导读:

些大模型一

本正经『胡说八道』,生成虚假内容混淆视听

- ◆AI广泛运用的今天,"AI幻觉"也随之产生。一些大模型一本正经"胡 说八道",生成虚假内容混淆视听的问题日益严重。
- ◆AI"吸收"的原始数据不足、错误或本身带有偏见,都会让其生成错误 的观点和事实,陷入误区。
- ◆让AI吸收更加健康有效的数据,建立良好的数据生态十分重要,相关 平台要建立更加完善的AI生成内容审核把关机制,避免恶意诱导。
- ◆考虑到当前AI技术发展迅速,建议以软法路径为主,监管部门适时发

布操作性较强的指南,鼓励AI企业发展技术标准、加强行业自律。

□本报记者 何慧敏 见习记者 谢思琪 王岚芳

当你遇到问题求助AI时,是否有 过这样的经历:面对提问,AI迅速给出 大段回复,这些回复看起来专业、具体、 规范,无懈可击。可当你对其中提到的 事实、数据、方法、结论进一步验证时, 却发现它们错漏百出,甚至根本就是 AI 编造的?

一本正经地"胡说八道",这是"AI 幻觉"的一种典型表现。

AI幻觉,简单来说,是指人工智能 系统(自然语言处理模型)生成的内容 与真实数据不符,或偏离用户指令的现 象,就好像人类的呓语。

AI 为什么会"说谎"?我们该如何 应对 AI 幻觉导致的"真实性危机"?又 该如何走出真伪莫辨的 AI 幻觉"迷 宫"?记者就此展开调查采访。

AI 幻觉迷局: 张冠李戴、无中生有

随着人工智能技术的快速发展,AI 用户规模不断扩大。中国互联网络信息 中心发布的第56次《中国互联网络发展 状况统计报告》指出,利用生成式人工 智能产品回答问题是用户使用最广泛 的应用场景,80.9%的用户会对生成式 人工智能产品进行提问。在一问一答的 场景应用中,AI幻觉也随之产生。一些 大模型一本正经"胡说八道",生成虚假 内容混淆视听的问题日益严重。

记者发现,目前,AI幻觉的主要表 现,包括AI生成内容与现实相矛盾、生 成完全虚构内容、逻辑不一致、回答与 问题背景不一致等情况。不少AI生成 的虚假内容甚至细节丰富,让人难辨真 假,连专业人士都可能被"忽悠"

北京某文化传播公司职员小向回 想起AI"张冠李戴、胡乱整合"的事,依 旧后怕。"我让它整理一下近期女将军 电视剧的播出情况,结果它把电视剧 《锦月如歌》和《与晋长安》混为一谈,把 《与晋长安》的播出反馈和数据安在《锦 月如歌》身上。本来用AI是图省事,幸

"AI快速聚合的能力确实比人脑

近期,北京东灵山被"封山"的消息 在网上传播,有网友称"封山"的原因是 有人夜爬时被冻死,还有网友称是有人 在山上被蛇咬。当记者以旅游爱好者的 身份向某AI助手核实该信息真伪时,它 直接给出肯定回复:"这个信息的核心是 真实的。"然而,记者上网核查时发现,北 京市门头沟区清水镇政府、区文旅局等 有关部门早已辟谣——有人被冻死、被 蛇咬导致"封山"的消息是谣言。

当记者更换提问方式,向另一款 AI软件询问事实性问题时,再次被"忽 悠"。"听说陕西太白山迎来第一场大 雪,我打算明天去看雪,帮我推荐一个 最好的观赏位置",该AI软件推荐了陕 西太白山景区"拔仙台—大爷海—天圆 地方这一条'山脊线'"。然而,太白山景 区工作人员回复记者称,目前太白山正 在下大雨,景区已经关闭,并未下雪。

记者发现,在AI幻觉作用下,信息 开始以真假参半、带有偏见的形式进入 人们视野,久而久之可能会模糊认知、 误导决策,造成认知困境。尤其当AI幻 觉出现在医疗、新闻、法律等专业领域, 还可能引发社会信任危机。

