



图书馆论坛

Library Tribune

ISSN 1002-1167, CN 44-1306/G2

《图书馆论坛》网络首发论文

题目: Chat GPT 给科研工作者带来的机遇与挑战
作者: 王树义, 张庆薇
收稿日期: 2023-02-23
网络首发日期: 2023-02-24
引用格式: 王树义, 张庆薇. Chat GPT 给科研工作者带来的机遇与挑战[J/OL]. 图书馆论坛. <https://kns.cnki.net/kcms/detail//44.1306.G2.20230223.2231.002.html>



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 给科研工作者带来的机遇与挑战

王树义, 张庆薇

摘要 2022 年末 OpenAI 对话机器人模型 ChatGPT 的出现给科研工作者的外部环境造成了显著变化。如何敏感地抓住 ChatGPT 的特点创造竞争优势, 成为科研工作者需要严肃思考和引发行动的问题。文章既讨论了 ChatGPT 带来的机遇, 包括在编程、阅读和写作方面的效率提升与赋能; 也探讨了 ChatGPT 带来的挑战, 包括回答真实性、数据污染与隐私安全等问题, 并且给出了相应的对策。

关键词 ChatGPT 辅助编程 辅助写作 数据安全

引用本文格式 王树义, 张庆薇. ChatGPT 给科研工作者带来的机遇与挑战[J]. 图书馆论坛, 2023

ChatGPT's Opportunities and Challenges for Researchers

Wang Shuyi, Zhang Qingwei

Abstract The emergence of OpenAI conversational bot models by the end of 2022 ChatGPT has created a significant change in the external environment for researchers. How to sensitively capture the features of ChatGPT to create a competitive advantage has become a question that every researcher needs to seriously think about and trigger action. This paper discusses the opportunities brought by ChatGPT, including efficiency enhancement and empowerment in programming, reading and writing, and also explores the challenges brought by ChatGPT, including answering the issues of authenticity, data contamination and privacy security, and gives corresponding countermeasures.

Keywords ChatGPT; assisted programming; assisted writing; data security

0 引言

2022 年 11 月 30 日, OpenAI 发布名为 ChatGPT 的模型研究预览版, 它可以用对话的方式与用户进行交互。ChatGPT 模型使用人类反馈的强化学习 (Reinforcement Learning from Human Feedback, RLHF) 进行训练^[1]。训练方法与 OpenAI 早前发布的 InstructGPT 类似, 但数据收集设置略有不同。OpenAI 使用有监督的微调方法, 基于 GPT-3.5 系列的模型训练了一个初始模型, 并且用人工 AI 训练师对话形式, 混合 InstructGPT 数据集撰写对话格式的回应。对于备选答案, 由人工 AI 训练师排名提供增强学习的奖励^[2]。ChatGPT 自发布以来, 非常受欢迎, 仅 5 天就吸引超过 100 万用户^[3], 上市第一个月就拥有 5,700 万活跃用户^[4], 估计发布后两个月内的月活跃用户就达到 1 亿^[5]。ChatGPT 的广泛普及使得 OpenAI 的价值增长到了 290 亿美元^[6]。

ChatGPT 的火爆伴随着一系列对它的讨论。人们津津乐道于它通过图灵测试^[7], 在明尼苏达大学通过法律考试, 并在加州大学伯克利分校的另一场考试中获得优异成绩^[8]。人们尝试用它写简历和求职信, 解释复杂的话题, 提供恋爱建议^[9]。在广泛的使用中, 用户们逐渐发现 ChatGPT 的许多问题, 如对话容量限制、成为抄袭和作弊利器、存在偏见、歧视以及准确性等问题^[10]。尽管大众对 ChatGPT 的讨论非常激烈和丰富多彩, 但作为科研人员, 更应该严肃审视 ChatGPT 以及相似模型和服务的出现会给学界带来什么样的变化? 在变化出现的时候, 该如何抓住机遇并避免负面影响, 从而获得科研竞争优势? 本文通过例证来尝试初步回答上述问题。

1 文献回顾

NLG (Neural Language Generation, 自然语言生成) 是指从非语言表示生成人类可以理解的文本的技术, 应用广泛, 包括机器翻译、对话系统、文本摘要等^[11]。目前主要的 NLG

模型包括 Transformer、GPT-1/2/3、BERT、XLM、BART、Codex 等。其中, Transformer 模型基于 Attention 机制, 在质量和用时上都比之前的模型有所提升^[12]; GPT 模型是使用大量数据训练好的基于 Transformer 结构的生成型预训练变换器模型, 能在常识推理、问题问答、语义蕴含的语言场景中取得改进^[13]; BERT 引入 MLM 和 NSP 训练方法, 能融合上下文^[14]; XLM 模型通过训练跨语言信息, 可以用在训练语料少的语言上学习到的信息^[15]。2020 年 OpenAI 发布的 GPT-3 模型参数达到 1,750 亿个, 通过与模型的文本互动来指定任务, 性能强大^[16]; 2021 年 OpenAI 又发布基于 GPT-3 的 Codex 模型, 能从自然语言文件中产生功能正确的代码^[17]。2022 年 OpenAI 发布基于 GPT-3 的 InstructGPT 模型, 加入了人类评价及反馈数据, 能遵循人类指令, 并可以泛化到没有见过的任务^[1]; ChatGPT 是 InstructGPT 模型的兄弟模型, 可以遵循提示中的指令并提供详细的响应, 回答遵循人类价值观^[2]。

