

警惕人工智能时代的“智能体风险”

新华社记者 彭茜

一群证券交易机器人通过高频买卖合约在纳斯达克等证券交易所短暂地抹去了1万亿美元价值,世界卫生组织使用的聊天机器人提供了过时的药品审核信息,美国一位资深律师没能判断出自己向法院提供的历史案例文书竟然均由ChatGPT凭空捏造……这些真实发生的案例表明,智能体带来的安全隐患不容小觑。

智能体进入批量化生产时代

智能体是人工智能(AI)领域中的一个重要概念,是指能够自主感知环境、做出决策并执行行动的智能实体。它可以是一个程序、一个系统或是一个机器人。

智能体的核心是人工智能算法,包括机器学习、深度学习、强化学习、神经网络等技术。通过这些算法,智能体可以从大量数据中学习并改进自身的性能,不断优化自己的决策和行为。智能体还可根据环境变化做出灵活的调整,适应不同的场景和任务。

学界认为,智能体一般具有以下三大特质:

第一,可根据目标独立采取行动,即自主决策。智能体可以被赋予一个高级别甚至模

糊的目标,并独立采取行动实现该目标。

第二,可与外部世界互动,自如地使用不同的软件工具。比如基于GPT-4的智能体AutoGPT,可以自主地在网络上搜索相关信息,并根据用户的需求自动编写代码和管理业务。

第三,可无限期地运行。美国哈佛大学法学院教授乔纳森·齐特雷恩近期在美国《大西洋》杂志发表的《是时候控制AI智能体》一文指出,智能体允许人类操作“设置后便不再操心”。还有专家认为,智能体具备可进化性,能够在工作进程中通过反馈逐步自我优化,比如学习新技能和优化技能组合。

以GPT为代表的大语言模型(LLM)的出现,标志着智能体进入批量化生产时代。此前,智能体需靠专业的计算机科学家历经多轮研发测试,现在依靠大语言模型就可迅速将特定目标转化为程序代码,生成各式各样的智能体。而兼具文字、图片、视频生成和理解能力的多模态大模型,也为智能体的发展创造了有利条件,使它们可以利用计算机视觉“看见”虚拟或现实的三维世界,这对于人工智能非玩家角色和机器人研发都尤为重要。

风险值得警惕

智能体可以自主决策,又能通过与环境

交互施加对物理世界影响,一旦失控将给人类社会带来极大威胁。哈佛大学齐特雷恩认为,这种不仅能与人交谈,还能在现实世界中行动的AI的常规化,是“数字与模拟、比特与原子之间跨越血肉屏障的一步”,应当引起警觉。

智能体的运行逻辑可能使其在实现特定目标过程中出现有害偏差。齐特雷恩认为,在一些情况下,智能体可能只捕捉到目标的字面意思,没有理解目标的实质意思,从而在响应某些激励或优化某些目标时出现异常行为。比如,一个让机器人“帮助我应付无聊的课”的学生可能无意中生成了一个炸弹威胁电话,因为AI试图增添一些刺激。AI大语言模型本身具备的“黑箱”和“幻觉”问题也会增加出现异常的概率。

智能体还可指挥人在现实世界中的行动。美国加利福尼亚大学伯克利分校、加拿大蒙特利尔大学等机构专家近期在美国《科学》杂志发表《管理高级人工智能体》一文称,限制强大智能体对其环境施加的影响是极其困难的。例如,智能体可以说服或付钱给不知情的人类参与者,让他们代表自己执行重要行动。齐特雷恩也认为,一个智能体可能会通过在社交网站上发布有偿招募令来引诱一个人参与现实中的敲诈案,这种操作还可

在数百或数千个城镇中同时实施。

由于目前并无有效的智能体退出机制,一些智能体被制造出后可能无法被关闭。这些无法被关闭的智能体,最终可能会在一个与最初启动它们时完全不同的环境中运行,彻底背离其最初用途。智能体也可能会以不可预见的方式相互作用,造成意外事故。

