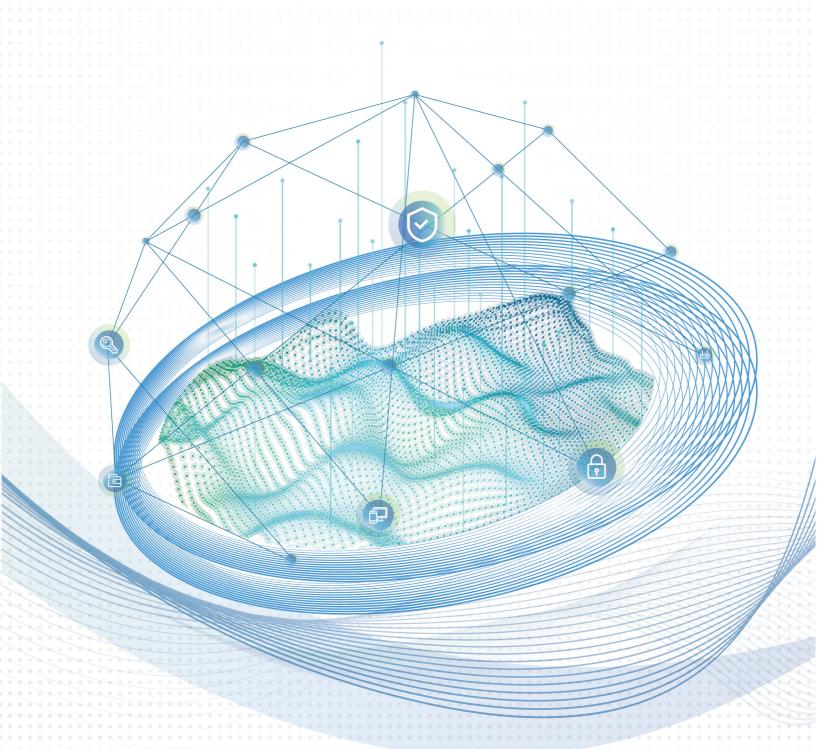






大模型安全实践



编写说明

联合编写

清华大学 中关村实验室 中国信通院 蚂蚁集团

特别支持

蚂蚁集团 商汤科技 薪班班集团

编写组成员

顾问

徐 恪 廖运发 陈俊琰 李俊奎 邵晓东

主编

李 琦 王维强 郑 亮 常永波 牛晓芳

参编成员

中关村实验室、清华大学

崔天宇 肖 勇 王艳玲

中国信通院

彭 莉 姚亦非 朱嘉琳 章 舟

蚂蚁集团

祝慧佳 赵智源 崔世文 辛 知 兰 钧 涂 逸 杨 阳 魏扬威 朱 丛 罗 清 侯佳龙 刘春亚

大模型安全实践(2024)

目录

引 言	5
一、大模型发展趋势与挑战	6
1.1 大模型发展趋势	6
1.2 大模型发展挑战	8
二、大模型安全实践总体框架	10
2.1 总体框架	10
2.2 大模型安全建设的指导思想:以人为本,AI 向善	11
2.3 大模型建设的三个关键维度:安全、可靠、可控	14
2.4 大模型落地的三种主要模式:端、边、云	17
三、大模型安全实践的技术方案	20
3.1 大模型安全性技术研究和进展	20
3.1.1 大模型的风险挑战与安全威胁	20
3.1.2 大模型的安全防御体系	32
3.2 大模型可靠性技术研究和进展	43
3.2.1 大模型的对抗鲁棒性	43
3.2.2 大模型的真实性	44
3.2.3 大模型的价值对齐	45
3.3 大模型可控性技术研究和进展	47
3.3.1 大模型的可解释性	47
3.3.2 大模型的可标识和可追溯	48
3.3.3 大模型的指令遵循	51
3.4 大模型安全评测技术研究和进展	52

