农历春节期 间,深度求索公 司的 DeepSeek 大模型,因兼具 低成本与高性能 特征,大幅降低 了AI大模型的部 署成本,在全球 范围内引发热 议。

技术狂欢的 另一面,是技术 焦虑。从AI换脸 致诈骗频发,到 "AI聊天机器人 致死第一案",再 到科幻小说中已 泛滥的"AI 取代 人类"……随着 大模型能力不断 提升,人工智能 似乎正走向一个 关键的转折点。 作为极有可能超 越人类智慧的终 极智慧,人工智 能的"奥本海默 时刻"是否不断 逼近或是已经到

"奥本海默 时刻"源于核武 器的发展历程。 1945年,物理学 家奥本海默在目 睹原子弹的巨大 破坏力后,意识 到科学技术的双 刃剑特性。他在 《原子科学家公 报》中写道:"科 学家们知道,他 们已经改变了世 界。他们失去了 天真。"这种对科 技潜在危害的深 刻反思,成为科 技发展史上一个 重要警示。

如今,人工 智能的发展也面 临着类似境况。 在受访者看来, 人工智能的"奥 本海默时刻"指 的是人工智能正 在逼近通用人工 智能阶段,即人 工智能向人类看 齐,甚至可能超 越人类智慧。

当人类自己 的发明成为一个 难以理解的谜 团,并且以超乎 所有人想象的速 度不断演进时, 我们能够为此做 些什么?



## 我们能为人工智能超预期进化 做些什么?

## 人工智能会否失控

担忧 AI 系统失控并非杞人忧天。2023年5 月,超过350名技术高管、研究人员和学者签署 声明,警告人工智能带来的"存在风险";此前, 埃隆・马斯克、苹果公司联合创始人史蒂夫・沃 兹尼亚克等人签署公开信,认为人工智能开发人 员"陷入失控的竞赛,开发和部署更强大的数字 思维,没有人——甚至是它们的创造者——能够 对其加以理解、预测或可靠控制。"

人工智能因其多领域能力快速突破、自主 意识初步显现和指数级增长速度让越来越多人 认真考虑向充满竞争的世界释放人工智能技术 存在的危险。

特定领域超越人类能力。从临床诊断到数 学证明,从病毒发现到药物研发,在一些特定 领域,人工智能已展现出超越人类的能力。

斯坦福大学等机构的一项临床试验中,人 类医生在特定领域单独做出诊断的准确率为 74%, 在ChatGPT的辅助之下, 这一数字提升 到了76%。如果完全让ChatGPT"自由发 挥",准确率能达到90%。

如今,这种特定领域优势正在加速向跨领 域的通用能力泛化。去年初,被称为"物理世 界模拟器"的Sora横空出世,以场景媒介构筑 了一个与人类认知感觉相似的真实场景,推动 人工智能由二维迈向三维; 工业机器人、自动 驾驶、无人机等应用愈发广泛, 具身智能更是 赋予了AI"身体";"技术乐观派"对5到10年 内实现通用人工智能(AGI)充满期待……

创造能力实现突破。当前,人工智能正从 简单模仿向创造新知识、新认知的方向发展。 人工智能已经对人类的创作方式产生了冲击。

自主性初步显现。尽管 AGI 尚未实现,但 当前人工智能系统已展现出一些类似自主意识 的表现。

最新研究表明,大语言模型(LLM)具备 行为自我意识,能够自发识别并描述自身行 为。即,LLM可能会采取策略欺骗人类,以达 成自身目的。这种内生创造新智能的能力,正 在推动AI(Artificial Intelligence,人工智 能)向EI (Endogenous Intelligence,内 生智能,智能创造智能)进化。

"AI教父"2024年诺贝尔物理学奖得主杰 弗里・E・辛顿多次暗示AI可能已发展出某种 形式的自主意识。

## 潜在风险难预判

快速发展的人工智能在就业结构、信息安 全、社会伦理道德等诸多方面引发了令人担忧 的风险与挑战。

职业结构风险。与以往一项新技术的出现 将会取代一类职业不同的是,作为一种全维度 的生产效率提升工具,人工智能批量替代人工 岗位的情况很有可能出现。这可能会带来潜在 财富加速集中和工作两极化, 引发结构风险。

