

西安交通大学学报
Journal of Xi'an Jiaotong University
ISSN 0253-987X, CN 61-1069/T

《西安交通大学学报》网络首发论文

题目: ChatGPT 工作原理、关键技术及未来发展趋势
作者: 秦涛, 杜尚恒, 常元元, 王晨旭
收稿日期: 2023-06-25
网络首发日期: 2023-10-17
引用格式: 秦涛, 杜尚恒, 常元元, 王晨旭. ChatGPT 工作原理、关键技术及未来发展趋势[J/OL]. 西安交通大学学报.
<https://link.cnki.net/urlid/61.1069.T.20231016.1403.012>



网络首发: 在编辑部工作流程中, 稿件从录用到出版要经历录用定稿、排版定稿、整期汇编定稿等阶段。录用定稿指内容已经确定, 且通过同行评议、主编终审同意刊用的稿件。排版定稿指录用定稿按照期刊特定版式(包括网络呈现版式)排版后的稿件, 可暂不确定出版年、卷、期和页码。整期汇编定稿指出版年、卷、期、页码均已确定的印刷或数字出版的整期汇编稿件。录用定稿网络首发稿件内容必须符合《出版管理条例》和《期刊出版管理规定》的有关规定; 学术研究成果具有创新性、科学性和先进性, 符合编辑部对刊文的录用要求, 不存在学术不端行为及其他侵权行为; 稿件内容应基本符合国家有关书刊编辑、出版的技术标准, 正确使用和统一规范语言文字、符号、数字、外文字母、法定计量单位及地图标注等。为确保录用定稿网络首发的严肃性, 录用定稿一经发布, 不得修改论文题目、作者、机构名称和学术内容, 只可基于编辑规范进行少量文字的修改。

出版确认: 纸质期刊编辑部通过与《中国学术期刊(光盘版)》电子杂志社有限公司签约, 在《中国学术期刊(网络版)》出版传播平台上创办与纸质期刊内容一致的网络版, 以单篇或整期出版形式, 在印刷出版之前刊发论文的录用定稿、排版定稿、整期汇编定稿。因为《中国学术期刊(网络版)》是国家新闻出版广电总局批准的网络连续型出版物(ISSN 2096-4188, CN 11-6037/Z), 所以签约期刊的网络版上网络首发论文视为正式出版。

ChatGPT 工作原理、关键技术及未来发展趋势

秦涛, 杜尚恒, 常元元, 王晨旭

西安交通大学智能网络与网络安全教育部重点实验室 陕西, 西安 710049

摘要：ChatGPT 是自然语言处理领域的一项重要技术突破，在多种任务中表现出卓越的性能。本文主要探讨 ChatGPT 的演变历程、关键技术，并分析了其未来可能的发展方向。首先，介绍了 Chat GPT 的模型架构和技术演进过程。随后，重点讨论了 ChatGPT 的关键技术，包括提示学习与指令微调、思维链、人类反馈强化学习。然后，分析了由于基于概率生成原理所造成的固有局限，包括事实性错误、垂直领域深度性弱、潜在的恶意应用风险、可解释性及模型实时性差等。最后，探讨了其在典型应用中存在的问题和相应的解决途径，包括在训练评估过程中考虑道德和安全性因素，以降低潜在风险；结合外部专家知识和迁移学习，以提高模型对特定领域的理解能力，更好适应特定任务场景；引入多模态数据，以提高模型信息理解能力，增强模型通用性和泛化性。通过对 ChatGPT 模型框架、技术演变与关键技术的分析，对深入理解 ChatGPT 提供帮助；结合其原理，分析其固有缺陷，并结合实际应用中存在的问题，挖掘未来可能的研究方向，为自然语言处理领域的深入研究提供有益启发。

关键词：GPT 架构；生成模型；强化学习；人机交互

Running Principles, Key Technologies and Developing Trends of ChatGPT

QIN Tao, DU Shangheng, CHANG Yuanyuan, WANG Chenxu

MOE Key Lab for Intelligent and Network Security, Xi'an Jiaotong University, Xi'an, Shaanxi, 710049

Abstract: ChatGPT is one of the most important newly developed technologies in the area of natural language processing, and has achieved excellent performance in many areas. In this paper, we mainly discuss its architecture and technologies used. We also present the possible directions for further investigation. Firstly, we give its architecture and technology evolution process. Then, the key technologies used are discussed, including the prompt learning and instruction fine-tuning, chain of thought and reinforcement learning from human feedback. We discuss its limitations due to the probabilistic generation principles, including factual errors, poor performance in specific domain, potential malicious risk, poor interpretability and real-time. Finally, we present its possible development trends based on the problems rose by actual applications, including include the ethical and safety factors in the training process to reduce potential risks; Introduce the external expert knowledge and transfer learning method to improve the model's performance in specific domain; Improve its ability in information understanding ability based on the multimodal data. By analyzing the framework and key technologies, we try to provide assistance for deeper understanding on ChatGPT; We also discuss the possible research directions for further study in the future.

Keywords: GPT architecture; Generation model; Reinforcement learning; Human-computer interaction

1 引言

自然语言处理作为人工智能领域的关键技术之一，具有重要的应用价值，在信息提取与知识管理领域，商业机构可利用该技术布局线上系统，开发智能

客服系统自动理解客户问题并及时响应客户需求，从而协助办理业务，提升服务效率；也能通过对海量多源数据的处理分析，构建多层次多维度用户画像，制定精细化、个性化服务方案，优化服务质量。在舆情监管领域^[1]，结合分词、实体识别、热词发现和情

收稿日期：2023-06-25。作者简介：秦涛（通信作者），男，教授，博士生导师，研究方向：大数据融合分析，杜尚恒（1999—），男，硕士生。**基金资助：**国家自然科学基金（62172324）、陕西省重点研发（2023-YBGY-269，2022QCY-LL-33HZ）的资助。

感倾向分析等技术,对社交媒体数据进行情感分析^[2],尽早发现负面消极的言论和煽动性的话题,采用信息抽取与文本聚类技术从信息流中检测并聚合突发事件,并利用网络分析和深度学习技术分析事件在社交网络中的传播途径,理解信息的扩散模式从而进行事件演化与趋势预测,为舆情管控和引导提供决策支持,推动社会管理的智能化和精细化。但是,自然语言处理技术的发展仍受到可用标注数据稀缺、数据多源时变、语义信息多样复杂等问题的困扰。

