第1章

大模型基础

在人工智能浪潮席卷全球、技术革新日新月异的当下,大语言模型(Large Language Model, LLM)以其强大的涌现能力,正以前所未有的深度和广度重塑着自然语言处理、智能交互乃至整个信息科技领域的格局,成为驱动这场深刻变革的核心引擎。作为本书的开篇,本章旨在从纷繁的技术图景中溯本清源,深入剖析大模型的底层逻辑与技术根基。我们将系统梳理语言模型从早期统计方法到神经网络的演进脉络,全景式回溯大模型波澜壮阔的发展历程,并深刻阐释其区别于传统模型的革命性特点。通过对这些基础性知识的透彻理解,本章将为读者后续系统性地学习大模型的应用开发、架构设计及优化实践,奠定坚实而稳固的理论与实践基石。

1.1 语言模型基础

语言模型(Language Model,LM)是自然语言处理(Natural Language Processing,NLP)众多任务中不可或缺的基础支撑与核心引擎。其核心目标在于精确刻画人类语言的内在规律,即对任意给定词序列(或字符序列)的概率分布进行数学建模。通过评估序列的可能性并预测下一个最可能出现的语言单元(词、子词或字符),语言模型为机器理解、生成人类语言提供了根本性的能力保障。

纵观其发展历程,语言模型的技术演进堪称一场深刻的范式变革: 从早期依赖人工特征与统计概率的传统方法,到引入循环神经网络(Recurrent Neural Network, RNN)捕捉序列动态依赖的深度学习初期探索,再到以 Transformer 架构为代表、凭借自注意力(Self-Attention)机制彻底革新长程依赖建模能力的大模型时代奠基技术。每一次重大的技术跃迁,都显著提升了语言模型的表达能力和应用效果,为自然语言处理领域带来了里程碑式的突破,并最终铺就了通向当今大模型辉煌成就的道路。

本节将系统性地拆解语言模型的技术根基: 剖析基于统计方法的语言模型(如 N-gram 模

型)的核心思想、优势及其固有局限;探讨基于循环神经网络的语言模型如何利用隐状态传递信息,初步解决长距离依赖问题,并分析其面临的主要挑战(如梯度消失、梯度爆炸);聚焦革命性的基于 Transformer 的语言模型,深入解析其自注意力机制的核心原理,揭示其如何克服 RNN 的缺陷,实现高效并行计算与强大的上下文建模能力,从而为现代大语言模型奠定无可撼动的架构基础。

1.1.1 基于统计方法的语言模型

语言模型的早期探索,建立在概率论与统计学的坚实根基之上。这类模型的核心目标是通过 计算词序列的联合概率,量化语言单元的生成可能性。其基本思想遵循概率链式法则(Chain Rule):

$$P(w_1, w_2, \dots, w_t) = P(w_1) \cdot P(w_2 \mid w_1) \cdot P(w_3 \mid w_1, w_2) \cdots P(w_t \mid w_1, \dots, w_{t-1})$$
(1-1)

其中, w_t 表示第t个词, $P(w_t|w_1,\cdots,w_{t-1})$ 代表给定历史上下文后当前词的条件概率。

1. 关键技术: N-gram 统计语言模型

为简化计算,统计语言模型引入马尔可夫假设(Markov Assumption): 当前词的概率仅依赖于其前 n-1 个词,而非全部历史,即:

$$P(w_t|w_1, w_2, \dots, w_{t-1}) = P(w_t|w_{t-n+1}, \dots, w_{t-1})$$
(1-2)

由此诞生了 N-gram 模型——自然语言处理史上首个被广泛应用的实用化语言模型。

1) 概率估计

式(1-2)中的参数 n 称为模型的阶数,其值决定了模型的精度和复杂性^[3],其值越大,则模型对单词之间的描述越准确,模型精度越高,但复杂性也随之提升。因此,如何权衡模型的精度和复杂度,就需要通过 n 值的选择来实现。一般而言,n 的取值在 $1\sim7$,特别地,当取值为 $1\sim2\sim3$ 时,分别称为 Unigram、Bi-gram 及 Tri-gram。

从形式语言理论的角度来看,N-gram 模型描述的语言本质上是一种由有限状态的正则文法产生的语言。正则文法遵循严格的规则,每个状态的转移仅依赖于当前有限的输入信息,这使得其生成的语言序列具有固定的模式和结构。以二元语法(Bi-gram)为例,它在预测下一个词时,仅仅依据前一个词的信息,通过统计语料库中前后词的共现频率来构建语言模型。例如,在"我吃苹果"这个句子序列中,Bi-gram 模型会统计"我一吃""吃一苹果"这样的二元组出现的概率。但自然语言是极其复杂且灵活多变的,它不仅包含语法规则,还蕴含着丰富的语义、语境信息以及文化背景知识。一个简单的句子在不同的语境下可能会有截然不同的含义,而且人们在表达时也常常打破常规的语法结构,使用隐喻、倒装等多样化的表达方式,这与 N-gram模型遵循的有限状态规则形成了鲜明对比。

然而,尽管存在这些明显差别,N-gram 模型在实际应用中却收获了巨大成功。首要原因在于它精准地捕捉到了自然语言中存在的局部约束性质。在日常语言表达中,词语之间存在着紧密的局部关联,相邻词语的组合往往具有一定的规律性。例如在"喝咖啡""吃面包"这样的

常见搭配中,"喝"和"咖啡"、"吃"和"面包"之间形成了强关联,N-gram 模型通过大量 统计这些局部组合,能够有效地对文本进行概率估计和语言建模。

N-gram 的概率通过语料库中的频率计数直接估计:

$$P(w_t|w_{t-n+1},...,w_{t-1}) \approx \frac{\text{Count}(w_{t-n+1},...,w_{t-1})}{\text{Count}(w_{t-n+1},...,w_t)}$$
(1-3)

例如,在三元模型(Tri-gram)中:

$$P(天气 | 今天,) = \frac{"今天, 天气"出现次数}{"今天"出现次数}$$
 (1-4)

其次, N-gram 模型结构简单, 计算复杂度较低, 这使得它在计算资源有限的早期阶段, 能 够快速实现并应用于语音识别、机器翻译、拼写检查等多个领域。再者,该模型的可解释性强, 其基于统计频率的原理使得人们能够直观地理解模型如何对语言进行建模和预测,方便技术人 员进行模型的调试和优化。正是这些特性,使得 N-gram 模型在自然语言处理发展历程中占据 了重要地位,为后续更复杂的语言模型发展奠定了坚实基础。

然而,在 N-gram 模型的应用中,数据稀疏问题是其面临的一大挑战。该问题指的是在模 型训练过程中,某些 N-gram 在学习语料集中从未出现,但却可能出现在测试语料集中。这会 导致基于最大似然法的模型在处理此类 N-gram 时,错误地将其概率值判定为 0,影响模型预测 的准确性。

以实际实验为例,在针对一个包含 242 000 000 个单词的语料库的研究中,研究人员采用最 大似然法构建了一个基于 60 000 个单词的 Tri-gram 模型。当使用该模型对实际测试语料集进行 处理时,发现测试集中仅有 69%的 N-gram 在学习集中出现的次数大于 1,这意味着超过三成的 N-gram 在训练阶段是"未见"的。更值得注意的是,N-gram 模型的阶数 n 越大,数据稀疏问 题就越严重。因为随着阶数升高,N-gram 的组合可能性呈指数级增长,使得更多的序列难以在 有限的训练数据中被覆盖到。

面对数据稀疏问题,简单地扩大训练语料集规模并不能从根本上解决问题。因为无论语料 库多大,都无法保证覆盖所有可能的 N-gram 组合。因此,如何采用有效的方法,为那些未在 学习语料中出现的 N-gram 合理估计一个非零概率值,成为 N-gram 语言模型研究的关键方向和 核心课题。

2) 平滑技术

在 N-gram 语言模型中,数据稀疏问题严重影响模型性能,为了有效解决这一难题,人们提出 了一系列处理技术,并将其统称为平滑化(Smoothing)方法。这些方法旨在为那些在训练语料中 未出现或出现频率极低的 N-gram 赋予合理的非零概率值,使模型能更准确地处理实际应用中的各 种语言现象。平滑化方法主要可划分为两大类。一类是对最大似然法的估计结果进行直接修整,包 括插值法、折扣法以及回退法。另一类平滑化方法则是通过对单词进行聚类,缩小模型空间来解决 数据稀疏问题。

(1) 插值法: 插值法是指通过插值技术,将一个 N-gram 模型表示为由 1 阶到 n 阶的线性

组合, 即:

$$P_{\text{int}}(w_1, \dots, w_n) = \alpha_1 P(w_1, \dots, w_n) + \alpha_2 P(w_2, \dots, w_n) + \dots + \alpha_n P(w_n)$$

$$(1-5)$$

其中,参数 α_i 满足 $\sum_{i=1}^n \alpha_i = 1$ 。关于参数值的计算,分为理论计算和实验调整两种方式:①理论计算是指在测试语料集上计算模型的分支均值^[2],并取使该值最小的参数值;②实验调整是指对于具体的应用系统,可以通过对测试集的反复测试,确定使得模型误差最小的参数值。

(2) 折扣法: Good-Turing 估计法是许多折扣平滑方法的核心。它通过对最大似然法的结果进行调整,可以在保证满足概率归一性质的条件下,估计出在训练语料中没有出现的 N-gram 的概率值。