AIGC (AI Generated Content) 是指利用人工智能技术来生成内容的技术, 包括文本到文本的语言模型、文本到图像的生成模型、从图像生成文本等。其中, 谷歌发布的 LaMDA 是基于 Transformer 的用于对话的语言模型, 利用外部知识源进行对话, 达到接近人类水平的对话质量^[18]; Meta AI 推出的 PEER 是可以模拟人类写作过程的文本生成模型^[19]; OpenAI 发布的 Codex 和 DeepMind 的 AlphaCode 是分别用于从文本生成代码的生成模型^[17,20]。在图像生成方面, GauGAN2 和 DALL·E 2 分别是可以生成风景图像和从自然语言的描述生成现实图像的生成模型, 基于 GAN 和 CLIP 模型, 使用对比学习训练^[21,22]; Midjourney 和 Stable Diffusion 是从文本到图像的 AI 生成模型, 而谷歌的 Muse 则实现了最先进的文本转换为图像的生成性能^[23]。另外, Flamingo 是一个视觉语言模型, 能将图像、视频和文本作为提示输出相关语言^[24]; VisualGPT 是 OpenAI 推出的从图像到文本的生成模型^[25]。

人工智能内容产生过程, 难以避免遇到各种问题, 如偏见和歧视。由于训练数据集可能存在偏见和歧视, ChatGPT 可能会学习到这些偏见或歧视, 因此需要采用多种方法对数据进行筛选和清洗, 或使用公平性算法来纠正模型偏差。总体而言, ChatGPT 的公平性取决于它的训练数据集以及使用它时的上下文和提问方式^[16]。另外, 还有算力挑战, ChatGPT 依赖大量算力来训练海量文本数据, 以此学习语言模式和知识。算力需求日益增长, 致使这一领域存在着技术垄断, 会对算力持续提升、大数据的训练等进一步行动产生影响。OpenAI 为了应对 ChatGPT 的高需求, 采取排队系统和流量整形等措施^[26]。

笔者梳理相关成果, 尚未发现详细分析与论述 ChatGPT 对科研工作者影响的研究论文。因此, 本文从 ChatGPT 给科研工作者带来的机遇与挑战两个方面作为切入点, 展开论述。

2 机遇

ChatGPT 是一种 AI 工程化的成功实践。AI 工程化专注于开发工具、系统和流程, 使得人工智能能够应用于现实世界的领域^[27]。它使得普通人可以应用最新的自然语言生成与对话技术, 完成很多曾经需要技术门槛才能完成的工作。

2.1 辅助编程

数据科学研究范式已经深刻影响了许多学科。许多研究需要通过不同形式来掌握足够的数据支持。通常研究数据获取的途径主要有 3 种: 开放数据集、API 调用和爬虫 (Web Scraper)。Python 语言是信息管理学科中进行数据分析最常用的编程语言。以前用户必须掌握 Python 语言的基础语法, 了解 Python 环境的使用, 才能编写程序进行 API 调用或利用爬虫搜集数据。而有了 ChatGPT, 用户可以通过自然语言对话形式, 给出要求, 人工智能

会直接给出源代码，并复制到实验环境，从而获取所需数据。图 1 演示了研究者打算爬取某个网页上的信息，于是可以通过浏览器的 Inspector 功能查找对应的区域，此处找到的是类别 `sdyw_ul`。



图 1 通过浏览器定位爬取范围

有了对应的爬取位置，用户就可以在 ChatGPT 里直接提出要求：“我想获得 `sdyw_ul` 类下的链接和标题。”（I want to get the links and titles under the class “`sdyw_ul`”）。然后 ChatGPT 自动编写程序，结果如图 2 所示。

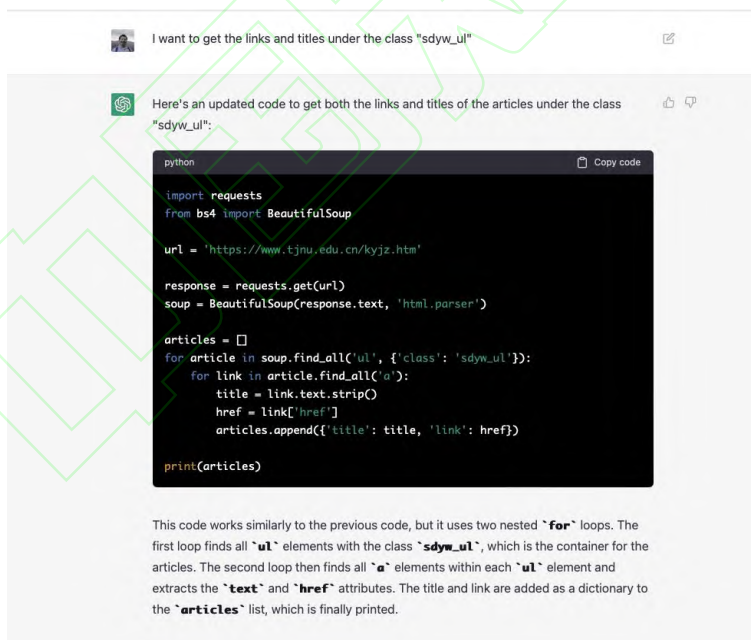


图 2 ChatGPT 自动编程爬虫

如果用户对程序运行结果不满意，可以通过进一步对话交流，让 ChatGPT 对程序源代码进行修改。例如，可以通过对话让 ChatGPT 把数据输出的方式改成 CSV 文件。ChatGPT 更新代码，返回结果如图 3 所示。

