

# 以現代化資料平台 加快深度學習

TABOR CUSTOM PUBLISHING 監製·合作單位：

The logo for 'datanami' features the word 'datanami' in a bold, lowercase, sans-serif font. Above the 'a' is a stylized icon of a hand with fingers spread, resembling a sunburst or a data visualization element.

• BIG DATA • BIG ANALYTICS • BIG INSIGHTS •

贊助單位：



**PURESTORAGE®**

# 人工智慧之社會影響

我們的生活周遭充斥著數不清的感應器，例如相機、智慧型手機與汽車等，這些感應器和我們的企業、教育系統等各種組織皆時時刻刻製造著數量龐大的資料。大數據的時代，我們能夠透過人工智慧（AI）、機器學習和深度學習技術，來挖掘出蘊藏在這些龐大資料中前所未見的深度洞見。

在華盛頓特區舉行的 2017 年美國網際網路協會年度盛會中，Amazon 執行長 Jeff Bezos 曾表示：「我們正進入一個文藝復興般的黃金年代。」 Bezos 表示：「我們現在能透過機器學習與 AI 解決許多問題，這在過去幾十年都還是只是科幻小說裡的內容。例如自然語言的理解，還有機器視覺的問題等，我們所處的真的是一個充滿驚奇的復興年代。機器學習與 AI 將大幅擴展未來的視野。它能夠推動改善，並使任何企業、政府組織、或慈善事業具備更強大的能力，基本上世界上沒有任何機構是無法從中受益的。」

Amazon、Apple、百度、Facebook、Google（Alphabet）、Microsoft 和 NVIDIA 等科技公司都有專門團隊負責 AI 的開發計畫，研究領域包含圖像識別、自然語言理解、視覺搜索、機器人技術、自動駕駛汽車、以及文字轉語音技術等。他們創新的人工智慧、機器學習與深度學習計畫包括以下幾個例子：

- ▶ **Amazon：** [Amazon 使用 AI 和複雜的學習演算法](#) 持續評估市場動態，以判定合適的推薦商品和選入黃金購物車（Buy Box）的商品。
- ▶ **Apple：** iPhone 與其他 Apple 產品上的 [Siri 虛擬助手](#) 以深度學習技術輔助搜尋，並透過互動式的語音介面提供相關解答。
- ▶ **百度：** 由百度開發的語音識別系統 [Deep Speech 2](#) 能夠輕易辨識英語及華語語音。在某些情況下，其轉譯的準確度甚至優於人類。

- ▶ **Facebook：** Facebook 的 [DeepMask 和 SharpMask](#) 軟體與 MultiPathNet 神經網絡共同運作，讓 Facebook 可以透過分析個別像素理解圖像內容。

- ▶ **Google（Alphabet）：** Google 執行長 Sundar Pichai 指出，用戶點擊 Google 地圖的畫面時，[Google 的街景服務](#) 會透過 AI 自動辨識路牌或店家招牌，以協助定義該地資訊。

- ▶ **Microsoft：** [Microsoft 的 AI](#) 能夠在 PowerPoint 中透過認知視覺系統分析照片並自動生成替代文字，或提供建議的圖表協助內容呈現。

- ▶ **NVIDIA：** [NVIDIA DRIVE™ PX](#) 是一個開放式的 AI 車輛專用運算平台，有利於汽車製造商和一級供應商加速自駕車的生產。

## 人工智慧的重要性

Forrester 在研究報告「人工智慧：2017 年對企業的可能性」（Artificial Intelligence: What's Possible for Enterprises in 2017）中指出，AI 已漸漸地成為現實。越來越多的公司組織、研究人員與教育機構已開始關注其潛能。該報告發現：「在 391 名受訪的商業和科技業專業人員中，只有 12% 已開始使用 AI 系統。然而，有 58% 的人已著手研究 AI 技術和企業的運用條件，另有 39% 已開始尋找並設計欲部署的 AI 功能。除此之外，這份 2016 年 11 月發佈的報告也發現，有 36% 的受訪者正著手進行相關的企業教育訓練，或建構實證 AI 潛力的企業案例。」