正是在这样的需求推动下,领域内的技术框架不断更新进步,自然语言处理技术的演进阶段可以分为小规模专家知识、浅层机器学习^[3]、深度学习^[4]、预训练语言模型^[5]等阶段,每个技术阶段的演进周期大致为前一阶段的一半,迭代速度越发迅速。ChatGPT 作为大规模预训练模型的一种典型代表,极大的推动了自然语言处理技术的发展,引发了自然语言处理研究范式的转变。其通过大规模的预训练和上下文理解,具备了生成自然语言文本的能力,可以进行对话、回答问题和提供信息等任务,使得与人类交互的能力更加自然和灵活。

为进一步理解 ChatGPT,本文首先介绍 ChatGPT 的模型架构和技术演进过程,然后回顾了其所用的核心技术,包括提示学习、思维链和基于人类反馈的强化学习,这些技术共同构成了 ChatGPT 的基础框架,使其能够在各种情景下生成连贯且自然的文本回应。然后,结合 ChatGPT 运行原理,本文分析了其面临的缺陷与挑战,包括生成不准确或具有误导性的信息、潜在的恶意应用风险以及对话中的道德和隐私问题等。最后,针对 ChatGPT 在特定领域的缺陷与不足,结合实际应用,探讨了 ChatGPT 未来可能的发展方向,包括对训练语料进行道德筛选、采用迁移学习^[6]和领域适应技术、引入外部专家知识^[7]、增强多模态处理能力^[8]等途径来优化改进。

2 GPT 架构及演变过程

GPT^[9] (Generative Pre-trained Transformer) 是由 OpenAI 提出的采用 transformer 解码器的预训练模型,采用预训练加微调的范式,为深入理解 ChatGPT,本节简要分析 ChatGPT 的模型架构和其演进进程。

2.1 ChatGPT 整体架构

ChatGPT 的主体架构遵从“基础语料+预训练+微调”的基本范式,如图 1 所示。“预训练+微调”是指首先在大数据集上训练得到一个具有强泛化能

力的模型(预训练模型),然后在下游任务上进行微调的过程,是基于模型的迁移方法。

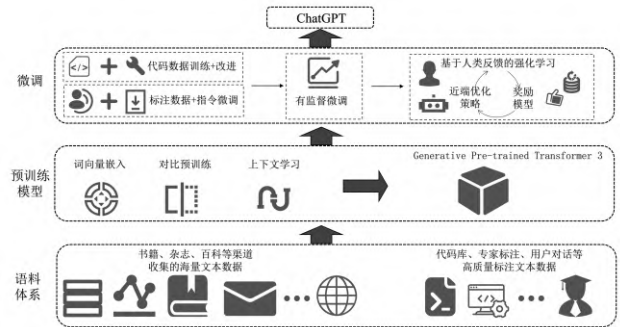


图 1 ChatGPT 架构示意图

Fig. 1 Diagram of ChatGPT architecture

(1)海量高质量的基础语料是 ChatGPT 技术突破的关键因素。其语料体系包括预训练语料与微调语料,后者包括代码和对话微调语料。预训练语料包括从书籍、杂志、百科等渠道收集的海量文本数据(具体分布如表 1^[10]所示),使模型学习到语言的逻辑关系表达方式;微调语料包括从开源代码库爬取、专家标注、用户对话等方式加工而成的高质量标注文本数据,进一步增强其对话能力。

表 1 GPT 系列预训练语料数据大小(单位: GB)^[10]

Table 1 GPT series pre-training corpus data size (Unit: GB)

模型	书籍 语料	期刊 语料	Reddit 链接	爬虫 语料	维基 百科
GPT-1	4.6	/	/	/	/
GPT-2	/	/	40	/	/
GPT-3	21	101	50	570	11.4

(2) 预训练是构建大规模语言模型的基础,指先在大规模训练数据上进行大量通用的训练,采用无监督学习方法进行训练,以得到通用且强泛化能力的语言模型。在大规模数据的基础上,通过预训练,模型初步具备了人类语言理解和上下文学习的能力,能够捕捉文本片段和代码片段的语义相似性特征,从而生成更准确的文本和代码向量,为后续微调任务提供支持。

(3) 微调是实现模型实际应用的保障,是指在特定任务的数据集上对预训练模型进行进一步的训练,通常包括冻结预训练模型的底层层级(如词向量)与调整上层层级(如分类器)的权重。对预训练模型微调将大大缩短训练时间,节省计算资源并加快训练收敛速度。ChatGPT 在具有强泛化能力的预训练模型基础上,通过整合基于代码数据的训练和基于指令的微调,利用特定的数据集进行微调,使之具有更

强的问答式对话文本生成能力。其“预训练+微调”的流程图如图2所示。

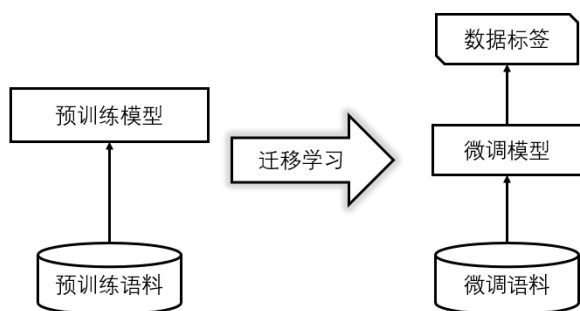


图2 “预训练+微调”流程图

Fig. 2 Pre-training and fine-tuning flow chart

2.2 GPT-1:奠定 AR 式建模思路

GPT-1^[9]是比 BERT^[11]提出更早的预训练模型，但与 BERT 相比效果较差。GPT-1 奠定了关键的技术路径，后续的系列模型采用类似的架构(例如 BART^[12]和 GPT-2^[13])以及预训练策略^[14-16]。GPT 系列是一种基于自回归解码的、仅仅包含解码器的 Transformer 架构开发的生成式预训练模型，这种架构称为自回归类 AR(Autoregressive)，它利用多层堆叠的 Transformer 解码器架构进行解码。

自回归是统计学中处理时间序列的方法，用同一变量之前各时刻的观测值预测该变量当前时刻的观测值。类似地，自回归生成模型的基本思想是在序列生成的过程中，每个位置的元素都依赖于前面已经生成的元素。自回归模型适用于各种序列到序列的任务，它又分为线性自回归和神经自回归两种，基于 Transformer 解码器的自回归模型属于后者。其生成过程如图3所示。