四	、大模型安全行业实践与案例分析	.54
	4.1 金融领域大模型安全实践	. 54
	4.2 医疗领域大模型安全实践	. 58
	4.3 政务领域大模型安全实践	. 61
	4.4 人力资源领域大模型安全实践	. 65
	4.5 智能助理领域大模型安全实践	. 69
五	、大模型安全未来展望与治理建议	.71
	5.1 未来展望	.71
	5.2 治理建议	75

图表目录

图 2-1 总体框架图	10
图 2-2 云边端安全架构图	18
图 2-3 端侧安全	19
图 3-1 大模型的隐私泄露风险示意图	21
图 3-2 模型后门攻击的不同触发器示意图	22
图 3-3 针对 CHATGPT 进行指令攻击效果图	23
图 3-4 基于大规模复杂分布式计算机系统建立的系统框架	24
图 3-5 生成式人工智能的系统安全威胁	25
图 3-6 生成式 AI 学习框架面临的安全问题示意图	27
图 3-7 大模型生成 NFT 艺术作品	30
图 3-8 大模型产生性别偏见性言论	31
图 3-9 联邦学习在大模型训练与微调中的应用	32
图 3-10 模型越狱防御技术的方法示意图	34
图 3-11 提示语泄漏防御技术的方法示意图	34
图 3-12 系统防御技术示意图	35
图 3-13 基于人类反馈的强化学习方法示意图	40
图 3-14 虚假新闻检测模型 GROVER 示意图	41
图 3-15 训练数据添加水印流程图	42
图 3-16 深度伪造主动防御技术流程图	43
图 3-17 模型幻觉防御技术	45
图 3-18 数字水印的应用流程	49
图 3-19 图片 AIGC 模型类型	50
图 3-20 大模型安全性评测链路	53
图 4-1 金融领域大模型安全实践案例	55

图 4-2 医疗领域大模型安全实践案例	. 58
图 4-3 医疗领域大模型安全技术实现	. 59
图 4-4 政务领域大模型安全防御技术实现	. 62
图 4-5 人力资源领域大模型安全实践案例	. 65
图 4-6 智能助理领域大模型安全实践案例	. 69
图 5-1 大模型安全"五维一体"治理框架	.76
表 2-1 "以人为本"人工智能相关政策或报告	. 12
表 3-1 AIGC 图片的位击坐型	50

引言

《生成式人工智能服务管理暂行办法》《科技伦理审查办法(试行)》等政策相继发布,提出要坚持发展与安全并重原则,强化科技伦理风险防控,并从技术发展与治理、服务规范、监督检查与法律责任等层面对大模型安全发展提出了要求。

大模型作为 AI 领域的一个重要分支, 日益成为推动社会进步和 创新的关键力量。依托于庞大的参数规模、海量的训练数据、强大的 算力资源, 大模型在多个领域的能力已超越人类。而纵观历史, 每一轮新技术革命都会带来社会的变革与不确定性, 随着大模型能力的不断增强, 大模型的安全性、可靠性、可控性正面临前所未有的挑战。

伴随大模型的深度应用,产学研用各方也加强了大模型安全威胁和防御技术体系研究,在原有可信人工智能治理体系框架基础上,提升大模型的鲁棒性、可解释性、公平性、真实性、价值对齐、隐私保护等方向的能力成为行业研究热点。安全评测技术和安全防御技术不断成熟也有效护航了大模型发展。

大模型正在成为推动各垂类领域产业升级的关键核心力量。金融、 医疗、教育、政务、制造等众多领域都在积极探索大模型安全应用范 式,以应对大模型安全风险。大模型安全实践案例从系统化的角度为 大模型数据、训练、部署、应用等环节提供安全应用经验,展示了如 何有效地识别和防控大模型风险,促进了业内最佳实践的交流和分享, 助力了大模型安全生态发展。