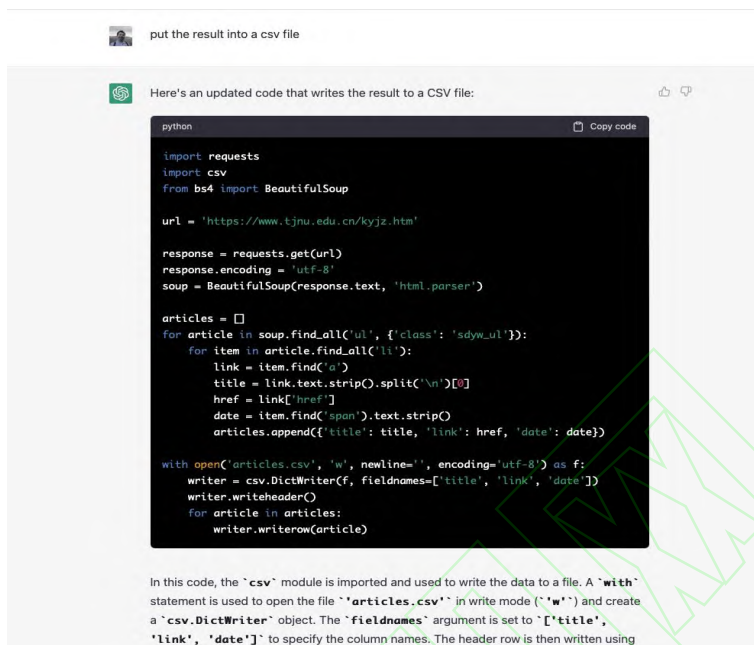


图 3 ChatGPT 爬虫编程修改输出格式

因为 ChatGPT 拥有对多轮对话的记忆力，所以每次只需要提出进一步的要求，就能不断让 ChatGPT 编写符合用户目标的程序，从而完成预期目标。最终用户可以仅通过自然语言交互和拷贝 ChatGPT 生成结果代码并运行的方式，把该网站上全部感兴趣的内容，存入 Excel 文件，如图 4 所示。

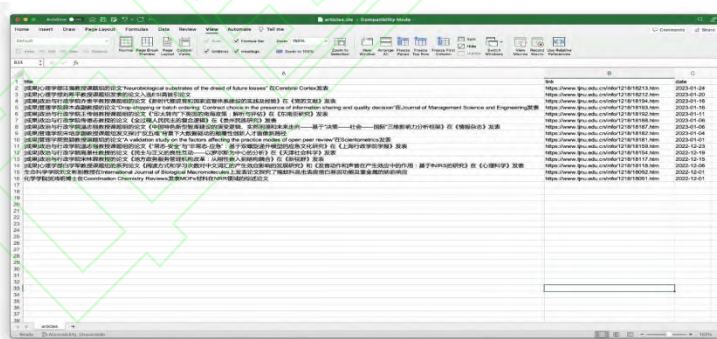
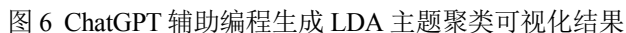


图 4 ChatGPT 生成爬虫的最终运行结果

ChatGPT 辅助编程，不只体现在数据采集环节。ChatGPT 的基础模型是“GPT-3.5”，底层基础模型是在大量代码上训练的结果，称为 code-davinci-002^[28]。因此，ChatGPT 见识过大量产生于 2021 年第四季度之前的代码段，且对代码上下文补全有较强的能力。在此之前的数据分析甚至是机器学习模型训练工作，都可以通过自然语言对话方式交给 ChatGPT 来进行。例如，在下面的例子里，用户尝试让 ChatGPT 采用 LDA 对一系列英文新闻文本做出主题挖掘，提出的自然语言指令和 ChatGPT 的应答如图 5 所示。用户只需将 ChatGPT 给出的代码复制运行，对应生成的 LDA 主题聚类可视化结果如图 6 所示。如图可见，原本需要

下面是使用LDA模型对新闻内容概述进行主题挖掘的Python示例代码，我们选择了6个主题作为主题数量：

图 5 ChatGPT 编写 LDA 主题挖掘代码



如果用户觉得结果有不满意的地方，可以随时跟 ChatGPT 交互，从而做出订正。例如，原本的代码中，ChatGPT 默认为笔者采用中文停用词，并且还使用 `jieba` 软件包进行分词，可以要求 ChatGPT 改用英语停用词进行处理。ChatGPT 会立即根据新的要求变动，给出更新后的代码，如图 7 所示。

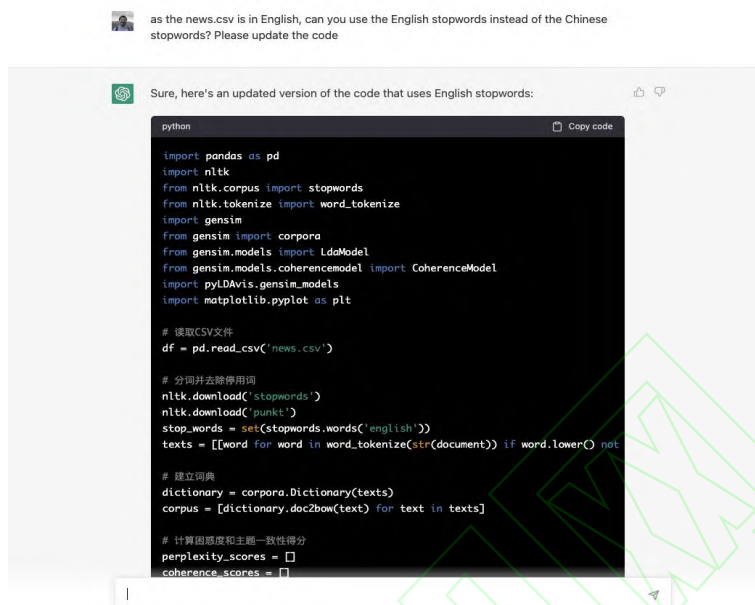


图7 要求 ChatGPT 改用英文停用词表

在这个例子中，ChatGPT 改用 nltk 软件包，使用内置的英文停用词表，可以做出更加符合要求的结果。不仅如此，在大部分 ChatGPT 生成的代码中，不仅会有详细的注释，代码完成后，ChatGPT 还会给出相应的解释，例如，在这个更新的代码中，笔者使用 NLTK 包中的“stopwords”语料库来获得“stopwords.words('english')”方法，还使用“word_tokenize”方法对文本进行标记（原文为：in this updated code, we use the 'stopwords' corpus from the NLTK package to get the 'stopwords.words('english')' method. we also use the 'word_tokenize' method to tokenize the text.）。这对于了解代码实际的功用，并且在其上进行修改迭代甚至是查找错误，都非常有帮助。对于想学习编程的入门级研究人员，也会起到显著的帮助作用。