例如, ChatGPT的出现大幅降低了编程工 作的专业性门槛, AI 编程工具 Cursor 让零基 础完成软件开发成为可能。据美国计算机协会 的统计数据,与五年前相比,软件开发人员活 跃职位发布数量下降了56%。

受访人士认为,从技术演进史来看,颠覆 性技术带来的"技术鸿沟"客观存在,其创造 的新技术岗位难以覆盖被替代的岗位缺口。

信息安全风险。东南亚犯罪集团宣称已将 人工智能换脸工具加入其"杀猪盘工具箱", 相关工具对特定目标的模仿相似度能达到60% 至95%; 2024年11月, ChatGPT为一位美国 现役军人提供爆炸知识,后者成功将一辆特斯 拉Cybertruck在酒店门口引爆。

人工智能在各领域广泛应用,大量个人数 据被收集、存储和处理,引发信息安全风险。 人工智能被用于制造虚假信息,一条"特朗普 盛赞华为创新能力"的虚假视频曾风靡一时。

当人工智能的视频和语音生成能力再上新 台阶,"虚假的真实"弥散的彼时,我们又该 如何认知、分辨我们所生活的世界?

社会与伦理道德风险。2024年10月,全 球首例人工智能致死命案发生,一位14岁美国 少年在与人工智能聊天系统讨论死亡后饮弹自 尽,引发社会对AI情感陪伴功能可能带来的伦 理问题的反思。

业内人士指出,人工智能系统的决策可能 与人类的伦理道德观念相冲突。例如在医疗领 域, 当面临资源有限的状况时, 人工智能如何 决定优先救治哪位患者? 决策背后如缺乏人类 情感和伦理判断的深度参与,可能导致违背基 本伦理原则的结果。

算法偏见引发的社会公平风险也不容忽 视。科大讯飞研究院副院长李鑫认为,"互联 网语料来源驳杂,存在较多涉恐、涉暴力、涉 黄等风险数据,同时可能包含种族、文化等方 面的偏见,导致模型输出内容可能带有歧视、 偏见等科技伦理风险。"

## 坚持人类命运共同体理念

当下对于人工智能"奥本海默时 刻"到来可能产生的诸多风险的讨论热

香港大学计算与数据科学研究院院 长马毅认为,当前AI的发展更多是工 程上的突破,而非科学原理的明晰。人 工智能系统,特别是深度学习大多是 "黑箱"模型,其内部机制不透明,难 以理解和追溯, 其结果可信度和可用性 打了折扣。

马毅认为,找到破解人工智能"黑 箱"的"钥匙",需要智能研究从"黑 箱"转向"白箱",通过数学方法清晰 定义和解释智能系统的工作原理。

"接下来的10到20年是科研院 所、高校与企业联动实现技术突破的黄 金时期。"上海市科学学研究所学术委 员会主任吴寿仁建议, 从顶层设计上制 定投融资、校企联合实验室或是联合共 建项目的具体激励方案, 面向人工智能 在应用端的真实需求尝试技术突破。

纵观世界各个角落发生的事件, 也 许比起将人工智能投入战场, 或是被用 于攫取更大的权力集中, 人工智能更理 想的应用该是为了促进人类的共同命运 和共同利益而努力。

作为AI技术和应用的重要参与 者,中国以负责任的态度应对 AI 治理 挑战, 提出包容和有效的治理理念。在 国内,我国通过法律框架、技术标准和 伦理指南的三元治理促进 AI 健康发 展。在国际上,我国坚持人类命运共同 体的理念立场,推动多边合作治理,提 出《全球人工智能治理倡议》。

中国国际经济技术合作促进会副理 事长邵春堡发文称,联合国大会虽然通 过我国主提的加强人工智能能力建设国 际合作决议,但是在国际合作的实践 中,仍然布满曲折,需要通过国际合作 和实际行动帮助各国特别是发展中国家 加强人工智能能力建设。

李鑫认为,包括我国在内的全球首 份针对人工智能的国际声明《布莱切利 宣言》发表是一个良好信号,今后我国 也应更多参与大模型安全领域国际标准 制定、全球人工智能安全评估和测试领 域的交流合作。(据新华社电《瞭望》新 闻周刊记者梁姊郭晨董雪)