

高通用性

高擴充性

成本低

雲端運算的架構

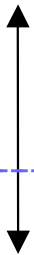
User Level



User-Level
Middleware



Core
Middleware



System Level

應用

Social Computing, Enterprise, ISV, ...

程式語言

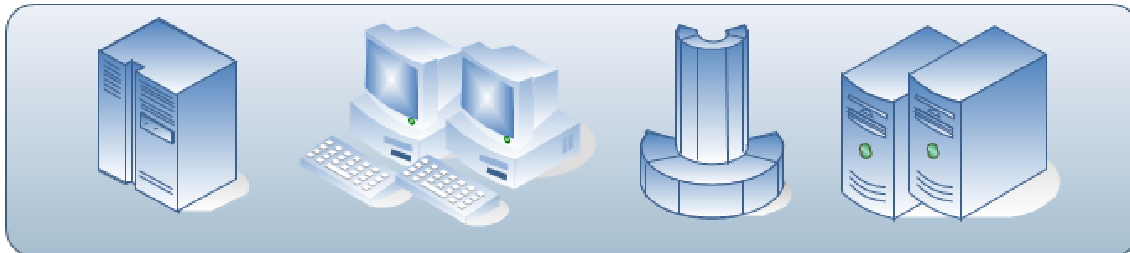
Web 2.0 介面, Mashups, Workflows, ...

控制

Qos Negotiation, Admission Control,
Pricing, SLA Management, Metering...

虛擬化

VM, VM management and Deployment



現有的雲端運算服務

- Windows
- Google
- Amazon
- Yahoo
-



Amazon : Web Service

- AWS
- 虛擬化的技術：Amazon EC2
 - Small (Default) \$0.10 per hour \$0.125 per hour
 - All Data Transfer \$0.10 per GB
- 儲存服務：Amazon S3
 - \$0.150 per GB – first 50 TB / month of storage used
 - \$0.100 per GB – all data transfer in
 - \$0.01 per 1,000 PUT, COPY, POST, or LIST requests
- 觀念：Paying for What You Use