# 人工智慧的大爆炸

人工智慧的出現始於三個關鍵技術的完美結合，亦稱為為人工智慧的大爆炸。這三個關鍵的驅動因素分別為深度學習演算法、以圖形處理器 (GPU) 為基礎的並行處理器、以及大數據的普及。

## 深度學習— 可自行撰寫軟體的新型態運算模型

傳統的電腦程式在設計上，是經由特定的指令代碼來循序運算資料。深度學習技術則讓電腦系統能夠藉由分析資料，來提供相關的深度見解及預測結果。機器學習指的是任何不經人工編寫即可自我學習的程式。深度學習（亦稱深度結構學習或分層學習）是使用人工神經網絡的一種機器學習類型。深度學習系統可以採用完全監督、半監督或無監督的形式進行。

## AI 實際運用： 衡量品牌影響力

SAP 的 Brand Impact 軟體是一個將深度學習用於分析品牌曝光度的例子。該 SAP 軟體可以在處理了數千筆的圖片及影片資料後，訓練出自圖片中自行辨識標誌與其

他品牌資訊的能力，無須由人工編寫明確的程式碼。許多品牌仰賴贊助轉播活動，業者一般會在活動後花費長達六週的人工流程製作報告，以評估在活動中展示商標對品牌影響力的投資報酬。調整品牌行銷的支出可能須要花一整個季度的時間才能完成。

SAP 的 Brand Impact 軟體，使用在 NVIDIA DGX-1 系統上訓練的深度神經網路。[SAP 的深度學習分析](#) 能夠立即並精準地辨識影片中的標誌。透過該 SAP 軟體，只要一天的時間就能產出可供審計的結果。圖 1 展示影片分析的案例。

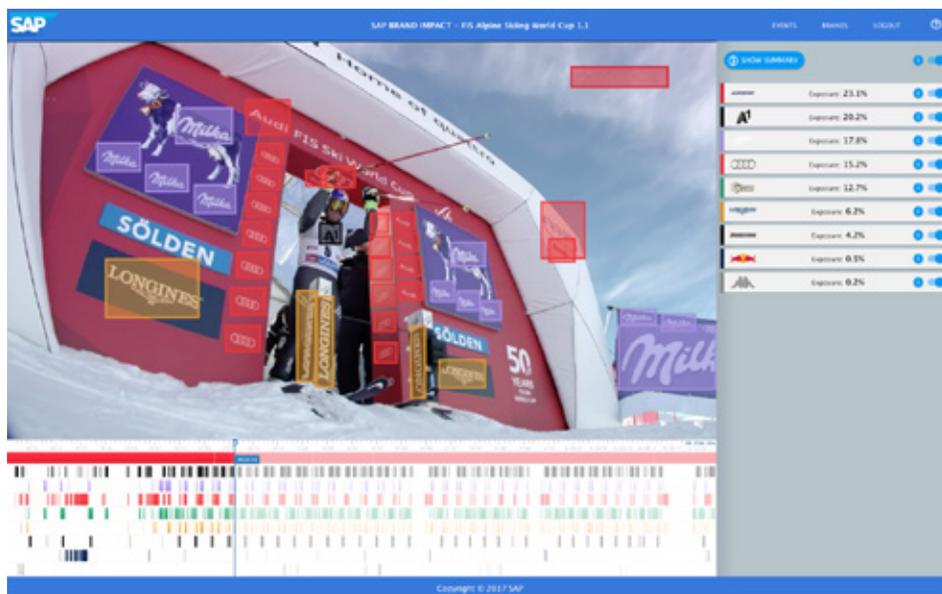


圖1. SAP Brand Impact - 以接近即時的速度辨識品牌標誌的位置。(來源：SAP)

## GPU：現代化的平行處理器

現代的運算工作，一般以多核 CPU 或 GPU 進行。如圖 2 所示，多達二十核心的 CPU 或上千個核心的 GPU 在今日並不少見。現代的 CPU 和 GPU 都屬於並行處理器，能夠同時運算一項以上的作業。

