GPU 在機器學習領域的應用現狀

眾所周知, GPU 已經成為支撐AI應用的 一種關鍵計算加速設備, GPU 的流多處 理器架構非常適合用來加快深度神經網路 應用中的大量矩陣運算過程。大量實測資 料表明,跟通用處理器相比,GPU 在運 行深度神經網路時具有顯著的效能優勢。 主流機器學習框架,如 TensorFlow 和 PyTorch 都支援使用 GPU 來加速深度神 經網路的訓練與推理計算。然而這些主流 深度學習框架只能將一個或多個 GPU 設 備分配給單個使用者或 AI 應用使用,這 給那些希望共用使用昂貴的 GPU 計算加 速設備的使用者帶來了很多困擾。此外由 於神經網路的訓練經常是一個需要反復調 整參數持續改進的過程,對於那些只需消 耗單個 GPU 部分計算資源的神經網路應 用,獨佔整個 GPU 會造成極大的資源浪 費,資料表明典型的 GPU 資源利用率只 有 20%到30%。

另外現有基於 GPU 的 AI 應用只能在具有 GPU 設備的電腦上運行·而在通用處理器占主導的邊緣計算與電信基礎架構網路等環境中·帶有 GPU 設備的電腦相對仍只是少數·如何能支援在任意計算節點來啟動和運行 GPU 加速的 AI 應用·目前急需解決方案。

綜上所述,現代 AI 應用呼喚一種能夠把 GPU 硬體拆分使用的加速器虛擬化技術 ,使之能夠靈活地滿足多種工作負載的需 求;同時通過在物理 GPU 和 AI 應用之 間實現一個虛擬抽象層,解耦 GPU 設備 與 AI 應用的位置綁定。

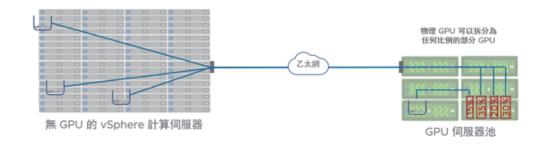
GPU 共用解決方案 Bitfusion

VMware vSphere 7整合Bitfusion功能·針對AI和ML應用提供彈性基礎架構 Bitfusion與VMware vSphere的整合將幫助企業節省成本·實現開箱即用式資源分享·並在正確的時間為正確的工作負載提供GPU等恰到好處的硬體加速器資源。

GPU 虛擬化技術 Bitfusion

Bitfusion 就是一種能夠滿足上述需求的 GPU 虛擬化技術·它能夠把 GPU 拆分成任意大小 (用百分比來指定)來分配給不同的工作負載或使用者使用;並且它不要求工作負載和 GPU 在同一台伺服器上·工作負載通過網路來遠端使用 GPU 資源。利用 Bitfusion 技術可以搭建一個 GPU 伺服器池·把所有的 GPU 資源集中在一起·然後再根據需求把 GPU 按比例拆分出一個小的部分 GPU (Partial GPU)來供工作負載使用。下圖就是一個典型的示意圖·我們把一塊 GPU 拆分成了4塊 (大小分別為 15%、35%、20%、30%)·分別提供給四個工作負載使用·其中三個工作負載運行在無 GPU 的 vSphere 計算伺服器上·一個工作負載運行在 GPU 伺服器池中的某一台伺服器上。

Bitfusion 的 GPU 資源池有點類似於存儲區域網路 SAN (Storage Attached Network),所以也有人把它叫作 GPU Attached Network。



Bitfusion 的 GPU 虛擬化技術可以解決以下問題:

- 解決 GPU 短缺問題:當多個用戶爭奪有限的 GPU 資源時 (尤其是高端的 GPU) · 我們可以把 GPU 拆分成多個部分來分配給多人使用 · 每個工作負載只用到部分的 GPU 。 開發和測試等一些試驗性的工作只需要驗證模型是否工作正常 · 對於 GPU 性能的要求不是那麼高 · 這一類工作可以使用更小的 GPU 分片。
- 提高 GPU 的利用率:把一個工作負載配置成使用 33%的 GPU · 我們同時運行這個工作負載的三個實例 · 我們就可以幾乎達到 100%的 GPU 利用率 · 機器學習的吞吐率也大大提高 (通常可以看到 2.5 倍的提高)。



Bitfusion 解決方案具有以下特點:

- 拆分成任意大小: Bitfusion 可以指定任意大小的拆分·例如 1%; 如前所述· 這特別適用 干開發測試等試驗性的應用場景。
- GPU 獨立性:拆分出的 GPU 相互獨立,各自運行不同的 AI 框架和模型,絕對不會相互影響。
- 大小可動態調整: 拆分出來的部分 GPU 可以動態調整大小‧例如從同一塊物理 GPU 拆分出來的兩塊部分 GPU 大小分別為 45% 和 55% · 55% 的 GPU 可以進一步拆分成更小的兩塊 35% 和 20% · 而不會影響另一塊 45% 部分 GPU 上工作負載的正常運行。
- 支援多個物理 GPU: 從不同物理 GPU 中拆分出來的多個部分 GPU 可以分配給同一個使用者和工作負載·這既可以提高整個 GPU 資源池的利用率·也有助於開發和調試多 GPU 工作負載應用。

GPU 共用解決方案 Bitfusion

VMware vSphere的Bitfusion功能在單獨下載後即可使用,不影響當前基礎架構,並且能與現有工作流程和生命週期無縫整合。 VMware去年收購Bitfusion,就希望將這項技術整合至VMware vSphere之中,Bitfusion所提供的軟體平台可將特定的物理資源,從環境中所連接的伺服器中解耦合,其中包括在虛擬化基礎架構中共用GPU,將其作為網路可訪問資源池,而非單個伺服器的孤立資源。