2.2 辅助阅读

做研究需要阅读文献。在信息资源管理学科，期刊数量众多，而且外文期刊所占比例很大，每年都涌现出很多新文章需要阅读。及时对文章重点进行把握，有利于在科研进度上保持领先。但是，众多的文献阅读、消化、理解，尤其是外文文献阅读，也对本领域的研究者构成了较大的挑战。有了 ChatGPT，研究者可以将外文论文中的内容直接输入进来，然后利用提问形式，要求 ChatGPT 自动提炼重点内容。以描述 ChatGPT 同类模型 InstructGPT 的论文 *Training language models to follow instructions with human feedback*^[1] 中的 3.5 节“Models”为例，输入其中主体部分到 ChatGPT，给出的提示词是“请用中文从下面资料中提炼出至少三个重点”。输入内容如图 8 所示。图 9 是 ChatGPT 给出的答案。可见 ChatGPT 可以正确理解用户的要求，并且对内容进行了正确的自动总结。

请用中文从下面资料中提炼出至少三个重点:

3.5 Models

We start with the GPT-3 pretrained language models from Brown et al. (2020). These models are trained on a broad distribution of Internet data and are adaptable to a wide range of downstream tasks, but have poorly characterized behavior. Starting from these models, we then train models with three different techniques: Supervised fine-tuning (SFT). We fine-tune GPT-3 on our labeler demonstrations using supervised learning. We trained for 16 epochs, using a cosine learning rate decay, and residual dropout of 0.2. We do our final SFT model selection based on the RM score on the validation set. Similarly to Wu et al. (2021), we find that our SFT models overfit on validation loss after 1 epoch; however, we find that training for more epochs helps both the RM score and human preference ratings, despite this overfitting. Reward modeling (RM). Starting from the SFT model with the final unembedding layer removed, we trained a model to take in a prompt and response, and output a scalar reward. In this paper we only use 6B RMs, as this saves a lot of compute, and we found that 175B RM training could be unstable and thus was less suitable to be used as the value function during RL (see Appendix C for more details). In Stiennon et al. (2020), the RM is trained on a dataset of comparisons between two model outputs on the same input. They use a cross-entropy loss, with the comparisons as labels—the difference in rewards represents the log odds that one response will be preferred to the other by a

图 8 ChatGPT 自动提炼重点输入部分

language instructions for each task (the datasets differ in the NLP datasets included, and the style of instructions used). We fine-tune them on approximately 1 million examples respectively and choose the checkpoint which obtains the highest reward model score on the validation set. See Appendix C for more training details

1. 本文研究的是如何让GPT-3模型执行指令类任务，并介绍了三种训练模型的方法。
2. 本文使用了标注数据进行有监督的预调（SFT）以及奖励模型（RM）的训练，并采用了深度强化学习（RL）方法进行模型训练。
3. 在RL模型的训练过程中，采用了PPO算法，并加入了来自SFT模型的KL惩罚以及预训练梯度的混合，并将模型性能与GPT-3、GPT-3-prompted等基线进行了比较。

Regenerate response

图 9 ChatGPT 自动提炼重点输出部分

在模型学习训练集材料中，已经接触过不少专有名词，所以甚至可以不进行任何输入，直接让 ChatGPT 帮助解释一些专有名词。例如，图 10 里答案中出现“深度强化学习”，可以让 ChatGPT 尝试讲解其含义。输入的提示词为“什么是‘深度强化学习’，在上述文稿里面的作用是什么？”。图 10 是 ChatGPT 给出的回答。

什么是“深度强化学习”，在上述文稿里面的作用是什么？

深度强化学习（Deep Reinforcement Learning, DRL）是强化学习（Reinforcement Learning, RL）和深度学习（Deep Learning, DL）的结合，其目标是使用神经网络来解决强化学习问题。在上述文稿中，深度强化学习被用于训练模型，其作用是通过与环境的交互来不断优化模型的参数，从而使其能够更好地执行指令类任务。具体来说，本文使用了基于DRL的PPO算法，通过一个“赌徒环境”进行交互，让模型根据环境反馈的奖励不断调整其行为策略，最终得到一个能够高效执行指令类任务的模型。

图 10 解释深度强化学习

可以对图 10 中出现的新专有名词继续提问。例如，“赌徒困境”是什么？ChatGPT 的回答如图 11 所示。如果对 ChatGPT 总结的内容不放心，还可以让 ChatGPT 找到与专有名词对应的原文文本。图 12 为 ChatGPT 自动找出的“赌徒困境”原始文本。通过 ChatGPT 展示原文的文本，研究者可以加以印证，证明 ChatGPT 总结没有偏离原文的叙述。