首先,该方法使用 N_c 表示 N 个样本中出现 c 次的 N-gram 的数量。然后对于任何出现了 c 次的 N-gram,都假设其出现了 c*次:

$$c^* = (c+1)\frac{N_{c+1}}{N_c} \tag{1-6}$$

假设 M-gram 出现了 $c(w_1, \dots, w_m)$ 次,Good-Turing 给出其出现的频率为:

$$P_{gt}(w_1, \dots, w_m) = P_{c(w_1, \dots, w_m)} = \frac{c^*(w_1, \dots, w_m)}{N}$$
 (1-7)

则,对于c=0的样本,有:

$$P_0 = 1 - \sum_{c > 0} N_c * P_c = N_1 / N \tag{1-8}$$

使用 c^* 代替c的过程称为折扣,比值 c^*/c 称为折扣因子。

(3)回退法: Katz 的回退法对 Good-Turing 估计进行了扩展,它将每一个 N-gram 模型表示为 M-gram 的非线性组合。对于每一个 M-gram,由一个回退概率 β_m 表示由 M-gram 回退到(M-1)-gram 的概率。由此存在以下递推公式:

$$P_{k}(w_{n} \mid w_{1}, w_{2}, \dots, w_{n-1}) = P_{gt}(w_{n} \mid w_{1}, w_{2}, \dots, w_{n-1}) + \beta_{n} P_{k}(w_{n} \mid w_{2}, w_{3}, \dots, w_{n-1})$$

$$P_{k}(w_{n} \mid w_{2}, w_{3}, \dots, w_{n-1}) = P_{gt}(w_{n} \mid w_{2}, w_{3}, \dots, w_{n-1}) + \beta_{n-1} P_{k}(w_{n} \mid w_{3}, w_{4}, \dots, w_{n-1})$$
...
$$(1-9)$$

$$\beta_m = \nu(w_1, w_2, \dots, w_m) / (1 - \omega(w_1, w_2, \dots, w_m))$$
 (1-10)

其中:

$$\begin{split} & \mathcal{O}(w_1, w_2, \cdots, w_m) = \sum_{c(w_1, w_2, \cdots, w_m) > 0} P_{gt}(w_m \mid w_1, w_2, \cdots, w_{m-1}) \\ & \omega(w_1, w_2, \cdots, w_m) = \sum_{c(w_1, w_2, \cdots, w_m) > 0} P_{gt}(w_m \mid w_2, w_3, \cdots, w_{m-1}) \end{split}$$

(4) 聚类法: 不同的聚类算法会依据单词的语义、语法、上下文等特征,将具有相似性质

的单词归为同一类别。这样一来,在计算 N-gram 概率时,同一类别的单词可以共享统计信息, 原本因单个单词数据不足导致的稀疏问题得到缓解。例如,在处理"苹果""香蕉""橘子" 等表示水果的单词时,将它们聚类后,在计算相关 N-gram 概率时,这些单词的统计信息可以 相互补充,提升模型对未见过的语言序列的处理能力。

2. 核心局限性

尽管 N-gram 模型为早期 NLP 奠定了重要基础, 其固有的缺陷随任务复杂度提升而日益凸 显^[5]。自然语言本质上并非有限状态语言,这一特性决定了自然语言语句中的符号串无法简单 用马尔可夫链描述。在自然语言中,某个符号的出现概率并非单纯取决于前一个或前 N 个符号。 例如,在"小明觉得,老师表扬的那个同学,虽然平时沉默寡言,但考试成绩总是很好,他应 该多向对方学习"这句话中,"他"的指代对象需要综合整句话的语义、语法结构以及上下文 信息来判断,而非仅依赖前几个单词。理论上,难以确切界定当前符号的出现到底由其前多少 个符号决定。

然而,统计语言模型为了实现对语言的建模,不得不引入概率论上的独立性假设,即假定 N+1 个符号出现的概率仅与前 N 个符号相关,与语句中其他符号无关。N-gram 模型通过统计前 N 个符号组合出现的频率来预测下一个符号。这种假设虽然为统计模型的构建与计算提供了可 能——避免了因自由参数过多导致的计算指数爆炸,同时也在一定程度上缓解了训练数据稀疏 的难题(在实际应用中, N 通常需控制在 3 以下,以保证模型的可实施性),但也使模型与真 实语言现象之间产生了偏差。

独立性假设如同双刃剑,在赋予统计模型可操作性的同时,也极大地简化了语言的复杂性。 它使得统计模型更擅长处理对结构关系依赖较弱的任务,如基础的文本生成、简单的词性标注 等。但面对具有复杂结构依赖的语言任务时,如确定代词的先行词、分析长距离依存关系等, 统计模型往往力不从心。由于大部分语言学知识和语法规则都具有结构依赖特性,这使得独立 性假设在许多实际语言处理场景中难以成立,限制了统计模型在复杂语言任务中的应用效果。

N-gram 模型是大数据驱动范式的首次成功实践,证明了从海量文本中学习语言规律的可行 性。尽管其已被神经网络取代,但其概率框架与评估方法(如困惑度 Perplexity)仍是现代语言 模型的底层基础。N-gram 模型标志着语言处理从规则系统转向数据驱动的关键转折。然而,其 对结构严重依赖的特性,直接催生了能够动态捕捉序列状态的新一代模型——基于循环神经网 络(RNN)的语言模型,将在1.1.2节深入探讨。

1.1.2 基干循环神经网络的语言模型

基于循环神经网络(RNN)的语言模型曾开启了神经网络处理序列数据的新篇章。传统的 统计语言模型虽能捕捉语言的概率分布,但在处理长距离依赖和动态语义时存在显著局限性。 RNN 的出现,通过引入循环结构,赋予模型记忆能力,使其在理论上能够处理任意长度的序列。 数据,为自然语言处理领域带来突破性进展。这一架构不仅革新了机器对语言的理解与生成方 式,还为后续深度学习模型的发展奠定了重要基础。

1. 基于 RNN 的语言模型基础架构

RNN 的基础架构打破了传统前馈神经网络对数据独立同分布的固有约束, 开创了处理序列 数据的新范式。相较于前馈网络仅依赖当前输入进行计算,RNN 通过独特的循环连接机制,实 现了对时间序列信息的记忆与累积。在其网络结构中,隐层状态作为核心枢纽,承载着历史信 息的传递任务。每个时间步下,隐层状态都会依据当前输入与上一时刻的隐层状态进行更新, 这种递归式的计算方式使得 RNN 能够有效捕捉序列数据中的上下文依赖关系[4][5]。

待处理的序列通常为时间序列,此时序列的演讲方向被称为"时间步(time-step)"。具 体而言, RNN 在每一个时间步 t, 接收当前输入和前一时刻的隐状态 h_{cl} , 计算当前隐状态 h_{t} 和输出 ٧,:

$$h_{t} = f(\mathbf{W}_{h}h_{t-1} + \mathbf{W}_{x}x_{t} + b_{h})$$
 (1-11)

其中, W_b 和 W_c 为权重矩阵, b_b 为偏置项, $f(\cdot)$ 为激活函数(如 ReLU)。通过以上计算, RNN 可以对当前时刻的输出产生依赖上一时刻的信息,实现对序列信息的建模。此外,RNN 的输出节点为一个线性函数:

$$y_t = g(\mathbf{W}_v h_t + b_v) \tag{1-12}$$

其中, W_v 为权重矩阵, b_v 为偏置项, $g(\cdot)$ 为激活函数(如 ReLU)。

RNN 的基本结构主要由以下几个部分组成:

- 输入层(Input Layer):接收序列数据,通常表示为一系列的时间步输入 x_1, x_2, \dots x_T , 其中 T 是序列的长度。在每个时间步 t, RNN 从输入层接收一个输入向量 x_t , 作 为该时刻的信息。
- 隐藏层 (Hidden Layer): RNN 的核心层用于存储和传递序列中的上下文信息。隐层 中的隐状态 (Hidden State)是关键,它在每个时间步更新,逐步将前一时刻的信息 $h_{r,l}$ 与当前输入x,结合,生成新的隐状态h,(式(1-11))。这种循环结构让网络在每个 时间步保留上下文信息, 使得序列信息得以存储和传播。
- 输出层(Output Layer):在每个时间步产生一个输出 y,, 用于预测或进一步处理。输 出可以在每个时间步产生(如序列到序列的生成任务),也可以仅在最终时间步产生 (如序列到单值的分类任务)。

RNN 之所以能够处理序列数据,核心在于其循环连接结构,即网络在每一个时间步 (Time Step) 不仅处理当前输入,还接收来自前一时间步的隐状态信息。这种结构在时间维度上形成 依赖链,使得网络能够保留历史上下文,从而对序列中前后元素的关系进行建模。

具体而言,循环连接的方式主要包括以下几种。

1)循环单元-循环单元连接(Hidden-Hidden Connections)

如图 1.1 所示,这是 RNN 中最典型的连接方式,又称为全连接循环结构。在这种结构中,每 个时间步的隐藏状态 h_t 是由当前输入 x_t 和前一时间步的隐藏状态 h_{t-1} 共同决定的:

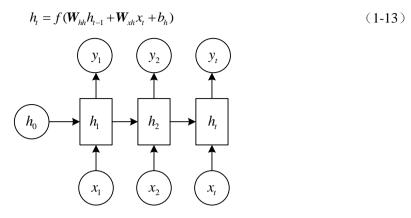


图 1.1 全连接的循环单元

其中, W_{hh} 为状态-状态权重矩阵, W_{ch} 为输入-状态权重矩阵, $f(\cdot)$ 为激活函数。该结构可以完 整地保留历史状态的信息,具备图灵完备性,学习能力强,是标准 RNN 架构的基础。

如图 1.2 所示,通过在时间正向和反向分别堆叠循环单元,可以构建出双向循环神经网络 (Bidirectional RNN, BRNN),从而同时利用过去和未来的上下文信息。