1997 年，NVIDIA 率先推出 GPU 加速運算技術，這是一種新型態的運算模式，能夠加快大規模並行的工作負載。NVIDIA 近期推出的 Tesla® V100 採用最新的 NVIDIA Volta™ 架構，只需要一顆 GPU 就能夠提供等同於 100 顆 CPU 的效能。

在今日，多核 CPU 和 GPU 應用於加快深度學習、分析工作及工程應用程式，令資料科學家、研究人員和工程師終於能夠應對以往無法克服的挑戰。新的深度學習演算法運用了由人腦作為啟發的大規模並行神經網絡。深度學習的模型能夠經由大量的學習案例編寫自身的軟體，不須由工程師手動編寫。這在面對圖片、影片和文字處理等一般工作時，更能夠提供超越人類的準確度。

圖 2 說明在並行運算與新型的深度學習大規模平行演算法結合之下，如何在圖像辨識方面提供了超越人類的準確率。

## 大數據

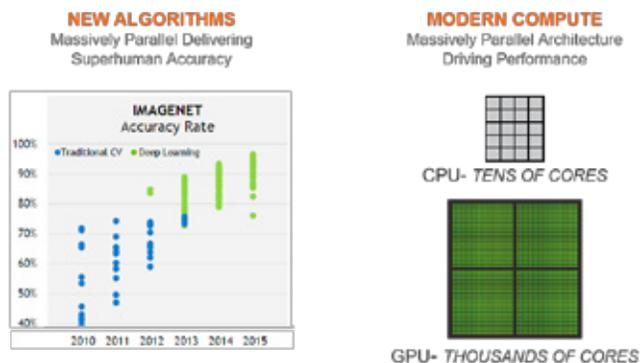


圖 2. 並行運算、新型演算法與大數據三方結合產生的智能大爆炸由 Pure Storage 提供。

資料是一個組織中最重要的資產。事實上，經濟學人雜誌在 2017 年 5 月的報導中聲稱資料已比石油更有價值。資料至今仍在不斷成長。IDC 的「2020 年數位宇宙 (The Digital Universe in 2020)」報告中寫道：「所謂的數位宇宙指的是一切新創資料的集合體，包含串流影片與其他儲存的數位資料。這些資料正以每兩年增加一倍的驚人速度不斷成長，2013 年的總資料量為 4.4 ZB，這數字預估到了 2020 年將超過 50 ZB。」

深度學習技術和神經網絡已問世很長一段時間了。那為什麼深度學習技術一直到現在才嶄露頭角呢？大數據又具有甚麼樣的價值？人工智慧領域的頂尖人物吳恩達先生在 2016 年的 Spark 研討會 中，描述大數據和深度學習的發展。吳恩達指出，如果採用邏輯迴歸等較傳統的學習演算法，並餵入大量資料，系統效能會遇到瓶頸，因為傳統的演算法並無法從更多的資料中擷取更深度的見解。吳恩達表示深度神經網路則是不同的。如圖 3 所示，餵入神經網路的訓練資料越多，就會變得越精準。多虧了創新的演算法、GPU 運算系統的效能大躍進、以及大數據的不斷成長，深度學習技術的普及率才得以快速發展。

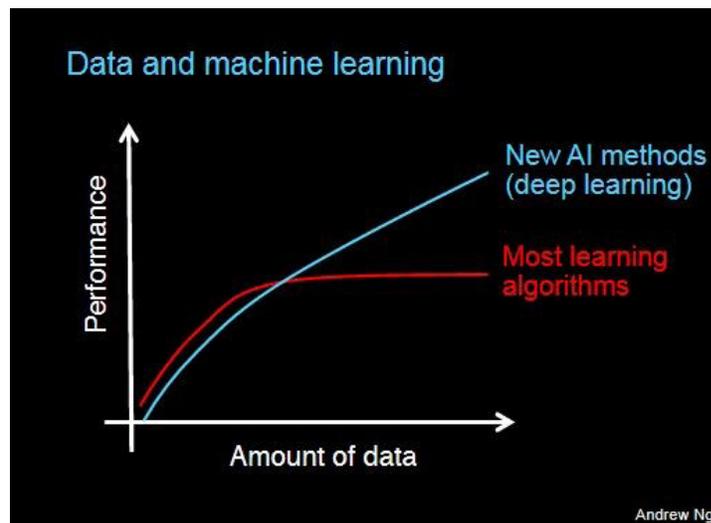


圖 3. 深度學習的效能會隨著數據量一同增長。(來源：吳恩達)

