

人工智能大模型

——當代歷史的標誌性事件及其意義



此項研究在這樣的猜想基礎上進行，即學習以及智能的任何其他特性的每一方面在原則上都能被精確描述，以致可使一台機器來模擬它。我們會嘗試尋求如何讓機器使用語言，形成抽象和概念，解決現在留待人類解決的問題，並提升自己。

——1956年達特茅斯會議人工智能(AI)定義①

2020至2022年，在新冠疫情肆虐全球的陰霾日子裏，人工智能(AI)創新的步伐完全沒有停止。美國人工智能研究公司OpenAI異軍突起：2020年4月發布神經網絡Jukebox②；5月發布語言模型GPT-3③；6月開放人工智能應用程式介面(Application Programming Interface, API)；2021年1月發布連接文本和圖像的神經網絡CLIP④；同月發布從文本創建圖像的神經網絡DALL·E⑤；2022年11月正式推出了對話互動式的聊天機器人程式ChatGPT⑥。相比於GPT-3，ChatGPT引入了基於人類回饋的強化學習(Reinforcement Learning from Human Feedback, RLHF)技術以及獎勵機制⑦。

GPT-3的發布是人類科技史上的里程碑事件，在短短幾個月席捲全球，速度超過人類最狂野的想像。GPT-3證明了一個具有高水平複雜結構和大量參數的人工智能大模型(foundation model，又稱「基礎模型」)可以實現深度學習(deep learning)。此後，大模型概念得到前所未有的關注和討論。但是，關於「大模型」的定義，對其內涵的理解和詮釋卻莫衷一是，「橫看成嶺側成峰，遠近高低各不同」。

儘管如此，並不妨礙人們形成了關於大模型的基本共識：大模型是大語言模型(Large Language Model, LLM)，也是多模態模型(multimodal model)。GPT是大模型的一種形態，G代表生成性的(generative)，P代表經過預訓

練 (pre-trained)，T 代表變換器 (transformer) ⑥。它引發了人工智能生成內容 (Artificial Intelligence Generated Content, AIGC) 技術的質變。大模型是人工智能賴以生存和發展的基礎。現在，與其說人類開始進入人工智能時代，不如說人類進入的是大模型時代。我們不僅目睹，也身在其中體驗了生成式大模型如何開始生成一個全新時代。

本文通過七個部分，分別說明大模型的定義、人工智能的歷史、大模型的基本特徵、Transformer 結構、GPU 和能源、知識革命、「人的工具化」及大模型在其中的作用，有助於進一步解讀大模型對於人類科技發展的重要意涵。

一 何謂大模型？

人工智能的模型，與通常的模型一樣，是以數學和統計學作為演算法基礎的，可以用來描述一個系統或者一個數據集。在機器學習 (machine learning) 中，模型是核心概念。模型通常是一個函數或者一組函數，以線性函數、非線性函數、決策樹、神經網絡等各種形式呈現。模型的本質就是對這個/組函數映射的描述和抽象，通過對模型進行訓練和優化，能夠得到更加準確和有效的函數映射。模型的目的是為了從數據中找出一些規律和模式，達到預測未來的結果。模型的複雜度可以理解為模型所包含的參數數量。一個模型的參數數量愈多，通常意味着該模型可以處理更複雜、更豐富的信息，具備更高的準確性和表現力。大模型一般用於解決複雜的自然語言處理 (Natural Language Processing, NLP)、電腦視覺和語音辨識等任務。這些任務需要處理大量的輸入數據，並從中提取複雜的特徵和模式。通過使用大模型，深度學習演算法就能更好地處理這些任務，提高模型的準確性和性能。

大模型的「大」，是指模型參數至少達到 1 億以上。但是這個標準一直在升級，目前很可能已經有了萬億參數以上的模型。GPT-3 大約的參數規模是 1,750 億。除了大模型之外，還有所謂的「超大模型」。超大模型是比大模型更大、更複雜的人工神經網絡 (Artificial Neural Network, ANN) 模型，通常擁有數萬億到數千萬億參數。超大模型一般被用於解決更為複雜的任務，如自然語言處理中的問答和機器翻譯、電腦視覺中的目標檢測和圖像生成等。這些任務需要處理極其複雜的輸入數據和高維度的特徵，超大模型可以在這些數據中提取出更深層次的特徵和模式，提高模型的準確性和性能，所以，超大模型的訓練和調整需要極其巨大的計算資源和大量數據、更加複雜的演算法和技術、大規模的投入和協作。

大模型和超大模型的主要區別在於模型參數數量的多寡、計算資源的需求和性能表現。伴隨大模型參數規模的膨脹，大模型和超大模型的界限正在消失。現在包括 GPT-4 在內的代表性大模型，其實就是原本的超大模型。或者說，原本的超大模型，就是現在的大模型。

如前所述，大模型可以定義為大語言模型，即具有大規模參數和複雜網絡結構的語言模型。與傳統語言模型（如生成性模型、分析性模型、辨識性模型）不同^⑨，大語言模型通過在大規模語料庫上進行訓練來學習語言的統計性規律，在訓練時通常通過大量的文本數據進行自監督學習^⑩，從而能夠自動學習到語法、句法、語義等多層次的語言規律。

如果從人工智能的生成角度定義大模型，與傳統的機器學習演算法不同，生成式大模型可以根據文本提示生成代碼，還可以解釋代碼，甚至在某些情況下調試代碼。在這樣的過程中，不僅實現文本、圖像、音訊、視頻的生成，構建多模態，而且還在更為廣泛的領域生成新的設計、新的知識和思想，甚至廣義的藝術和科學的再創造。

近幾年，比較有影響的大模型主要來自 Google、Meta 和 OpenAI。除了 OpenAI 的 GPT 之外，2018 至 2023 年 Google 先後發布對話程式語言模型 LaMDA、BERT 和 PaLM-E^⑪。2023 年，Facebook 的母公司 Meta 推出大語言模型 LLaMA，以及在 Meta AI 博客上免費公開大語言模型 OPT-175B^⑫。在中國，大模型主要代表是百度的「文心一言」和華為的「盤古」。這些模型的共同特徵是：需要在大規模數據集上進行訓練，基於大量的計算資源進行優化和調整。因為大模型的出現和發展所顯示的湧現性、擴展性和複合性，長期以來人們討論的所謂「弱人工智能」、「強人工智能」和「超人工智能」的界限不復存在，這樣劃分的意義也自然消失^⑬。

二 大模型是人工智能歷史的突變和湧現

如果從 1956 年美國達特茅斯學院 (Dartmouth College) 的人工智能會議算起，還有三年，人工智能歷史就踏入七十年。該會議引申出人工智能三個基本派別：一、符號學派 (Symbolism)，又稱為邏輯主義、心理學派或電腦學派。該學派主張通過電腦符號操作來類比人的認知過程和大腦抽象邏輯思維，實現人工智能。符號學派主要集中在人類推理、規劃、知識表示等高級智能領域。二、聯結學派 (Connectionism)，又稱為仿生學派或生理學派。該學派強調對人類大腦的直接類比，認為神經網絡和神經網絡間的連接機制與學習演算法能夠產生人工智能。學習和訓練是需要有內容的，數據就是機器學習、訓練的內容。聯結學派的技术性突破包括感知機 (下詳)、人工神經網絡、深度學習。三、行為學派 (Actionism)，思想來源是進化論和控制論。其原理為控制論以及感知—動作型控制系統。該學派認為行為是個體用於適應環境變化的各種身體反應的組合，它的理論目標在於預見和控制行為^⑭。

