

巨量資料基本概念

學習目標 >>>

- 瞭解何謂巨量資料、巨量資料的發展以及主要巨量資料種類
- ★ 瞭解巨量資料的問題與挑戰



3-1 巨量資料的基本觀念

根據 Google 前執行長 Eric Emerson Schmidt 的說法，人類在 2003 年之後每年產生的資料量，是人類歷史活動總合之前一年的資料。換句話說，每年產生出資料的數量是成指數成長的，如此一來，資料「放在哪」、「放得下」、「不會掉」，這三點非常重要。以每年產生 10 ZB (2.5 EB × 365 天) 的資料量，又以指數方式成長的速度來看，人類的資料量在 2020 年前會到達 YB 等級。

● 表 3-1 位元組的次方單位

位元組的次方單位		
名字	縮寫	次方
kilobyte	KB	10^3
megabyte	MB	10^6
gigabyte	GB	10^9
terabyte	TB	10^{12}
petabyte	PB	10^{15}
exabyte	EB	10^{18}
zettabyte	ZB	10^{21}
yottabyte	YB	10^{24}

以目前大家常用的硬碟為 1TB 容量，那麼 10ZB 就是 100 億顆這個容量的硬碟，而且是每年 100 億顆。這麼多硬碟、資料及備份，且隨時可存取，還要從這麼多資料中找出有意義的資訊，這些都是巨量資料中重要的課題。

3-1-1 什麼是巨量資料（Big Data）？

什麼是巨量資料？多大的資料量可以稱為巨量資料？於 20 世紀 80 年代早期，指的是資料量大到需要儲存在數千萬個磁帶中的資料；到了 90 年代，指的又是資料量超過單一桌上型電腦儲存能力的資料，不同的年代有著不同的答案。如今，巨量資料

指的是那些關聯式資料庫難以儲存、單機資料分析統計工具無法處理的資料，這些資料需要儲存在擁有數千萬台機器的大規模平行系統上。巨量資料出現在日常生活和科學研究的各個領域，資料的持續增長使人們不得不重新考慮資料的儲存和管理。

隨著網路的興起，漸漸的人們也開始習慣在網路上分享和交流資訊。舉例來說，社交網路 Facebook 擁有龐大的使用者群，且還在不斷增加中。這些使用者每天所發出的文章及對話記錄更是不計其數，其資料量已經達到 PB 級別，傳統的解決方案已經不可能去最佳化地處理這些資料。Facebook 自己開發了 Cassandra 系統，現在又採用 Hbase，這些針對巨量資料的管理系統能夠提供給使用者較好地服務，而且具有可擴充性和容錯性，這兩點是巨量資料問題所需要的效能。微博服務商 Twitter 也面臨巨量資料的挑戰，訊息的發送量達到每天數億條，而查詢量則達到每天數十億次，這要求儲存管理系統不僅能夠儲存大規模資料，而且能夠提供大量的讀寫服務。Twitter 原先使用 MySQL 資料庫，之後由於使用者暴增，便將資料移轉到 NoSQL 系統上，儘管 NoSQL 還未成熟，但卻是解決巨量資料較有效的方案。其他的網際網路公司同樣面臨著巨量資料帶來的問題，如 Google 搜尋引擎需要處理大規模的網頁資訊，YouTube 則需要儲存和提供使用者分享的視訊資料，維基百科傳送使用者分享的知識等，這些都有關大規模資料資訊儲存與管理。

隨著電子商務的發展，越來越多的人在網上選購商品，電子商務網站需要儲存大量的商品資訊和使用者交易資訊相關的大規模資料，同時網站需要提供迅速的請求反應，以加強使用者體驗來吸引客戶。而且網站還要對這些巨量資料進行處理和分析，以確實針對使用者之個人偏好推薦商品，使巨量資料成為系統建構和業務成敗的關鍵因素。中國商業網站淘寶使用 Hbase 來儲存資料，同時不斷探索自己的解決之路，開發了支援巨量資料的資料庫系統 OceanBase 來實現部分線上應用。全球最大的線上拍賣和購物網站 eBay 也積極尋求巨量資料的解決方案，它以 Hadoop 建立了自己為基礎的叢集系統 Athena 來處理大規模資料，同時開發了自己的開放原始碼雲端平台專案 Turmeric 來更進一步地開發和管理各種服務。同時，各大零售公司無論是線上銷售還是實體銷售，都會注意收集客戶的消費資訊，以便能更個人化地提供服務或推薦商品給客戶，這些都是有關大規模資料的應用。