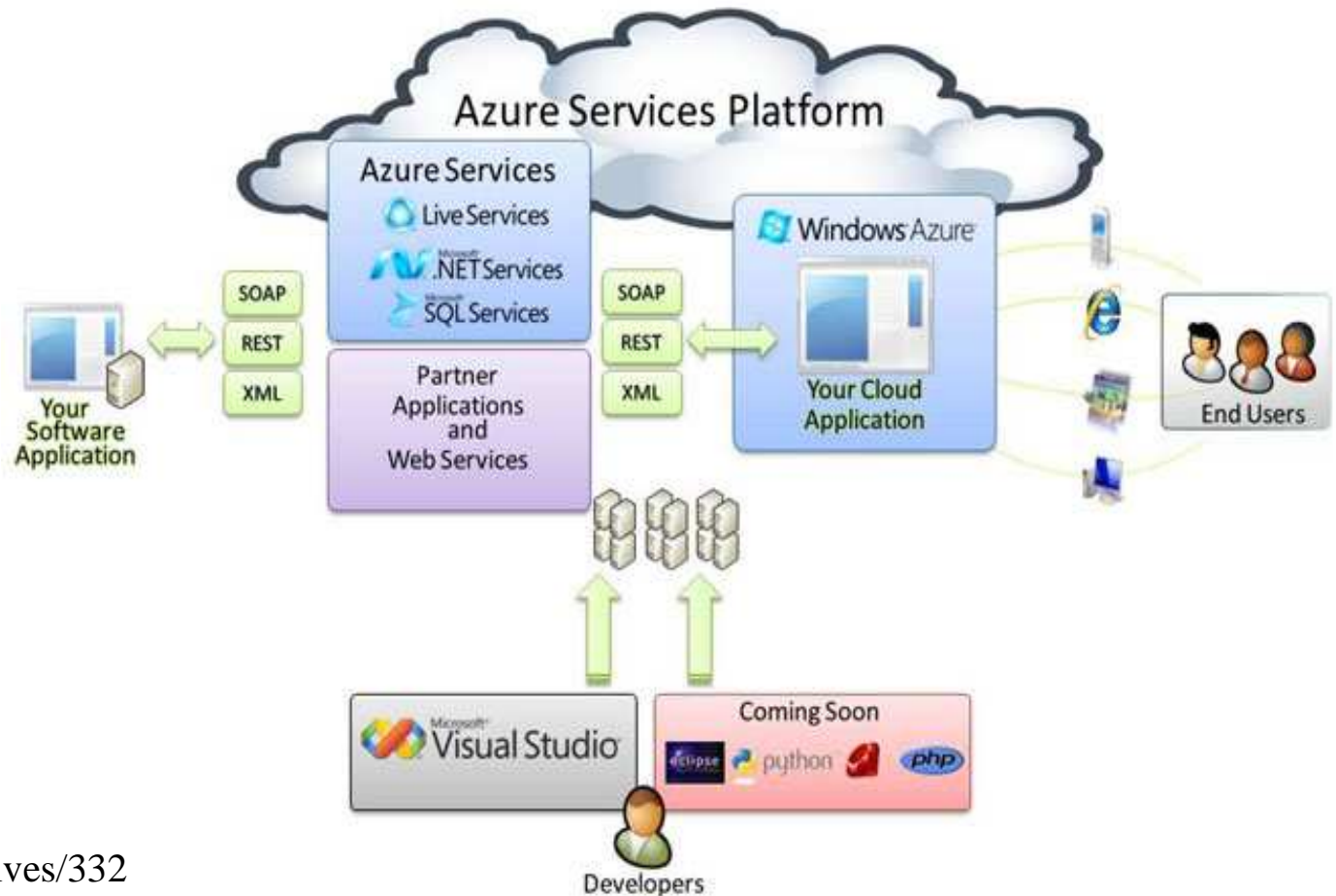
Google : App Engine

- 網路平台，讓開發者可自行建立網路應用程式於 google 平台中。
- 提供：
 - 500MB of storage
 - up to 5 million page views a month
 - 10 applications per developer account
- Limit：
 - Language: Python、Java
 - web applications



Windows : Azure

- Windows Azure 是一套雲端服務作業系統。作為 Azure 服務平台的開發、服務代管及服務管理環境。
- .Net services
- SQL services
- Live services



Yahoo : Hadoop

- Apache 項目，Yahoo 資助、開發與運用
 - 2006年開始參與開源的雲端運算框架Hadoop，並將其使用在內部服務中。
 - 2008年2：目前最大的Hadoop應用
 - 2千臺伺服器
 - 執行超過1萬個Hadoop虛擬機器
 - 5 Petabytes的網頁內容
 - 分析1兆個網路連結



雲端運算產業類型

SaaS

Software as a Service

PaaS

Platform as a Service

IaaS

Infrastructure as a Service

雲端運算產業

架構即服務

- 提供了核心計算資源和網絡架構的服務
- infrastructure stack:
 - Full OS access
 - Firewalls
 - Routers
 - Load balancing

IaaS

雲端運算產業

Examples

- Flexiscale
- AWS: EC2 (Amazon Elastic Compute Cloud)

IaaS



雲端運算產業

平台即服務

- 提供平台給系統管理員和開發人員，以為它構建、測試及部署定製應用程序
- 管理系統的成本昂貴
- Popular services
 - Storage
 - Database
 - Scalability

PaaS

IaaS

雲端運算產業

Examples

- Google App Engine
- AWS: S3 (Simple Storage Service)
- Microsoft Azure

PaaS

IaaS

雲端運算產業

SaaS

PaaS

IaaS

軟體即服務

- 通過Internet提供軟體的模式，用戶向提供商租用基於Web的軟體，來管理企業經營活動，且無需對軟體進行維護，服務提供商會全權管理和維護軟體

雲端運算產業

SaaS

軟體即服務

- 不用管理硬體與軟體
- 操作簡單 (瀏覽器)
- Pay per use
- Instant Scalability
- Security
- Reliability

PaaS

IaaS

雲端運算產業

SaaS

PaaS

IaaS

Examples

- Google Docs
- CRM
- Financial Planning
- Human Resources
- Word processing
- Salesforce.com

比較表

服務 屬性	Amazon EC2	Google App Engine	Microsoft Azure	Yahoo Hadoop
架構	Iaas/Paas	Paas	Paas	Software
服務型態	Compute/ Storage	Web application	Web and non- web	Software
管理技術	OS on Xen hypervisor	Application container	OS through Fabric controller	Map / Reduce Architecture
使用者介面	EC2 Command-line tools	Web-based Administration console	Windows Azure portal	Command line and web
APIs	yes	yes	yes	yes
收費	yes	maybe	yes	no
程式語言	AMI (Amazon Machine Image)	Python	.NET framework	Java,



財團法人國家實驗研究院

國家高速網路與計算中心
NATIONAL CENTER FOR HIGH-PERFORMANCE COMPUTING

Hadoop 簡介

王耀聰 陳威宇

Jazz@nchc.org.tw

waue@nchc.org.tw

2008. 04 . 27-28

國家高速網路與計算中心(NCHC)

Outline

- 什麼是 Hadoop ?
- 有什麼特色 ?
- 怎麼來的呢 ?
- 有誰在用 ?
- 有實用案例嗎 ?