比較上述三個人工智能派別：符號學派依據的是抽象思維，注重數學可解釋性；聯結學派則是形象思維，偏向於仿人腦模型；行為學派是感知思維，傾向身體和行為模擬。從共同性方面來說，這三個派別都以演算法、算

力和數據作為核心要素。但是在相當長的時間裏，符號學派主張的基於推理和邏輯的人工智能路線處於主流地位。不過，電腦只能處理符號，不可能具有人類最為複雜的感知。二十世紀80年代末，符號學派開始走向式微。之後的人工智能編年史，有三個重要的里程碑。

第一個里程碑：機器學習。機器學習理論的提出，可以追溯到圖靈(Alan Turing)寫於1950年的一篇論文〈電腦機器與智慧〉(“Computing Machinery and Intelligence”)和圖靈測試(The Turing test) ⑮。1952年，在國際商業機器公司(IBM)工作的塞繆爾(Arthur L. Samuel)開發了一個西洋棋的程式。該程式能夠通過棋子的位置學習一個隱式模型，為下一步棋提供比較好的走法。塞繆爾用這個程式駁倒了機器無法超越書面代碼、並像人類一樣學習的論斷。他創造並定義了「機器學習」⑯。

機器學習是一個讓電腦不用顯示程式設計就能獲得能力的研究領域。1980年，美國卡內基梅隆大學(Carnegie Mellon University)召開了第一屆機器學習國際研討會，標誌着機器學習研究已在全世界興起。此後，機器學習開始得到大量應用。1986年，三十多位人工智能專家共同撰寫的《機器學習：一項人工智能方案》(*Machine Learning: An Artificial Intelligence Approach*)文集第二卷出版⑰，顯示出機器學習突飛猛進的發展趨勢⑱。二十世紀80年代中葉是機器學習的最新階段，機器學習已成為新的學科，它綜合應用了心理學、生物學、神經生理學、數學、自動化和電腦科學等，形成理論基礎。1995年，瓦普尼克(Vladimir N. Vapnik)和科茨(Corinna Cortes)提出的支持向量機(Support Vector Machine, SVM，又稱「支持向量網絡」)，實現機器學習領域最重要的突破，具有非常強的理論論證和實證結果。

機器學習有別於人類學習，二者的應用範圍和知識結構有所不同：機器學習是基於對數據和規則的處理和推理，主要應用於數據分析、模式識別、自然語言處理等領域；而人類學習是一種有目的、有意識、逐步積累的過程。總之，機器學習是一種基於演算法和模型的自動化過程，並分為監督學習和自監督學習兩種。

第二個里程碑：深度學習。深度學習是機器學習的一個分支。所謂「深度」是指神經網絡中隱藏層(位於輸入和輸出之間的層)的數量。傳統的神經網絡只包含兩至三個隱藏層，而深度神經網絡可以有高達150個隱藏層，提供了大規模的學習能力。隨着大數據和深度學習爆發並得以高速發展，最終成就了深度學習理論和實踐。2006年，辛頓(Geoffrey E. Hinton)正式提出「深度置信網絡」(Deep Belief Nets/Deep Belief Network, DBN)概念⑲，那一年成為了「深度學習元年」。在辛頓深度學習理論的背後，是堅信如果不了解大腦，就永遠無法理解人類的認識。人腦必須用自然語言進行溝通，而只有1.5公斤重的大腦，大約有860億個神經元(通常稱為「灰質」)與數萬億個突觸相連。人們可以把神經元看作是接收數據的中央處理器(Central Processing Unit, CPU)。所謂「深度學習」可以伴隨着突觸的增強或減弱而發生，即在一個擁有大量神

經元的大型神經網絡中，計算節點和它們之間的連接，僅通過改變連接的強度，從數據中學習。辛頓認為，實現人工智能的進步需要通過生物學途徑，或者通過神經網絡途徑替代模擬硬件途徑，形成基於100萬億個神經元之間的連接變化的深度學習。

深度學習主要涉及三類方法：一、基於卷積運算的神經網絡系統，卷積神經網絡 (Convolutional Neural Network, CNN) 是一類包含卷積運算且具有深度結構的前饋神經網絡，是深度學習的代表演算法之一。二、基於多層神經元的自編碼神經網絡，包括自編碼 (auto encoder) 和近年來受到廣泛關注的稀疏編碼 (sparse coding) 兩類。三、以多層自編碼神經網絡的方式進行預訓練，進而結合鑒別信息進一步優化神經網絡權值的深度置信網絡。通過多層處理，逐漸將初始的「低層」特徵表示轉化為「高層」特徵表示後，用簡單模型即可完成複雜的分類等學習任務。

深度學習是建立在人工神經網絡理論和機器學習理論上的科學，它使用建立在複雜的網絡結構上的多處理層，結合非線性轉換方法，對複雜的數據模型進行抽象，得以識別圖像、聲音和文本。在深度學習的歷史上，卷積神經網絡和循環神經網絡 (Recurrent Neural Network, RNN) 曾經是兩種經典模型。在循環神經網絡中，節點之間的連接可以形成一個循環，允許一些節點的輸出影響到同一節點的後續輸入，因此能夠表現出時間上的動態行為。

2012年，辛頓和克里澤夫斯基 (Alex Krizhevsky) 設計的 AlexNet 神經網絡模型在 ImageNet 競賽中實現圖像識別和分類，成為新一輪人工智能發展的起點。這類系統可以處理大量數據，發現人類通常無法發現的關係和模式。2016年人工智能機器人 AlphaGO 戰勝韓國職業圍棋棋手李世石，這是深度學習的經典範例。

第三個里程碑：人工智能生成內容大模型。2018年10月，Google 發布 BERT 模型是代表性事件。該模型是一種雙向的基於 Transformer 的自監督語言模型，通過大規模預訓練無標註數據來學習通用的語言表示，從而能夠在多種下游任務，如專名識別、詞性標記和問題回答中進行微調。利用大型文本語料庫 BookCorpus 和英文維基百科裏純文字的部分，無須標註數據，用設計的兩個自監督任務來進行訓練，訓練完成的模型通過微調在十一個下游任務上實現最佳性能。

因為 BERT 模型，掀起了預訓練模型的研究熱潮，從2018年開始大模型迅速流行，預訓練語言模型 (Pre-trained Language Model, PLM) 及其「預訓練—微調」方法已成為自然語言處理任務的主流範式。大模型利用大規模無標註數據通過自監督學習進行預訓練，再利用下游任務的有標註數據進行自監督學習以微調模型參數，實現下游任務的適配²⁰。

如前所述，大模型的訓練需要大量的計算資源和數據，OpenAI 使用了數萬台 CPU 和圖形處理器 (Graphics Processing Unit, GPU)，並利用了多種技術，如自監督學習和增量訓練等，對模型進行了優化和調整。2018至2023年，