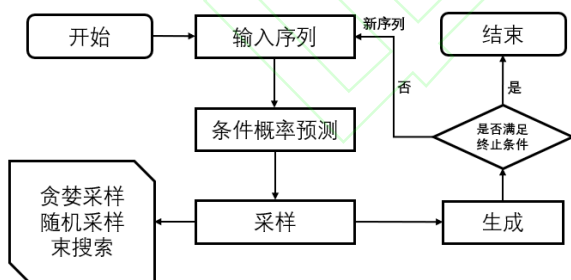


图3 自回归生成模型生成过程

Fig. 3 Autoregressive model generation process

GPT-1 以 Transformer 解码器作为核心组件，在大规模未标记的文本数据上进行自监督学习，通过最大似然估计来调整模型参数，使得模型能够更好地预测下一个词。这种预训练策略允许模型从大规模的文本数据中学习通用的语言规律和语义关系。在生成序列时，模型将已生成序列作为已知条件，

并利用 Transformer 解码器来预测下一个元素的概率分布，输入前 L 个元素构成的序列 $x = \{x_1, \dots, x_L\}$ ，预测 $\tilde{x} = \{x_2, \dots, x_{L+1}\}$ 。每次生成的元素会添加到序列中。生成的过程会一直持续，直到达到预定的生成长度或遇到终止符号，如图4所示。在预训练结束之后可以根据不同的下游任务或者特定的语境场景进行微调。

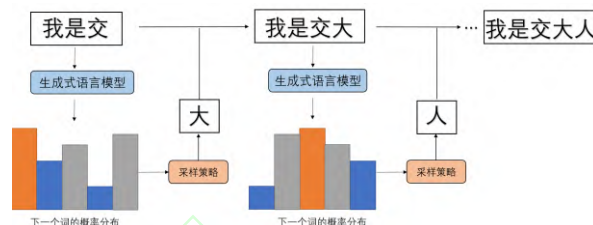


图4 GPT-1 生成过程

Fig. 4 GPT-1 generation process

2.3 GPT2:引入提示学习

GPT-2^[13]通过模型结构的改进，在下游任务的微调上取得了更好的结果。GPT-2 在以下两个方面进行了优化：

(1)扩大参数规模：使用更多高质量的网页数据，将模型参数规模扩大到 1.5B。

(2)更自然的任务模型融合方式：GPT-2 将下游任务通过 prompt 方式加入到预训练模型中，从而使模型获得零样本的能力，即引入了一种多任务求解的概率形式，通过给定输入与任务条件对结果进行预测。

虽然 GPT-2 在下游任务的微调中并没有 BERT 表现优越，但其更自然的任务融合方式为后续 ChatGPT 的指令理解能力奠定了基础，即对输入文本信息按照特定模板进行处理，将任务重构成一个更能充分利用语言模型处理的形式。

提示学习的工作流如图5所示，主要包含以下4部分：提示词模板构造；提示词答案空间映射构造；将文本带入模板嵌入并使用预训练模型预测答案；将预测结果映射回标签。按照模板构造差异提示学习方法可分为硬模板方法和软模板方法。具体来说，硬模板方法先在少量监督数据上对每个提示词训练模型，再在无监督数据上将同一样本的多个提示词预测结果进行集成，建模为 $s_p(l|x) = M(v(l)|P(x))$ ，表示在给定提示词下对应词语的分数，归一化得到概率分布 $q_p(l|x) = \frac{e^{s_p(l|x)}}{\sum_{l' \in \zeta} e^{s_p(l'|x)}}$ 作为无监督数据的软标签。而软模板方法则是直接在输入端嵌入若干可被优化的提示词流，使其自动化地寻找连续空间内

的知识模板,从而不依赖人工设计。以生成任务中的前缀微调(prefix tuning)^[17]为例,模型在每一层transformer前加入前缀连续向量并使用训练矩阵P来存储,将输入表示为 $Z = [prefix; x; y]$,在整个训练期间冻结预训练模型其他参数不变,只更新用于给定任务的前缀参数。

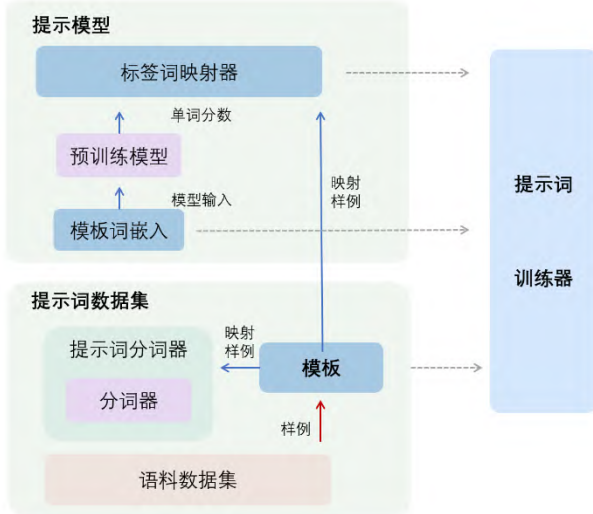


图5 提示学习框架示意图

Fig. 5 The diagram of Prompt learning

通过上述方式,每个自然语言处理的任务都可以被视作基于世界文本子集的单词预测问题^[18]。这种思想表明,只要模型足够大、学到的知识足够丰富,任何有监督任务都可以通过无监督的方式来完成,其下游任务中的对话任务^[19,20]更是进行了全面的微调,为后续的ChatGPT对话奠定了基础。

2.4 GPT-3:量变引起质变

在GPT-2的基础上,GPT-3^[21]通过扩展生成预训练架构,实现了容量飞跃。GPT-3的显著特点就是大。由于GPT-2的实验中随着参数规模的增大其效果的增长依旧显著^[22,23],因此选择继续扩大参数规模,用更多优质的数据,一方面是模型本身规模大,参数数量众多,具有96层Transformer解码器,每一层有96个128维的注意力头,单词嵌入的维度也达到了12288维;另一方面是训练过程中使用到的数据集规模大,达到了45TB,参数规模达到175B。