3-1-2 巨量資料的定義及特性

巨量資料由巨型資料集組成，這些資料集大小常超出人類在可接受時間下的收集、使用、管理和處理能力，巨量資料必須藉由電腦對資料進行統計、比對、解析，方能得出客觀結果。美國在 2012 年就開始著手大資料，歐巴馬更在同年投入 2 億美金在大資料的開發中，更強調大資料會是之後的未來石油。

巨量資料，談的不僅僅是資料量（Volume），還包含了時效性（Velocity）、多樣性（Variety）及可疑性（Veracity）：

- ➁ Volume：資料的大量產生、處理及保存，談的就是巨量資料，也有人稱為海量資料。
- ➂ Velocity：時效性這個詞，有多種解釋，但我們認為用 IBM 的解釋來說是比較恰當的，也就是處理的時效，既然前頭提到巨量資料其中一個用途是做市場預測，那處理的時效如果太長就失去了預測的意義，所以處理的時效對巨量資料來說也是非常關鍵的，500 萬筆資料的深入分析，可能只能花 5 分鐘的時間。
- ➃ Variety：多樣性指的是資料的形態，包含文字、影音、網頁、串流等結構性及非結構性的資料。
- ➄ Veracity：可疑性指的是當資料的來源變得更多元時，這些資料本身的可靠度、品質是否足夠，若資料本身就是有問題的，那分析後的結果也不會是正確的。

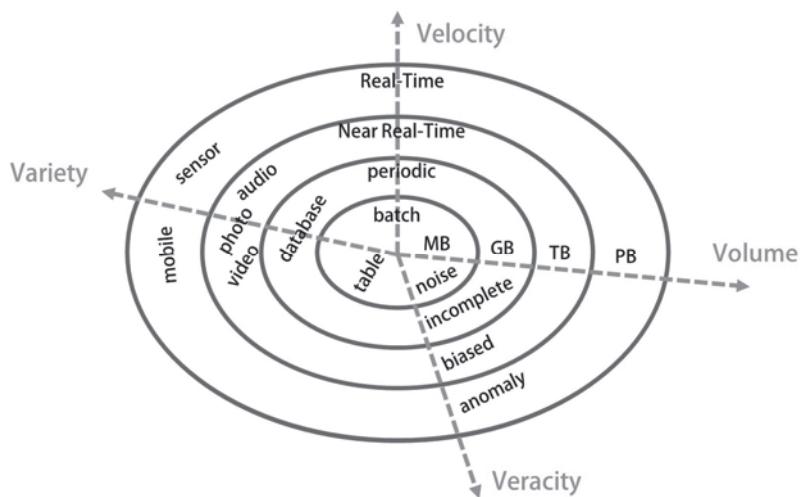


圖 3-1-2 巨量資料的 4Vs 定義

◎ 範例解析 ◎

Q：下列何者為巨量資料的定義？

- (A) 巨量資料除了指資料量的龐大之外，資料特性包含變化速度快及多樣性
- (B) 巨量資料中有 80% 都為結構化資料
- (C) 儲存的資料內容不包含影片及電子郵件
- (D) 巨量資料強調資料的數量，能夠為企業帶來商業機會

解答：(A)

解析：(B) 巨量資料具有很高的非結構化資料，此為巨量資料的多樣性（Variety）

- (C) 巨量資料含有多樣性特質，儲存的資料不再是單純的結構化資料，影片和電子郵件等非結構化資料也屬於巨量資料的內容中
- (D) 巨量資料的資料量（Volume）並沒有強調數量要多少才達標準，或是此數量能否為企業帶來商機

3-2 巨量資料的應用

各個領域的科學研究同樣面臨巨量資料的挑戰，從生物基因到天文氣象，從物理實驗到臨床醫學，致力於測量技術和裝置的發展，這些領域在實驗或實作上產生了大量的資料，而人們需要對這些資料進行處理分析，並進一步挖掘出有價值的資訊，但這不是十分容易的事情。