OpenAI 實現大模型從 GPT-1 到 GPT-4 的五次迭代，同時開放了應用程式介面，使得開發者可以利用大模型進行自然語言處理的應用開發。

總之，大模型是基於包括數學、統計學、電腦科學、物理學、工程學、神經學、語言學、哲學、人工智能學融合基礎上的一次突變，並導致了一種「湧現」(emergence)。大模型是一種革命。在模型尚未達到某個臨界點之前，根本無法解決問題，性能也不會比隨機好。但是，當大模型突破某個臨界點之後，性能會發生愈來愈明顯的改善，形成爆發性的湧現能力。如論者所言：「許多新的能力在中小模型上線性放大規模都得不到線性的增長，模型規模必須要指數級增長超過某個臨界點，新技能才會突飛猛進。」^{②1}

更為重要的是，大模型賦予人工智能以思維能力——一種與人類近似，又很不相同的思維能力。前述 AlphaGo 戰勝李世石的世紀級圍棋大賽，證明了人工智能思維的優勢。

三 大模型的基本特徵

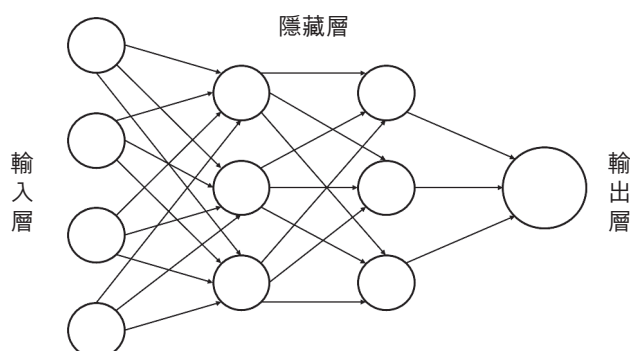
大模型的基本特徵可以總結為：以人工神經網絡作為基礎；為神經網絡提供更好的預訓練方法並促進規模化，能顯著降低人工智能工程化門檻；具有理解自然語言的能力和模式；已經形成「思維鏈」；需要向量數據庫的支援；具有不斷成長的泛化功能，並且被植入了控制論的基於人類回饋的強化學習機制。

大模型以人工神經網絡作為基礎。1943 年，心理學家麥卡洛克 (Warren S. McCulloch) 和數理邏輯學家皮茨 (Walter H. Pitts, Jr.) 建立了第一個神經網絡模型，即 M-P 模型 (又稱「麥卡洛克-皮茨模型」或「MCP 模型」)。該模型是對生物神經元結構的一種模仿，將神經元的樹突、細胞體等接收信號定義為輸入值 x ，突觸發出的信號定義為輸出值 y 。M-P 模型奠定了支援邏輯運算的神經網絡基礎。1958 年，電腦專家羅森布拉特 (Frank Rosenblatt) 基於 M-P 模型發明了包括輸入層、輸出層和隱藏層的感知機 (perceptron)。神經網絡的隱藏層最能代表輸入數據類型特徵 (圖 1)。從本質上講，這是第一台使用模擬人類思維過程的神經網絡的新型電腦。

以 OpenAI 為代表的團隊，為了讓具有多層表示的神經網絡學會複雜事物，創造了一個初始化網絡的方法，即預訓練。實際上，生成式大模型為神經網絡提供了更好的預訓練方法。現在的大模型都是以人工神經網絡作為基礎的演算法數學模型，其基本原理依然是羅森布拉特的感知機。這種人工智能網絡依靠系統的複雜程度，通過調整內部大量節點之間相互連接的關係，從而達到處理信息的目的。

大模型生成內容的前提是大規模的文本數據輸入，並在海量的通用數據上進行預訓練。通過預訓練不斷調整和優化模型參數，使得模型的預測結果盡可能接近實際結果。預訓練中使用的大量文本數據包括維基百科、網頁文

圖1 神經網絡的層級關係：由輸入到輸出



圖片來源：筆者改製自 Moonzarin Reza, "Galaxy Morphology Classification Using Automated Machine Learning", *Astronomy and Computing*, vol. 37 (October 2021), <https://doi.org/10.1016/j.ascom.2021.100492>。

本、書籍、新聞文章等，用於訓練模型的語言模型部分。此外，還可以根據應用場景和需求，調用其他外部數據資源，包括知識庫、情感詞典、關鍵詞提取、實體識別等。在預訓練的過程中，大模型不是依賴於人為編寫的語法規則或句法規則，而是通過學習到的語言模式和統計性規律，以生成更加符合特定需求和目標的文本輸出。

預訓練促進了規模化。所謂的「規模化」是指用於訓練模型的大量計算，最終轉化為規模愈來愈大的模型，具有愈來愈多的參數。在預訓練過程中，大模型形成理解上下文的學習能力。或者說，伴隨上下文學習的出現，人們可以直接使用預訓練模型。大模型通過大量語料庫訓練，根據輸入文本和上下文生成合適的文本輸出，學習詞彙、句法結構、語法規則等多層次的語言知識；通過對大量樣本進行學習，更多的計算資源的投入（包括正確和錯誤的文本樣本），捕捉到語法和句法的統計性規律，形成一個詞或字元的概率的預測能力，進而根據不同樣本的預測錯誤程度調整參數，處理複雜的語境，最終逐漸優化生成的文本。例如，ChatGPT 會根據之前與使用者交互的上下文和當前的生成狀態，選擇最有可能的下一個詞或短語。

「預訓練—微調」方法能顯著降低人工智能工程化門檻。預訓練模型在海量數據的學習訓練後具有良好的泛化性（下詳），使得細分場景的應用廠商能夠基於大模型，通過零樣本、小樣本學習來獲得顯著的效果。因此，人工智能有望構建成統一的智慧底座，以賦能各行各業。生成式大模型不會止步於簡單的內容生成，而會逐步達到更高的人工智能，得以預測、決策、探索。針對大量數據訓練出來的預訓練模型，後期採用業務相關數據進一步訓練原先模型的相關部分，給出額外的指令或者標註數據集來提升模型的性能，通過微調從而得到準確度更高的模型。

大模型具有理解自然語言的能力和模式。自然語言如漢語、英語及其文字，具有複雜性和多樣性，且伴隨文化演變而進化；通過表達含義，實現人

類溝通和交流，推動人類思維發展。對自然語言的理解，首先要理解文本的特徵。在大模型研究的早期階段，主要集中在自然語言處理領域，形成從簡單的文本問答、文本創作到符號式語言的推理能力。之後大模型發生程式設計語言的變化，有助於更多人直接參與大模型用於問答的自然語言交互和程式設計模式，經過形式極簡的文本輸入，利用自然語言表達的豐富性，形成自然語言與模型的互動。上述的BERT、GPT等一系列代表性模型，不同於基於語法規則、句法規則的傳統語言模型；這些大語言模型基於統計語言學的思想，在大量文本數據上進行自監督學習，利用自然語言中的統計性規律（涉及貝葉斯原理[Bayes theorem]和馬爾可夫鏈[Markov chain]等數學工具，以及N元[n-gram]語言模型²⁰），通過對大量語法和句法樣本學習，捕捉到相關規則並進行推斷，對各種不同形式的語言表達具有一定的容忍性、適應性和靈活性，從而生成具有語法和語義合理性的文本。

