



解決方案簡介

解決大數據分析的難題

Intel 和 SAP 的解決方案如何共同運作，提供動態分層的資料儲存與分析功能

數十年間，企業都仰賴著儲存於關聯式資料庫的結構化資料，以帶動商業決策。在與日俱增的處理與儲存能力之外，全新的資料收集方法越來越多，因而開啟了收集與分析大規模資料集的新管道。企業可運用比以前更大規模也更快速的方法收集資料，並使用資料進行關鍵的即時決策。

無論如何，在老舊的分析方法中，資料儲存在傳統硬碟式關聯式資料庫內，這會拖慢即時敏捷度。隨著資料量增加，傳統的分析工具會拖慢資料庫速度，可能導致仰賴資料庫的應用程式回應速度降低。企業需要可以擴展的解決方案，並針對大小超過數百 TB 甚至 PB 的資料集提供報告和近乎即時的分析結果。

為了解決這些挑戰，Intel 與 SAP 攜手合作，打造創新的概念性驗證系統，以儲存與分析大規模、快速的資料集。此解決方案簡介說明多層架構方法的設計理念與優勢、系統的軟硬體元件，以及效能測試的結果。

大數據分析的挑戰

現代的企業可以從包括社群媒體網站到製造工廠感應器等數千、數百萬種來源擷取大量的即時資料。如此廣泛的收集來源可產生各種不同的資料類型，包括相片、影片，以及文字等。傳統硬碟式關聯式資料庫管理系統 (RDBMS) 需要嚴格的配置，可能不是最適合收集或分析未結構化資料的解決方案。在大型資料庫表格上進行複雜即時查詢的速度通常十分緩慢，導致使用者不滿，並影響到決策程序。此外，低落的應用程式效能可能會影響使用者處理執行報告或收集數據等日常任務，最終可能會影響到資料安全性。隨著企業將更多資料推送到傳統資料庫中，受挫的用戶可能會轉向採用資料擷取方法，將資料由受保護的主要資料庫抽取至小型的本機資料庫中，這可能會導致安全性與隱私權受到侵害。

企業所面對的另一個挑戰是將新舊資料儲存在相同的資料庫中。資料的存取頻率會隨著時間減少。廠房管理員可能會想要檢視即時

的廠房生產數據，稽查人員則可能需要檢視每季的廠房生產數據。在這兩個情況中，新舊資料始終都互相結合，資料庫引擎必須在資料庫內查詢不斷增加的資料。這可能導致效能低落，限制重要的即時分析功能。

Intel 與 SAP 的目標是打造一個分析系統，解決擷取、儲存，以及組織大型未結構化資料集的效能挑戰，同時維持查詢效能。

設計理念：動態資料分層

為了滿足平衡效能與速度、數量，以及多樣性的挑戰，Intel 與 SAP 將多溫度動態資料分層概念作為系統設計的基礎。此設計理念是將經常存取的資料放置於快速的隨機存取記憶體 (RAM) 中，並將存取頻率較低的資料移動到較慢、成本較低的儲存空間中。

多溫度資料分為三種類別：

- 「熱資料」是組織經常存取的資料，舉例來說，這可能包括用於即時預測與回報，或是即時查詢的資料。熱資料通常必須儲存在快速的 RAM 中，而不是以硬碟為基礎，甚至以固態硬碟 (SSD) 為基礎的儲存空間內，以提供即時的成果。
- 「溫資料」不會經常存取，但還是經常使用，此資料仍然需要儲存在以硬碟或 Intel® SSD DC P3700 系列等 SSD 為基礎的儲存空間中，搭配快速的分析引擎。舉例來說，客服中心可能需要執行每月的來電狀態報告，但這並不需要即時的結果。資料在封存前，通常都保留在此狀態中。
- 「冷資料」僅會偶爾存取，並已經考慮將其封存。例如，製造廠房可能不會經常需要擷取數年前的資料，但他們可能會想要保存這樣的資料，以展示由改善程序帶來的每年效能效益。

針對較小型的資料集，三種溫度的資料皆可放置於單一系統或小型系統叢集中。但持續增加的成长率與使用率可能需要更為複雜的分層方法，運用多種層級的運算、記憶體，以及儲存設備。