3-2-1 巨量資料的挑戰

隨著下一代基因定序技術的發展，基因中所蘊含的資訊逐漸被人們所發掘，使人們獲得更多更準確的基因資料，但是如何比對基因資料及如何從這些資料中挖掘出所需要的資訊，這是生物資訊學遇到的新挑戰。在環境氣象研究中，科學家已經收集了數十年甚至上百年的氣象環境資料，在這些資料中分析氣候的變化需要巨量資料處理技術的支援。而在醫學藥物研究中蒐集大量病人的生理資料和藥物測試的資料，這些資料的規模也很大，需要從中分析出有用的資訊。在人文社會科學中，社會學家開始注意網際網路之社交網路上的人際交往和社會關係，其有關的資料量也是非常龐大

的，從巨量資料中找出社會學家有興趣的內容是富有挑戰性的。人工智慧研究方面，人們希望電腦擁有人類的學習能力和邏輯推理能力，這需要機器儲存大量的經驗資料和知識資料，還需要從這些大量資料中迅速獲得所需要的內容，並進行分析處理，進一步做出正確且有效的判斷。

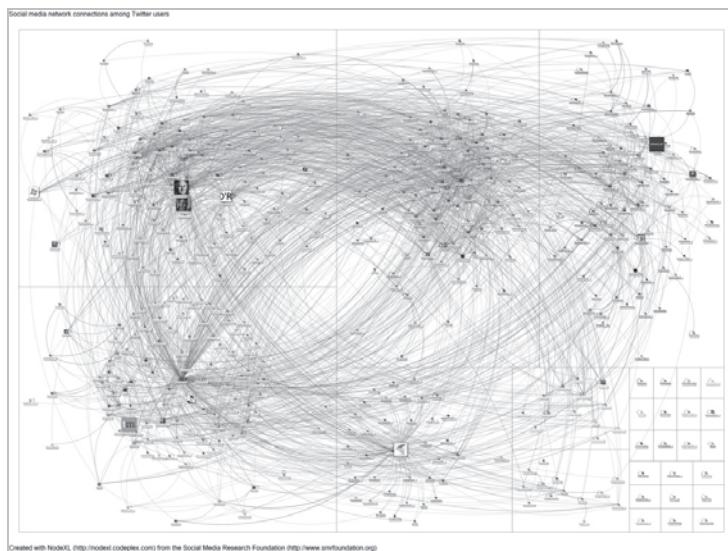


圖 3-2 透過巨量資料分析社交關係
(http://www.flickr.com/photos/marc_smith/6934127903/)

3-2-2 巨量資料的應用範圍

如今感測器的廣泛使用使資料獲取更加方便，這些感測器會連續地產生資料，如：即時監控系統、網路流量監測等。除了感測器源源不斷地產生資料外，許多領域也會產生類似的大量資料，如經濟金融領域中股票價格和交易資料、零售業中的交易資料、通訊領域中的資料等都是，這些資料最大的特點就是巨量，因為他們每時每刻連續不斷地產生，但與其他的巨量資料不同，此種資料連續有序，變化迅速，而且對處理分析的反應度要求較高，因此對於資料的處理和挖掘常常採用不同的方法。經濟金融領域各個方面都產生巨量資料，如證券價格變化和股票交易形成的資料、企業或個人各種經濟活動而產生的資料等。現代經濟已經步入巨量資料時代，在新時代下可以帶來創新和生產率增長，並可能出現新的商業模式。妥善利用經濟生活產生的巨量資料，可以發揮重要的經濟作用，不僅有利於企業的商業活動，也有利於國民經濟，加

強國家的競爭力。以大規模為目的的經濟資料，人們除了需要加強取得、儲存和分析資料的能力，同時需要保障資料的安全和隱私，但這仍然是極大的挑戰。

傳統的關聯式資料庫並不能夠極佳地解決巨量資料帶來的問題，單機的統計和視覺化工具也變得力不從心，一些新的資料管理系統和平行資料庫、網格資料庫、分散式資料庫、雲端平台、可擴充式資料庫等孕育而生，它們為解決巨量資料提供了多種選擇。

3-2-3 巨量資料的挖掘

資料挖掘（Data Mining）更是巨量資料應用中重要的一環，古人云：「物以類聚，人以群分」。這句話正是描述了資料挖掘的一種演算法—集群演算法（Clustering Algorithms）。要看一個人是怎樣的，只需要看他周圍都有什麼樣的朋友，而從資料挖掘的角度來說，用集群演算法預測一個物件的特徵，只需要看他周圍物件的特徵。