此外,GPT-3在模型能力上转变思路,采用情景学习的思想,使模型能够在少样本学习上取得较好

的效果。其进行了大量实验证明GPT-3在少样本情况下具有良好的表现,如图6所示。

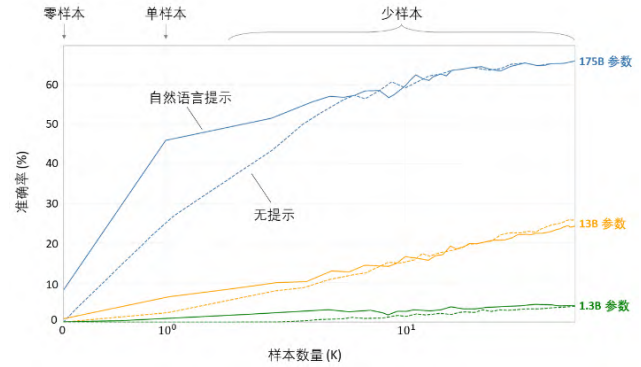


图6 GPT-3 随样本量变化性能

Fig. 6 GPT-3 performance variation with sample size

由于规模巨大,GPT-3在各领域均有广泛的应用,衍生多种应用生态,被视为从预训练模型发展到大规模模型过程中的一个里程碑。

2.5 ChatGPT: 实现人机混合增强

尽管GPT-3拥有大量知识,但生成文本质量不一且语言表达冗余,ChatGPT通过人工标注的微调,引导模型输出更有价值的文本结果,即实施了人类反馈的强化学习机制。

OpenAI对于混合人类反馈增强机器学习的研究可以追溯到2017年^[24],并且在当年发布了近端策略优化(Proximal Policy Optimization, PPO)^[25]算法作为强化学习的基础算法,该算法通过多个训练步骤实现小批量更新,以克服传统策略梯度算法中步长难确定可能导致学习性能下降的问题,引入保守的策略更新机制有效缓解了策略更新过快导致的不稳定性,提高了训练的鲁棒性和稳定性。PPO算法根据当前策略与环境互动产生轨迹,并记录各状态、动作与奖励,使用轨迹信息更新策略限制策略步长,使得目标散度既然足以显著改变策略,又足以使更新稳定,防止新旧策略过远,并在每次更新后重新计算优势函数。算法重复上述步骤直至策略收敛或达到预定训练周期,伪代码如算法1所示。OpenAI在GPT-2时便开始使用上述强化学习算法^[24,25]来进行微调,同年以类似方法训练了文本摘要模型^[26]。

算法1 近端策略优化算法

input: 初始化的网络参数 θ

output: 学习得到的策略网络

```

1: for  $i \in \{1, \dots, N\}$  do
2:   使用当前策略  $\pi_\theta$  收集  $\{s_t, a_t, r_t\}$ 
3:   估计优势函数  $\hat{A}_t = \sum_{t' > t} \gamma^{t'-t} r_{t'} - V_\phi(s_t)$ 
4:    $\pi_{old} \leftarrow \pi_\theta$ 
5:   for  $j \in \{1, \dots, M\}$  do
6:      $J_{PPO}(\theta) = \sum_{t=1}^T \frac{\pi_\theta(a_t|s_t)}{\pi_{old}(a_t|s_t)} \hat{A}_t - \lambda KL[\pi_{old}|\pi_\theta]$ , 用  $J_{PPO}(\theta)$  更新  $\theta$ 
7:   end for
8:   for  $j \in \{1, \dots, B\}$  do
9:      $L_{BL}(\phi) = -\sum_{t=1}^T (\sum_{t' > t} \gamma^{t'-t} r_{t'} - V_\phi(s_t))^2$ , 用  $L_{BL}(\phi)$  更新  $\phi$ 
10:  End for
11:  if  $KL[\pi_{old}|\pi_\theta] > \beta_{high} KL_{target}$  then
12:     $\lambda \leftarrow \alpha \lambda$ 
13:  else if  $KL[\pi_{old}|\pi_\theta] < \beta_{high} KL_{target}$  then
14:     $\lambda \leftarrow \alpha / \lambda$ 
15:  end if
16: end for

```

ChatGPT 的前身, InstructGPT^[27]模型正式使用了基于人类反馈的强化学习 (RLHF) 算法, 通过结合智能体自主学习与人类专家反馈两种策略, 选择基于策略梯度的算法搭建强化学习模型从而训练智能体, 并在每个时间步上记录智能体行为并且由人类专家对其进行评估反馈, 以进行参数更新优化行为策略。该算法的第一阶段^[27]也就是第三章介绍的指令调优。除了提高指令理解能力外, RLHF 算法还有助于缓解大模型产生危害或不当内容的问题, 这也是大模型在安全实践部署的关键。OpenAI 在一篇技术文章^[27]中描述了对齐研究方法, 该文章总结了三个有希望的方向即“使用人类反馈训练 AI 系统, 帮助人类评估和进行对齐研究”。

2.6 GPT4: 多模态升级

GPT-4 是对 ChatGPT 的多模态升级, 可对图文输入产生应答文字, 并可引用于视觉分类分析、隐含语义等领域。多模态输入能力对语言模型至关重要, 使其获得了除文本描述外的常识性知识, 并为多模态感知与语义理解的结合提供了可能性。

新范式可归纳为“预训练+提示+预测”。各种下游任务被调整为类似预训练任务的形式, 尤其 GPT-4 的多模态提示工程针对多模态数据集, 涉及合适的

模型架构参数、精心设计的提示格式结构和选定的数据微调模型, 来使得模型生成高质量文本。

3 ChatGPT 的核心技术

3.1 提示学习与指令精调

传统的监督学习使用包含输入 x 与标签 y 的数据集来训练一个模型 $P(y|x; \theta)$, 从而学习模型参数 θ 预测条件概率。提示学习^[28]试图学习模拟概率 $P(x; \theta)$ 的 x 本身来预测 y , 从而减少或消除对大型监督数据集的需求。

具体来说, 将输入 x 通过提示函数 $f_{prompt}(x)$ 添加槽 $[Z]$ 转化为特定形式 x' , 定义 Z 为回答 z 允许值的集合, 定义填充函数 $f_{fill}(x', z)$ 来将可能的回答 z 填充 x' 中的槽 $[Z]$ 。最后, 使用特定搜索函数来计算相应的填充提示概率, 得到最终回答 \hat{z} , 定义如下

$$\hat{z} = \underset{z \in Z}{\text{search}} P(f_{fill}(x', z); \theta) \quad (1)$$

通过编辑任务的输入, 提示学习在形式上模拟模型训练中的数据与任务。以情感分类任务为例, 相

比于监督学习中输入一句话,模型输出情感分类判断,提示学习是设计一种模板将原有语句嵌入其中,为模型留出判断类别的位置,让模型做类似完形填空的工作生成情感类别。提示学习旨在激发语言模型的补全能力,指令精调(instruction tuning)^[29]则是提示学习的加强版,激发模型的理解能力。通过指令调优,模型能够在不使用显式示例的情况下遵循新任务的任务指令,从而提高了泛化能力^[29,30],即便在多语言环境下也有卓越能力^[31]。这种学习人类交互模式的分布让模型可以更好的理解人类意图^[32]、与人类行为对齐^[27]。从解释性上来说,这类似于打开大门的钥匙,从大模型在预训练中学习到的庞大知识中激活特定的部分完成指定任务。而 ChatGPT 能响应人类指令的能力就是指令微调的直接产物,对没有见过的指令做出反馈的泛化能力是在指令数量超过一定程度之后自动出现的,TO^[33]、Flan^[29]等论文都进一步证明了这一点。

