

無限融合引領AI突破

理大研究大模型訓練

中國AI晶片進口受限，科學家另闢蹊徑，DeepSeek大幅減少AI算力要求後，本港學者亦開拓研究，讓更多人參與AI。

上星期，DeepSeek再在arxiv發表論文，提出注意力機制NSA，減少超長文本訓練和推論資源算力，引起關注。另一群AI學者在arxiv則提出嶄新大模型融合策略InfiFusion，大幅減少算力需求，文章第一通訊作者為香港理工大學楊紅霞教授。

楊紅霞是少數具大模型訓練實戰經驗專家，從美國回流，預見AI模型面臨轉捩點，研發新一代模型訓練，推動Model over models (MoM) 範式帶動AI升級。

節省成本驚人

InfiFusion融合數個不同模型專長，集合變成稱為「Pivot Model」大模型，合併不同模型結構詞彙，保留各自能力，增強跨模型的推理能力，真正是強強聯手，同時大幅減少整體算力的要求。小模型訓練成本相對低，讓不同領域專家參與，各自訓練小模型，通過InfiFusion融合成更強大模型。

從論文的數據，InfiFusion節省的成本驚人，性能卻比相同大小模型更好。InfiFusion提出兩種融合策略：Pairwise Fusion (成對融合) 和 Unified Fusion (統一融合)，前者是主模型與來源模型分別配對，逐對找到最佳權重，以優化模型特殊性能。統一融合則是所有模型同時訓練主模型，目標

是找到全域最佳權重，整個模型性能臻至最佳。

目前AI訓練由少數人操縱，許多人無從參與，長遠空虛通用AI出現。

香港理工大學楊紅霞教授

以成對融合訓練出的InfiFusion & TA，140億參數規模模型為例，InfiFusion & TA利用450小時GPU訓練，各種能力平均得分79.96分。透過統一融合，訓練相同規模的模型InfiFusion u僅160小時，竟獲79.92分，性價比極高。市場上相同大模型，以qwen2.5-14b-instruct為例，140億參數平均獲75.65分，卻需180萬小時訓練。從訓練結果，可見融合訓練之後，模型能力並沒損失，整體能力還提高了。

統一融合InfiFusion u適合一般用途，成對融合適合特定任務模型組合。論文通過較小模型，證明InfiFusion較傳統蒸餾(Distillation)融合的方法，成本大幅降低，更重要是可無限融合，過程不流失模型能力，打破了規模局限。蒸餾是從較大教師模型，轉移知識到較小學生模型，融合多個模型較難，InfiFusion卻是無限融合。

InfiFusion低成本高效益，亦可取代以微調(Finetuning)和RAG (Retrieval Augmented Generation) 等為模型加入專有知識技術的技巧。

InfiFusion通過更高效融合，保留模型專業知識和能力，可細緻調整和優化，以應對不同任務。楊紅霞相信，模型訓練趨向以「模型為中心」(Model centric)，以解決更多專門領域難題，可擴展至新知識領域，不止限於服務一般領域。

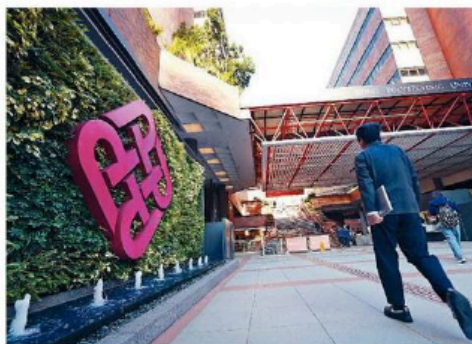
論文證明InfiFusion融合，模型保留各自優勢。例如一個模型擅於推理，另一個具醫學知識，兩者互補迸發不同優勢，最終在更多場景表現更好。據楊紅霞解釋，InfiFusion長遠意義，在於全球知識領域盈千上萬，散布不同專家手上，目前AI訓練由少數人操縱，許多人無從參與，長遠空虛通用AI出現。

預訓練無以為繼

「專業知識分拆變成小模型後，再在模型層面上融合，獲得基礎大模型能力，減少了訓練成本，更多人貢獻知識，模型能力不遜直接訓練。」

楊紅霞提到，建立基礎模型過程，包括了預訓練(Pretraining)和對齊(Alignment)兩階段，預訓練消耗大量數據和算力，通過scaling law壓縮大量知識，學習數據模式、特徵和知識，掌握語言結構和語義資訊。

楊紅霞說，網上知識即將耗盡，預訓練無以為繼。一年前，美國AI科學家Ilya Sutskever洞悉預訓練時代結束；結果DeepSeek改在對齊階段發力，以微調(Fine-tuning)和強化學習(Reinforcement Learning)推進模型能力，透過調整模型行為，實現如深理解「有限數據」，改進推理和決策。



楊紅霞說，DeepSeek在對齊階段，提高實際應用性能和推理決策，原理有如教導孩子成長。「以教小孩為比喻，預訓練好比孩子天生智商和知識量，對齊是後天教導，所以同樣智商小孩，不同教導方式下，才能有不同發揮。」DeepSeek從教導入手，增強模型推理能力。

性能上超越微調及RAG

即使模型推理能力增強，始終只是學習網上原始數據，離不開以「數據為中心」(Data Centric)，只能原地踏步，一味比拼數學和編程等，展示網上掌握技能，卻缺乏專業解難能力。

InfiFusion讓模型能力泛化至尖端領域，性能上大幅超越微調和RAG，具靈活性和泛化能力，有助發展科學研究，新材料發現和醫學診斷。楊紅霞說，InfiFusion作為一種新範式具深

DeepSeek大幅減少AI算力要求後，包括在內的本港學者亦開拓研究，讓更多人參與AI。



楊紅霞說，DeepSeek在對齊階段，提高實際應用性能和推理決策，原理有如教導孩子成長。

厚潛力，吸引行業專家參與AI，集思廣益，海納百川建立通用型的AI，不再只靠堆算力，涉及億計美元成本，大部分人卻摒諸門外。

她指出，本港醫學知識非常豐富，可讓本港研究人員投入大模型訓練，匯聚醫學知識，憑香港多年累積知識，通過AI造福更多病人。