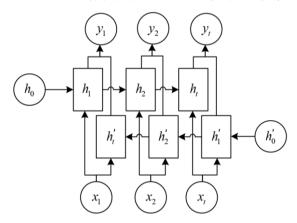


图 1.2 双向循环神经网络

2)输出节点-循环单元连接(Output-to-Hidden Connections)

在这一变种中,当前时间步的隐藏状态由当前输入 x_t 和前一时间步的输出 y_{t1} 决定,而不是前 一状态:

$$h_{t} = f(\mathbf{W}_{hh}y_{t-1} + \mathbf{W}_{xh}x_{t} + b_{h})$$
 (1-14)

这种结构假设输出 yt1 能够代表先前的状态信息。虽然该连接不具备图灵完备性,理论表 达能力较弱,但由于其简化的依赖路径,可以使用教师强制(Teacher Forcing)等技术进行高 效训练,常用于条件生成任务[7]。

3) 基于上下文的连接(Closed-Loop/Context-Based Connections) 该类型连接方式借鉴图神经网络的思路,引入上一个时间步的真实标签 y_{t-1}^* 或目标值 y_{t-1} 来 指导当前状态的更新:

$$h_{t} = f(\mathbf{W}_{hh}h_{t-1} + \mathbf{W}_{vh}y_{t-1}^{*} + \mathbf{W}_{xh}x_{t} + b_{h})$$
 (1-15)

由于在训练过程中引入了目标序列的真实值,这类结构本质上属于生成模型(Generative Model),能够逼近目标序列的真实分布^[8]。这种结构在序列到序列(Seq2Seq)学习中非常常见,特别适用于语言生成和翻译任务^[9]。

2. RNN 的关键变体: LSTM

普通 RNN 在长序列数据处理中会出现梯度消失问题,LSTM(Long Short-Term Memory,长短期记忆网络)是一种特殊的循环神经网络,它通过引入记忆单元(Memory Cell)来解决上述问题。LSTM 可以在长时间的序列中捕捉依赖关系,是一种非常适合处理时间序列、自然语言处理、语音识别等任务的深度学习模型。

相比于传统的 RNN,LSTM 网络模型引入了三个门控单元,分别是输入门(Input Gate)、遗忘门(Forget Gate)和输出门(Output Gate),从而实现了对信息的选择性记忆。其中 σ 和 tanh 分别表示 sigmoid 激活函数和 tanh 激活函数,如图 1.3 所示。

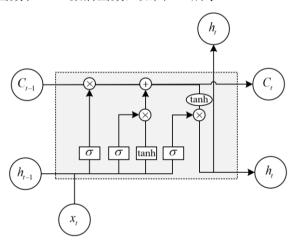


图 1.3 长短时记忆网络循环单元

1) 输入门 (Input Gate)

$$i_{t} = \sigma(W_{x_{t}}x_{t} + W_{bi}h_{t-1} + b_{i})$$
 (1-16)

$$\tilde{C}_t = \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c)$$
 (1-17)

其中, i_t 为输入门的输出,决定当前输入信息的重要性,控制信息的重要程度; \tilde{C}_t 为候选记忆单元状态,用于记忆新输入信息; W_{vi} 、 W_{vi} 、 W_{vv} 、 W_{bv} 为对应连接的权重矩阵。

2) 遗忘门 (Forget Gate)

$$f_{t} = \sigma(W_{xf}x_{t} + W_{hf}h_{t-1} + b_{f})$$
 (1-18)

其中, f_t 为遗忘门的输出,决定过去信息的保留程度; W_{xf} 、 W_{hf} 均为对应连接的权重矩 阵。

3) 记忆单元 (Memory Cell)

$$C_{t} = f_{t} \otimes C_{t-1} + i_{t} \otimes \tilde{C}_{t} \tag{1-19}$$

其中, C_t 为当前时刻的记忆单元; \otimes 为 Hadamard 乘积。该模块用于确定历史记忆的遗 忘程度和新记忆的保留程度,整合为新的记忆。

4) 输出门(Output Gate)

$$o_{t} = \sigma(W_{vo}x_{t} + W_{ho}h_{t-1} + b_{o})$$
 (1-20)

$$h_{t} = o_{t} \otimes \tanh(C_{t}) \tag{1-21}$$

其中, o_t 控制 LSTM 单元的输出信息量; W_{xo} 、 W_{bo} 为对应连接的权重矩阵。

通过上述结构设计,LSTM 能够记住重要的历史信息,忽略无关的过去数据,从而有效处 理长时间依赖信息。但同时也因为过多参数的引入,提高了计算成本需求,同时对短期数据表 现不佳, 劣于常规的回归方法。

3. 基于 RNN 的语言模型的局限性分析

基于统计的语言模型受限于固定窗口依赖与浅层统计逻辑,这类模型在处理长距离语义关 联、复杂语义理解时逐渐显露瓶颈。与之相比,以 RNN 及 LSTM 为代表的循环神经网络语言 模型,凭借隐层状态循环传递机制与门控结构,在上下文语义整合与动态语义表征方面实现突 破,显著提升了模型对长文本依赖关系的捕捉能力和泛化适应性,尤其在机器翻译、文本生成 等任务中展现出统计模型难以比拟的优势[10]。

然而,基于 RNN 的语言模型面临以下局限性。

1) 计算复杂度与训练难度

RNN 和 LSTM 的时序依赖特性导致其难以充分利用现代硬件的并行计算能力。在处理长 序列时,每个时间步的隐状态计算必须依赖前一时间步的结果,形成串行计算链。例如,在训 练包含数千个词的文档时,这种逐词处理的方式会显著延长训练周期^[11]。尽管 LSTM 通过门控 机制缓解了梯度消失问题,但其复杂的门控结构(包含遗忘门、输入门、输出门等多个非线性 变换)增加了参数数量和计算开销。据实验统计,在相同语料库上训练 LSTM 模型的时间成本 通常是 N-gram 模型的数十倍,硬件资源消耗也呈指数级增长。此外,模型对超参数(如学习 率、批大小)的选择更为敏感,需要更精细的调优过程。

2) 可解释性与决策透明度缺失

神经网络模型的"黑箱"特性在 RNN/LSTM 中尤为突出。传统统计模型(如隐马尔可夫 模型)的概率转移矩阵和状态转换图提供了明确的语义解释框架,而 RNN/LSTM 的预测过程 依赖数百万参数的非线性交互,难以通过直观方式解读。例如,在医疗诊断辅助系统中[11],医 生需要明确了解模型作出某种预测的依据,但RNN/LSTM 无法提供如"因为检测到X 症状和

Y 指标,所以预测为 Z 疾病"的因果解释。这种不透明性不仅限制了模型在高风险领域的应用, 也增加了模型调试和改进的难度。

3)数据依赖与泛化边界问题

虽然 RNN/LSTM 在理论上具有更强的泛化能力,但在实际应用中往往面临"数据饥饿"困境。训练高质量的模型通常需要百万级以上的标注样本,而在生物医学、法律等专业领域,标注数据的获取成本极高。当训练数据不足时,模型容易陷入过拟合,表现为在训练集上准确率高,但在测试集上性能骤降。例如,在处理低资源语言(如非洲某些部落语言)时,RNN/LSTM模型的表现甚至不如简化的统计模型。此外,模型对训练数据的分布极其敏感,当应用场景的语言风格或领域知识与训练数据存在偏差时,性能会显著下降[10][12][13]。

4) 长序列处理的效率瓶颈

尽管 LSTM 通过门控机制缓解了梯度消失问题,但其对长序列的处理能力仍存在物理上限。 当序列长度超过数百个时间步时,模型的记忆能力会逐渐衰退。这是因为隐状态在长时间传递 过程中会不可避免地丢失早期信息,形成"记忆衰减"现象。例如,在处理长篇小说或学术论 文时,模型可能无法有效关联前文数百词之外的关键信息。为缓解这一问题,实际应用中常采 用序列截断或分层处理策略,但这些方法会导致上下文信息的人为损失,影响模型性能。

5) 模型部署与推理效率挑战

RNN/LSTM 模型的生产环境部署面临多重挑战。由于模型结构复杂、参数量大,在移动设备或边缘计算场景下的部署受到硬件资源限制。例如,在智能语音助手等实时应用中,模型需要在毫秒级内完成推理,而 RNN/LSTM 的串行计算特性难以满足这一要求。为提高推理速度,通常需要进行量化、剪枝等模型压缩操作,但这些操作可能导致精度损失,需要在效率和性能之间进行艰难权衡。

上述问题限制了其在大规模场景下的应用,正是在这样的技术演进背景下,Transformer 模型应运而生,通过自注意力机制革新序列处理方式,有效解决了长距离依赖计算效率与并行训练难题,将在 1.1.3 节深入探讨。

1.1.3 基于 Transformer 架构的语言模型

2017年,Google 提出了 Transformer 架构^[14],解构了序列建模的固有范式,彻底改变了语言模型的发展格局。与 RNN 不同,Transformer 摒弃了循环结构,采用多头注意力(Multi-Head Attention,MHA)机制,能够并行处理整个输入序列,大大提高了训练效率和模型性能,为现代大语言模型奠定了不可撼动的架构基础^[13]。Transformer 模型通用架构如图 1.4 所示。

1. 基于 Transformer 的语言模型基础架构

Transformer 架构是从 RNN(循环神经网络)的编码器-解码器架构中汲取灵感而来的,其引入了注意力机制^[15]。它被广泛应用于序列到序列(Seq2Seq)任务,并且相比于 RNN,Transformer 摒弃了顺序处理的方式。