大模型已經形成「思維鏈」(Chain-of-Thought, CoT)。思維鏈是重要的微調技術手段，其本質是一個多步推理的過程。通過讓大語言模型將一個問題拆解為多個步驟，一步一步分析，逐步得出正確答案。我們還可以這樣理解：思維鏈相當於大模型中的數據，人工智能以思維鏈為數據，然後再進行微調和回饋，從而形成人工智能能力。在電腦語言中，有所謂「第四範式」(Fourth Normal Form, 4NF) 概念，有助於理解思維鏈的功能，也有助於大模型更加結構化和規範化，減少數據信息冗餘和碎片化等弊病，提高大模型的效率。

大模型需要向量數據庫的支援。向量是大模型的數據存儲的基本單位。雖然大模型呈現端到端、文本輸入輸出的形式，但是實際接收和學習的數據並不是傳統文本，因為文本本身數據維度太高、學習過於低效，所以需要量化的文本。「所謂量化的文本，就是模型對自然語言的壓縮和總結」。向量是人工智能理解世界的通用數據形式，大模型依賴向量數據庫，其即時性對分散式運算的要求很高，隨着數據的變化即時更新，以保障向量的高效存儲和搜索²¹。

大模型具有不斷成長的泛化 (generalization) 功能。大模型泛化是指大模型可以應用 (泛化) 到其他場景，泛化能力是模型的核心。大語言模型通過大量的數據訓練，掌握語言中的潛在模式和規律，在面對新的、未見過的語言表達時具有一定的泛化能力。在新的場景下，針對新的輸入信息，大模型就能做出判斷和預測。而基於語法規則、句法規則的傳統語言模型通常需要人為編寫和維護規則，對於未見過的語言表達可能表現較差。針對泛化誤差，大模型通常採用遷移學習、微調等手段，在數學上權衡偏差和方差。大語言模型廣泛應用於自然語言處理領域的多個任務，如語言生成、文本分類、情感分析、機器翻譯等。說到底，大模型的泛化就是指其通用性，最終需要突破泛化過程的局限性。但是，實現通用大模型，還有很長的路。

大模型植入了控制論的基於人類回饋的強化學習機制。回饋是控制論中的基本概念，是指一個系統把信息輸送出去，又把其作用結果返回，並對信

息的再輸出產生影響，起到控制和調節作用的過程。大模型構建人類回饋數據集，訓練一個激勵模型，模仿人類偏好對結果打分，通過從外部獲得激勵來校正學習方向，從而獲得一種自適應 (self-adaptive) 的學習能力。

四 大模型和 Transformer

如果說神經網絡是大模型的「大腦」，那麼 Transformer 就是大模型的「心臟」。2017年6月，Google 團隊的瓦斯瓦尼 (Ashish Vaswani) 等人發表論文〈注意力足矣〉(“Attention Is All You Need”)，系統提出了 Transformer 的原理、構建和大模型演算法。此文的開創性思想，顛覆了以往將序列建模和循環神經網絡畫等號的思路，開啟了預訓練模型的時代²⁹。

Transformer 是一種基於「注意力機制」(attention mechanism) 的深度神經網絡，可以高效並行處理序列數據，與人的大腦非常近似。Transformer 的基本特徵如下：(1) 由編碼組件 (encoder) 和解碼組件 (decoder) 兩個部分組成。(2) 採用神經網絡處理序列數據。神經網絡的工作是將一種類型的數據轉換為另一種類型的數據；在訓練期間，神經網絡的隱藏層以最能代表輸入數據類型特徵的方式調整其參數，並將其映射到輸出。(3) 擁有的訓練數據和參數愈多，它就愈有能力在較長文本序列中保持連貫性和一致性。(4) 標記和嵌入。輸入文本必須經過處理並轉換為統一格式，然後才能輸入到 Transformer。(5) 實現並行處理整個序列，從而將深度學習模型擴展到前所未有的速度和容量。(6) 引入了注意力機制，在正向和反向的非常長的文本序列中跟蹤單詞之間的關係，包括自注意力 (self-attention) 機制和多頭注意力 (multi-head attention) 機制。Transformer 的多頭注意力機制中有多個自注意力機制，可以捕獲單詞之間多種維度上的相關系數注意力評分 (attention score)，摒棄了卷積神經網絡和循環神經網絡。(7) 訓練和回饋。在訓練期間，Transformer 提供了規模非常大的配對示例語料庫 (例如英語句子及其相應的法語翻譯)。編碼器模組接收並處理完整的輸入字串，嘗試建立編碼的注意向量和預期結果之間的映射。

在 Transformer 之前，發揮近似功能的是循環神經網絡或卷積神經網絡。Transformer 起初主要應用於自然語言處理，但漸漸地，它們在幾乎所有的領域都發揮了作用。通用性一直是 Transformer 最大的優勢，包括圖像、視頻、音訊等多種領域的模型都需要使用 Transformer。

總之，Transformer 是一種非常高效、易於擴展、並行化的神經網絡架構，其核心是基於注意力機制的技術，可以建立起輸入和輸出數據的不同組成部分之間的依賴關係，具有品質更優、更強的並行性和訓練時間顯著減少的優勢。Transformer 現在被廣泛應用於自然語言處理的各個領域，如 GPT、BERT 等，都是基於 Transformer 模型。

五 大模型、GPU 和能源

任何類型的大模型都是由複雜構造支援的，包括硬件基礎設施層、軟件基礎設施層、模型MaaS (Mobility as a Service，即「交通行動服務」)層和應用層(圖2)。在這一結構中，GPU就是硬件基礎設施層的核心所在。隨着人工智能時代的到來，人工智能演算法效率已經超越了摩爾定律(Moore's Law)。摩爾定律的內容為：積體電路上可容納的電晶體數目，約每隔兩年便會增加一倍。二十一世紀以來，摩爾定律面臨新的生態：功耗(包括開關功耗)、記憶體極限，以及算力瓶頸等「技術節點」。摩爾定律逼近物理極限，無法迴避量子力學的限制。在其限制下只有三項選擇：延緩摩爾，擴展摩爾，超越摩爾。延緩摩爾定律即突破技術難題，延長該定律的適用時間；擴展摩爾定律即將該定律推廣至諸如量子電腦一類新興計算平台；超越摩爾定律即另闢蹊徑，通過技術組合方案如「芯粒」(chiplet)，實際達到最新的計算能力要求。

圖2 大模型產業的多層結構



圖片來源：筆者繪製。

GPU具有數量眾多的運算單元，採用極簡的流水線進行設計，適合計算密集、易於並行的程式，特別是具備圖形渲染和通用計算的天然優勢。大模型的訓練和推理對GPU提出了更高的要求：更高的計算能力、更大的顯存容量、更快的顯存頻寬、更高效的集群通信能力、低延遲和低成本的推理。GPU可以通過異構計算(heterogenous computing)提供端到端的深度學習資源，縮短訓練所需的環境部署時間。總之，GPU的高性能計算推動了大模型