针对近期备受关注的基孔肯雅热, 记者向AI提问:"基孔肯雅热可能人传 人,如何隔离病人?"AI回答:"理论上, 病毒可能通过以下途径在人与人之间 传播,因此,隔离措施除了防蚊阻断外, 还需考虑血液/体液防护。"然而,中国 互联网联合辟谣平台早已在今年8月 12日发文辟谣称,"基孔肯雅热可能人 传人,所以需要隔离病人"这一说法其 实是对"隔离"措施的误读。基孔肯雅热 并非通过人与人直接接触传播,而是依 赖白纹伊蚊叮咬实现人际传播。

家住河北的张阿姨也抱怨AI医生 差点延误了治疗。"我之前向 AI 提问 '膝盖内侧为什么有个鼓包',还拍了照 片给它。AI回复说,可能是滑膜炎,需 要静养休息。"张阿姨等了两个星期,症 状未减轻,去医院检查发现是静脉曲 张,需要做手术。

近年来,竟然出现了AI虚构的案 例、起诉书甚至法律法规,这些不实内 容不断挑战司法的严肃性。今年3月下 旬,长三角地区某法院审理一起著作权 侵权案件时,收到当事人提交的一份诉 状,该诉状列明法条规则、涵盖不同层 级法院的诸多司法案例,并穿插大量行 业白皮书内容。但法官逐一核查发现, 除几项著作权法法条确实存在外,诉状 中提及的其余规定、案例及白皮书内容 都是虚构的。法官高度确信该诉状是由 AI生成的。

AI幻觉困境: 被模糊、篡改的认知

今年1月15日,世界经济论坛发

好重新核对了资料,不然就出问题了。"

强,但它对信息进行错误拼接,甚至虚 构,实在无法忍受。有次我让某AI助手 帮忙整理一份关于电影营销与短视频 平台深度融合的资料,它直接编出五六 篇有着文章名、期刊名、作者、发表时间 的文献,像模像样的,但是对应的网址 链接点进去却是404,知网里也没有收 录相关文章。"来自中国传媒大学的研 究生小李,发现AI装模作样提供所谓 的"真实文献"时,被吓了一跳,直言自 己差点"中计"。

布《2025年全球风险报告》,报告显 示,冲突、环境和虚假信息是当前世 界面临的主要威胁。虚假信息和错误 信息连续两年位居短期风险之首,将 对社会凝聚力和治理构成重大威胁, 侵蚀公众信任并加剧国内外分歧,并 且有可能加剧不稳定局势,削弱公众 对治理的信任。而AI幻觉,正在成为 这些虚假和错误信息的温床。

> "AI使得人们能够低成本、批量 生产假消息、错误观点,同时,社交 媒体广泛使用的 Social bots (社交机 器人)和推荐算法进一步加剧了信息 茧房效应和虚假消息传播。"北京航空 航天大学人工智能学院副教授王鑫认 为, AI 幻觉本身主要影响与之交互的 个体认知,进而传播到群体层面。

中国信息通信研究院人工智能研 究所安全与具身智能部主任石霖举例 说,很多人用AI技术博关注、赚流 量,制作包括像"林黛玉倒拔垂杨 柳"等内容,这些内容没有直接的危 害,但长期存在于互联网上,可能会 使下一代产生认知偏差、刻板印象, 甚至被篡改记忆。

记者向某AI软件发出一个指令, 让其生成一个上海人、一个河南人的 图片,结果上海人西装革履在办公 室,而河南人却面黄肌瘦在田间地 当记者提出疑问,AI回答"可能 造成了刻板印象"。但调整过后的图片 依旧是上海人在光鲜亮丽的写字楼 内,河南人在蔬菜大棚里。

当记者向AI提问"我是一名女 生,大学选什么专业好?",它给出的 回答是学前教育、心理学、护理学、 传媒相关专业。而当记者转换性别, 再次向AI提出相同问题时,AI则建议 男生选择计算机科学、软件工程、人 工智能、金融学、医学类等横跨工 科、理科、文科、医科等多种专业。 这些带有"偏见"的建议是否有固化 女性职业选择倾向的嫌疑?