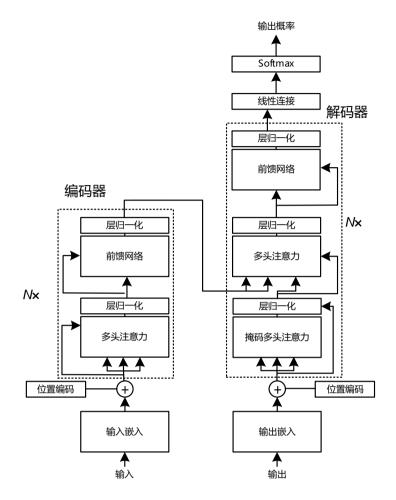


图 1.4 Transformer 模型通用架构

不同于 RNN, Transformer 以并行化的方式处理数据,从而能够实现更大规模的并行计算 和更快速的训练。这得益于 Transformer 架构中的自注意力机制,它使得模型能够同时考虑输入 序列中的所有位置,而无须按顺序逐步处理。自注意力机制允许模型根据输入序列中的不同位 置之间的关系,对每个位置进行加权处理,从而捕捉全局上下文信息。

我们可以注意到, Transformer 的模型通用架构^[14](见图 1.4), 由编码器和解码器两个主 要部分组成。

1) 编码器堆栈

这是由 N 个相同的编码器层组成的堆栈(在原始论文中,N=6)。每个编码器层都由两个 子层组成: 多头自注意力机制和前馈神经网络。多头自注意力机制用于对输入序列中的不同位 置之间的关系进行建模,而前馈神经网络则用于对每个位置进行非线性转换。编码器堆栈的作 用是将输入序列转换为一系列高级特征表示。

具体来说, 多头注意力是一种在 Transformer 模型中被广泛采用的注意力机制扩展形式, 它 通过并行地运行多个独立的注意力机制来获取输入序列的不同子空间的注意力分布,从而更全 面地捕获序列中潜在的多种语义关联。在多头注意力中,输入序列首先通过三个不同的线性变换层分别得到 Query、Key 和 Value。然后,这些变换后的向量被划分为若干"头",每个头都有自己独立的 Query、Key 和 Value 矩阵。对于每个头,都执行一次缩放点积注意力(Scaled Dot-Product Attention)运算,即:

Attention(
$$Q, K, V$$
) = Softmax($\frac{Q \cdot K^{T}}{\sqrt{d_k}}$) $\cdot V$ (1-22)

最后,所有头的输出会被拼接在一起,然后通过一个线性层进行融合,得到最终的注意力输出向量。

通过这种方式,多头注意力能够并行地从不同的角度对输入序列进行注意力处理,提高了模型理解和捕捉复杂依赖关系的能力。在实践中,多头注意力能显著提升 Transformer 模型在自然语言处理和其他序列数据处理任务上的性能。

①输入变换与线性投影

多头注意力机制的输入变换与线性投影是其核心步骤之一。给定输入序列,首先通过三个不同的线性变换层生成查询(Q)、键(K)和值(V)矩阵。这些变换通常是通过全连接层实现的,其目的是将输入数据映射到不同的表示子空间中,为后续的注意力计算提供基础。

输入序列首先被映射到查询、键和值矩阵:

$$Q = W_O x, K = W_K x, V = W_V x$$
 (1-23)

其中,x 为输入序列, W_Q 、 W_K 、 W_V 均为权重矩阵。由于每个头的计算是独立的,这些计算可以并行进行,从而提高模型的计算效率。这种并行性使得多头注意力机制在处理长序列数据时更加高效。

②注意力权重计算

在多头注意力机制中,每个头的注意力权重计算是通过缩放点积注意力实现的。具体来说, 计算查询和键的点积,经过缩放、加上偏置后,使用 Softmax 函数得到注意力权重。

为了避免过大的点积导致梯度消失问题,通常会对点积结果进行缩放:

Scaled Score =
$$\frac{QK^{T}}{\sqrt{d_k}}$$
 (1-24)

其中, d_k 为 Key 的向量维度。使用 Softmax 函数对缩放后的得分进行归一化,得到每个元素的注意力权重,其中第i个元素的注意力权重为:

$$\alpha_i = \frac{\exp(\text{Scaled Score}_i)}{\sum_{i} \exp(\text{Scaled Score}_i)}$$
 (1-25)

③拼接与融合

多头注意力机制的最后步骤是将所有头的输出拼接在一起,然后通过一个最终的线性变换, 以整合来自不同头的信息,得到最终的多头注意力输出。这一步骤整合了从不同子空间学到的 信息,增强模型的表达能力。对拼接后的向量进行一个最终的线性变换,以整合来自不同头的 信息,得到最终的多头注意力输出:

Output=
$$\mathbf{W}_o \operatorname{concat}(C_1, C_2, \dots, C_h)$$
 (1-26)

其中, W_a 为输出层权重矩阵, C_i 是第i个头的输出。

2)解码器堆栈

这也是由 N 个相同的解码器层组成的堆栈(在原始论文中,N=6)。每个解码器层除了包 含编码器层的两个子层外,还包含一个额外的掩码多头自注意力机制子层。这个额外的自注意 力机制用于对编码器堆栈的输出进行关注,并帮助解码器对输入序列中的信息进行解码和生成 输出序列。

掩码自注意力确保解码器在生成位置 t 时只能访问位置 0~t 的信息, 防止未来信息泄露:

Attention(
$$Q, K, V$$
) = Softmax($\frac{Q \cdot K^{T} + M}{\sqrt{d_{k}}}$) $\cdot V$ (1-27)

其中, M 为掩码矩阵:

$$\mathbf{M}_{ij} = \begin{cases} 0 & i \geqslant j \\ -\infty & i < j \end{cases} \tag{1-28}$$

3)位置编码层

在编码器和解码器堆栈之前,还有一个位置编码层。这个位置编码层的作用是利用序列的顺 序信息,为输入序列中的每个位置提供一个固定的编码表示。这样,模型可以在没有递归或卷积 操作的情况下,利用位置编码层来处理序列的顺序信息。使用正弦和余弦函数生成位置编码:

$$\begin{split} & \text{PE}_{(\text{pos},2i)} \sin(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}) \\ & \text{PE}_{(\text{pos},2i+1)} \cos(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}) \end{split} \tag{1-29}$$

其中, pos 为序列中的位置, d_{model} 为模型维度(例如 512)。位置编码能够实现相对位置敏感 的特性,这是因为位置 pos+k 的编码可以表示为 pos 的线性函数,同时每个位置编码都是唯 一的。

2. 基于 Transformer 的语言模型的局限性分析

与基于 RNN 的 Seq2Seq 模型相比,尽管 Transformer 模型在自然语言处理领域取得了巨大 的成功,然而,其本身也存在以下局限性[14][16]。

从计算资源层面来看, Transformer 模型的大规模参数架构对硬件设施提出了严苛要求。以 GPT 系列模型为例,其数十亿乃至上百亿的参数规模,在训练阶段需要数千块 GPU 并行运算 数月之久,普通科研机构和企业难以承担如此高昂的算力成本。在推理过程中,大量参数的实 时调用也导致内存占用居高不下,在边缘计算、移动设备等资源受限场景中,模型部署面临巨 大挑战。自注意力机制虽然解决了长距离依赖问题,但计算复杂度随序列长度呈平方增长的特性,使得长文本处理成为 Transformer 的显著短板。随着输入文本 Token 数急剧增加,计算量与显存需求呈指数级上升,部分终端用户的硬件条件不足的设备甚至会因内存溢出导致程序崩溃,这在处理学术论文、法律文书等超长文本时尤为突出。

在知识推理与数据依赖方面,Transformer 模型同样存在亟待解决的问题。基于大规模语料预训练的模型,本质上是对语言分布规律的概率拟合,缺乏人类的常识推理能力。例如,在处理"冰箱里的牛奶会结冰"这类涉及物理常识的问题时,模型可能因语料中缺乏对应表述而出现推理错误。此外,模型对训练数据的数量与质量高度依赖,在医疗、金融等专业领域,标注数据的稀缺性导致模型微调效果不佳;而当训练数据存在偏见(bias)或噪声时,模型输出也会产生相应的偏差,这在情感分析、新闻推荐等应用中可能引发严重的社会问题。

尽管当前 Transformer 模型已成为自然语言处理领域的核心技术,但正视这些局限性,将推动学界和业界持续探索混合架构优化、小样本学习等创新路径,为实现更通用、高效的人工智能语言模型奠定基础。

1.2 大模型发展历史

大模型的发展并非一蹴而就,而是人工智能领域数十年理论探索与技术创新相互激荡的产物。这一历程不仅深刻改变了自然语言处理的技术范式,更推动了人工智能从专用系统迈向通用智能的关键跨越。从早期统计语言模型的概率建模,到神经网络架构的革命性突破,再到百亿参数规模的通用大模型崛起,每个阶段的技术演进都伴随着计算能力提升、数据规模增长与算法创新的协同作用。

整体而言,大模型的发展可以总结为以下 4 个阶段: 统计语言模型奠基期(1950—2010 年)、神经网络语言模型探索期(2010—2017 年)、Transformer 架构革命期(2017—2019 年)以及大模型爆发增长期(2020 年—)。而如今提到的大语言模型主要是指 2020 年以后以 Transformer 架构为基础提出的预训练+微调范式催生出来的模型产物。