# 傳統的儲存裝置為何不能滿足深度學習的需求

儘管並行運算和演算法都有了重大的進展，但負責儲存和傳送大數據的科技，卻大多仍以序列時代所設計的老舊技術為基礎。新的運算模式勢必要搭配新型態的儲存系統才能滿足大量的資料需求。

光是在過去數年內，深度學習所需的運算量和 GPU 能提供的運算量就躍升了高達十倍以上。但傳統硬碟和固態硬碟的效能同樣的時間內卻無所成長。隨著非結構化資料的數量爆發性地成長，傳統的儲存裝置已難以滿足各種新興大數據運算程式的需求。

現今大多數的伺服器採用直連儲存架構 (DAS) 或分散式直連儲存架構 (DDAS)，資料集以分散的方式儲存於伺服器中的多個硬碟內。DDAS 讓資料科學家得以使用現成的系統或零組件產品來進行分析作業，例如 X86 架構處理器和標準的硬碟。然而，這樣的作法隱藏著許多問題。在現代化的資料分析技術開發之初，任何一個儲存平台的容量和速度都不足以應付如此龐大的資料量以及大數據軟體對大量頻寬的需求。

現代化分析的目標是透過分析從資料中擷取深度見解。這些資訊往往是來自於軟體日誌與物聯網 (IoT) 裝置，並以非結構化的資料形式存在。較舊的系統只能處理高度正規化的資料，無法有效分析人工智慧和深度學習工作常用的半結構化或非結構化資料。因此，傳統的儲存裝置已經成為了應用程式的主要瓶頸。數十年舊的序列技術無法有效處理非結構化的資料，效能也因而大受限制。如果資料是第四次產業革命中新貨幣，為什麼儲存產業的技術還停留在序列的時代呢？

## 運用深度學習創造實際的商業利益

公司機構在投資了人工智慧或深度學習技術後，體驗到了這些好處：

**Forrester 報告結果**：25% 的受訪者表示他們使商業流程的效率獲得提升，21% 表示客戶滿意度獲得了改善，18% 則表示節省了成本。

**醫療機構**：部署了靈活的預測分析功能後發現，由於慢性疾病與其他照護管理問題，10% 的員工共耗費了 70% 的資源。

**資料中心**：DeepMind AI 將資料中心的散熱成本降低了 40%。

**降低運算晶片的測試成本**：一間主要的晶片製造商運用預測式分析輔助晶片測試，並省下了共三百萬美元的製造成本。

**環法自行車賽採用機器學習**：2017 年的環法自行車賽使用機器學習進行預測，並將歷史資料與 2017 年的比賽內容合併，進而提供比以往更詳盡的比賽資訊。

**智慧電網**：麥肯錫全球研究院的研究發現，AI 能夠透過感應器和機器學習技術每分鐘持續調整電網，使電網更加智慧化，以實現最大的發電效率。

**提高工作績效**：工廠使用 UpSkill 的 Skylight 平台與擴增實境 (AR) 眼鏡，以更高的效率和更低的錯誤率提供操作員工作中所需的資訊。員工的平均績效最高提升了 32%。

# 專為深入學習設計的 FLASHBLADE

在大數據的新時代，為了滿足資料分析、科學發明，或是電腦繪圖等各種用途，應用程式運用了功能強大的大型伺服器農場以及極高速的網路來存取數以 PB 計的大量資料。這些新的應用程式需要更快更有效的儲存裝置，傳統的解決方案早已不堪負荷。

現代需要的是一個全新的創新儲存架構，以呼應先進的應用程式，並同時在每秒讀寫次數 (IOPS)、傳輸量、延遲與容量等各個面向提供業界最佳的效能與突破性的儲存密度。PureStorage® 全新的 FlashBlade™ 快閃記憶體儲存裝置能夠滿足以上所有需求。FlashBlade 能夠有效應付大數據與並行工作，可助您帶動未來的創新、探索與真知灼見。