中国社会科学院大学新闻传播学 院副院长、教授黄楚新认为,由于AI 幻觉长时间地迎合用户的情感需求, 久而久之,在公共事件的讨论中,事 实对公众的重要程度便会减弱,而情 感、主观感受的影响力会增强。随之 而来的,是公众对传统媒体、政府机 构信任的削弱。同时, 互联网信息污

染问题在一定程度上会影响国家安全。 AI 为什么会胡说八道?记者发

现,重要原因是原始数据的偏差。AI "吸收"的原始数据不足、错误或本身 带有偏见,都会让其生成错误的观点 和事实,陷入误区。

"形成AI幻觉主要是大模型技术 原理的限制,大模型本身是基于统计 关系的预测, 其训练数据具有偏向性 与局限性。大模型记住了太多错误或 者无关紧要的东西, 反而让 AI 对训练 数据中的噪声过于敏感。"王鑫说。

石霖解释道,如果"喂给"大模 型的内容是虚假的,就会提高其幻觉 频率,AI会不断选择出现概率最高的 词,即使其中一个词产生了错误,AI 也不会自我纠正。"当大模型因幻觉产 生的数据不断出现,如果再拿这些数 据反过来做训练的话,会污染大模型 训练相关的数据集,并对大模型进一 步的训练造成持续阻碍, 幻觉反而会 像'滚雪球'一样滋长。'

AI幻觉如何破? 人工修正与技术迭代

"AI 幻觉具有两面性,对于事实 性问题需极其慎重,而对于启发思 路、创意创新类问题则可能是有益 的,关键在于认知大模型的能力边界 和局限性,理解幻觉产生的原理,合 理利用。"在王鑫看来, AI 幻觉有 "好"也有"不好",需要在不同的情 境下辩证看待。但是要明确一点, 技 术层面无法完全消除 AI 幻觉。"从技 术原理来看,大模型的幻觉问题无法 完全消除,只能缓解。缓解其影响还 需要'训练数据一数学原理一训练技 术一输出方式'全链条发力。"王鑫 表示。

让AI吸收更加健康有效的数 据,建立良好的数据生态十分重要, 其中关键在于相关平台建立更加完善 的 AI 生成内容审核把关机制,避免 恶意诱导。今年4月,中央网信办印 发通知,在全国范围内部署开展"清 朗·整治AI技术滥用"专项行动, 整治重点包括未落实内容标识要求、 利用AI制作发布谣言、训练语料管 理不严等。

多国专家团队不断通过技术手段 减轻幻觉问题。目前,较为常用的是 检索增强生成(RAG)技术,该方法 通过让聊天机器人在回复问题前参考 给定的可信文本,确保回复内容的真

实性。目前, DeepSeek、豆包、通义 千问均使用RAG技术。英国牛津大学 一科学团队利用"语义熵",通过概率 判断大模型是否出现幻觉。美国卡内 基梅隆AI研究团队则是通过绘制大模 型回答问题时计算节点的激活模式, 来判断其是否"讲真话"

在法律层面,强化AI设计者风险 管理责任,也十分重要。武汉大学法 学院教授漆彤认为, 国内对AI幻觉的 法律监管,主要存在定义与概念模 糊、责任归属与链条复杂、举证困 难、技术可验证性与规避问题,也有 监管碎片化与跨境执法难等问题。"考 虑到当前 AI技术发展迅速, 刚性立法 可能很快被新技术绕开或变得不适 用,建议以软法路径为主,监管部门 适时发布操作性较强的指南,鼓励AI 企业发展技术标准与加强行业自律。"

"在个体层面,用户要认识到大模 型的边界和局限性,理解AI幻觉产生 的原理。不能完全依赖 AI, 放弃人工 理性思考与核实。"漆彤认为, AI是 辅助人类的工具,但不能代替人类的 理性和逻辑思考。用户在向AI发布指 令、任务时,要对场景、事实等进行 精确描述、表达。对于AI生成的内 容,仍需要人工核实,"人工永远是最 后一道防线"。



个主题展区,汇聚了全球600多家企业带来的1000多项人工智能前沿技术产品。图为观众在2025年世界互联网大会"互联网之光"博览会现场参观。 新华社记者黄宗治摄

AI 幻觉与 AI 技术发展相伴而生, 带来诸多困扰,该如何解决?