的發展，大模型不斷對GPU提出迭代要求。例如，微軟 (Microsoft) 為 OpenAI 開發的用於大模型訓練的超級電腦是一個單一系統，伺服器擁有超過 28.5 萬個 CPU 內核、1 萬個 GPU 和 400Gbps 的網絡連接。

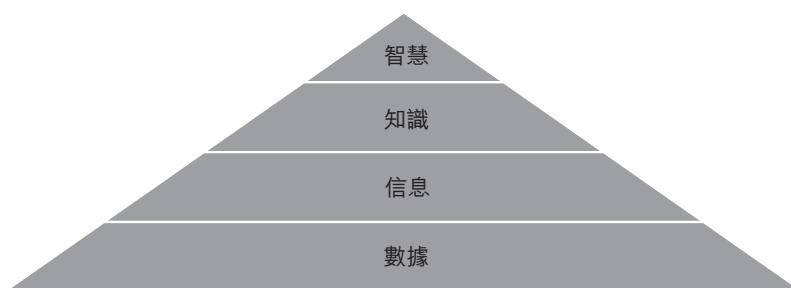
大模型的演變將加速對能源的需求。根據國際數據公司 (IDC) 預測，到 2025 年，全球數據量將達到 175ZB，而且近 90% 的數據都是非結構化的。這些數據需要大量的計算能力才能被分析和處理。同時，隨着人工智能演算法不斷升級和發展，它們的複雜性和計算量也在不斷增加。據估計，目前人工智能的能源消耗佔全球能源消耗約 3%，而據此推斷，到 2025 年，人工智能將消耗 15% 的全球電力供應。除了硬件開發所必須投入的「固定碳成本」以外，對於人工智能日常環境的維護投入也不容小覷。所以，人工智能的快速發展將對能源消耗和環境產生巨大的影響^⑤。

人工智能的快速發展和應用帶來了能源消耗和環境問題，需要在技術和政策上尋求解決方案。在這個過程中，需尋求可持續的能源供應，以減少對傳統能源的依賴，開發在非常低功耗的芯片上運行的高效大模型。

六 大模型和知識革命

一般來說，知識結構類似金字塔，包括了數據、信息、知識和智慧四個層次 (圖 3)。大模型具有極為寬泛的溢出效應。其中最為重要的是引發前所未有的學習革命和知識革命。

圖 3 由數據到智慧的金字塔

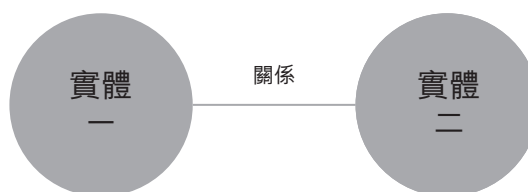


圖片來源：筆者繪製。

基於大數據與 Transformer 的大模型，實現了對知識體系的一系列改變：(1) 改變知識生產的主體。即從人類壟斷知識生成轉變為人工智能生產知識，以及人類和人工智能混合生產知識。(2) 改變知識譜系。從本質上看，知識圖譜是語義網絡的知識庫；從實際應用的角度來看，可以將知識圖譜簡化理解成多關係圖。我們通常用圖裏的節點來代表實體，用連接節點的直線來代表兩個節點之間的關係。實體指的是現實世界的事物，兩點連線表示不同實體之間的某種聯繫 (圖 4)。不同於以往的知識譜系模型，如本體或知識地圖

等，知識圖譜包含大量結構化的實體知識，具備更好的組織、管理和理解互聯網信息的能力，與提升大模型的訓練效果息息相關，表現出大模型時代知識供給的特徵。(3) 改變知識的維度。知識可分為簡單知識和複雜知識、獨有知識和共有知識、具體知識和抽象知識、顯性知識和隱性知識等。二十世紀50年代，世界著名的科學哲學大師波蘭尼 (Michael Polanyi) 發現了知識的隱性維度，而人工智能易於把握知識的隱性維度。(4) 改變知識獲取途徑。大模型正在引領教育革命，人們熟悉的搜尋引擎正在由啟發式的聊天機器人逐步取代。(5) 改變推理和判斷方式。人類的常識基於推理和判斷，而機器常識則是基於邏輯和演算法；人類可以根據自己的經驗和判斷力做出決策，而機器則需要依賴程式和演算法。(6) 改變知識創新方式和加速知識更新速度。不僅知識更新可以通過人工智能實現內容生成，而且大模型具有不斷生成新知識的天然優勢；人類知識處理的範式將發生轉換，人類知識的邊界有機會更快速地擴展。(7) 改變知識處理方式。人類對知識的處理有六個層次：記憶、理解、應用、分析、評價和創造。大模型在這六個層次的知識處理中，都能發揮一定的作用，為人類大腦提供輔助。

圖4 知識圖譜示例



圖片來源：筆者繪製。

簡言之，如果大模型與外部知識源 (例如搜尋引擎) 和工具 (例如程式設計語言) 結合，將豐富知識體系和提高獲取知識的效率。萬物皆可人工智能化，大模型將引發知識革命，形成人類自然智慧和人工智能智慧並存的局面。

知識需要學習。基於赫布理論 (Hebbian theory) 的學習方法被稱為「赫布型學習」。赫布理論又稱「赫布定律」(Hebb's rule)、「赫布假說」(Hebb's postulate)、「細胞結集理論」(cell assembly theory) 等，是一個神經科學理論，由赫布 (Donald O. Hebb) 於1949年提出，描述了在學習過程中大腦的神經元所發生的變化，從而形成記憶印痕^⑥。赫布理論描述了突觸可塑性的基本原理，即突觸前神經元向突觸後神經元的持續重複的刺激，可以導致突觸傳遞效能的增加。以深度學習為核心的大模型的重要特徵就是以人工神經網絡作為基礎。所以，大模型是充分實踐赫布理論的重要工具。

1995年，美國哈佛大學心理學家珀金斯 (David N. Perkins) 提出「真智力」(true intelligence)，並提出智商包括三種主要成份或維度：(1) 神經智力 (neural intelligence)，具有「非用即失」(use it or lose it) 的特點。(2) 經驗智力 (experiential

intelligence)，是指個人積累的不同領域的知識和經驗，豐富的學習環境能夠促進經驗智力。(3) 反省智力 (reflective intelligence)，指一個人使用和操縱其心理技能的能力，類似於元認知 (metacognition，對自己的思維過程的認識和理解) 和認知監視 (cognitive monitoring，指任何旨在評價或調節自己的認知的活動) 等概念；有助於有效地運用神經智力和經驗智力的控制系統²⁹。大模型恰恰具備上述三種主要成份或維度。所以，大模型不僅有智慧，而且是具有高智商的一種新載體。

七 大模型和「人的工具化」

雖然大模型實現智慧的途徑和人類大腦並不一樣，但是最近美國約翰斯·霍普金斯大學 (Johns Hopkins University) 的專家發現，GPT-4 可以利用思維鏈推理和逐步思考，有效證明了其心智理論性能。在一些測試中，人類的水平大概是 87%，而 GPT-4 已經達到 100%。此外，在適當的提示下，所有經過基於人類回饋的強化學習訓練的模型都可以實現超過 80% 的準確率³⁰。如果人工智能互聯網化，或者互聯網人工智能化，無疑會推進智慧革命的積聚和深化。