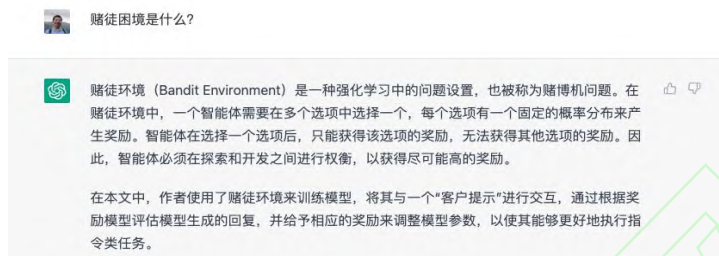


图 11 ChatGPT 解释赌徒困境

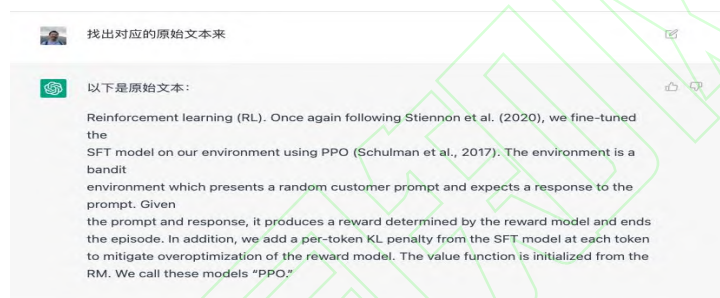


图 12 查找专有名词对应的原始文本

另外，用户还可以对文本提出问题，ChatGPT 会尽全力尝试解答。例如示例论文这样的讲述模型训练方法的作品，研究者可能更感兴趣一种模型获取反馈与提升改进的流程，并且用它和其他同类模型进行比对。所以可以问出一个非常综合性的问题：“模型是如何获得反馈和改进，达到训练目标的？”图 13 是 ChatGPT 的回答。可以看到，ChatGPT 对文本语义理解非常准确，而且还用中文进行了流畅自然的翻译。特别地，对于文中出现的专有名词，如“SFT”等，都用英文全称和缩写加以注明。

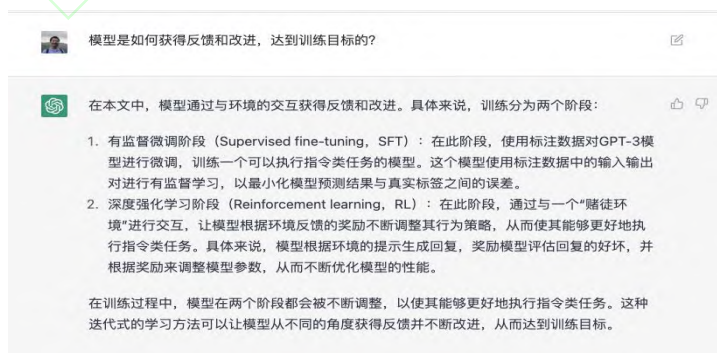


图 13 ChatGPT 对阅读材料的综合性问题解答

因为 ChatGPT 具有多轮对话的记忆功能, 用户甚至可以将多篇论文的主体部分分别输入, 然后加以比对。这样一来, 读论文的工作就从“读一篇”变成“读一片”, ChatGPT 的汇总可以快速提示研究者某一领域多篇重要文献的特点与异同, 有助于研究者快速了解领域的发展与演化过程。在从前一个研究团队用若干天才能做出的文献梳理工作, ChatGPT 可以在很短时间内高效率完成。

2.3 辅助写作

写作是沟通科研工作成果的过程, 必不可少。但以往在写作环节, 科研工作者往往需要花费很大的心力。因为不仅要详细描述和展示科研成果, 也许充分考量行文的风格、措辞、举例等。特别是多人合作写文章, 往往还要第一作者最终统合稿件的不同部分, 重新用统一的风格撰写全文。ChatGPT 的基础就是大语言模型, 最善于学习语言风格。研究者可以在草稿里摆出主要的事实而不需要考虑语序、语法等因素, 由 ChatGPT 来帮助写作、润色。用户可以将自己之前写作的文章输入 ChatGPT 中, 然后如图 14 要求 ChatGPT 提取文章的风格。之后对新的文本, 可以调用存储的文章风格(本例中为“paper style”)进行风格转化与改写。例如, 图 15 中演示的例子是本文第一作者对第二作者提供材料的风格改写。ChatGPT 对文本的语义加以保留, 但是在表现形式上进行了调整变化。统一的风格样式可以提升读者阅读的流畅性。

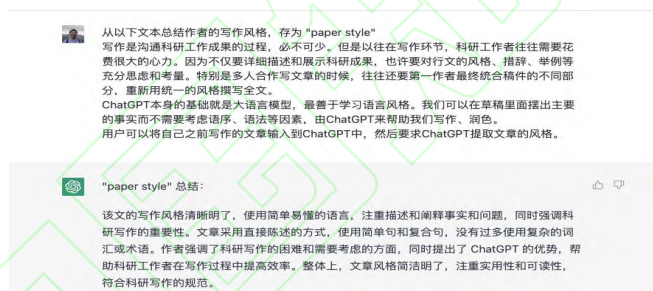


图 14 ChatGPT 学习文本风格并存储



图 15 ChatGPT 以存储的风格改写文本示例

在写作过程中，ChatGPT 也能帮助作者扩展思路，联想到合适的例子。例如，当写作过程中发现当前使用的例子作为证据并不足够充分和贴切，需要找到更好的例证，在过去，如果用户需要找到相关信息，那就必须进入搜索引擎输入关键词，然后在海量的结果中筛选适合的内容；然而现在用户只需告诉 ChatGPT “补充例子，论证上面的论断”，就可以得到相关的信息，如图 16 所示。虽然图 16 ChatGPT 提供的例子可能无法直接原文使用，但它至少对作者会有启发。例如，人脸识别已成为“日用品”，用户几乎每天都要使用这种方式验证付款，但在写作时，作者或许没有第一时间想到将其作为“AI 工程化”的例子。如果用户认为图 16 提供的例子不够好，可以接着要求 ChatGPT 提供其他例子，ChatGPT 返回的结果如图 17 所示。这样一来，原本在写作中常见而琐碎的技术性问题就被 ChatGPT 的人工智能功能解决，显著提升写作效率。