巨量資料的挖掘是指對於資料進行處理以及研究，並從資料中分析有用的資訊和發現知識的過程。資料分析和處理資料是我們常用的詞彙，那麼資料分析和資料挖掘有什麼樣的區別呢？

從本質上的角度來說，資料分析和資料挖掘都是為了從收集來的資料中找尋有用的資訊、發現知識並對資料加以詳細研究和概括歸納的過程。將純粹的資料賦與其意義，成為一項有意義的資訊（information）。在不少場景中，資料分析和資料挖掘這兩種概念是可以互換的，而它們之間最大的區別是資料本身的不同，這主要表現在以下兩個方面上：

1. 資料量的不同，資料分析通常是儲存在資料庫或檔案中，一個應用的資料量級別在 MB 或是 GB，而資料挖掘的應用資料，動輒 TB 甚至 PB（可參照表 3-1）。
2. 資料型態的不同，資料挖掘的物件不僅僅是文字，還有音訊、視訊和圖片資料，並且不僅是結構化之後的資料，還有半結構化和不規則的資料在其中，資料挖掘的對象都落在這些不規則的資料上。

舉例來說，兩者之間的差別就像是淘金者和礦山主，不同點在於淘金者只在一條小溪上工作，甚至幾十個人共用一條小溪，通常只能透過手動作業用沙漏過濾出沙裡的黃金。而礦山主則佔有整座極大的礦山，由於礦山擁有成分複雜的礦石和數量繁多且不同種類的礦物，這時礦山主就不能僅僅依靠手動作業，而需要建立一個以機器為工作的主力，來達到大幅的效率產出。

查詢、報表、連線應用分析以上這些屬於傳統的資料分析，本質上和資料挖掘的區別是在沒有明確假設的前提下挖掘資訊、發現知識。資料挖掘所得出的資訊通常具有先前未知性、有效性和可實用性三種特徵；而資料分析主要是一個假設檢驗的過程。就如同上述中的例子，資料挖掘如同在礦山中不知該從何下手找到礦物，沒有事先的明確方向，處處都充滿著未知。

資料挖掘建立在擁有大量資料，並且能夠讓機器方便讀取的資料儲存之上，大部分採用機器學習（Machine Learning）的演算法，是自動發掘知識的過程。然而這些不表示資料分析可以完全被取代。就像現在大工業只是取代了手動生產的生產組織形式，而手動生產中的方法、技能等都被現代大工業吸收進來，重新指定了新的意義。同樣的巨量資料的挖掘也需要資料分析的演算法和想法，只是用新的方法重新組織施行。資料挖掘並不是一種嶄新的學科，而是綜合了統計分析、機器學習、人工智慧、資料庫等諸多方面的研究成果的邊緣學科（由多門學科交叉產生出來的學科）。

◎ 範例解析 ◎

Q : 巨量資料的運用包含以下哪幾項服務？

- (A) 結合生理預測與事件預測，分析資料來做預防醫學
- (B) 分析資料以提供企業潛在商機智慧
- (C) 運用 CRM、上網、APP 等交叉分析，找出電信退租的共同特徵
- (D) 利用網站點擊次數統計，來規劃旅遊預算

解答： (A)、(B)、(C)

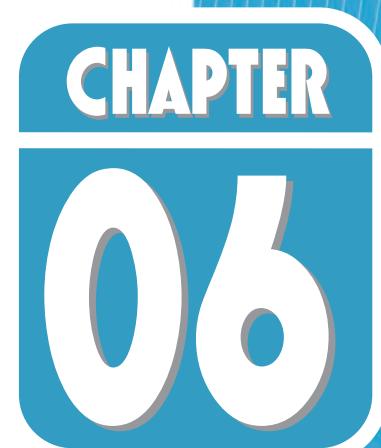
解析： (A) 單人的資料做生理預測及事件預測就擁有龐大的資料量，資料分析當然需要巨量資料技術的運用