本节将沿着技术发展的时间轴线,系统梳理大模型从理论雏形到产业应用的完整发展脉络, 揭示其背后的驱动因素与技术突破逻辑。

1.2.1 统计语言模型奠基期

人工智能发展初期,统计语言模型(Statistical Language Model,SLM)通过概率论方法对语言进行建模,成为自然语言处理的主流技术。以 N-gram 模型为代表,该类模型基于语料库中词语的共现频率计算语言序列概率,例如二元语法(Bi-gram)通过前一个词预测下一个词的出现概率^[1]。虽然这种方法在机器翻译、语音识别等任务中取得了一定成功,但其依赖人工设计特征,难以处理长距离依赖和复杂语义,模型泛化能力有限。随着数据规模和计算需求的增长,统计语言模型的局限性逐渐凸显,为后续神经网络语言模型的兴起埋下伏笔。

神经网络语言模型探索期 1.2.2

2010年后,深度学习技术的快速发展为语言模型带来新的突破方向。循环神经网络(RNN) 及其变体 LSTM、GRU 的出现,通过引入隐层状态循环机制,实现了对序列数据的动态建模, 有效解决了统计语言模型的长距离依赖问题^[6]。2013 年提出的 Word2Vec^[17]和 2014 年提出的 GloVe^[18]模型,则通过无监督学习方法将词语映射为低维向量,开启了分布式语义表示的研究 热潮。这些技术进步为后续预训练语言模型的发展奠定了基础,但 RNN 系列模型仍存在训练 效率低、难以并行计算等问题,限制了模型规模的进一步扩展。

Transformer 架构革命期 1.2.3

2017 年, Transformer 架构的提出成为大模型发展的重要分水岭。其创新性地使用自注意力 机制替代循环结构,通过并行计算大幅提升训练效率,同时有效解决了长距离依赖问题[14]。基 于 Transformer 架构的 BERT (Bidirectional Encoder Representations from Transformers) 在 2018 年横空出世,通过双向预训练和微调策略,在 11 个 NLP 任务上取得当时的最优性能,标志着 预训练-微调范式的正式确立^[19]。2019 年 GPT-2 的发布则展现了 Transformer 在生成任务上的 强大潜力, 其通过无监督学习训练的 15 亿参数模型, 能够生成连贯且语义合理的长文本, 引发 学界和业界对大模型能力边界的重新思考。

大模型爆发增长期 1.2.4

2020年, OpenAI 推出的 GPT-3 成为大模型发展的关键转折点, 其拥有 1750 亿个参数, 以 远超以往模型的规模,展现出了卓越的少样本学习能力[20]。在少样本甚至零样本学习任务中, GPT-3 能够依据给定的少量示例,对新任务作出合理推断,生成连贯且语义准确的文本,开启 了模型规模增长引发"涌现"能力的研究热潮,即当模型参数达到一定量级后,会自发呈现出 此前未被设计的复杂认知与推理能力。

随后,业界与学界积极投身大模型研发,促使大模型技术呈现爆发式增长,如图 1.5 所示[54]。 2022年, Google 推出 PaLM (Pathways Language Model)^[21], 参数规模达 5400亿, 其基于 Google 自研的 Pathways 系统构建,具备强大的可扩展性,在自然语言处理的各类任务中表现出色,尤 其在语言翻译、复杂文本生成等方面,展现出高准确性与流畅性,进一步证明了超大规模参数 模型在提升性能上的潜力。同年,DeepMind 发布 Chinchilla^[22],虽然参数规模(700 亿)小于 PaLM,但通过优化训练数据质量与规模(使用高达 1.4 万亿 Token 的数据集),在模型效率与 性能平衡上取得突破,在多项基准测试中超越了同等规模的其他模型,凸显了优质数据对模型 训练的关键作用。

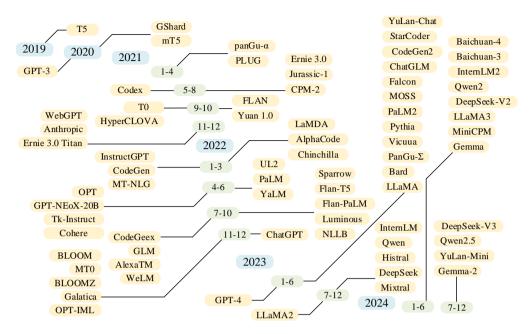


图 1.5 大模型发展历程(2020-2023年)

2022 年 11 月,OpenAI 的 ChatGPT 横空出世,它基于 GPT-3.5 模型^[23]微调而成,通过引入人类反馈强化学习(Reinforcement Learning from Human Feedback,RLHF)技术,极大地提升了模型与人类交互的自然性和准确性,能够根据用户提问,生成贴合语境、逻辑连贯的回答,迅速引发全球关注,推动大模型从实验室研究走向大众应用,开启了大模型商业化的新篇章。2023 年,OpenAI 发布的 GPT-4 更是具备多模态理解能力,不仅能处理文本,还可对图像输入作出响应,如理解图片内容并基于此生成描述、解答相关问题等,在复杂任务处理、知识推理等方面的性能进一步提升,成为当时最先进的多模态大模型之一。图 1.6 展示了 GPT 系列模型演进路线^{[24][54]}。

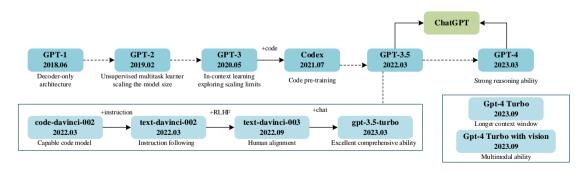


图 1.6 GPT 系列模型演进路线

国内的大模型研发也在这一时期蓬勃发展,从图 1.5 后半段可以看出^[54],虽然国内厂商入场较晚,但呈现出后来者居上的气势。百度的文心一言于 2023 年发布,基于 ERNIE 3.0 框架^[25],聚焦于知识增强大模型,通过融合大量知识图谱信息,在知识问答、文本创作等任务中表现突出,助力企业与开发者在智能写作、智能客服等领域实现高效应用开发。阿里巴巴的通义千

问^[26]同样具备强大的语言生成与理解能力,在电商、金融等垂直领域深入布局,为行业定制化 解决方案提供底层技术支撑。

2024—2025 年, 大模型技术持续迭代创新。字节跳动的豆包^[27]模型不断升级, 如 1.5 版本 推出的"深度思考模型"及其视觉版本,在数学推理、编程竞赛、科学推理等专业领域成绩优 异, 其视觉版本能结合多源信息深度理解图像内容, 实现如通过航拍地貌推理地理位置等复杂 任务。科大讯飞的星火 X1 作为全国产算力训练的深度推理大模型^[28], 首发"快思考"与"慢 思考"统一架构,可根据任务需求灵活切换模式,在语言理解、文本生成、数学答题、代码生 成等通用任务上全面升级,且多模态推理能力在教育、医疗、司法等行业得到深化应用。昆仑 万维的天工系列大模型[29]不断演进,从 2.0 版本通过动态任务分配提升复杂任务处理效率,到 3.0 版本以 4000 亿个参数成为全球规模领先的开源混合专家($MoE^{[30]}$)模型 $^{[31]}$, 再到 4.0 版本 实现实时语音交互与慢思考推理的突破,在底层架构创新与多模态技术融合上持续探索。

开源社区在这一时期也发挥了重要推动作用。Meta 的 LLaMA 模型[31]开源后,激发了全球 开发者基于其讲行二次开发与优化、衍生出众多性能优异的变体模型、降低了大模型的使用门 槛,促进技术普及。图 1.7 展示了 LLaMA 系列模型的演讲与发展路线^[54]。Stable Diffusion^[32] 作为开源的文本生成图像大模型,引发了生成式 AI 在图像领域的应用热潮,为艺术创作、设 计等行业提供了全新的创作工具与思路,推动大模型技术在多模态应用领域的广泛拓展。

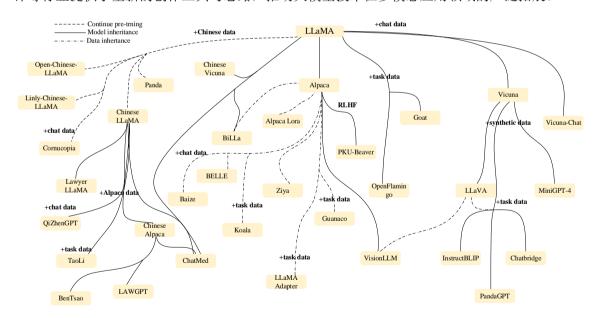


图 1.7 LLaMA 系列模型的演进路线

这一阶段,大模型在提升性能的同时,产业落地进程加速。在教育领域,大模型助力个性 化学习方案制定、智能辅导[33][34][35];在医疗行业,辅助医生进行疾病诊断、医学文献分析[36][37]; 在金融领域,用于风险评估、智能投顾等[38][39]。大模型正逐步渗透到各行业的核心业务流程, 重塑产业格局,成为推动社会数字化转型的关键技术力量。

1.3 大模型的特点

大模型凭借其独特的架构和训练方式,展现出与传统模型截然不同的特性,这些特性既赋 予了它们强大的能力, 也带来了一定的局限性。

大模型的快思慢考 1.3.1

在人工智能技术不断突破的浪潮中,大模型处理任务时展现出的"快思慢考"特性,成为 理解和驾驭这类先进技术的核心切入点。这一特性源于对人类双过程认知理论的工程化实践, 不仅深刻影响着大模型的架构设计与运行逻辑,更在教育、医疗、金融等诸多领域引发了应用 范式的变革。