## PURE STORAGE

有鑑於 Pure Storage 在全快閃記憶體資料儲存裝置的創新發展，Gartner 已連續四年在旗下的[固態陣列魔力象限 \(Magic Quadrant for Solid State Arrays\)](#) 報告中將 Pure Storage 評為業界的領導者。自從 Pure Storage 首次推出橫向擴充的儲存平台 FlashBlade 以來，該公司針對即時大數據分析、財務分析與生產工作所提供的儲存裝置已佔據了大幅市場。

FlashBlade 的架構是從頭至尾專為現代化的分析工作所設計，可提供高效能、符合成本效益、並且能夠輕鬆持有、簡單使用的橫向擴充儲存裝置，足以應付 PB 規模的營運資料。**FlashBlade 是專為快閃記憶體以及深度學習所需的大規模並行工作所設計，架構中已不含任何供機械式硬碟使用的空間。**

FlashBlade 的關鍵特性是能夠提供具擴充性的彈性效能，亦即應客戶當下的需求提升效能、容量、連結性與功能的能力。之所以辦得到，是因為 FlashBlade 不論是軟體或是硬體都從根本上具備了大規模並行特性。憑藉其獨特的 [Evergreen™ 商業模式](#)，客戶永遠不須要重新購買任何已擁有的儲存空間，並且可以隨著技術發展而升級，過程中對運行的服務、系統效能、以及資料的完整性都不會有影響。圖 4 為 FlashBlade 機箱範例。



圖4. FlashBlade 機箱：1.6 PB 4U 規格 由 Pure Storage 提供。

# 提供 AI 所需的資料傳輸量

深度學習系統通常使用小型的檔案讓訓練機保持繁忙。圖 5 中的範例透過 NVIDIA DGX-1 伺服器 and FlashBlade 資料儲存平台運行深度學習的訓練工作。在這個例子中，每一套 DGX-1 皆透過 AlexNet 模型與微軟的 CNTK 框架處理每秒 13,000 張的圖像。

該訓練模型必須隨機存取小型文件，傳統的儲存系統並無法有效處理這樣的任務。在這個例子中，FlashBlade 提供的資料傳輸量足以讓多個 DGX-1 系統發揮最大的訓練效能。