(B) 經濟金融領域中股票價格和交易資料，也都具有巨量資料的特質，這樣的資料分析也須運用巨量資料的技術

雲端服務規劃與建置

學習目標 >>>

- ★ 瞭解雲端服務規劃
- ★ 瞭解雲端服務建置



6-1 雲端服務設計概念

6-1-1 公有雲及私有雲運用評估

美國國家標準與技術研究所對於雲端運算的佈署模式區分為四種：公有雲、私有雲、社群雲、混合雲。其中雲端運算最熱門的議題就是混合雲，顧名思義就是混合公有雲以及私有雲的一種佈署方式。在第一章中已經對於這四種佈署模式有過介紹，而這邊再以簡單的說法介紹這四種佈署模式，並對於企業規劃做運用評估。

對於企業來說，建置私有雲是使用虛擬化技術的方式來提升對於硬體資源的使用率，採用公有雲服務則可以降低設備維運與管理成本，上述中提到的都是它們各自的優點，透過這兩者不同的優點，視應用系統的特性彈性調配選擇公有雲或私有雲，創造最大效益，這就是混合雲的優勢所在。社群雲則是因為架設私有雲所需要的費用高，以及私密資料不希望放置在公有雲上，所以集合數個相互信任的組織，建立一個共同的私有雲。

混合雲的佈署方式通常有兩種，第一個是先建置私有雲，日後再逐漸將工作負載逐步移往公有雲，另一種則是相反過來，先使用公有雲，再將工作負載移動至私有雲上。不管是哪一種方式，都必須有一套客觀的評估標準來決定哪些系統或資料適合放上公有雲，這可以分為經營面和技術面來評估。

➤ 經營面

- **法規面**：主要是指企業所在國家或所屬產業是否有特殊規範，導致應用系統不適合移往公有雲，例如：某些資料必須放在境內…等。
- **敏感性資料**：企業上一定會有屬於商業機密性的資料，可能也不適合移往公有雲，雖然這和公有雲的安全性無關，而是一種信任以及情感上的原因，畢竟大家都比較信賴自身所能掌握的東西。
- **系統關鍵**：企業對系統的營運都會有系統停機的忍受度，多數提供公有雲服務的廠商都會有著 99.95% 可用性保證，這表示一個月停機時間不超過 21 分鐘，倘若能容忍的停機時間不超過 10 分鐘，自然就得放在私有雲由企業自行掌控。

 技術面

- ◆ **系統相依性**：往往企業所擁有的應用系統數量隨著業務量以及公司規模不斷增加，系統的生成時間和方式都不一樣，但彼此間都有相依性，也就是說當移動某一應用系統時，受影響的不是只有該應用系統，因此評估及顧慮將系統放上公有雲是非常重要的。
- ◆ **網路基礎設施**：應用系統在企業內部中所使用的是區域網路，不需要考慮到頻寬問題，但當它放上公有雲時就不一樣了，連線方式從內部區域網路到對外有線網路，網路基礎設施是否也要隨之調整，也是需要思考的問題。
- ◆ **資安政策**：當應用系統放上公有雲後，要如何確保符合既有的資安政策。

6-1-2 多租戶的設計概念

雲端運算的 SaaS 是基於 Web，通常會提供給成千上萬的用戶共同使用，同時具有良好的擴充性，且在運作過程當中，用戶之間彼此並不會互相干擾。對於用戶來說，不但不需要擔心軟體運行與硬體資源的維護問題，而且只需要負擔少量的租賃費用即可享受和以往相同或更好的服務；對於軟體廠商來說，每個用戶都共用一整套完整的軟硬體設施，在開發、維護與擴充上只需要專注於同一個實體，不須擔負額外的成本以提供不同的實體給不同的用戶所使用。多租戶技術就是軟體廠商目前得以滿足上述需求的關鍵技術，使大量用戶共同使用相同的軟硬體資源，每個用戶可以依需求進行客製化設定，且不影響其他用戶的使用。

多租戶的核心是同一個應用程式的實例（Instance），可以處理多個用戶的請求（Request），在多租戶的領域上，每個用戶被稱為租戶。多租戶平台如圖 6-1 所示：