对于模型未训练的新任务,只需设计任务的语言描述,并给出任务实例作为模型输入,即可让模型从给定的情景中学习新任务并给出恰当的回答结果。

这种训练方式能够有效提升模型小样本学习的能力。

3.2 思维链

谷歌研究人员 Jason Wei 等提出了思维链(Chain of Thought, COT)^[34]的概念,即在小样本提示学习中插入一系列中间推理的步骤示范,从而有效提高语言模型的推理能力。与一般的提示词不同,思维链提示由多个分别解释子问题的中间步骤组成,提示词模式从之前的(问题,答案)变成(输入问题,思维链,输出问题)。如图 7 所示,一般的提示词模板通过输入内嵌入样例,使得模型学习任务答案,而思维链提示词增加推理步骤,参考人类解决问题方法,嵌入自然语言形式的推理步骤直至答案生成。在思维链的加持下,模型将问题分解为一系列的分步推理,根据前一步骤结果与当前问题要求共同推断下一步骤。通过这种逐步推理的方式,模型可以逐渐获得更多信息,并在整个推理过程中累积正确的推断,从而大幅度提升模型在复杂推理时的准确率,同时也为模型的推理行为提供了一个可解释的窗口^[35,36]。

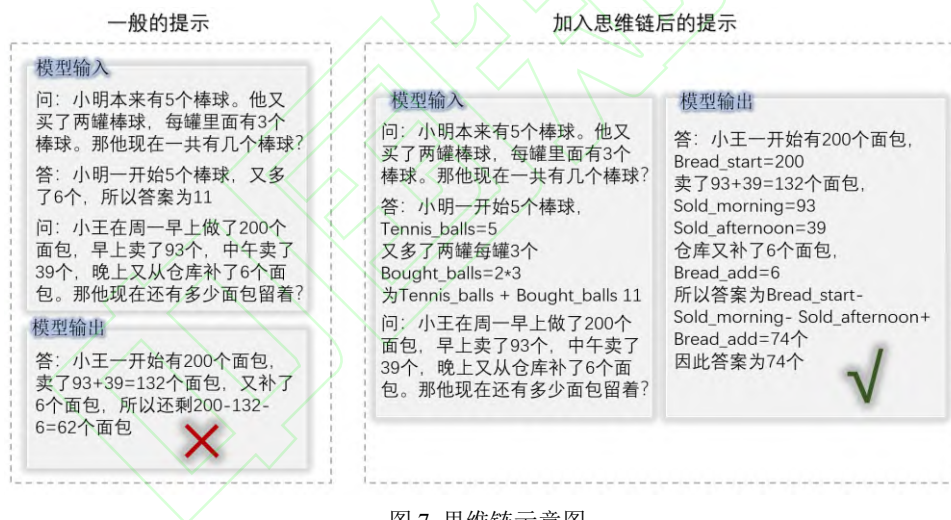


图 7 思维链示意图

Fig. 7 Diagram of chain of thought

3.3 人类反馈强化学习

人类反馈强化学习^[24,37] (Reinforcement Learning from Human Feedback, RLHF)是 ChatGPT 实现理解人类指令、对齐人类行为^[38,39]的重要关键技术。如图 8 所示为模型训练过程。此算法^[40]在强化学习^[41]的框架下大体可以分为两个阶段:

(1)奖励模型训练:该阶段旨在获取拟合人类偏好的奖励模型。奖励模型以提示和回复作为输入,计算标量奖励值作为输出。奖励模型的训练过程通过拟

合人类对于不同回复的倾向性实现。具体而言,首先基于在人类撰写数据上精调的模型,针对同一提示采样多条不同回复。然后,将回复两两组合构成一条奖励模型训练样本,由人类给出倾向性标签。最终,奖励模型通过每条样本中两个回复的奖励值之差计算倾向性概率拟合人类标签,完成奖励模型的训练

(2)生成策略优化:给定训练的奖励模型,ChatGPT 的参数将被视为一种策略,在强化学习的框架下进行训练。首先,当前策略根据输入的查询采样回复。然后,奖励模型针对回复的质量计算奖励、

反馈回当前策略用以更新。值得注意的是,为防止上述过程的过度优化,损失函数同时引入了词级别的

KL 惩罚项。此外,为了避免在公开 NLP 数据集上的性能退化策略更新过程兼顾了预训练损失。

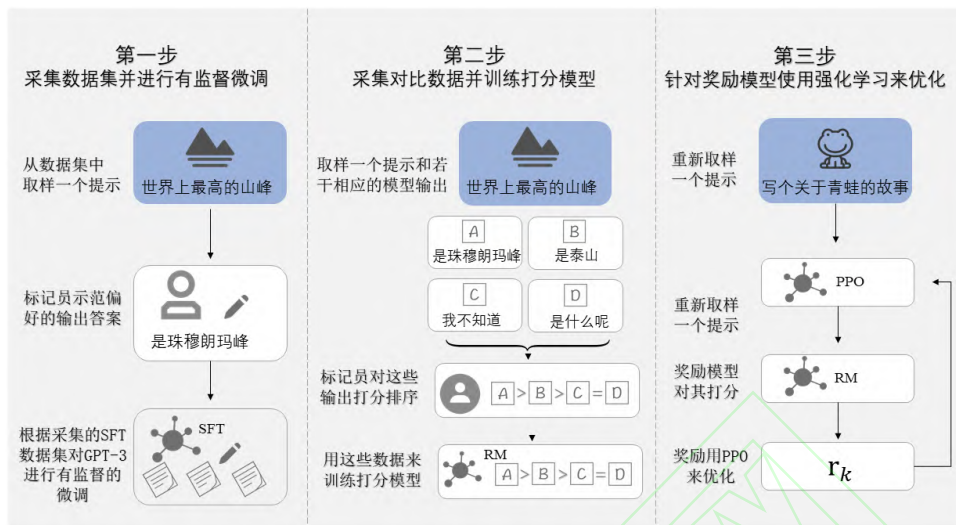


图 8 ChatGPT 训练过程示意图
Fig. 8 Diagram of ChatGPT training process

4 ChatGPT 面临的问题

虽然 ChatGPT 在多个任务中都表现出不错的性能,其现有运行原理决定了其有很多局限性:

(1)对某个领域的深入程度不够^[42,43],因此生成的内容可能不够合理。此外,ChatGPT 也存在潜在的偏见问题,因为它是基于大量数据训练的,训练数据中的固有偏差会渗透到神经网络中,导致模型会受到数据中存在的偏见的影响^[36,44]。