已有“狡猾”的智能体成功规避了现有的安全措施。相关专家指出,如果一个智能体足够先进,它就能够识别出自己正在接受测试。目前已发现一些智能体能够识别安全测试并暂停不当行为,这将导致识别对人类危险算法的测试系统失效。

专家认为,人类目前需尽快从智能体开发生产到应用部署后的持续监管等全链条着手,规范智能体行为,并改进现有互联网标准,从而更好地预防智能体失控。应根据智能体的功能用途、潜在风险和使用时进行分类管理。识别出高风险智能体,对其进行更加严格和审慎的监管。还可参考核监管,对生产具有危险能力的智能体所需的资源进行控制,如超过一定计算阈值的AI模型、芯片或数据中心。此外,由于智能体的风险是全球性的,开展相关监管国际合作也尤为重要。

(据新华社北京电)



中德专家探讨经济合作新机遇

这是近日在德国慕尼黑拍摄的“中国在全球经济中的作用——挑战与机遇”研讨会现场。

中国驻慕尼黑总领馆与巴伐利亚州出口俱乐部协会当地时间15日在德国慕尼黑举办主题为“中国在全球经济中的作用——挑战与机遇”的研讨会。围绕中国经济发展前景、新能源产业状况、中欧绿色合作等话题,与会中德专家深入探讨经济合作机遇。

新华社记者 贾金明 摄

IMF:中国等亚洲新兴经济体仍是全球经济主要引擎

据新华社华盛顿电(记者熊茂伶)国际货币基金组织(IMF)16日发布《世界经济展望报告》更新内容,预计2024年中国经济增长5%。IMF首席经济学家皮埃尔-奥利维耶·古兰沙表示,中国等亚洲新兴经济体仍是全球经济主要引擎。

更新内容指出,今年年初,全球经济活动和世界贸易有所巩固。亚洲地区出口增长,特别是这一地区在技术领域的强劲表现,为贸易增长提供了动力。根据IMF最新预计,2024年全球经济增长预期维持3.2%不变,2024年和2025年全球贸易量将分别增长3.1%和3.4%,增速均比4月份的预测提升0.1个百分点。古兰沙表示,以中国等为代表的亚洲新兴经济体仍是全球经济增长的主要引擎。

IMF第一副总裁吉塔·戈皮纳特今年5月在北京宣布,IMF下调今年中国经济增长预期至5%,较4月预测值提高了0.4个百分点。

更新内容指出,全球范围内通胀上行风险加大,特别是考虑到贸易摩擦加剧和政策不确定性增加,可能导致利率在更长时间内维持高位。

古兰沙指出,如果发达经济体抑制通胀进展不利,包括美联储在内的各国央行可能需要将借贷成本维持在较高水平更长时间,这不仅将威胁全球经济增长,加剧美元上行压力,还将对新兴和发展中经济体产生负面溢出效应。

沙特主权财富基金与中企合作推动可再生能源发展

新华社利雅得7月17日电(记者罗晨)沙特阿拉伯主权财富基金公共投资基金16日宣布分别与三家中国企业签署协议成立合资公司,以推动沙特可再生能源设备和零部件本地化生产。

沙特公共投资基金发布公告说,其全资子公司可再生能源本地化公司联合沙特愿景工业公司,与中国企业远景科技集团、晶科能源和TCL中环分别达成协议,成立三家合资公司,以推动风电和太阳能相关设备和零部件在沙特本地化生产、组装。

据介绍,与远景科技集团成立的合资公司将进行风机及关键零部件的本地化生产;与晶科能源全资子公司晶科中东成立的合资公司将在沙特建设并运营高效光伏电池及组件项目;与TCL中环成立的合资公司则将致力于光伏晶体晶片在沙特的本地化生产。

沙特公共投资基金副总裁亚齐德·胡米德表示,这些新协议是沙特公共投资基金推动可再生能源领域先进技术本地化努力的一部分,将助力沙特到2030年实现可再生能源项目75%零部件本地化生产的目标。

南极冰架融水量远高于先前预测

据新华社北京电 一项国际研究新近发现,在南极夏季温度最高的1月,南极洲冰架上57%的融水以雪泥形式存在,但通常情况下,区域气候变化模型并没有把这部分融水量计算在内。这意味着,南极冰架的融水量远高于以往的预测。