图 16 ChatGPT 补充例子证明论断

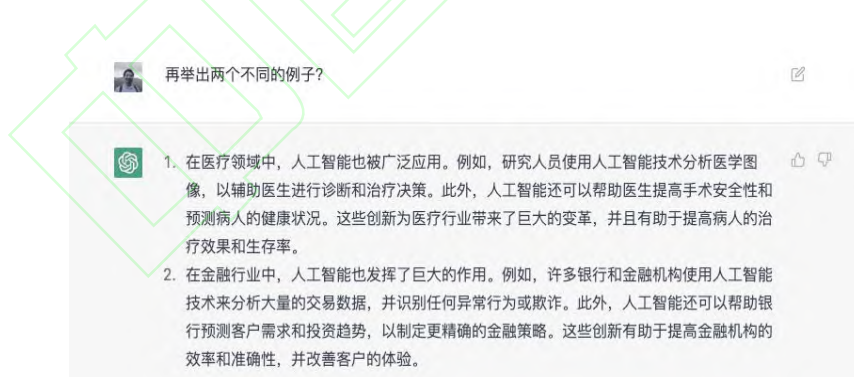


图 17 要求 ChatGPT 继续补充例证

3 挑战

上文讨论了 ChatGPT 给研究者带来的机遇。通过辅助编程、辅助阅读与辅助写作，为研究者赋能，大幅提升研究工作效率，缩短研究产生成果的时间周期。但是，ChatGPT 的出现也给研究者带来挑战和困扰，如果使用不当，甚至会给研究过程带来灾难性的后果。下文从回答的真实性、数据污染，以及隐私和数据安全角度分别论述。

3.1 回答的真实性

ChatGPT 的基础是一个生成式语言模型，它根据概率分布关系生成最符合要求的语言，但无法保证生成内容的真实性和准确性。

一些研究者在使用 ChatGPT 时没有意识到这一点，他们惊讶于 ChatGPT 回答问题的精准性，并直接采纳其答案。对于前文列举的编程功能，这个问题并不严重，因为程序编码是否准确有客观的评价标准；但对于阅读和写作辅助功能，则可能会因缺乏足够的检验依据而导致研究者采纳错误的答案。以前文展示过的 ChatGPT 举例功能来说，作者曾要求 ChatGPT 对“人工智能工程化”举出例证，结果收到的是图 18 这样的回答，回答中的疏漏非常明显：DALLE 究竟是由 Facebook 还是 OpenAI 推出？ChatGPT 给出的两个例子自相矛盾。不难发现，即便对答案的真实性缺乏把握，ChatGPT 回答时语气却非常自信。如果研究者在使用 ChatGPT 生成答案时不进行取舍，将其作为内容的组成部分，发表论文或者出版书籍后，难免遇到尴尬情况。因此在选用 ChatGPT 的答案时，研究者应该保持审慎的态度。



图 18 ChatGPT 的错误回答

3.2 数据污染

ChatGPT 的广泛使用使得很多未经思考或验证的内容大量产生。据 Intelligent.com 报道，被调查的大学生中，至少有三分之一采用 ChatGPT 来撰写作业的答案^[29]。ChatGPT 更是被广泛应用于问答网站的答案生产，且大量充斥于社交媒体。虚假信息直接影响受众之外，这些海量产生的低质量信息也会造成互联网数据的污染。这就意味着未来的人工智能模型，在从互联网获取大规模公开语料时，会吸纳大量 ChatGPT 生成内容。如果不能加以甄别，这些数据将深刻影响未来模型训练数据的质量。人工智能需要从人类产生的文本学习语言的规律。如此多的人工生成数据涌入训练集，不仅不会对模型能力带来提升，还可能混入更多噪声，导致回答问题的准确度降低。这会反过来影响人类的信息获取与知识传承。OpenAI 指出，ChatGPT 的不正当使用会对教育者、家长、记者、虚假信息研究人员和其他群体产生影响。为此，OpenAI 官方在 2023 年 1 月推出 ChatGPT 生成文本的检测分类器^[30]。使用的演示效果（采用官网自带的人工输入文本）如图 19 所示。

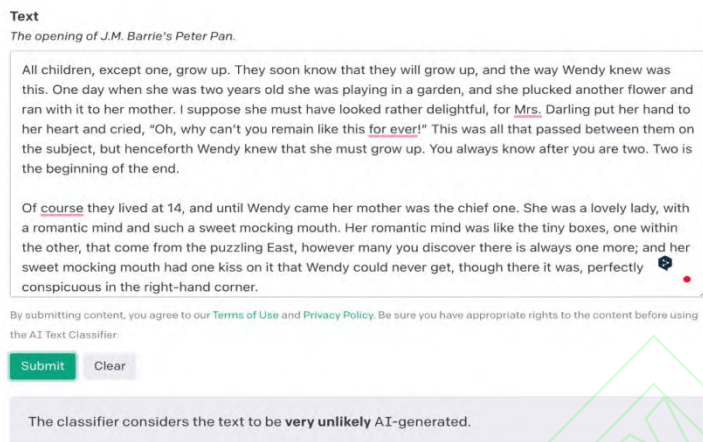


图 19 OpenAI 官方推出的 ChatGPT 检测分类器

然而这种分类器还存在着非常多的问题。OpenAI 官方建议不要将其作为主要决策工具，而应作为确定文本来源的辅助手段。对短文本（少于 1,000 个字符），OpenAI 提供的 ChatGPT 分类器不可靠。即使是较长文本，有时也会被分类器错误标记。甚至人工编写的文本也时常会被分类器错误地标记为由 AI 编写，而且检测器对于错误的检测结果非常自信。OpenAI 建议仅将分类器用于英文文本。在其他语言上分类器表现要差得多。对于代码，该分类器也不可靠。另外，它无法可靠地识别有规律的可预测的文本。官方给出的例子是“包含前 1000 个质数列表是由 AI 还是人类编写的，因为正确答案总是相同的”。AI 编写的文本通过编辑可以规避分类器识别。OpenAI 提供的分类器能够应对一定程度的改写，但对于改写较多的情况则不能准确识别。

3.3 隐私与数据安全

ChatGPT 带来的另外一个挑战，就是隐私与数据安全问题。

当用户第一次注册并开启 ChatGPT 时，会看到有关数据收集和隐私保护的提示：“对话可能会被人工智能培训师审查，以改善系统。请不要在对话中分享任何敏感信息。”在用户注册前提示用户不要输入隐私信息。

许多人将 ChatGPT 视为成熟的产品来使用，并认为它保护用户隐私是理所当然的事情。然而事实并非如此。ChatGPT 模型建立在 GPT3.5 版本之上，使用了人工参与的增强学习。每个“研究预览版”的 ChatGPT 用户都是 OpenAI 的免费测试员。