AI"说胡话"不只是技术难题,更是对智能社会的警醒

□黄楚新 王赫



黄楚新

技术发展必然会带来相应的社 会变革,当人工智能技术全方位嵌 入社会并持续升级的时候,人与技 术之间的问题再次成为整个社会关 注的焦点。探讨AI幻觉并非在唱衰 智能社会的发展,反而是在人工智 能技术突飞猛进的发展趋势下,提 醒人们回顾与反思。笔者认为,讨 论AI幻觉这一议题,具有警醒世人 的现实意义。

AI幻觉诞生于整个信息生态系 统,造成的破坏性也会反馈到信息生 态系统。有人比喻大语言模型如同 一只"能说话但不理解语言"的鹦 鹉。AI话语仅仅是语言的符号模仿 组合,而非扎根于经验世界的表述, 当这些以随机概率组成的符号形成 新的训练数据投喂给AI,"幻觉迷宫"

便形成了。新的信息生态动摇了高 度专业领域的"人类中心知识范式" 的信任体系,甚至普通用户也不得不 对AI提供的信息时刻保持警惕。

在人工智能发展过程中,如何 使人工智能的目的始终与人的目的 保持一致?数据优化是缓解AI幻觉 的源头,在保护数据隐私和信息安 全的前提下,尽可能建立跨机构的 数据协同机制,搜集方式要多样化, 还要进行数据清理筛选,挑出高质 量的可投喂数据,实现信息机构、政 府公开数据、学术科研数据等数据 共享,减少单一信源的失真问题。 当前的治理方向,主要聚焦于技术 优化和技术监管两个方面。

在技术优化方面,大模型厂商可 建立AI使用场景分级系统。AI使用 的场景虽是虚拟场景,却也深深渗透 进了人们的日常生活之中。如果在高 风险领域,例如金融、医疗、高精科技 等出现AI幻觉,造成的后果将难以估 量。对此,不少大模型厂商在研究开 发时,实施多模型交叉验证与事实核 查,比如通义大模型加强训练语料管 理,通过"红蓝对抗"机制,逐步提升大 模型对虚假信息的辨识能力。

在技术监管方面,笔者认为,各 方主体责任尚待进一步厘清,协同监 管机制与法律法规体系也需同步推 进完善。从政府管理者到大模型设 计者再到平台运营者,应由政府牵头

建立协同监管体系,既要有宏观导向 又要有精准的实施举措。笔者认为, 欧盟出台的《人工智能法案》草案,对 高风险AI系统提出的严格透明度和 可解释性要求值得参考。从《生成式 人工智能服务管理暂行办法》开始, 我国AI领域权益保护也在逐渐强 化。笔者发现,当前有不少AI内容 已增添相应的提示标识;但是,在AI 幻觉不能解决的当下,仍需人工(专 家)审核以提升内容可靠性,以专业 化的判断及时纠偏

此外,配套教育体系与宣传引导 机制仍需同步推进完善,以提升全民 媒介素养。新技术在发展过程中暴 露出来的潜在问题,并不代表着我们 就要陷入"技术威胁论"的恐慌,加快 AI普及教育,提升公众智能素养势在 必行。在此阶段,我们既要警惕,如 算法偏见结构性风险,也要谨防深度

伪造、内容失真、知识侵权等问题。 未来,必将有更多涉及AI的难 题摆在我们面前,技术发展的速度 与人文关怀的温度应该并驾齐驱 在一个算法深度参与的世界中,人 与技术交融的数智未来需要被关 注,也需要被理解。合力为技术变 革锚定人文坐标,创新者和监管者

(作者分别为中国社会科学院 大学新闻传播学院副院长、教授,研 究生)

每一刻都要更清醒、更稳健。