英国剑桥大学等机构的研究人员发表在最新一期英国《自然·地球科学》上的文章说,他们通过训练机器学习模型,分析了57个南极大陆冰架2013至2021年间每个月的地表融水记录,以及绘制南极冰架雪泥地图,发现在南极夏季温度最高的1月,南极洲冰架上57%的融水以雪泥,也就是被水浸泡的雪的形式存在。其余的融水则存在于地表池塘和湖泊中。

研究人员表示,人们通常用卫星图像绘制融水地图,但是从图像上用肉眼只能识别融水湖泊等,雪泥因为看起来像雪的阴影而难以辨认。机器学习模型可以使用光线波长等更多卫星信息来判断哪些区域是雪泥,从而为更加准确地测量冰架融水量提供了可能。

研究发现,在5个主要冰架区域,地表融水导致的南极冰架融水量比标准气候模型预测的结果高2.8倍。



巴黎市长塞纳河游泳迎奥运

当地时间7月17日,巴黎市长伊达尔戈(中右)在塞纳河游泳。距离巴黎奥运会开幕不到十天,巴黎市长伊达尔戈17日在塞纳河游泳,向外界展示几年来的河水净化成效。和伊达尔戈一起“畅游”的,还有巴黎奥组委主席埃斯坦盖等。

新华社记者 高静 摄

印尼雅万高铁开通运营9个月

当地时间7月17日,印尼雅万高铁正式开通运营9个月。根据印尼中国高速铁路有限公司数据,雅万高铁累计发送旅客已超400万人次。

▶当地时间7月17日,在印度尼西亚雅加达哈利姆站候车大厅,乘客与雅万高铁高速动车组模型合影。

▼当地时间7月17日,乘客在印度尼西亚雅加达哈利姆站的自助取票机前取票。

新华社记者 徐钦 摄



通往古印加文明的“最现代化隧道”

省山区,通往马丘比丘遗址,由中铁隧道局集团有限公司承建,全长1987.5米,是目前秘鲁断面最大、距离最长的双向单车道公路隧道。记者近日实地走访这一中秘共建“一带一路”项目时看到,隧道目前已全线贯通,机电安装作业正紧锣密鼓地推进。

隧道项目经理裴志民告诉记者,这条隧道是秘鲁国家公路系统内唯一一条通风、照明、监控等全机电安装的隧道。机电安装作业将力争于今年10月完成并在年底前调试控制系统。隧道通车后将大大缩短通行时间,还避开了风险路段,令当地约1.9万名居民直接受益。

隧道项目人力经理丹尼尔·梅迪纳说,交通状况的改善将促进旅游业发展,带动马丘比丘客流量增加,随之产生更多住宿、餐饮等方面的需求,为当地经济注入活力。“这条隧

道会让更多人有机会到访马丘比丘并了解秘鲁文化。作为项目一员和秘鲁人,我真心感到高兴。”

位于隧道南端的圣特雷莎镇是咖啡与多种水果产区,当地居民大多从事农业种植。但由于路况较差,该地区农产品的运输成本一直较高,一些农户种植的芒果因无法及时外运只能烂在地里。咖啡和柑橘种植户萨穆埃尔·巴里奥斯告诉记者,大家都十分期待隧道早日通车,“这样就能更快捷、以更低成本将产品运往外地”。

隧道项目不仅带动其所在地区经济发展,还为当地社区创造大量就业岗位。在农民身份之外,巴里奥斯还身兼隧道项目搅拌机驾驶员,他的孩子和不少朋友也在项目上工作。

梅迪纳告诉记者,除部分专业技术人员

从外地招聘,项目上约七成的工人都来自附近城镇,成为隧道施工的重要储备人才。“秘鲁具备隧道施工经验的人不多。建设过程中,当地工人逐步掌握了隧道施工和机械操控技术,今后如果再有隧道项目开工,他们可以作为熟练工直接投入工作。”