Hadoop ?

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

Hadoop

- 以Java開發
- 自由軟體
- 上千個節點
- Petabyte等級的資料量
- 創始者 Doug Cutting
- 為Apache 軟體基金會的 top level project

特色

- 巨量
 - 擁有儲存與處理大量資料的能力
- 經濟
 - 可以用在由一般PC所架設的叢集環境內
- 效率
 - 藉由平行分散檔案的處理以致得到快速的回應
- 可靠
 - 當某節點發生錯誤，系統能即時自動的取得備份資料以及佈署運算資源

起源:2002-2004

- Lucene
 - 用Java設計的高效能文件索引引擎API
 - 索引文件中的每一字，讓搜尋的效率比傳統逐字比較還要高的多
- Nutch
 - nutch是基於開放原始碼所開發的web search
 - 利用Lucene函式庫開發

起源：Google論文

- Google File System
 - 可擴充的分散式檔案系統
 - 設計目的在於可以給大量的用戶提供總體性能較高的服務
 - 適用於分散式、對大量資訊進行存取的應用
 - 可運作在一般的普通主機上，且提供錯誤容忍的能力
- “The Google File System “發表於SOSP'03 October，並將設計的概念公開

起源：Google論文

- Google's GFS & MapReduce papers published:
 - SOSP 2003 : “The Google File System”
 - OSDI 2004 : “MapReduce : Simplified Data Processing on Large Cluster”
 - OSDI 2006 : “Bigtable: A Distributed Storage System for Structured Data”
- directly address Nutch's scaling issues

起源:2004~

- Dong Cutting 開始參考論文來實做
- Added DFS & MapReduce implement to Nutch
- Nutch 0.8版之後，Hadoop為獨立項目
- Yahoo 於2006年僱用Dong Cutting 組隊專職開發
 - Team member = 14 (engineers, clusters, users, etc.)

誰在用 Hadoop

- Yahoo 為最大的贊助商
- IBM 與 Google 在大學開授雲端課程的主要內容
- Hadoop on Amazon Ec2/S3
- More…:

- A9.com
- ADSDAQ by Contextweb
- EHarmony
- Facebook
- Fox Interactive Media

- IBM
- ImageShack
- ISI
- Joost
- Last.fm

- Powerset
- The New York Times
- Rackspace
- Veoh
- Metaweb

Hadoop於yahoo的運作資訊

年份	日期	節點數	耗時 (小時)
2006	四月	188	47.9
2006	五月	500	42
2006	十一月	20	1.8
2006	十一月	100	3.3
2006	十一月	500	5.2
2006	十一月	900	7.8
2007	七月	20	1.2
2007	七月	100	1.3
2007	七月	500	2
2007	七月	900	2.5

Sort benchmark, every nodes with terabytes data.

Hadoop於yahoo的部屬情形

資料標題：Yahoo! Launches World's Largest Hadoop
Production Application

資料日期：February 19, 2008

Number of links between pages in the index	roughly 1 trillion links
Size of output	over 300 TB, compressed!
Number of cores used to run single Map-Reduce job	over 10,000
Raw disk used in the production cluster	over 5 Petabytes

Hadoop於yahoo的部屬情形

資料標題：Scaling Hadoop to 4000 nodes at Yahoo!

資料日期：September 30, 2008

Total Nodes	4000
Total cores	30000
Data	16PB

	500-node cluster		4000-node cluster	
	write	read	write	read
number of files	990	990	14,000	14,000
file size (MB)	320	320	360	360
total MB processes	316,800	316,800	5,040,000	5,040,000
tasks per node	2	2	4	4
avg. throughput (MB/s)	5.8	18	40	66

瞭解
更多

Hadoop 與 google 的對應

Develop Group	Google	Apache
Sponsor	Google	Yahoo, Amazon
Algorithm Method	MapReduce	Hadoop
Resource	open document	open source
File System (MapReduce)	GFS	HDFS
Storage System (for structure data)	big-table	Hbase
Search Engine	Google	nutch
OS	Linux	Linux / GPL