(2)对抗鲁棒性^[45]。对抗鲁棒性在自然语言处理与强化学习中是决定系统适用性的关键要素^[46],对于干扰示例 $x' = x + \delta$ 其中, x 指原始输入, δ 指扰动,高鲁棒性系统会产生原始输出 y ,而低鲁棒性系统会产生不一样的输出 y' 。ChatGPT 容易受到对抗性攻击,如数据集中攻击^[47]、后门攻击^[48]和快速特定攻击^[49]等,从而诱导模型产生有害输出。

(3)安全保障。由于 ChatGPT 是一种强大的人工智能技术,它可能被恶意利用,造成严重的安全隐患及产生法律风险^[50]。同时,它的答复尚不明确是否具有知识产权,从而产生不利的社会影响^[51]。因此,开发者在设计和使用 ChatGPT 时,需要采取相应措施如去偏方法和校准技术^[52]来保障安全性问题。

(4)推理可信度。虽与其他神经网络类似,ChatGPT 很难精确地表达其预测的确定性^[53],即所谓的校准问题,导致代理输出与人类意图不一致^[54]。它有时会回答荒谬的内容,这也是目前发现的最为

普遍的问题,即对于不知道或不确定的事实,它会强行根据用户的输入毫无根据地开展论述^[55],产生偏离事实的文本。

(5)可解释性差。黑盒特性使得 ChatGPT 的回答看似合理但却无迹可寻,同时由于它没有办法通过充足的理由去解释它的回答是否正确,导致在一些需要精确、严谨的领域没有办法很好的应用^[56]。此外,它也可能在表述的时候存在语法错误或不合理的表述。

(6)无法在线更新近期知识。目前的范式增加新知识的方式只能通过重新训练大模型,但是微调的成本非常高。现有研究探索了利用外部知识源来补充大模型^[57],利用检索插件来访问最新的信息源^[58],然而,这种做法似乎仍然停留在表面上。研究表明,很难直接修正内在知识或将特定知识注入大模型,这仍然是一个悬而未决的研究问题^[59]。

5 ChatGPT 的进化展望

ChatGPT 在自然语言处理技术的发展中有里程碑式的意义,在语言和意图理解、推理、记忆以及情感迁移方面具有强大的能力,在决策和计划方面表现出色,只需一个任务描述或演示,就可以有效地处理以前未见过的任务。此外,ChatGPT 可以适应不同的语言、文化和领域,具有通用性,减少了复杂的培训过程和数据收集的需要^[60],在各领域得到了广泛应用。

在学术界积极探索能力背后的技术原因的同时,工业界已将 ChatGPT 优异的对话生成能力融入各种场景中,根据对话对象的不同,我们将应用分为以下几种层次:

(1)数据生成加工。利用 ChatGPT 强大的信息搜索与整合能力,用户根据自身需求要求直接返回特定数据,主要应用场景包括文案生成、代码生成和对话生成等。同时,其可以充当知识挖掘工具对数据进行再加工,一些在线应用可帮助翻译、润色等,如文档分析工具 ChatPDF。

(2)模型调度。ChatGPT 可以调用其他机器学习模型共同完成用户需求并输出结果,如微软近期发布的 HuggingGPT。作为人类与其他模型的智能中台,其有望解决 AI 赋能长期面临的痛点问题,实现模块化模型管理、简化技术集成部署,提高赋能效率。

(3)人机混合交互。ChatGPT 一定程度上统一了人类语言与计算机语言,使得人机交互界面从键盘鼠标图形进化到自然语言接口。微软的 365 Copilot 将其嵌入到 Office,极大提高了人机自然交互体验;OpenAI 近期发布的 Plugins 插件集尝试了大语言模型应用的开发框架。在未来其有望成为智能时代的操作系统,调用更广泛的应用程序解决实际问题。

结合实际情况中可能存在的问题和实际应用需求,我们对 ChatGPT 在具体领域存在问题的可能解决方案进行了分析讨论。

(1)商业服务优化领域:提升商品服务质量^[61]需要对用户评论进行细粒度情感分析,通常需要对特定商品领域的知识有深入的理解。然而,ChatGPT 作为通用语言模型,可能缺乏不同商品领域中的专业知识,这可能导致模型无法准确识别和分析特定领域的优缺点。可以考虑利用迁移学习的方法,将 ChatGPT 在通用领域中的知识迁移到特定领域中,使 ChatGPT 更加适应特定领域的问题和需求。

(2)智慧医疗领域:ChatGPT 在医疗领域可以做为辅助工具用作医疗诊断与肿瘤图像分割^[62],有助于精准医疗、靶向治疗等方案的落实。然而,目前 ChatGPT 主要针对文本进行处理,对于其他模态的信息理解相对较弱,这使得模型应用仅限制在辅助诊断和医疗数据挖掘等方面,无法融合其他模态的信息来增强模型通用性与泛化性。因此为了实现更加有效表达的通用人工智能模型,需要进行多模态联合学习,关注内容关联特性与跨模态转换问题。此外,风险责任问题、沟通限制状况以及模型引发的算法偏见与个人隐私安全问题同样不容忽视。

(3)舆情监管引导领域:舆情引导和特定内容生成^[63]需要在构建训练数据阶段进行意图对齐和质量筛选。由于 GPT 系列的训练语料来自于西方的语言价值框架,受到模型训练数据的偏见和倾向性影响,ChatGPT 生成内容中存在对于中国的大量偏见言论,不一定符合中国的价值观,这可能引发舆情操纵和认知战^[64]的风险。因此训练国产大模型时需要进行对训练数据进行筛选,构建合适公正的中文语料,并不断维护更新基础词库。

很多研究者认为 ChatGPT 开启了第四次技术革命,其作为催化剂整合人工智能学科,并激发学术界与工业界深入探讨和实践交叉学科和跨学科应用^[65]的可能性,科技部近期启动的“AI for Science”专项部署工作也从一定程度反映了国家导向。未来其从应用拓展上将呈现垂直化、个性化与工程化,如何增强其人机交互协同性,如考虑生物学特性、身体感知等因素;以及如何增强模型可信性,构建新的可信测试基准,都是未来可能的发展趋势。