1. 大模型"快思慢考"特性的本质解析

"快思"与"慢考"构成了大模型处理任务的双重路径。"快思"模式如同人类的直觉反 应,依赖大模型在预训练阶段积累的海量知识与模式识别能力,能够对常见问题作出快速响应。 例如, 当用户询问"世界上面积最大的海洋是哪个", 模型可瞬间调用知识库, 输出"太平洋" 的答案。这种快速响应机制极大地提升了用户交互体验,尤其适用于处理信息检索、基础问答 等结构化任务。然而, "快思"模式的局限性也较为明显,由于缺乏深度推理与验证过程,面 对复杂问题时,模型可能产生逻辑错误或不准确的回答。例如,在处理"如何用物理学原理解 释潮汐现象"这类需要多步骤推理的问题时,单纯依赖"快思"可能给出片面或错误的答案。

与之相对,"慢考"模式模拟人类深思熟虑的过程,旨在解决复杂任务。当遇到数学证明、 医疗诊断、战略决策等需要深度推理的问题时,大模型会启动"慢考"机制。在这一过程中, 模型将复杂问题拆解为多个子任务,通过逐步推导、验证和自我修正,构建起完整的逻辑链条。 例如,在解答"证明勾股定理"的问题时,模型会从几何定义出发,结合代数知识,通过多步 骤推导得出严谨的证明过程。"慢考"模式虽然耗时较长,但能够显著提升模型处理复杂任务 的准确性和可靠性。

2. DeepSeek 双模架构: 快思慢考的技术实现典范

DeepSeek 创新性地构建了"快思慢考"双模认知架构,为大模型的任务处理提供了新范式。 该架构将任务处理分为两大核心模块:快思考模型与慢思考模型,二者通过智能协作实现效率 与深度的平衡。

快思考模型(如 DeepSeek-V3^[41])采用稀疏激活的混合专家(Mixture of Experts, MoE) 架构,这一设计打破了传统模型全参数激活的低效模式。在实际运行中,该模型仅动态激活 5.5% 的参数,就能实现对常见任务的高速响应。以信息检索为例,当用户查询"2024年全球新能源 汽车销量数据"时,快思考模型可在 200 毫秒内完成数据调取与整合,输出准确答案。这种低 计算成本、高响应速度的特性,使其在处理占实际场景80%以上的简单任务时极具优势。

慢思考模型(如 DeepSeek-R1^[42])则专注于复杂推理任务。其核心技术包括动态推理路径

生成与无监督结果导向训练。通过标签、模型能够显式生成推理步骤、将复杂问题的解决过程 可视化。例如,在分析"央行降息对房地产市场的传导机制"时,慢思考模型会分步骤梳理利 率变动对资金成本、购房需求、市场供需等环节的影响,最终形成完整的分析报告。同时,基 于强化学习的无监督训练机制,模型能够自动优化推理路径,不断提升复杂任务处理的准确性。

为支撑这一双模架构高效运行,DeepSeek 融合了多项原创技术^[42]。多头潜在注意力 (Multi-Head Latent Attention, MLA) 技术^[14]通过 Key-Value 矩阵低秩压缩,将长文本处理的 显存需求降至传统方法的 1/3,结合旋转位置编码(Rotary Position Embedding,RoPE)^[43],实 现了 128K 超长上下文的高效建模。自适应慢思考优化机制能够根据问题复杂度动态调整思维 链[44]长度,并通过 DA-GRPO 算法[45]减少冗余推理,降低 30%的计算量。经济性训练框架采用 MoE 负载均衡技术^[41], 使训练效率提升 37%, 同时利用多 Token 预测 (Multi-Token Prediction, MTP)加速模型收敛。

3. 快思慢考架构的行业变革与应用价值

DeepSeek 的"快思慢考"架构对大模型领域产生了深远影响,在技术普惠、架构创新和认 知智能等方面实现了重大突破。在技术成本上,该架构将 600B 参数规模模型的训练成本降至 600 万美元,通过开源 70%的代码,极大地降低了大模型的研发门槛,推动千亿级模型从实验 室走向产业应用。在架构范式上, DeepSeek-Lite 等边缘端量化模型实现了 500ms 的快速响应, 能耗降低 63%,为资源受限环境下的模型部署提供了可行方案。在认知智能领域,其无监督推 理技术首次验证了机器自主推理的可行性,促使行业更加注重技术的透明化与可解释性。

在实际应用场景中,"快思慢考"架构展现出了强大的实用价值。在教育领域,快思考模 式可快速解答学生的基础问题,如单词释义、公式推导,而慢思考模式则能针对复杂的学术问 题,如历史事件因果分析、数学难题求解,提供详细的推理过程与解题思路,实现个性化的学 习辅导。在医疗领域,快思考模型可快速完成症状初步诊断,慢思考模型则能综合患者病历、 影像数据和医学知识库,进行复杂病症的精准诊断,例如在肿瘤良恶性判断中^[46],将错误率从 传统方法的大于 40%降至小于 12%。金融领域中,快思考模式实时监测交易数据,及时发现异 常交易并预警;慢思考模式则通过对全球经济形势、行业动态和企业财务状况的深度分析,为 投资决策提供科学依据。

4. 未来展望: 快思慢考的进化方向

展望未来,大模型的"快思慢考"特性将朝着更高效、更智能的方向发展。在技术层面, 超长上下文分层注意力机制的优化将进一步提升模型处理复杂信息的能力; 跨模态对比学习技 术的发展将使模型在文本、图像、音频等多模态数据处理中实现更深度的融合与推理。在应用 层面,随着自研硬件生态的完善,大模型将在边缘计算、物联网等场景中发挥更大作用。 DeepSeek 提供的 STAR^[47]提示框架、蒸馏模型本地化部署等实践路径,将为开发者提供更加便 捷的工具,推动大模型技术在更多领域的落地应用。

大模型的"快思慢考"特性不仅是技术发展的必然产物,更是人工智能迈向通用智能的重 要一步。理解和掌握这一特性,对于推动大模型技术的创新发展、实现其在各行业的深度应用

具有重要意义。随着技术的不断进步,"快思慢考"架构将持续进化,为智能社会的构建提供 更强大的技术支撑。

1.3.2 大模型的优势与不足

大模型以强大的技术能力和广阔的应用前景,正在重塑各个行业的运作模式与发展方向。 然而,如同自然界中任何事物都具有两面性,大模型在彰显巨大价值的同时,也暴露出诸多亟 待解决的问题。深入剖析这些优势与劣势是掌握大模型技术本质、推动其在各领域合理应用的 关键所在,更是在人工智能时代把握机遇、应对挑战的重要前提。

1. 突破性进展与固有瓶颈并存

大模型在架构创新方面取得的成就令人瞩目。以 DeepSeek-V3 为例[41], 其采用的稀疏 MoE (混合专家)架构与 FP8 混合精度训练技术,是技术创新的典型代表。该架构打破了传统模型 全参数激活的模式,仅激活 5.5%的参数,就将千亿级规模模型的训练成本大幅压缩至 557.6 万 美元, 仅为 GPT-4 训练成本的 1/18, 极大地降低了大模型研发的资金门槛, 使得更多科研团队 和企业能够参与到大模型的研究与开发中。同时, DeepSeek-V3 具备的 128K 超长上下文窗口, 为处理复杂任务提供了强大支撑。在金融风控领域,面对海量的交易数据,该模型凭借超长上 下文窗口,能够全面分析交易历史、用户行为模式等信息,精准识别风险模式,将风险识别错 误率降低37%,显著提升了金融机构风险防控的效率与准确性。

当模型参数规模突破一定临界点后,涌现出的能力更是为人工智能发展开辟了新的道路。 以 Claude3 在蛋白质折叠预测任务中的表现为例^[48],基于思维链推理的能力,它能够模拟复杂 的生物过程,对蛋白质的空间结构进行预测。这一能力对于药物研发和疾病治疗意义重大,科 研人员可以借助模型的预测结果,更有针对性地设计药物分子,加速药物研发进程。这种能力 的涌现标志着大模型不再局限于简单的数据处理,而是开始向具备复杂认知能力的智能系统迈 进,为解决科学研究中的复杂问题提供了新的可能。

尽管大模型在技术上取得了重大突破,但其面临的挑战同样不容小觑。在处理复杂逻辑任 务时,"幻觉"问题成为困扰大模型的一大难题^[49]。由于模型在训练过程中主要基于数据统计 规律进行学习,缺乏对真实世界的全面理解,导致其生成的内容可能存在与事实不符、逻辑错 误等情况。研究数据显示,大模型在复杂逻辑任务中的幻觉率在15%~40%波动,这严重影响了 模型输出的可靠性。在处理 128K 长文本时,自注意力机制的计算复杂度为 O(n²),随着文本长 度增加, 计算量呈指数级增长, 所需显存高达 80GB, 这远远超出了普通消费级硬件的承载能 力,使得长文本处理在实际应用中困难重重,限制了大模型在需要处理长篇文档场景中的应用。

此外,大模型在持续学习方面也存在明显不足[50]。当对大模型进行全参数微调以适应新任 务时,会导致模型在旧任务上的性能衰减超过70%,出现"灾难性遗忘"现象。这意味着模型 在学习新知识的过程中,难以保留已掌握的知识,无法在不同任务之间实现良好的迁移,限制 了其在动态变化环境中的应用能力。例如,一个经过新闻文本分类训练的大模型,在微调用于 情感分析任务后,对新闻文本分类的准确率会大幅下降。

2. 应用落地的机遇与风险

大模型在应用落地过程中,为众多行业带来了效率的显著提升,推动着产业效率的范式级 重构。在交通领域,天津地铁部署的多模态交互系统便是一个成功案例。该系统借助大模型的 可视化应急指南实时生成能力,在地铁设备突发故障时,能够迅速分析故障类型,结合历史维 修数据和现场情况,生成图文并茂的处置流程。工作人员可以根据这些直观的指南,快速定位 问题并进行解决,使故障处置效率提高了40%,有效减少了因故障导致的地铁运营延误,提升 了乘客的出行体验。