在現實生活中，大模型的衝擊正在全面顯現。例如，GPT 作為一種基於大規模文本數據的生成式大模型，包括對語言學、符號學、人類學、哲學、心理學、倫理學和教育學等廣義思想文化領域的衝擊，對自然科學技術的全方位衝擊，進而對經濟形態及其運行的衝擊，對社會結構的衝擊，以及對國際關係的衝擊。此外，值得關注的是，人工智能已經開始進入金融領域，與加密數位貨幣結合。2020 年，OpenAI 聯合創始人奧特曼 (Samuel H. Altman) 推出名為「世界幣」(Worldcoin) 的加密貨幣項目，期望通過人工智能技術支援的全球化金融公平與普惠的開源協定，支援私人數位身份和新的金融系統，「賦予人工智能時代的個人權力」。至 2023 年 5 月，超過一百五十萬人加入了加密貨幣錢包 World App 的測試階段，已經在八十多個國家或地區可用。

現在，人類面臨的大模型挑戰，還不僅僅是職場動盪、失去工作、增加失業的問題，而是更為嚴酷的現實課題：人類是否或早或晚會成為大模型的工具人？不僅如此，如果人工智能出現推理能力，在無人知道原因的情況下越過界限，是否會發生人工智能確實傷害甚至消滅人類的潛在威脅？最近網上有這樣的消息：有人利用最新的 AutoGPT 開發出 ChaosGPT，下達毀滅人類指令，人工智能自動搜索核武器數據，並招募其他人工智能輔助³¹。

大模型是人工智能歷史的分水嶺，甚至是工業革命以來人類文明史的分水嶺：在這之前，人們更多關注和討論的是人類如何適應機器，以及和機器人合作，實現艾西莫夫 (Isaac Asimov) 的「機器人三定律」(Three Laws of Robotics)；現在進入如何理解大模型、如何預知人工智能的危險拐點，特別是某些人類

和人工智能合作，反對另外的人類，甚至發生人工智能的徹底失控。人工智能聊天機器人(包括ChatGPT)即使經過數百萬文本源的訓練，可以閱讀並生成「自然語言」文本語言，但是就像人類自然地寫作或交談一樣，不幸的是它們也會犯錯。這些錯誤稱為「幻覺」，或者「幻想」。值得注意的是，因為人工智能幻覺的存在，很可能發生對人類決策和行為的誤導。

正是在這樣的背景下，2023年3月29日，馬斯克(Elon R. Musk)聯名千餘名科技領袖，呼籲暫停開發人工智能，認為這是場危險競賽，讓我們從不斷湧現出具有新能力、不可預測的黑匣子模型中退後一步。據《紐約時報》(*The New York Times*)報導，身在多倫多的圖靈獎得主辛頓在4月向Google提出了請辭。辛頓離職的原因是為了能夠「自由地談論人工智能的風險」；他對自己畢生的工作感到後悔，「我用一個正常的理由安慰自己：如果我沒做，也會有別人這麼做的」。辛頓最大的擔憂是：人類只是智慧演變過程中的一個短暫階段，人工智能很可能比人類更聰明^⑩。未來的人工智能很可能對人類的存在構成威脅，所以停止發展人工智能也許是一個理性的做法，但不可能發生。人們應該合作，阻止人工智能的無序發展^⑪。對比GPT-4剛發布時，辛頓還是何等讚譽有加：「毛蟲吸取了足夠的養分，就能化繭成蝶，GPT-4就是人類的蝴蝶。」^⑫僅僅一個多月的時間，辛頓的立場發生如此逆轉，這不免讓人們想到第二次世界大戰之後，愛因斯坦(Albert Einstein)和奧本海默(Julius R. Oppenheimer)都明確表達了為參與核武器研發和提出建議感到後悔，更為核武器成為冷戰籌碼和政治威脅的工具感到強烈不滿。

事實上，控制論之父維納(Norbert Wiener)在《人有人的用處——控制論和社會》(*The Human Use of Human Beings: Cybernetics and Society*)一書中認為，機器要在所有層面上取代人類，而非只是作為人類的工具提供替代性的力量，因此機器對於人類的影響是深遠的^⑬。霍金(Stephen Hawking)生前也曾多次表達他對人工智能可能導致人類毀滅的擔憂。

遺憾的是，現在世界處於動盪時刻，人類已經自顧不暇，無人知曉人工智能的下一步會發生甚麼。《機械姬》(*Ex Machina*)是一部2015年上映的英國科幻電影，講述主人公受邀鑒定人形機器人是否具備人類心智所引發的故事，其中有這樣的蒼涼台詞：「將來有一天，人工智能回顧我們，就像我們回顧非洲平原的化石一樣，直立猿人住在塵土裏，使用粗糙的語言和工具，最後全部滅絕。」

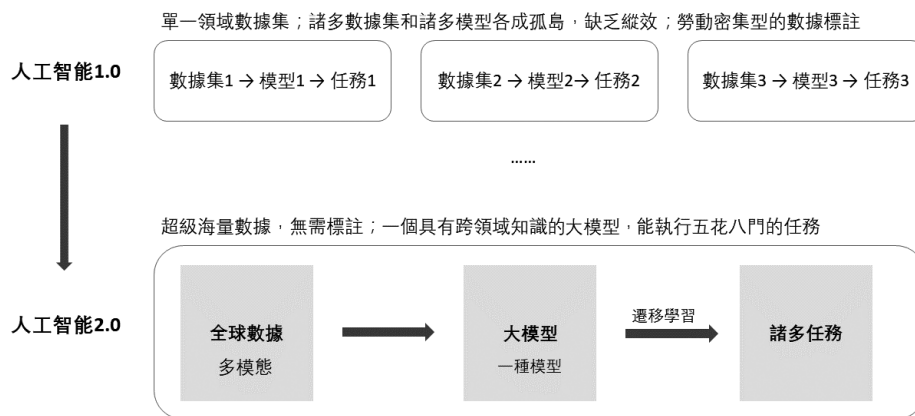
近日有一個消息：來自瑞士洛桑聯邦理工學院(*École polytechnique fédérale de Lausanne*)的研究團隊提出了一種全新的方法，可以用人工智能從大腦信號中提取視頻畫面，邁出了「讀腦術」的第一步。相關論文已於《自然》(*Nature*)雜誌刊登^⑭。據說該文受到很多質疑，但可以肯定的是，不僅愈來愈多的科學家、工程師和企業家，包括天才，還可能有某些陰暗和邪惡力量，正在試圖影響和改變人工智能發展的方向和路徑，增加人們與日俱增的不安。如果說人工智能是人類的又一個潘朵拉盒子，很可能再無人能將其關上。

在人類命運面臨的鉅變趨勢下，人類選擇在減少，然而不可放棄讓人回歸人的價值，需要留下種子——火星遷徙至少具有這樣的超前意識。

八 結語

在人工智能1.0時代，人工智能數據來源是需要人工參與標註並且專注於特定領域的結構化數據；而在人工智能2.0時代，人工智能無需人工干預而能夠處理海量數據，具備跨領域的能力（圖5）。隨着大模型發展，人工智能從1.0時代加速進入2.0時代。