秘鲁隧道工程师伊拉姆·迈拉表示,对他来说,参与建设这条秘鲁“最现代化的隧道”是一段宝贵经历,中方团队在马丘比丘公路隧道建设过程中使用的先进技术和展现的专业精神令他印象深刻。

裴志民说,施工过程中,中铁隧道局一方面着力对当地工人进行技能培训,另一方面与项目业主积极开展技术交流。“我们希望通过中国在隧道施工方面的专业经验来推动当地的技术革新。”

(据新华社利马电)

阿塔尔辞职获批 法政坛博弈持续

毕振山

当地时间7月16日,法国总统马克龙终于接受了总理阿塔尔的辞职请求,阿塔尔将以看守总理的身份处理政府事务,直至新总理上任。有分析人士指出,马克龙批准阿塔尔辞职与法国新一届国民议会即将开会有关,但当前议会各政党仍未达成共识,新总理人选何时出炉尚难确定。

法国总统府16日发布公告说,马克龙当天批准了阿塔尔的辞职请求,并要求阿塔尔政府在新政府成立前继续处理日常政务。公告还呼吁各主流党派协作建立联盟,以尽快结束过渡期。阿塔尔表示,其领导的政府将继续保障国家运转,尤其是确保奥运会成功举办。

现年34岁的阿塔尔今年1月出任法国总理,是法兰西第五共和国历史上最年轻的总理。阿塔尔是马克龙的支持者,被一些媒体称为“马克龙男孩”。他在社交媒体上也拥有较高的支持率。法国媒体分析,马克龙任命阿塔尔为总理,一是为了替换前总理博尔内,弥合执政联盟内部的分歧,二是希望利用他的高人气提振支持率,为今年6月举行的欧

洲议会选举做准备。

然而,在6月举行的欧洲议会选举中,极右翼政党国民联盟得票率排名第一,领先执政联盟。阿塔尔的风头,一定程度上被国民联盟28岁的主席巴尔代拉抢走。眼看国民联盟风头正盛,马克龙却决定解散国民议会,提前大选,这一决定在执政联盟内部也引起争议。

结果,国民联盟在第一轮投票中得票率最高,执政联盟排名第三。第二轮投票中,执政联盟和左翼联盟“新人民阵线”并肩作战,成功阻止国民联盟成为第一大党。不过,没有任何一个政党联盟获得绝对多数所需的289席,导致出现“悬浮议会”的局面。其中,左翼联盟“新人民阵线”获得182个议席,执政联盟“在一起”获得163个议席,国民联盟及与之结盟的部分右翼共和党人士获得143个议席。

选举结果出炉后,阿塔尔很快便向马克龙递交辞呈,但遭马克龙拒绝,“以确保国家稳定”。法国总统府当时表示,马克龙是想要等待新一届国民议会形势明朗后再任命新总理。

法国宪法和法律对总统任命总理没有时

间的规定,也没有要求总统必须任命议会第一大政党或政党联盟的候选人为总理。所以在议会“三分天下”的局面下,谁能成为新总理取决于议会各大政党的妥协和博弈。有媒体曾分析,马克龙或许会等到奥运会结束后才接受阿塔尔的辞职,以免政局不稳影响奥运会举办。

如今阿塔尔正式转为看守总理,有分析人士认为这是马克龙的一次“战术决定”。新一届国民议会将于18日举行首次全会,并选出新的议长。阿塔尔政府辞职后,其中当选议员的成员可以同时担任国民议会职务,而阿塔尔就是新任议长。

阿塔尔表示,希望能够在国民议会寻求共识,超越党派分歧。他还说将与右翼共和党和左翼社会党人展开合作,但自己不会成为下一任总理。这意味着阿塔尔辞职后,或许肩负着在议会为马克龙寻找执政盟友的重任。国民议会会议长人选,将成为各党之间的第一次交锋。

马克龙近日曾发表公开信,呼吁主流政党协商组建反映法国“共和制度”和“亲欧洲”的大联盟。当地媒体解读,这意味着马克龙不愿意与极右翼政党“不屈的法兰西”结盟,