專案目標

專案目標包括每小時擷取數億列資料、支援跨數百萬列資料進行即時查詢、提供探索式與預測式分析工具，以及動態管理資料在跨多個系統層級的移動過程。

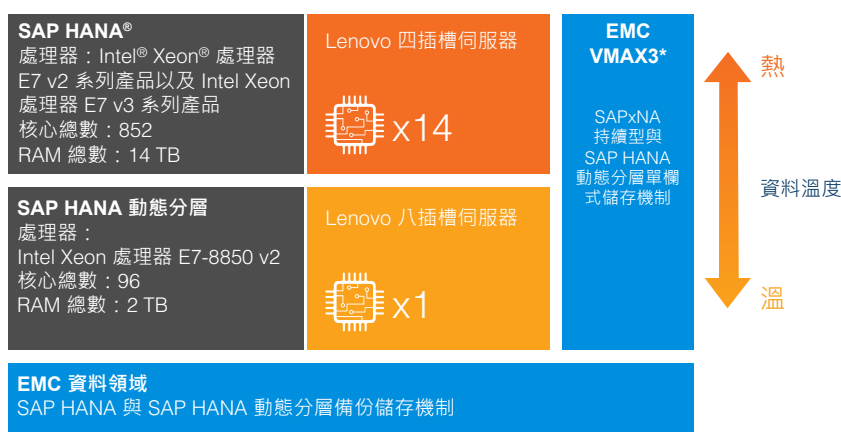


圖 1. Intel 與 SAP 分析解決方案架構概覽，資料根據預先定義的到期標準流經熱資料與溫資料狀態

雖然 Intel 與 SAP 著眼於將多溫度設計理念作為他們系統架構的基礎，但僅實作並測試了「熱資料」與「溫資料」層級。

架構概覽

Intel 與 SAP 設計並打造了多層的資料分析與儲存系統，可處理全球數以千計的使用者，以及數百 TB 的資料。為了驗證概念而設定的特定效能目標包括：

- 每小時由數千個資料收集裝置擷取數億列資料，並使新資料在 15 分鐘內備妥以供查詢
- 支援跨全球成千上萬名使用者的數億列資料進行即時查詢與報告，且不會對效能有明顯的影響
- 提供整合式探索與預測式分析工具，為過去的活動與未來的趨勢提供更深入的瞭解
- 透過溫度，動態管理資料在跨多個儲存設備與運算系統層級的移動過程

此系統結合兩個使用資料溫度模型定義的資料管理與分析層級，包括：

1. 「熱資料」階層建立於 SAP HANA® 與搭載 Intel® Xeon® 處理器 E7 v2 系列產品與 Intel Xeon 處理器 E7 v3 系列產品的伺服器上¹。此階層使用 EMC VMAX3*，提供持續型儲存。

2. 「溫資料」階層建立於 SAP HANA 動態分層與搭載 Intel Xeon 處理器 E7 v2 系列產品的伺服器上。此階層使用相同的 EMC VMAX3 儲存系統，以進行長期資料儲存。

熱資料管理：SAP HANA®

為了達到每小時擷取數億列資料的目標，Intel 與 SAP 在搭載 Intel Xeon 處理器 E7 v2 系列產品與 Intel Xeon 處理器 E7 v3 系列產品的伺服器叢集上導入 SAP HANA。SAP HANA 是一種記憶體內資料庫平台，與傳統的硬碟式關聯式資料庫不同。所有資料都保留在系統記憶體中，提供比硬碟式資料庫高出許多層次的存取速度。SAP HANA 的記憶體內技術可為大型、高速的資料集提供即時分析能力，同時 SAP HANA 叢集組態有助於達成擷取並即時分析大量資料的目標。

溫資料管理：SAP HANA 動態分層

熱資料用於即時分析與高速資料，溫資料層級需要在儲存於成本較低的儲存空間中，相對更大的資料集上進行預測式分析。Intel 與 SAP 導入了 SAP HANA 的動態分層選項，運行於搭載 Intel Xeon 處理器 E7-8850 v2 系列產品的伺服器上，以達成這些目標。

使用 SAP HANA 的動態分層選項，可以在達到預設的過期門檻時，將資料移出系統記憶體，並移到成本較低的硬碟或 SSD 儲存設備中。SAP HANA 動態分層運用 SAP HANA Data Warehouse Foundation 的資料生命週期管理員功能，將過期的資料移動到由 SAP® IQ 管理的近線儲存中。

SAP HANA 動態分層可為跨結構化、半結構化，以及未結構化的資料集，提供大規模的進階分析功能。這為溫資料提供快速大量載入功能，以及更加符合成本效益的近線儲存。因此，SAP HANA 動態分層可為溫資料分析提供強大且符合成本效益的中繼儲存與分析平台。

表 1. 熱資料層級伺服器配置

熱資料層級伺服器配置	
伺服器	14 台四插槽 Lenovo 伺服器搭載 Intel® Xeon® 處理器 E7 系列產品 ²
核心總數	852
總記憶體	14 TB

表 2. 熱資料層級網路配置

熱資料層級網路配置	
內部網路	10 GB 乙太網路 (GbE) 連線，使用 Lenovo RackSwitch G8272*
用戶端網路	10 GbE 連線，使用 Lenovo RackSwitch G8272
管理網路	1 GbE 連線，使用 Lenovo RackSwitch G8052
儲存設備網路	8 GB 光纖通道連線，連接 EMC VMAX3* 系列產品儲存系統與 SAP HANA® 伺服器叢集