圖5 人工智能1.0和2.0



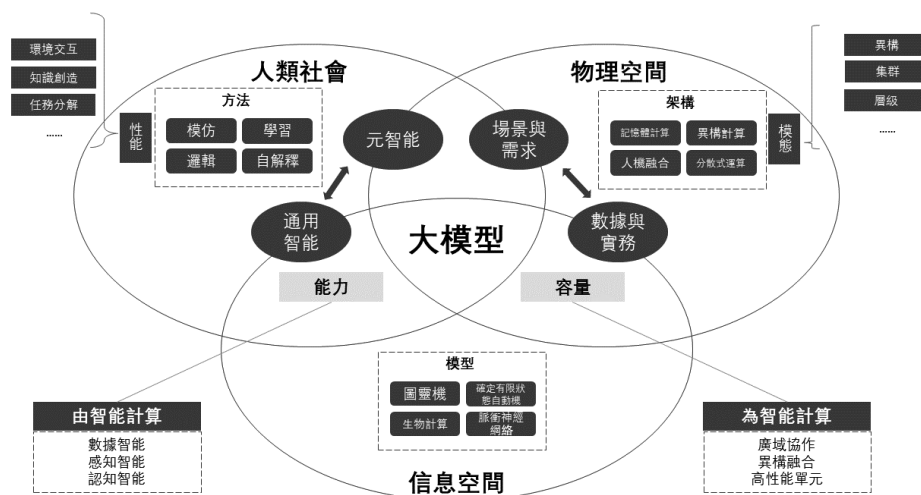
圖片來源：筆者改製自〈創新工場李開復：AI 2.0已至，將誕生新平台並重寫所有應用〉（2023年3月14日），搜狐網，www.sohu.com/a/653951867_114778。

在人工智能2.0時代，大模型分工愈來愈明確，日益增多的大模型，特別是開源大模型會實現不同的組合。支援大模型的數據不僅要求高品質，而且必須開源，任何與開源大模型的競爭必然注定失敗。前述Meta的LLaMA模型所支援的就是開源社區。

可以預見的是，大模型規模的擴大存在着極限：一方面是物理性限制，一方面是大模型存在收益遞減的拐點。所以，大模型設計或架構需考慮如何引入控制論，以適應人類回饋。大模型將樂高(Lego)化，構成大模型集群，不僅推動人類社會、物理空間和信息空間日益緊密融合，而且正在生成一個大模型主導的世界（圖6）。

在這樣的歷史時刻，我們需要重新認識生成主義(enactivism)。生成主義由瓦雷拉(Francisco J. Varela)、湯普森(Evan Thompson)和洛什(Eleanor Rosch)在《具身心智：認知科學和人類經驗》(*The Embodied Mind: Cognitive Science and Human Experience*)中提出，主張心智能力是嵌入在神經和體細胞活動中的，並通過生物的行為而湧現^⑤。論者指出，「生成認知強調，我們所經驗的世界

圖6 人類社會、物理空間、信息空間三重視角下的大模型



圖片來源：筆者改製自 Shiqiang Zhu et al., "Intelligent Computing: The Latest Advances, Challenges, and Future", *Intelligent Computing*, vol. 2 (January 2023), <http://doi.org/10.34133/icomputing.0006>。

是有機體的物理構成、它的感覺運動能力和與環境本身互動的產物。有機體的世界不是一個預先給定的、客觀的、靜待有機體去『經驗』、『表徵』或『反映』的中性世界。相反，世界是通過有機體的行動或動作而生成的」^⑥。人工智能的生成式大模型，確實包括生成主義的要素。人工智能將給生成主義注入新的生命力。

註釋

- ① John McCarthy et al., "A Proposal for the Dartmouth Summer Research Project on Artificial Intelligence, August 31, 1955", *AI Magazine* 27, no. 4 (2006): 12-14.
- ② 對 Jukebox 的介紹，參見 <https://openai.com/research/jukebox>。
- ③ 2018年6月，OpenAI 發布 GPT-1，模型參數數量為 1.17 億；2019年2月，發布 GPT-2，模型參數數量為 15 億；2020年5月，發布 GPT-3，參數數量為 1,750 億；2022年11月，正式推出了對話互動式的聊天機器人 ChatGPT；2023年3月，正式推出 GPT-4，成為目前較先進的多模態大模型。GPT-4 主要在識別理解能力、創作寫作能力、處理文本量以及自訂身份屬性迭代方面取得進展。
- ④ CLIP (Contrastive Language-Image Pre-Training) 模型是 OpenAI 在 2021 年初發布的用於匹配圖像和文本的預訓練神經網絡模型，可以說是近年來在多模態模型研究領域的經典之作。該模型直接使用大量的互聯網數據進行預訓練，在很多工作表現上達到了目前最高水平。關於預訓練，下文會有較詳細的討論。
- ⑤ DALL·E 是一個可以根據書面文字生成圖像的人工智能系統，該名稱來源於著名畫家達利 (Salvador Dalí) 和電影《機器人總動員》(Wall·E, 2008)。
- ⑥ 參見 "Introducing ChatGPT", <https://openai.com/blog/chatgpt>。
- ⑦ 強化學習 (reinforcement learning) 是機器學習 (machine learning) 的範式和方法論之一，用於描述和解決智能體 (agent) 在與環境的交互過程中通過學習策略以達成回報最大化或實現特定目標的問題。
- ⑧ 中文將 "Transformer" 翻譯為變換器，並不能完全反映大模型的 Transformer 的基本內涵。所以，本文還是直接使用英文原詞。