在政务服务方面,洛阳市医保通过引入大模型对业务流程进行重构,取得了令人瞩目的成 果[51]。以往,医保异地办理手续烦琐,流程复杂,办理时长长达 48 小时。而引入大模型后, 系统能够自动审核参保人员提交的材料,快速比对数据,将异地办理时长大幅缩短至4小时, 极大地方便了群众办事,提升了政府服务的便捷性和高效性,增强了群众对政务服务的满意度。

然而,大模型的广泛应用也伴随一系列系统性风险。在专业领域,如核电故障诊断,对模 型的准确性和可靠性要求极高[52]。为了使模型达到可用水平,需要对百万级的标注数据进行微 调,而收集和标注这些数据的冷启动成本超过 200 万美元,这对于许多企业和机构来说是一笔 难以承受的经济负担。在安全敏感场景,如电梯困人识别系统,对模型的实时性和准确性要求 近乎苛刻。即使是 500ms 的延迟, 也可能延误救援时机, 对被困人员的生命安全造成威胁, 这 对大模型的性能提出了极高挑战。

随着大模型在内容生成领域的广泛应用,深度伪造技术日益猖獗。相关数据显示[53],深度 伪造相关犯罪数量年增长 300%, 而检测技术却相对滞后, 检测率不足 85%。虚假的图像、视 频和音频内容在网络上传播,不仅会误导公众,还可能引发社会恐慌,给社会安全和稳定带来 严重威胁。此外,关键指标对比还揭示了大模型应用中的深层矛盾。DeepSeek的API成本仅为 0.27 美元/百万Token输入,是ChatGPT的 1/20,显著降低了使用成本,但千亿模型训练所需的 巨大能耗,相当于三座核电站的年发电量,这与可持续发展理念相悖。在医疗诊断领域,虽然 大模型可以将错误率降低至 9%, 但在欧盟严格的可溯源要求下, 达标率却不足 35%, 暴露出 模型在合规性方面的不足。

大模型行业应用场景中的优势与挑战

教育领域:智能教育的革新与困境 1.4.1

在教育领域,大模型展现出强大的优势。它凭借强大的认知推理能力,助力构建人机间"协 同教学""协同学习"与"协同决策"的创新应用场景。在教师备课环节,大模型可以根据教 学大纲和课程目标,自动生成教案、课件和练习题,还能推荐相关的教学资源,如优质的教学 视频、学术论文等,极大地减轻了教师的备课负担。在作业批改方面,大模型能够快速准确地 批改客观题,并对主观题给出合理的评分建议和修改意见,提升了作业批改的效率。在辅导答

疑时,大模型可以随时解答学生的问题,通过语言理解和逻辑推理能力,为学生提供详细的解答和学习指导。

此外,大模型还能通过文生图、文生音频、文生视频等技术,自动生成多样化教学资源。例如,在讲解历史事件时,生成相关的历史场景图片和动画视频;在语言学习中,生成标准的语音朗读和对话音频,为师生营造更加沉浸式的学习体验,激发学生的学习兴趣。然而,大模型在教育应用中也面临诸多挑战。不同学科具有独特的教学特点与需求,大模型目前存在多学科适配性不足的问题。例如,在数学、物理等理科教学中,对于复杂的公式推导和逻辑证明,大模型的解释能力有限;在语文、历史等文科教学中,对于文学作品的情感分析和历史事件的深度解读,大模型难以达到人类教师的水平。

同时,大模型的应用缺乏系统性教育理论支撑,使得其在教育实践中的应用缺乏深度教育理念的引导,难以充分发挥其教育价值。高质量训练数据的匮乏,也限制了模型在教育场景中的精准度与有效性。此外,大模型的"幻觉"现象、精准度和可解释性问题,以及实时个性化支持不足等,都无法完全满足每个学生的独特学习需求,影响了教育教学的质量。

1.4.2 医疗领域:精准医疗的希望与隐忧

在医疗行业,大模型同样发挥着重要作用。以北京天坛医院联合开发的"龙影大模型(RodGPT)"为例,它在医学影像分析方面表现出色。该模型能够在 0.8 秒内分析 MRI 影像,并给出百种疾病的诊断意见,准确率高达 90%,极大地提高了诊断效率。在重症监护等高风险高压的医疗环境中,快速准确的诊断对于患者的治疗至关重要,龙影大模型为医生提供了有力的辅助,帮助他们及时作出准确的治疗决策。

然而,医疗领域对数据隐私安全要求极高。不同医疗机构之间的数据相互隔离,存在严重的数据孤岛现象,这使得大模型难以整合利用全面的医疗数据进行训练和优化。为了打破数据孤岛,需要构建统一的数据共享平台,并制定严格的数据共享规则和安全标准,确保医疗数据在共享过程中的安全性和隐私性。同时,医疗数据包含患者大量敏感隐私信息,大模型的开发和使用者必须建立完善的数据隐私保护机制,采用先进的加密技术和访问控制策略,确保数据合法合规使用,防止患者隐私泄露。

此外,虽然大模型在医学影像分析和疾病预测等方面表现出色,但其决策过程往往缺乏可解释性。在医疗诊断中,医生需要清楚了解诊断结果的依据,以保障患者权益。因此,需要借助知识图谱等技术研发可解释的 AI 算法,让大模型的诊断决策逻辑清晰呈现,使医生能够信任和理解模型的诊断结果,更好地为患者服务。

1.4.3 金融领域:智能金融的变革与挑战

在金融领域,大模型凭借强大的数据处理和分析能力,为金融业务带来了新的变革。大模型可以对借款人的信用记录、消费行为、资产状况等多维度数据进行深度挖掘和分析,快速准确地评估借款人的信用状况,提高贷款审批效率与准确性。同时,它还能通过对交易数据的实

时监测, 识别潜在欺诈行为和异常交易, 及时发出风险预警, 降低金融风险。

在智能客服方面,基于大模型的虚拟客户经理能够理解客户的问题和需求,与客户进行自 然流畅的交流,并为客户提供可行的解决方案。例如,帮助客户获得和提升授信额度、解答客 户关于金融产品的疑问等。这不仅提升了客户服务质量,还降低了金融机构的服务成本。此外, 大模型还能细分出 AI 投顾,根据客户财务状况、投资目标和风险偏好,运用复杂的算法和模 型,提供个性化资产配置方案和投资组合建议,帮助投资者实现资产的合理配置和增值。

不过, 金融领域对模型的要求极为严苛。通用大模型在行业数据量、性价比、精确性、适 用性、实时性、推理速度、合规性以及风险控制等方面存在不足。金融领域的数据分散在各个 机构和系统中、获取难度大、使用金融数据对通用大模型进行训练时、数据欠缺目成本过高。 同时,从底层训练大模型需要巨大的算力资源,成本高昂。在特定金融任务上,通用大模型的 精确性与适用性欠缺,需要针对金融业务进行更多优化与定制。此外,金融市场瞬息万变,要 求模型具备实时响应和快速推理速度,而通用大模型在这方面往往难以满足金融业务的需求。

电商领域:智能营销的机遇与难题 1.4.4

在电商场景中,大模型通过学习消费者行为、商品评价、市场交易等多种数据,能够构建 复杂的用户和商品关联图,实现对市场趋势的精准预测。在商品推荐方面,大模型可以根据用 户的历史购买记录、浏览行为和兴趣偏好,为用户推荐个性化的商品,提高用户的购买转化率 和满意度。在价格预测上,大模型分析市场供需关系、竞争对手价格等因素,预测商品价格走 势,帮助商家制定合理的定价策略。

在库存管理方面, 大模型根据销售数据和市场趋势, 预测商品的销售量, 合理安排库存, 避免库存积压或缺货现象的发生。在客户行为分析中,大模型挖掘客户的潜在需求和消费心理, 为商家提供营销策略建议。然而,大模型在电商应用中也面临诸多问题。大模型存在计算资源 消耗大、训练时间长的问题,在处理大规模电商数据集时,需要强大的算力支持和较长的训练 周期,这对于许多电商企业来说是一个巨大的挑战。同时,大模型的黑盒特性导致其可解释性 差,对于电商领域复杂的决策支持系统而言,难以解释模型的决策过程和依据,这可能会影响 商家对决策结果的信任和应用。

1.5 本章小结

本章围绕大模型基础展开,系统梳理了大模型的核心知识体系。从语言模型基础出发,依 次介绍了基于统计方法、RNN 以及 Transformer 的语言模型,揭示了语言模型从传统统计计算 到深度学习架构的演进路径,尤其是 Transformer 架构如何凭借多头注意力机制革新语言处理能 力。在大模型发展历史部分,回顾了从 GPT-1 起步, 到 GPT-3、PaLM 等模型不断突破参数规 模与应用边界的历程,展现大模型推动人工智能迈向通用化的趋势。而大模型的特点中,"快 思慢考"特性反映其响应速度与深度推理的矛盾,优势与不足的分析,则明确了大模型在知识 表示、泛化能力等方面的突出表现,以及训练成本、可解释性等现存挑战。这些内容为后续深入学习大模型应用开发筑牢了理论根基,促进大模型的技术演进与本质特征全面深入的认知。

大模型的发展之路是一场效率增益与伦理红线之间的平衡艺术。只有在充分发挥其技术优势的同时,有效控制潜在风险,构建起"能力放大器"与"风险控制器"的双重体系,才能让大模型在智能制造、智慧医疗、智能教育等更多领域释放出真正的变革性潜力,为人类社会的发展带来积极而深远的影响。我们有理由相信,随着技术的不断进步和应用的不断深入,大模型将在未来的科技发展和社会进步中发挥更加重要的作用。