表 3. 溫資料層級伺服器與儲存設備配置

溫資料層級伺服器與儲存設備配置	
伺服器	一個八插槽 Lenovo 伺服器 ³
CPU	Intel® Xeon® 處理器 E7-8850 v2
核心總數	96
總記憶體	2 TB
儲存空間	EMC VMAX3* 系列產品儲存系統

表 4. 溫資料層級網路配置

溫資料層級網路配置	
用戶端網路	四個 10-GbE 連線彙整為邏輯 40-GB 連線，使用 Lenovo RackSwitch G8272*
管理網路	1 GbE 連線，使用 Lenovo RackSwitch G8052
儲存設備網路	兩個 8-GB 光纖通道連線，連接 EMC VMAX3* 系列產品儲存系統與 SAP HANA® 資料分層伺服器

概念性驗證效能成果

在許多方面，這個概念性驗證系統都超過了原始設計目標⁴。

- 此系統可達到每小時 240 億列資料的載入速度，遠遠超過每小時擷取數億列資料的原始目標。
- 將資料載入溫資料層級時僅使用 264 GB 記憶體，以及在 Intel Xeon 處理器上產生百分之 14 的高峰使用量，這能讓系統載入資料時同時快速回應即時查詢。此結果符合為使用者提供即時查詢服務的設計目標，且不會在載入資料時發生明顯的效能降低。
- 在熱資料層級中，僅需要 88 秒即可將 18 億列資料備妥以供分析。在溫資料層級，使用 SAP HANA 預測式分析程式庫 (PAL)，僅需不到兩分鐘即可彙整與分析 18 億列資料。為 20 億列資料進行擷取、轉換，與載入 (ETL) 動作需要 4 分鐘 48 秒，且僅需要 88 秒即可備妥以供分析。整個程序僅需要 6 分鐘，超越原始設計目標的 15 分鐘。
- 總共載入 8,000 億列資料，相當於 51 TB 的未壓縮資料。EMC VMAX3 儲存系列產品可擴充到 4 PB 的容量，並可達到每秒 55,643.78 MB 的輸送量⁵。
- SAP HANA 資料生命週期管理員功能可在 2.5 小時內將 4.29 億列資料從熱資料層級移動到溫資料層級，符合資料行動力的設計目標。

總結

精密的資料收集與分析可協助組織在現代的資料導向商業世界中獲得競爭優勢。資料的速度、數量，以及多樣性與日俱增，加上分析工具日益精密複雜，就算是最強大的關聯式資料庫也可能會受到影響。為了克服傳統資料分析與儲存方式的限制，Intel 與 SAP 攜手合作，設計並建立了一個分析系統，可提供分層資料儲存與對大型資料集的分析能力，更可進行快速又高效率的查詢。此系統在資料載入速度、載入資料時的系統可用度，以及資料在載入後到可以分析的速度等方面，都超過原始概念性驗證的目標。

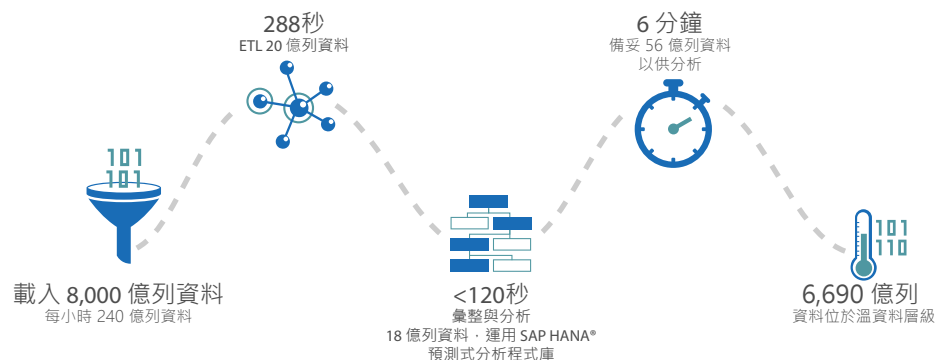


圖 2. Intel 與 SAP 協助加速資料收集與分析程序，同時最佳化長期資料儲存

Intel 與 SAP 結合硬體與軟體，為擷取與分析大規模、高速的資料集提供強大的基礎。SAP HANA 運行於搭載 Intel Xeon 處理器的伺服器叢集上，提供了大幅超過預期的

資料擷取量。此外，此概念性驗證展現了執行於 Intel Xeon 處理器上的 SAP HANA 動態分層，足以彙整與分析大型資料集，對使用者的即時查詢幾乎毫無影響。