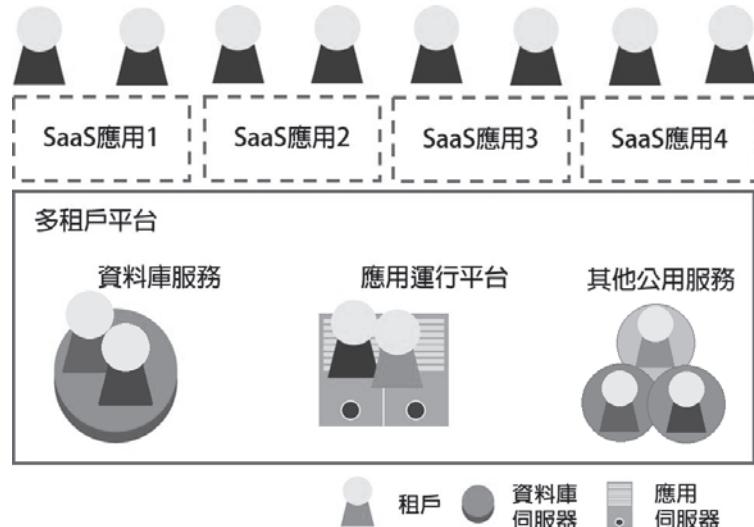


圖 6-1 SaaS 的多租戶概念

在多租戶應用程式的設計和實作上需要考慮三個關鍵：

1. **資源共享**：減少每個租戶的軟硬體與管理成本。
2. **資料隔離**：防止租戶之間非法存取與干擾。
3. **客製設定**：透過配置的方法支持所有與租戶相關的 UI (User Interface)、存取控制與資料模型 (Data Model)。

雖然由傳統的架構轉換為雲端運算多租戶的架構所能夠獲得的好處非常可觀，絕大多數的軟體廠商都知道這是必然的趨勢，但這個轉換過程所必須投入的人力、時間、金錢等現實成本，以及要使用哪種多租戶架構、該如何轉換、轉換後的品質、同一架構如何兼顧傳統佈署與雲端佈署等因素都是必須考慮清楚的。

6-1-3 水平擴充的設計概念

雲端運算中擴充 (Scale) 是很重要的考量。擴充的方式有很多種，這邊介紹垂直擴充 (Scale-up) 和水平擴充 (Scale-out) 兩種較為常見的擴充。

➤ 垂直擴充 (Scale-up)

垂直擴充稱作 Vertical scale。在儲存架構中垂直擴充的方法，是把現有的硬體“換掉”，換成容量更大或速度更快的硬體。例如：現在有的硬碟 500GB 容量不夠用，就把 500GB 的硬碟換成 1TB 的硬碟。

➤ 水平擴充 (Scale-out)

水平擴充稱作 Horizontal scale。水平擴充的方法則是保留原有的硬體，然後再接上一個硬體。例如現有的硬碟 500GB 不夠用時，再接上一個 500GB 的硬碟，這樣就有 1TB 的容量空間可以使用。

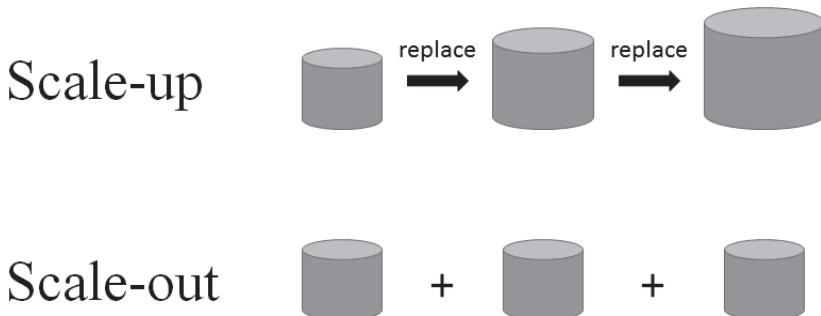


圖 6-2 儲存架構擴充方式

在雲端服務中，如果是採用垂直擴充資源的話，可以讓舊有的運算環境直接搬移至新環境上，並達到運算能力或使用資源的擴充，但這樣沒有有效的利用到每一個可用資源，所以大部分採用的是水平擴充這個方式，當儲存或是運算資源不足時透過添加來擴充，這樣才是個有延伸性的雲端服務。

6-1-4 負載平衡的設計概念

負載平衡 (Load Balance) 是一種網路技術，用於多個電腦間的資源分配（例如：網路流量、CPU 資源和儲存資源等），以達到優化資源使用、縮短回應時間，更重要的是避免過載的發生。