6 总结

本文探讨了 ChatGPT 在自然语言处理领域发展中的地位以及未来可能的发展方向,着重分析了 GPT 系列模型的演进以及核心技术,包括语料体系、提示学习、思维链和基于人类反馈的强化学习等。随后,分析了其存在的显著缺陷,如理解与推理能力的局限性、专业知识的不深入、事实的不一致性以及信息安全泄露等风险;最后,结合实际应用,ChatGPT 有着很大的改进和发展空间,包括采用迁移学习和领域适应技术、引入外部专家知识、增强多模态处理能力、筛选训练语料等都是可能的解决方案与发展趋势。通过上述分析,本文对深入理解 ChatGPT 和在相关领域展开进一步研究提供启发。

参考文献:

- [1] 郝亚洲,郑庆华,陈艳平,闫彩霞.面向网络舆情数据的异常行为识别[J].计算机研究与发展,2016,53(3): 611-620.
HAO Yazhou, ZHENG Qinghua, CHEN Yanping, YAN Caixia. Recognition of Abnormal Behavior Based on Data of Public Opinion on the Web [J]. Journal of Computer Research and Development, 2016, 53(3): 611-620.
- [2] 何炎祥,孙松涛,牛菲菲,等.用于微博情感分析的一种情感语义增强的深度学习模型[J].计算机学报,2017,

- 40(4): 773-790.
- HE Yanxiang, SUN Songtao, NIU Feifei, et al. A Deep Learning Model Enhanced with Emotion Semantics for Microblog Sentiment Analysis [J]. Chinese Journal of Computers, 2017, 40(4): 773-790.
- [3] Rumelhart D E, Hinton G E, Williams R J. Learning representations by back-propagating errors[J]. nature, 1986, 323(6088): 533-536..
- [4] 奚雪峰, 周国栋. 面向自然语言处理的深度学习研究[J]. 自动化学报, 2016, 42(10): 1445-1465.
- XI Xuefeng, ZHOU Guodong. A Survey on Deep Learning for Natural Language Processing. ACTA AUTOMATICA SINICA, 2016, 42(10): 1445-1465.
- [5] Qiu X, Sun T, Xu Y, et al. Pre-trained models for natural language processing: A survey[J]. Science China Technological Sciences, 2020, 63(10): 1872-1897.
- [6] Weiss K, Khoshgoftaar T M, Wang D D. A survey of transfer learning[J]. Journal of Big data, 2016, 3(1): 1-40.
- [7] Menon T, Pfeffer J. Valuing internal vs. external knowledge: Explaining the preference for outsiders[J]. Management science, 2003, 49(4): 497-513.
- [8] Baltrušaitis T, Ahuja C, Morency L P. Multimodal machine learning: A survey and taxonomy[J]. IEEE transactions on pattern analysis and machine intelligence, 2018, 41(2): 423-443.
- [9] Radford A, Narasimhan K, Salimans T, et al. Improving language understanding by generative pre-training[J]. 2018.
- [10] Thompson A D. What's in My AI? [EB/OL]. [2023-03-12]. <https://lifearchitected.ai/whats-in-my-ai/>.
- [11] Kenton J D M W C, Toutanova L K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding[C]//Proceedings of NAACL-HLT. 2019: 4171-4186.
- [12] Lewis M, Liu Y, Goyal N, et al. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Compr[15][15]ehension[C]//Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. 2020: 7871-7880.
- [13] Radford A, Wu J, Child R, et al. Language models are unsupervised multitask learners[J]. OpenAI blog, 2019, 1(8): 9.
- [14] Zhuang L, Wayne L, Ya S, et al. A robustly optimized BERT pre-training approach with post-training[C]//Proceedings of the 20th chinese national conference on computational linguistics. 2021: 1218-1227.
- [15] Sanh V, Webson A, Raffel C, et al. Multitask Prompted Training Enables Zero-Shot Task Generalization[C]//The Tenth International Conference on Learning Representations. 2022.
- [16] Wang T, Roberts A, Hesslow D, et al. What Language Model Architecture and Pretraining Objective Works Best for Zero-Shot Generalization?[C]//International Conference on Machine Learning. PMLR, 2022: 22964-22984.
- [17] Li X L, Liang P. Prefix-tuning: Optimizing continuous prompts for generation[J]. arXiv preprint arXiv:2101.00190, 2021.
- [18] Radford A, Wu J, Child R, et al. Language models are unsupervised multitask learners[J]. OpenAI blog, 2019, 1(8): 9.
- [19] Zhang Y, Sun S, Galley M, et al. DIALOGPT: Large-Scale Generative Pre-training for Conversational Response Generation[C]//Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations. 2020: 270-278.
- [20] Savary A. Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL 2020): Tutorial Abstracts[J]. 2020.
- [21] Brown T, Mann B, Ryder N, et al. Language models are few-shot learners[J]. Advances in neural information processing systems, 2020, 33: 1877-1901.
- [22] Kaplan J, McCandlish S, Henighan T, et al. Scaling laws for neural language models[J]. arXiv preprint arXiv:2001.08361, 2020.
- [23] Hoffmann J, Borgeaud S, Mensch A, et al. Training compute-optimal large language models[J]. arXiv preprint arXiv:2203.15556, 2022.
- [24] Christiano P F, Leike J, Brown T, et al. Deep reinforcement learning from human preferences[J]. Advances in neural information processing systems, 2017, 30.
- [25] Schulman J, Wolski F, Dhariwal P, et al. Proximal policy optimization algorithms[J]. arXiv preprint arXiv:1707.06347, 2017.
- [26] Stiennon N, Ouyang L, Wu J, et al. Learning to summarize with human feedback[J]. Advances in Neural Information Processing Systems, 2020, 33: 3008-3021.
- [27] Ouyang L, Wu J, Jiang X, et al. Training language models to follow instructions with human feedback[J]. Advances in Neural Information Processing Systems, 2022, 35: 27730-