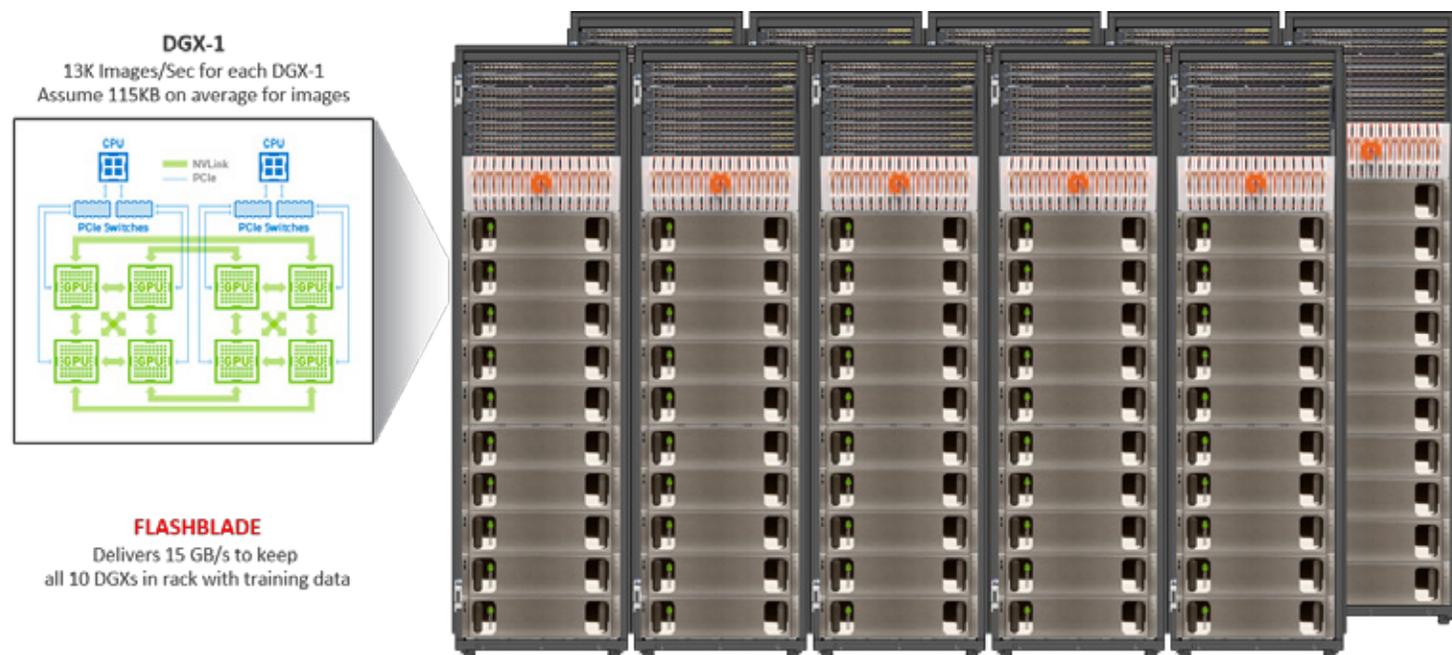


圖5. FlashBlade 如何提供 AI 所需的資料傳輸量。由 Pure Storage 提供。

# 總結

資料正以驚人的速度不斷增長，而這樣的快速增長將持續進行下去。AI、機器學習、以及深度學習等資料處理與分析的新技術讓經過特殊設計的應用程式不但能分析資料，還能從分析過程中學習，並提供近一步的預測資訊。

運算這些資料須仰賴多核心 CPU 與 GPU 的平行運算能力，以及極度快速的神經網路。然而，傳統的儲存解決方案是以數十年舊，且不具擴充性的架構為基礎，無法提供機器學習所需的大規模並行性。傳統的儲存裝置正逐漸成為大數據運算工作中的瓶頸，業界勢必需要一個全新的儲存技術才能滿足資料分析工作的效能需求。

Pure Storage 的 FlashBlade 全快閃記憶體儲存陣列即是設計來滿足這些需求。FlashBlade 的效能可隨著資料量提升線性成長。無論檔案大小，FlashBlade 都能提供真正線性的容量與效能擴充能力，因此非常適合 AI 和深度學習技術的現代化分析工作。

Pure Storage 的工程 VP Par Botes 表示：「現代化的運算框架造就了複雜度越來越高，且效能越來越強的分析作業與珍貴資料」「有了 FlashBlade，我們的使命是透過龐大、快速、且部署簡易的全快閃記憶體儲存平台將大數據轉變為快數據，並且為所有的產業與市場提供價值。」

## 關於 PURE STORAGE

Pure Storage <sup>iii</sup> (NYSE:PSTG) 協助企業推動新的可能性。Pure 的端對端資料平台 (包括 FlashArray、FlashBlade、以及與思科合作開發的融合式方案 FlashStack) 皆是由創新的軟體提供驅動。透過雲端連結，用戶可以行動裝置隨時隨地進行管理，此外亦支援 Pure 的 Evergreen 商業模式。Pure 的全快閃記憶體技術與便於

顧客的商業模式結合，其簡單、高效、且長盛不衰的產品方案有助於企業推動商業與 IT 轉型。Pure 的客戶包含各種規模，並來自於各式各樣不斷豐富化的產業。Satmetrix 認證 NPS 83.7 的高分證明了 Pure 的客戶滿意度為業界首屈一指。

i NVIDIA 和 SAP 合作創造新一波的 AI 商業應用程式 <https://blogs.nvidia.com/blog/2017/05/10/nvidia-sap-partner/>。

ii AI：新世代的電力，吳恩達，2016 Spark 高峰會，<https://www.youtube.com/watch?v=4eJhcxYR4I>。

iii Pure Storage、FlashBlade 和「P」標誌是 Pure Storage 在美國和其他國家的商標或註冊商標。其他標誌為其所屬公司之商標。