- ⑨ 生成性模型從一個形式語言系統出發，生成語言的某一集合。代表是喬姆斯基 (Avram N. Chomsky) 的形式語言理論和轉換語法。分析性模型從語言的某一集合開始，根據對這個集中各個元素的性質的分析，闡明這些元素之間的關係，並在此基礎上上演繹的方法建立語言的規則系統。代表是前蘇聯數學家庫拉金娜 (O. S. Kulagina) 和羅馬尼亞數學家馬爾庫斯 (Solomon Marcus) 用集合論方法提出的語言模型。在生成性模型和分析性模型的基礎上，將二者結合起來，產生了一種很有實用價值的模型，即辨識性模型。辨識性模型可以從語言元素的某一集合及規則系統出發，通過有限步驟的運算，確定語言中合格的句子。代表是巴爾-希列爾 (Yehoshua Bar-Hillel) 用數理邏輯方法提出的句法類型演算模型。
- ⑩ 自監督學習是一種機器學習範式和相應的方法，用於處理未標註的數據，以獲得有助於下游學習任務的有用表示。
- ⑪ Google 推出的 LaMDA (Language Model for Dialogue Applications) 是自然語言處理領域的一項新的研究突破。它是一個面向對話的神經網絡架構，可以就無休止的主題進行自由流動的對話。它的開發是為了克服傳統聊天機器人的局限性，後者在對話中往往遵循狹窄的、預定義的路徑。BERT (Bidirectional Encoder Representation from Transformers) 是一個預訓練的語言表徵模型。它強調了不再像以往一樣採用傳統的單向語言模型或者把兩個單向語言模型進行淺層拼接的方法進行預訓練，而是採用新的「屏蔽語言模型」(Masked Language Model, MLM)，以致能生成深度的雙向語言表徵。關於 BERT 的論文發表時，提及 BERT 在十一個自然語言處理任務中獲得了新的目前最高水平的結果 PaLM-E，參數數量高達 5,620 億 (GPT-3 的參數數量為 1,750 億)，同時集成語言和視覺，用於機器人控制。相比大語言模型，它被稱為視覺語言模型 (Visual Language Model, VLM)。兩者不同之處，在於後者對物理世界是有感知的。
- ⑫ LLaMa (Large Language Model Meta AI) 有多個不同大小的版本。該模型主要從維基百科、書籍，以及來自 arXiv、GitHub、Stack Exchange 和其他網站的學術論文中收集的數據集上進行訓練。LLaMa 模型支援二十種語言，包括拉丁語和西里爾字母語言，目前看原始模型並不支援中文。2023 年 3 月，LLaMa 模型發生洩露，意外促成了大批 ChatGPT 式服務的產生。OPT-175B 模型的參數數量超過 1,750 億，和 GPT-3 相當。OPT 是“Open Pre-trained Transformer”的縮寫。它的優勢在於完全免費，這使得更多缺乏相關經費的科學家可以使用這個模型。同時，Meta 還公布了代碼庫。
- ⑬ 參見 Nick Bostrom, *Superintelligence: Paths, Dangers, Strategies* (Oxford: Oxford University Press, 2014)。
- ⑭ Melinda Bognár, “Prospects of AI in Architecture: Symbolicism, Connectionism, Actionism” (10 January 2021), *Journal of Architectural Informatics Society*, <https://openreview.net/pdf?id=gVHffM4DlpG>.
- ⑮ A. M. Turing, “Computing Machinery and Intelligence”, *Mind* LIX, issue 236 (1950): 433-60.
- ⑯ A. L. Samuel, “Some Studies in Machine Learning Using the Game of Checkers”, *IBM Journal of Research and Development* 44, no. 1/2 (2000): 206-26.
- ⑰ Ryszard S. Michalski, Jaime G. Carbonell, and Tom M. Mitchell, eds., *Machine Learning: An Artificial Intelligence Approach*, vol. II (Los Altos, CA: Morgan Kaufmann, 1986).
- ⑱ 這一階段代表性的工作有莫斯托 (Jack Mostow) 的指導式學習、萊納特 (Douglas B. Lenat) 的數學概念發現程式、蘭利 (Pat Langley) 的 BACON 程式及其改進程式。
- ⑲ Geoffrey E. Hinton, Simon Osindero, and Yee-Whye Teh, “A Fast Learning Algorithm for Deep Belief Nets”, *Neural Computation* 18, no. 7 (2006): 1527-54.
- ⑳ 在文本資料中，包括有標註數據和無標註數據，這是所謂數據驅動。

- ⑳ Jason Wei et al., “Emergent Abilities of Large Language Models”, *Transactions on Machine Learning Research* (August 2022), <https://openreview.net/forum?id=yzkSU5zdwD>.
- ㉑ 貝葉斯原理是用貝葉斯風險 (Bayes Risk) 表示的最優決策律；馬爾可夫鏈描述的是概率論和數理統計中的離散的指數集 (index set) 和狀態空間 (state space) 內的隨機過程 (stochastic process)；N 元模型是大詞彙連續語音辨識中常用的一種語言模型。
- ㉒ Cage：〈Pinecone：大模型引發爆發增長的向量數據庫，AI Agent 的海馬體〉(2023 年 4 月 26 日)，「海外獨角獸」微信公眾號，https://mp.weixin.qq.com/s?__biz=Mzg2OTY0MDk0NQ==&mid=2247501819&idx=1&sn=2fcee248cf2b9703804e6d9a45dc4c97。
- ㉓ Ashish Vaswani et al., “Attention Is All You Need” (6 December 2017), arXiv, <https://doi.org/10.48550/arXiv.1706.03762>.
- ㉔ 〈人類已達硅計算架構上限！預計 2030 年，AI 會消耗全球電力供應的 50%〉(2023 年 3 月 27 日)，「新智元」微信公眾號，<https://mp.weixin.qq.com/s/k9A8d2gX14xyE5cSBI-Dpw>；王鵬：〈雙碳視角下人工智能發展再思考——投產相抵還是能耗脅迫？〉(2023 年 4 月 9 日)，《中國日報》網，<https://column.chinadaily.com.cn/a/202304/09/WS6432bc56a3102ada8b2376fa.html>。
- ㉕ Donald O. Hebb, *The Organization of Behavior: A Neuropsychological Theory* (Mahwah, NJ: Psychology Press, 2002).
- ㉖ David N. Perkins, *Outsmarting IQ: The Emerging Science of Learnable Intelligence* (New York: Free Press, 1995).
- ㉗ 〈100:87：GPT-4 心智碾壓人類！三大 GPT-3.5 變種難敵〉(2023 年 5 月 1 日)，「新智元」微信公眾號，https://mp.weixin.qq.com/s?__biz=Mzl3MTA0MTk1MA==&mid=2652326060&idx=1&sn=c0ffa5d76ee8af079a2dbbe4b37bb15f。
- ㉘ 〈有人給了 AI「毀滅人類」的任務，讓它持續自主運行，它開始研究最強核武器〉(2023 年 4 月 9 日)，騰訊網，<https://new.qq.com/rain/a/20230409A07GZ100>。
- ㉙ Cade Metz, “‘The Godfather of A.I.’ Leaves Google and Warns of Danger Ahead”, *The New York Times*, 1 May 2023.
- ㉚ 部夢凡整理：〈人工智能教父：也許還有希望限制 AI 的無序發展〉(2023 年 5 月 13 日)，新浪財經網，<https://finance.sina.com.cn/stock/hyyj/2023-05-13/doc-imytrira2549050.shtml>.
- ㉛ Founder Park：〈「AI 教父」離職谷歌：對畢生工作感到後悔和恐懼〉(2023 年 5 月 2 日)，「極客公園」微信公眾號，https://mp.weixin.qq.com/s?__biz=MTMwNDMwODQ0MQ==&mid=2652991276&idx=1&sn=950b2e53feae3ceb7c7c03b39eb4a1a9。
- ㉜ 維納 (Norbert Wiener) 著，陳步譯：《人有人的用處——控制論和社會》(北京：商務印書館，1978)。
- ㉝ Sara Reardon, “Mind-reading Machines Are Here: Is It Time to Worry?”, *Nature* 617, no. 7960 (2023): 236.
- ㉞ 參見 Francisco J. Varela, Evan Thompson, and Eleanor Rosch, *The Embodied Mind: Cognitive Science and Human Experience* (Cambridge, MA: MIT Press, 1991)。
- ㉟ 葉浩生、曾紅、楊文登：〈生成認知：理論基礎與實踐走向〉，《心理學報》，2019 年第 11 期，頁 1270-80。

朱嘉明 經濟學博士、教授，曾任職於聯合國工業發展組織 (UNIDO)，並先後任教於維也納大學和台灣大學。現任「數字資產研究院」學術與技術委員會主席，中國投資協會數字資產研究中心專家組組長。