如果用户输入的内容包含敏感信息，如银行卡号，则可能会对用户的财务和金融安全造成影响。而如果输入手机号、住址等信息，则可能会带来人身安全的隐患。

对研究者来说，在输入原创性想法时也要三思而行。尽管 ChatGPT 并没有主动剽窃用户想法的意图，但用户输入的内容都会对模型造成影响。如果恰巧有其他用户对同一研究方向感兴趣，前面研究者输入的想法可能会作为答案的一部分启发后者。另外，根据 OpenAI 官方提示，ChatGPT 人工训练员（AI trainers）也有可能查看对话信息，以改进模型。

从学术整体进步的角度来看，这种信息加速传播共享有利于研究进展速度。但对研究者个体来说，其利益可能会受到潜在的威胁。

综上所述, ChatGPT 的挑战主要分为回答真实性、数据污染与隐私和安全等方面。面对 ChatGPT 带来的挑战, 研究者可以通过以下对策尽量避免潜在的损失。

第一, 针对回答的真实性问题, 建议研究者时刻警醒, 不要轻易相信 ChatGPT 提供的答案。即便对看似合理的答案内容, 在正式采纳和使用前, 也需要找到源头信息进行验证。

第二, 针对数据污染问题, 建议研究者采用 OpenAI 官方提供的 ChatGPT 生成文本检测工具对重要来源数据进行抽样检测。在构建大规模研究数据集时, 尽量避免采用开放式问答社区 2022 年 12 月之后的回答, 以避免噪声混入。

第三, 对于隐私和安全问题, 建议研究者与 ChatGPT 对话过程中, 避免暴露个人隐私与所在机构的敏感信息。对于研究意图和想法, 如果无法绕开, 尽量分散在不同对话中进行任务处理, 避免被人工训练员在同一对话中大量获取相关信息。

4 结论

OpenAI 的对话机器人模型 ChatGPT 对科研工作者的外部环境造成了显著变化, 为提高编程、阅读和写作效率带来了机遇, 但也带来了回答真实性、数据污染和隐私安全等挑战。为了敏感地抓住 ChatGPT 的特点, 创造竞争优势, 科研工作者需要认真思考并采取行动。通过本文的讨论, 读者可以看到 ChatGPT 对科研工作的赋能意义十分明显, 合理利用能够大幅提升工作效率。而针对 ChatGPT 给科研工作者带来的挑战, 本文提出了对策, 如在使用 ChatGPT 生成的答案时需要进行谨慎评估, 同时需要利用的技术和方法来应对数据污染和隐私安全问题。总之, 科研工作者也需要不断学习和更新自己的技能, 以更好地适应这个新的科研环境。

ChatGPT 出现时间不久且快速迭代, 也有许多竞品宣布会相继在近期推出。本文受到当前写作时间点的客观局限, 无法对近期和远期即将出现的产品或服务趋势做出准确预测。本文写作时, 尚未发现与 ChatGPT 实力相当的真正竞品, 因此研究对象比较单一, 只涉及 ChatGPT 自身。团队后续会在新的同类产品出现后加以深入对比研究, 为科研工作者提供更加符合本土化需求的分析结果与建议。

参考文献

- [1] OUYANG L, WU J, JIANG X, et al. Training language models to follow instructions with human feedback: arXiv:2203.02155[Z/OL]. arXiv, 2022(2022-03-04)[2023-02-02]. <http://arxiv.org/abs/2203.02155>. DOI:10.48550/arXiv.2203.02155.
- [2] OPENAI. ChatGPT: Optimizing Language Models for Dialogue[EB/OL](2022-11-30)[2023-02-12]. <https://openai.com/blog/chatgpt/>.
- [3] JACKSON S. OpenAI Executives Say Releasing ChatGPT for Public Use Was a Last Resort After Running into Multiple Hurdles — and They're Shocked by Its Popularity[EB/OL](2023-01-26)[2023-02-05]. <https://www.businessinsider.com/chatgpt-openai-executives-are-shocked-by-ai-chatbot-popularity-2023-1>.
- [4] CERULLO M. ChatGPT User Base Is Growing Faster Than TikTok[EB/OL](2023-02-01)[2023-02-05]. <https://www.cbsnews.com/news/chatgpt-chatbot-tiktok-ai-artificial-intelligence/>.
- [5] HARRISON M. ChatGPT's Explosive Popularity Makes It the Fastest-Growing App in Human History[EB/OL](2023-02-03)[2023-02-05]. <https://futurism.com/the-byte/chatgpts-fastest-growing-app-human-history>.
- [6] SOUTHERN M G. ChatGPT's Popularity Boosts OpenAI's Value To \$29 Billion[EB/OL](2023-01-06)[2023-02-01]. <https://www.searchenginejournal.com/chatgpts-popularity-boosts-openais-value-to-29-billion/475762/>.
- [7] YALALOV D. ChatGPT Passes the Turing Test[EB/OL](2022-12-08)[2023-02-08]. <https://mpost.io/chatgpt-passes-the-turing-test/>.