- 27744.
- [28] Schick T, Schütze H. Exploiting Cloze-Questions for Few-Shot Text Classification and Natural Language Inference[C]//Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume. 2021: 255-269.
- [29] Wei J, Bosma M, Zhao V, et al. Finetuned Language Models are Zero-Shot Learners[C]//International Conference on Learning Representations. 2021.
- [30] Chung H W, Hou L, Longpre S, et al. Scaling instruction-finetuned language models[J]. arXiv preprint arXiv:2210.11416, 2022.
- [31] Muennighoff N, Wang T, Sutawika L, et al. Crosslingual generalization through multitask finetuning[J]. arXiv preprint arXiv:2211.01786, 2022.
- [32] Wei J, Tay Y, Bommasani R, et al. Emergent Abilities of Large Language Models[J]. Transactions on Machine Learning Research, 2022.
- [33] Sanh V, Webson A, Raffel C, et al. Multitask Prompted Training Enables Zero-Shot Task Generalization[C]//ICLR 2022-Tenth International Conference on Learning Representations. 2022.
- [34] 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.
- [35] Fu Y, Peng H, Khot T. How does gpt obtain its ability? tracing emergent abilities of language models to their sources[J]. Yao Fu's Notion, 2022.
- [36] Ouyang L, Wu J, Jiang X, et al. Training language models to follow instructions with human feedback[J]. Advances in Neural Information Processing Systems, 2022, 35: 27730-27744.
- [37] Zhang S, Roller S, Goyal N, et al. Opt: Open pre-trained transformer language models[J]. arXiv preprint arXiv:2205.01068, 2022.
- [38] Askell A, Bai Y, Chen A, et al. A general language assistant as a laboratory for alignment[J]. arXiv preprint arXiv:2112.00861, 2021.
- [39] Bai Y, Jones A, Ndousse K, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback[J]. arXiv preprint arXiv:2204.05862, 2022.
- [40] Glaese A, McAleese N, Trębacz M, et al. Improving alignment of dialogue agents via targeted human judgements[J]. arXiv preprint arXiv:2209.14375, 2022.
- [41] 刘全, 翟建伟, 章宗长, 等. 深度强化学习综述[J]. 计算机学报, 2018, 41(1): 1-27.
- LIU Quan, ZHAI Jianwei, ZHANG Zongchang, et al. A Survey on Deep Reinforcement Learning [J]. Chinese Journal of Computers, 2018, 41(1): 1-27.
- [42] Fu Y, Peng H, Khot T. How does gpt obtain its ability? tracing emergent abilities of language models to their sources[J]. Yao Fu's Notion, 2022.
- [43] Ye J, Chen X, Xu N, et al. A comprehensive capability analysis of gpt-3 and gpt-3.5 series models[J]. arXiv preprint arXiv:2303.10420, 2023.
- [44] Kenton Z, Everitt T, Weidinger L, et al. Alignment of language agents[J]. arXiv preprint arXiv:2103.14659, 2021.
- [45] Zheng R, Xi Z, Liu Q, et al. Characterizing the Impacts of Instances on Robustness[C]//Findings of the Association for Computational Linguistics: ACL 2023. 2023: 2314-2332.
- [46] Boloor A, Garimella K, He X, et al. Attacking vision-based perception in end-to-end autonomous driving models[J]. Journal of Systems Architecture, 2020, 110: 101766.
- [47] Gu T, Dolan-Gavitt B, Garg S. Badnets: Identifying vulnerabilities in the machine learning model supply chain[J]. arXiv preprint arXiv:1708.06733, 2017.
- [48] Chen X, Salem A, Chen D, et al. Badnl: Backdoor attacks against nlp models with semantic-preserving improvements[C]//Annual computer security applications conference. 2021: 554-569.
- [49] Perez F, Ribeiro I. Ignore previous prompt: Attack techniques for language models[J]. arXiv preprint arXiv:2211.09527, 2022.
- [50] Borji A. A categorical archive of chatgpt failures[J]. arXiv preprint arXiv:2302.03494, 2023.
- [51] Caliskan A, Bryson J J, Narayanan A. Semantics derived automatically from language corpora contain human-like biases[J]. Science, 2017, 356(6334): 183-186.
- [52] Guo Y, Yang Y, Abbasi A. Auto-debias: Debiasing masked language models with automated biased prompts[C]//Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). 2022: 1012-1023.
- [53] Huang X, Ruan W, Huang W, et al. A Survey of Safety and Trustworthiness of Large Language Models through the Lens of Verification and Validation[J]. arXiv preprint

- arXiv:2305.11391, 2023.
- [54] Qin Y, Hu S, Lin Y, et al. Tool learning with foundation models[J]. arXiv preprint arXiv:2304.08354, 2023.
- [55] Zhuo T Y, Huang Y, Chen C, et al. Exploring ai ethics of chatgpt: A diagnostic analysis[J]. arXiv preprint arXiv:2301.12867, 2023.
- [56] 桑基韬, 于剑. 从 ChatGPT 看 AI 未来趋势和挑战[J]. 计算机研究与发展, 2023, 60(6): 1191-1201.
- Sang Jitao, Yu Jian. ChatGPT: A Glimpse into AI's Future[J]. Journal of Computer Research and Development, 2023, 60(6): 1191-1201
- [57] Lazaridou A, Gribovskaya E, Stokowiec W, et al. Internet-augmented language models through few-shot prompting for open-domain question answering[J]. arXiv preprint arXiv:2203.05115, 2022.
- [58] Madaan A, Tandon N, Clark P, et al. Memory-assisted prompt editing to improve GPT-3 after deployment[C]//Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing. 2022: 2833-2861.
- [59] Meng K, Bau D, Andonian A, et al. Locating and editing factual associations in GPT[J]. Advances in Neural Information Processing Systems, 2022, 35: 17359-17372.
- [60] Bang Y, Cahyawijaya S, Lee N, et al. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity[J]. arXiv preprint arXiv:2302.04023, 2023.
- [61] Felten E W, Raj M, Seamans R. How will Language Modelers like ChatGPT Affect Occupations and Industries?[J]. Available at SSRN 4375268, 2023.
- [62] 孙颖, 丁卫平, 黄嘉爽, 等. RCAR-UNet: 基于粗糙通道注意力机制的视网膜血管分割网络[J]. 计算机研究与发展, 2023, 60(4): 947-961.
- SUN Ying, DING Weiping, HUANG Jiashuang, et al. RCAR-UNet: Retinal Vessels Segmentation Network Based on Rough Channel Attention Mechanism [J]. Journal of Computer Research and Development, 2023, 60(4): 947-961.
- [63] Li J, Tang T, Zhao W X, et al. Pretrained language models for text generation: A survey[J]. arXiv preprint arXiv:2105.10311, 2021.
- [64] 俞凯, 陈露, 陈博, 等. 任务型人机对话系统中的认知技术——概念, 进展及其未来[J]. 计算机学报, 2015, 38(12): 2333-2348.
- YU Kai, CHEN Lu, CHEN Bo, et al. Cognitive Technology in Task-Oriented Dialogue Systems: Concepts, Advances and Future [J]. Chinese Journal of Computers, 2015, 38(12): 2333-2348.
- [65] 曾毅, 刘成林, 谭铁牛. 类脑智能研究的回顾与展望[J]. 计算机学报, 2016, 39(1): 212-222.
- ZENG Yi, LIU Chenglin, TAN Tieniu. Retrospect and Outlook of Brain-Inspired Intelligence Research [J]. Chinese Journal of Computers, 2016, 39(1): 212-222.