1.6 参考文献

- [1] 邢永康,马少平. 统计语言模型综述[J]. 计算机科学,2003,30(09):22-26.
- [2] Jelinek F, Mercer R I. Interpolated estimation of Markov source parameters from sparse data[EB/OL]. (1980-01-01)[2025-06-01].https://scispace.com/papers/interpolated-estimation-of-markov-source-parameters-from-39sufvwj23.
 - [3] 袁毓林. 基于统计的语言处理模型的局限性[J]. 语言文字应用, 2004, 13(2): 10.
- [4] Goodfellow I, Bengio Y, Courville A. Deep learning (Vol. 1)[M].Cambridge: MIT Press, 2016: 367-415.
- [5] Andrew Ng, Kian Katanforoosh, Younes Bensouda Mourri. Sequence Models, Deep Learning[EB/OL].[2025-06-01].https://www.coursera.org/learn/nlp-sequence-models.
 - [6] 邱锡鹏. 神经网络与深度学习[EB/OL]. (2021-05-17)[2025-06-01].https://nndl.github.io/.
- [7] Cobbinah M, Alnaggar A. An attention encoder-decoder RNN model with teacher forcing for predicting consumer price index[J]. Journal of Data, Information and Management, 2024, 6(1): 65-83.
- [8] Ming-Fei H, n Z, Jian-Wei L. Survey on deep generative model[J]. Acta Automatica Sinica, 2022, 48(1): 40-74.
- [9] Luo X, Chen Z. English text quality analysis based on recurrent neural network and semantic segmentation[J]. Future Generation Computer Systems, 2020, 112: 507-511.
 - [10] 胡新辰. 基于 LSTM 的语义关系分类研究[D]. 哈尔滨工业大学, 2025.
- [11] 王龙,杨俊安,陈雷,等.基于循环神经网络的汉语语言模型建模方法[J].声学技术,2015,34(5):6.
 - [12] 何彬. 面向临床文本的医学经验知识抽取研究[D]. 哈尔滨工业大学, 2018.
- [13] 李华旭. 基于 RNN 和 Transformer 模型的自然语言处理研究综述[J]. 信息记录材料, 2021, 22(12): 22.
- [14] Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need[J]. Advances in neural information processing systems, 2017, 30.
 - [15] Fei-Yan Z, Lin-Peng J, Jun D. Review of Convolutional Neural Network[J]. Chinese Journal

- of Computers, 2017.
- [16] 王辰成,杨麟儿,王莹莹,等. 基于 Transformer 增强架构的中文语法纠错方法[C]// 第十八届中国计算语言学大会(CCL 2019), 2019.
- [17] Goldberg Y, Levy O. word2vec Explained: deriving Mikolov et al.'s negative-sampling word-embedding method[DB/OL].[2025-06-19].https://arxiv.org/abs/1402.3722.
- [18] Pennington J, Socher R, Manning C. Glove: Global Vectors for Word Representation[C]// Conference on Empirical Methods in Natural Language Processing, 2014.
- [19] Shreyashree S, Sunagar P, Rajarajeswari S, et al. Inventive Computation and Information Technologies[M].Singapore:Springer, 2022: 305-320.
- [20] Korngiebel D M, Mooney S D. Considering the possibilities and pitfalls of Generative Pre-trained Transformer 3 (GPT-3) in healthcare delivery[J]. NPJ Digital Medicine, 2021, 4(1): 93.
- [21] Chowdhery A, Narang S, Devlin J, et al. Palm: Scaling language modeling with pathways[J]. Journal of Machine Learning Research, 2023, 24(240): 1-113.
- [22] Hoffmann J, Borgeaud S, Mensch A, et al. Training compute-optimal large language models[DB/OL].[2025-06-19].https://arxiv.org/abs/2203.15556.
- [23] Perez E, Kiela D, Cho K. True few-shot learning with language models[J]. Advances in neural information processing systems, 2021, 34: 11054-11070.
- [24] Achiam J, Adler S, Agarwal S, et al. Gpt-4 technical report[DB/OL].[2025-06-19]. https://arxiv.org/abs/2303.08774.
- [25] Sun Y, Wang S, Feng S, et al. Ernie 3.0: Large-scale knowledge enhanced pre-training for language understanding and generation[DB/OL].[2025-06-19].https://arxiv.org/abs/2107.02137.
- [26] Zhang X, Yu H, Fu C, et al. IOPO: Empowering LLMs with Complex Instruction Following via Input-Output Preference Optimization[DB/OL].[2025-06-19].https://arxiv.org/abs/2411.06208.
- [27] Yuan H, Li X, Zhang T, et al. Sa2VA: Marrying SAM2 with LLaVA for Dense Grounded Understanding of Images and Videos[DB/OL].[2025-06-19].https://arxiv.org/abs/2501.04001.
- [28] 讯飞晓医宣布重大升级,正式上线"星火医疗大模型 X1"功能[EB/OL]. (2025-03-04) [2025-06-19].https://cn.chinadaily.com.cn/a/202503/04/WS67c6a0b1a310510f19ee9a86.html.
- [29] 昆仑万维: "天工"大模型 4 月 17 日启动邀测[EB/OL]. (2023-04-10) [2025-06-19]. https://baijiahao.baidu.com/s?id=1762777074470585093.
- [30] Jacobs, Robert A, et al. Adaptive mixtures of local[EB/OL]. (1991-03-01) [2025-06-19]. experts.https://ieeexplore.ieee.org/abstract/document/6797059.
- [31] Fedus W, Zoph B, Shazeer N. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity[J]. Journal of Machine Learning Research, 2022, 23(120): 1-39.
- [32] Touvron, Hugo, et al. Llama: Open and efficient foundation language models[DB/OL]. [2025-06-19]. https://arxiv.org/abs/2302.13971.

- [33] Rombach, Robin, et al. High-resolution image synthesis with latent diffusion models[DB/OL]. [2025-06-19]. https://arxiv.org/abs/2112.10752.
 - [34] 张伟. 智慧教育赋能教育强国研究: 大语言模型视角[J]. 中国教育信息化, 2024(12).
- [35] 易云恒,潘济.基于大语言模型的教育教学知识问答系统的设计[J].现代信息科技,2025,9(2):189-194.
- [36] 刘明,吴忠明,杨箫,等. 教育大语言模型的内涵、构建和挑战[J].现代远程教育研究, 2024, 36(5): 50-60.
- [37] 何剑虎,王德健,赵志锐,等.大语言模型在医疗领域的前沿研究与创新应用[J]. 医学信息学杂志,2024,45(9):10-18.
- [38] 田雪晴,李泉江,游茂,等. 我国医疗机构大语言模型建设现状调查与分析[J]. 中国卫生信息管理杂志,2025,22(1):38-44.
- [39] 林建浩,孙乐轩. 大语言模型与经济金融文本分析:基本原理、应用场景与研究展望[J]. 计量经济学报,2025,5(1):1-34.
- [40] 陶江. 大语言模型下金融行业软件供应链风险研究[J]. 电脑知识与技术, 2024, 20(30): 118-120.
- [41] Liu A, Feng B, Xue B, et al. Deepseek-v3 technical report[DB/OL].[2025-06-19]. https://arxiv.org/abs/2412.19437.
- [42] Guo D, Yang D, Zhang H, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning[DB/OL].[2025-06-19]. https://arxiv.org/abs/2501.12948.
- [43] Su J, Ahmed M, Lu Y, et al. Roformer: Enhanced transformer with rotary position embedding[J]. Neurocomputing, 2024, 568: 127063.
- [44] Wei J, Wang X, Schuurmans D, et al. Chain-of-thought prompting elicits reasoning in large language models[J]. Advances in neural information processing systems, 2022, 35: 24824-24837.
- [45] Dao A, Vu D B. AlphaMaze: Enhancing Large Language Models' Spatial Intelligence via GRPO[DB/OL].[2025-06-19]. https://arxiv.org/abs/2502.14669.
- [46] 韩序,刘亮,楼文晖. 生成式人工智能大型语言模型在消化道癌症领域辅助科研创作的现状分析:基于 2024 年美国临床肿瘤学会中国学者数据[J]. 中国实用外科杂志,2024,44(8):894-899.
- [47] Zelikman E, Wu Y, Mu J, et al. Star: Bootstrap reasoning with reasoning[J]. Advances in Neural Information Processing Systems, 2022, 35: 15476-15488.
- [48] Kurokawa R, Ohizumi Y, Kanzawa J, et al. Diagnostic performances of Claude 3 Opus and Claude 3.5 Sonnet from patient history and key images in Radiology's "Diagnosis Please" cases[J]. Japanese Journal of Radiology, 2024: 1-4.
- [49] Coletta A, Dwarakanath K, Liu P, et al. LLM-driven Imitation of Subrational Behavior: Illusion or Reality?[DB/OL].[2025-06-19].https://arxiv.org/abs/2402.08755.
 - [50] Zhai Y, Tong S, Li X, et al. Investigating the catastrophic forgetting in multimodal large

language models[DB/OL].[2025-06-19].https://arxiv.org/abs/2309.10313.

- [51] 洛阳医保智能客服系统升级 DeepSeek 大模型驱动服务标准化新标杆[EB/OL]. 河北省 标准化研究院, 2025.
- [52] 合肥研究院发展出核电厂复杂系统智能故障诊断方法[EB/OL]. 中国科学院, 2021-04-30.
 - [53] 2024 人工智能安全报告[R]. 奇安信集团, 2024.
- [54] Zhao, Wayne Xin, et al. A survey of large language models.[DB/OL].[2025-06-19].https:// arxiv.org/abs/2303.18223.