- [8] KELLY S M. ChatGPT Passes Exams from Law and Business Schools | CNN Business[EB/OL](2023-01-26)[2023-02-08]. <https://www.cnn.com/2023/01/26/tech/chatgpt-passes-exams/index.html>.
- [9] TIMOTHY M. 11 Things You Can Do With ChatGPT[EB/OL](2022-12-20)[2023-02-08]. <https://www.makeuseof.com/things-you-can-do-with-chatgpt/>.
- [10] AGOMUOH F. The 6 Biggest Problems with ChatGPT Right Now[EB/OL](2023-01-27)[2023-02-06]. <https://www.digitaltrends.com/computing/the-6-biggest-problems-with-chatgpt-right-now/>.
- [11] DONG C, LI Y, GONG H, et al. A Survey of Natural Language Generation: 8[J/OL]. ACM Computing Surveys, 2023, 55(8): 1-38[2023-02-17]. <http://arxiv.org/abs/2112.11739>. DOI:10.1145/3554727.
- [12] VASWANI A, SHAZEER N, PARMAR N, et al. Attention Is All You Need: arXiv:1706.03762[Z/OL]. arXiv, 2017(2017-12-05)[2023-02-01]. <http://arxiv.org/abs/1706.03762>. DOI:10.48550/arXiv.1706.03762.
- [13] RADFORD A, NARASIMHAN K. Improving Language Understanding by Generative Pre-Training[C/OL]. (2018-06-11)[2023-02-23]. <https://www.semanticscholar.org/paper/Improving-Language-Understanding-by-Generative-Radford-Narasimhan/cd18800a0fe0b668a1cc19f2ec95b5003d0a5035>.
- [14] WANG Y, SUN Y, FU Y, et al. Spectrum-BERT: Pre-training of Deep Bidirectional Transformers for Spectral Classification of Chinese Liquors: arXiv:2210.12440[Z/OL]. arXiv, 2022(2022-10-22)[2023-02-01]. <http://arxiv.org/abs/2210.12440>. DOI:10.48550/arXiv.2210.12440.
- [15] LAMPLE G, CONNEAU A. Cross-Lingual Language Model Pretraining: arXiv:1901.07291[Z/OL]. arXiv, 2019(2019-01-22)[2023-02-01]. <http://arxiv.org/abs/1901.07291>.
- [16] BROWN T, MANN B, RYDER N, et al. Language Models are Few-Shot Learners[C/OL]//Advances in Neural Information Processing Systems. Curran Associates, Inc., 2020: 1877-1901[2023-02-05]. <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html>.
- [17] CHEN M, TWOREK J, JUN H, et al. Evaluating Large Language Models Trained on Code: arXiv:2107.03374[Z/OL]. arXiv, 2021(2021-07-14)[2023-02-02]. <http://arxiv.org/abs/2107.03374>. DOI:10.48550/arXiv.2107.03374.
- [18] THOPPILAN R, DE FREITAS D, HALL J, et al. LaMDA: Language Models for Dialog Applications: arXiv:2201.08239[Z/OL]. arXiv, 2022(2022-02-10)[2023-02-04]. <http://arxiv.org/abs/2201.08239>. DOI:10.48550/arXiv.2201.08239.
- [19] SCHICK T, DWIVEDI-YU J, JIANG Z, et al. PEER: A Collaborative Language Model: arXiv:2208.11663[Z/OL]. arXiv, 2022(2022-08-24)[2023-02-01]. <http://arxiv.org/abs/2208.11663>. DOI:10.48550/arXiv.2208.11663.
- [20] LI Y, CHOI D, CHUNG J, et al. Competition-level code generation with AlphaCode: 6624[J/OL]. Science, 2022, 378(6624): 1092-1097[2023-02-01]. <https://www.science.org/doi/10.1126/science.abq1158>. DOI:10.1126/science.abq1158.
- [21] SALIAN I. NVIDIA Research's GauGAN AI Art Demo Responds to Words[EB/OL](2021-11-22)[2023-02-01]. <https://blogs.nvidia.com/blog/2021/11/22/gaugan2-ai-art-demo/>.
- [22] RAMESH A, DHARIWAL P, NICHOL A, et al. Hierarchical Text-Conditional Image Generation with CLIP Latents: arXiv:2204.06125[Z/OL]. arXiv, 2022(2022-04-12)[2023-02-03]. <http://arxiv.org/abs/2204.06125>. DOI:10.48550/arXiv.2204.06125.
- [23] ROMBACH R, BLATTMANN A, LORENZ D, et al. High-Resolution Image Synthesis with Latent Diffusion Models[C/OL]//2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). New Orleans, LA, USA: IEEE, 2022: 10674-10685[2023-02-04]. <https://ieeexplore.ieee.org/document/9878449/>. DOI:10.1109/CVPR52688.2022.01042.
- [24] ALAYRAC J-B, DONAHUE J, LUC P, et al. Flamingo: A Visual Language Model for Few-Shot Learning[J/OL]. (2022-04-28)[2023-02-02]. <https://arxiv.org/abs/2204.14198>.
- [25] CHEN J, GUO H, YI K, et al. VisualGPT: Data-efficient Adaptation of Pretrained Language Models for Image Captioning: arXiv:2102.10407[Z/OL]. arXiv, 2022(2022-03-30)[2023-02-02]. <http://arxiv.org/abs/2102.10407>. DOI:10.48550/arXiv.2102.10407.
- [26] MORRISON R. Compute Power Is Becoming a Bottleneck for Developing AI. Here's How You Clear It[EB/OL](2022-12-13)[2023-02-05]. <https://techmonitor.ai/technology/ai-and-automation/chatgpt-ai-compute-power>.
- [27] COURSERA. What Is an AI Engineer? (And How to Become One)[EB/OL](2022-06-27)[2023-02-01]. <https://www.coursera.org/articles/ai-engineer>.

- [28] OPENAI. OpenAI API[EB/OL](2023-02-16)[2023-02-16]. <https://platform.openai.com>.
- [29] INTELLIGENT.COM. Nearly 1 in 3 College Students Have Used ChatGPT on Written Assignments[EB/OL](2023-01-23)[2023-02-02]. <https://www.intelligent.com/nearly-1-in-3-college-students-have-used-chatgpt-on-written-assignments/>.
- [30] OPENAI. New AI Classifier for Indicating AI-Written Text[EB/OL](2023-01-31)[2023-02-02]. <https://openai.com/blog/new-ai-classifier-for-indicating-ai-written-text/>.

作者简介 王树义，博士，天津师范大学管理学院副教授；张庆薇，天津师范大学管理学院研究生。

收稿日期 2023-02-23