¹ 雖然 Intel 與 SAP 環境中有一些伺服器採用 Intel® Xeon® 處理器 E7 v2 系列產品，但 Intel 和 SAP 建議在新配置上採用搭載 Intel Xeon 處理器 E7 v3 系列產品，因為 Intel® Transactional Synchronization Extensions (Intel® TSX) 可將效能改善高達 6 倍。使用新的 Intel TSX 能使交換式工作負載效能改善最高 6 倍，係根據 SAP® 內部的 OLTP insert 及 select 測試，在 SUSE® Linux® Enterprise Server 11 SP3 測量每分鐘交易數 (TPM)。

組態：

- 基準 1.0：4S Intel Xeon 處理器 E7-4890 v2，512 GB 記憶體，SUSE Linux Enterprise Server 11 SP3，SAP HANA® 1 SP8，分數 14,327 TPM。
- 最多提升 1.8 倍 TPM：4S Intel Xeon 處理器 E7-4890 v2，512 GB 記憶體，SUSE Linux Enterprise Server 11 SP3，SAP HANA 1 SP9 分數，26,139 TPM。
- 最多提升 2.7 倍 TPM：4S Intel Xeon 處理器 E7-8890 v3，512 GB 記憶體，SUSE Linux Enterprise Server 11 SP3，SAP HANA 1 SP9，停用 Intel TSX，分數 39,330 TPM。
- 最多提升 6 倍 TPM：4S Intel Xeon 處理器 E7-8890 v3，512 GB 記憶體，SUSE Linux Enterprise Server 11 SP3，SAP HANA 1 SP9，啟用 Intel TSX，分數 89,619 TPM。

如需完整資訊，請參閱 <http://www.intel.com/performance/datacenter>。

² 熱資料層級包含 14 個伺服器，採用下列 CPU，每個伺服器具有 1 TB RAM 和四插槽：兩個伺服器包含 Intel® Xeon® 處理器 E7-8890 v3、四個伺服器包含 Intel Xeon 處理器 E7-8890 v2、五個伺服器包含 Intel Xeon 處理器 E7-8880 v2、兩個伺服器包含 Intel Xeon 處理器 E7-4890 v2，以及一個伺服器包含 Intel Xeon 處理器 E7-4860 v2。熱資料層級的核心總數為 852，RAM 總計 14 TB。

³ 溫資料層級包含一個八插槽伺服器，包含 Intel® Xeon® 處理器 E7-8850 v2。核心總數為 96，RAM 總計 2 TB。

⁴ 總共載入 8,000 億列資料，大約為 51 TB 的未壓縮資料。EMC VMAX3® 系統的速度可達到每秒 800 MB (約為可用容量的百分之 2)，同時每小時將 240 億列寫入溫資料層級。資料載入結合平行載入 100 個資料檔案，使用溫資料層級 192 個核心之中的 100 個，剩餘的核心用於即時查詢。每個資料檔案的大小為 5 GB，包含約 7,700 萬列資料。請注意，並未使用業界標準效能標竿測量效能。

⁵ 如需完整的 EMC VMAX3® 效能標竿資訊，請造訪 http://www.storageperformance.org/benchmark_results_files/SPC-2/EMC/B00073_EMC_VMAX-400K/B00073_EMC_VMAX-400K_SPC-2_executive-summary.pdf。

效能測試中使用的軟體與工作負載可能僅針對 Intel 微處理器進行最佳化。效能測試 (例如 SYSmark® 與 MobileMark*) 係使用特定的電腦系統、元件、軟體、作業及功能進行測量。這些因素若有任何異動，均可能導致測得結果產生變化。建議您參考其他資訊與效能測試數據，協助您充分評估欲購買產品的性能，包括該產品在搭配其他產品運作時的效能。

Intel 技術的功能與優勢取決於系統配置，而且可能需要支援的硬體、軟體或服務啟動。實際效能會依系統組態而異。沒有電腦系統能提供絕對的安全性。請諮詢當地的系統製造商或零售商，或造訪 intel.com 進一步瞭解。

本文件並未透過禁反言或任何其他方式，明示或暗示授予任何智慧財產權。

此處所述的產品可能包含設計瑕疵或錯誤，稱為勘誤內容，可能導致產品不符合先前發佈的規格。若有需要，可向本公司索取最新的勘誤表。

Intel 並不控制或稽核本文件提及的第三方效能標竿資料或網站。您應造訪該網站並確認本文件提及的資料是否正確。

Intel、Intel 圖誌、Intel Inside、Intel Inside 圖誌，以及 Xeon 是 Intel Corporation 在美國及其他國家的商標。

Copyright © 2016 Intel Corporation。保留一切權利。

*其他品牌和名稱為其所屬公司的資產。