雲端運算環境為了提供大量的使用者存取雲端服務且不使伺服器發生過載的情形，經常需要運用大量的雲端伺服器，但服務並非隨時都被大量的存取，在低使用率

的期間，若同樣保持大量的雲端伺服器開機，就會造成不必要的浪費。如上所述，雲端服務中不僅在使用量高的時候要能夠平衡，在使用量低的時候也要降低資源上的使用。

這樣的需求目前常見的解決辦法是動態資源分配，透過線上遷移 (live migration) 的動作，在節點負載過高時，將虛擬機線上遷移至其他負載量較低的節點上；或者是負載量較低時，將零星的虛擬機集中在同一個節點上降低資源的浪費。

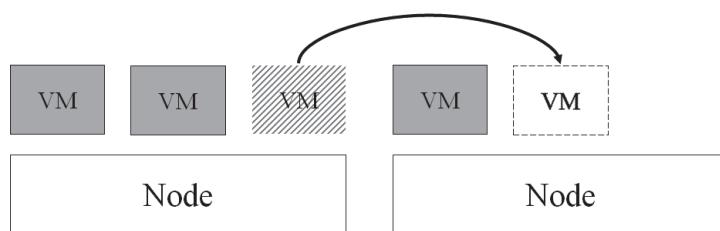


圖 6-3 動態遷移技術

6-1-5 容錯的設計概念

什麼是容錯 (fault-tolerant)？降低風險發生的可能性，這是容錯最主要的目的，這並不是一個新穎的名詞，早在自動化系統中就有這樣的需求存在，爲了能夠讓自動化系統在某個元件損壞時不造成整個系統停擺。當然在雲端服務中這樣的需求也是必要的，雲端服務中所強調的高可用性 (High Availability) 就是容錯重要的設計考量。

在雲端服務中軟硬體達到容許錯誤發生時，仍然可以正常持續運作，實現了高可用性。一般這些錯誤可能發生在節點 (Nodes)、應用程式 (Application)、網路 (Networks)、硬碟 (Disks) 四種狀況上。

上述的四種狀況上，設計概念大致都是一樣的，爲了達到盡可能不中斷的服務，都會有預留 Standby 的資源，以節點故障爲例：兩個節點分別爲 A 和 B，A 是正在運行服務的節點，而 B 是處於 Standby 狀態的節點，當容錯機制偵測到 A 節點故障時，會立刻將處於 Standby 的 B 切換爲運行並接替 A 的位置，此時 A 的故障會告知系統管理員，讓系統管理員重新將 A 維修後啓用，並使它設定在 Standby 的狀態。

企業在採用雲端服務時都會有服務水準協議（Service Level Agreement）的需求，在雲端服務中為了有一個明確的規範，將可用性化作一個百分率，百分率即是一年中的正常運行時間，下表是幾個常見的比率：

● 表 6-1 不同的可用性對照以年（365 天）與以月（30 天）的停機時間

可用性%	每年停機時間	每月停機時間
90% ("one nine")	36.5 天	72 小時
99% ("two nines")	3.65 天	7.20 小時
99.9% ("three nines")	8.76 小時	43.8 分鐘
99.95% ("three nines five")	4.38 小時	21.56 分鐘
99.99% ("four nines")	52.56 分鐘	4.32 分鐘

◎ 範例解析 ◎

Q：下列關於 Scale Out（水平擴充）敘述，哪些正確？

- (A) 靠增加處理器來提升運算能力和增加獨立伺服器來提升運算能力
- (B) 增強處理器等運算資源進行升級以獲得對應用性能的要求
- (C) 依靠多部伺服器、儲存協同平行運算，藉負載平衡及容錯等功能來提高運算能力及可靠度
- (D) 目前在開源軟體方面，也有許多水平擴充的系統，比如集群文件系統、分散式文件系統

解答：(A)(D)

解析：(A) 以增加的方式擴充運算能力及是 Scale Out（水平擴充）

- (B) 所提到的是透過升級來擴充效能，這樣的擴充方是屬於 Scale Up（垂直擴充）
- (C) 採用多個運算資源協同的平行運算，這是增加運算資源的方式，而非增強運算資源
- (D) 集群文件系統、分散式文件系統，這些系統都可以透過額外增加儲存節點加強系統的儲存資源