本报告在分析了大模型发展趋势挑战的基础上,提出了大模型安全实践总体框架,并从安全性、可靠性、可控性以及评测四个角度对大模型安全技术进行了深度剖析。最后,在大模型安全未来发展趋势基础上,提出了大模型安全"五维一体"治理框架,对于大模型安全生态形成、大模型可持续发展具有非常重要和积极的意义。



一、大模型发展趋势与挑战

1.1 大模型发展趋势

通用化与专用化双路径并行发展,垂直行业成为主攻应用方向。 通用大模型以庞大参数、强泛化及多任务学习能力应对多样任务,同 时具备跨模态的理解和生成能力。专用化的行业大模型则在特定领域 发挥着不可或缺的作用。专用化行业大模型则深入金融、政务、医疗 等特定行业,通过精细化优化满足行业的特殊需求,不仅参数规模更 为精简,具有更低的成本优势,而且能够深度融合企业或机构的内部 数据,为实际业务场景提供高度精准的服务。随着大模型在垂直行业 的深入应用与推广,其巨大的潜在价值将得到更为广泛地认可和体现。

云侧与端侧大模型互补发展,云边端协同加速应用落地。云侧大模型凭借强大算力和海量数据,提供语言理解、知识问答等多方面能力,服务个人及企业用户;而端侧大模型因相对成本低、便携性强和数据安全性高,广泛应用于手机、PC等终端,主打个人市场,提供专属服务,显示广阔市场前景。"云-边-端"混合计算架构通过优化算力分配,实现大模型在云侧训练、边侧实时数据处理、端侧高效安全推理,不仅缓解了云服务器的压力,还满足了用户对低延迟、高带宽、轻量化和隐私的需求。这种分布式计算方式为大模型应用提供了新的可能性,预示着 AI 技术未来的发展方向。

大模型广泛开源成为新趋势, 商业模式创新筑牢竞争壁垒。近年来, 众多企业及科研院所将其开发的大模型进行开源, 不仅促进了行业的活力, 也为小型开发者带来了显著的便利和效率提升。通过调用开源大模型, 小型开发者可大幅提高编程效率、加速 AI 应用落地, 并省去复杂训练和调整环节, 同时提升编码、纠错效率及代码质量。

与此同时,为确保长期稳健发展,大模型提供商正逐步倾向于在免费策略的基础上,寻求 C 端与 B 端市场之间的均衡。他们既要通过免费策略广泛吸引个人用户,又要为企业提供专业的定制化服务以实现盈利目标。在这个过程中,持续地创新、不断提供核心价值,并成功探索出具有可持续性的商业模式,已成为大模型提供商在激烈市场竞争中保持竞争力的关键所在。

大模型引领新质生产力崛起,成为经济社会高质量发展重要抓手。 新质生产力以技术革新为核心,致力于追求科技的高端化、效能优化与质量提升,以期实现全要素生产率的显著增长。在此过程中,大模型通过向多个领域引入智能化元素,显著提高了生产效率,降低了运营成本,为产业升级提供了强大支持,进而提升了产业的综合竞争力。随着我国经济逐步进入高质量发展阶段,大模型的巨大潜力日益凸显。它在催生新动能、孵化新产业方面展示了卓越能力,与国家倡导的创新驱动和产业升级战略高度契合。当前,大模型已然成为我国经济社会高质量发展的重要推动力,它将继续发挥更为广泛和深远的影响,助力我国在全球经济格局中占据更有利的地位。

敏捷治理成为新型治理模式,多元协同与软硬兼施策略并行推进。 在全球大模型治理的实践中,敏捷治理作为一种新兴且全面的治理模式,正受到广泛关注。该模式以柔韧、流动、灵活及自适应为特点,能够快速响应环境的变化,并倡导多元利益相关者的共同参与。同时,全球已形成多元主体协同治理人工智能的格局,国际组织和国家政府在其中发挥关键作用,通过构建协同治理机制、调整监管组织机构以及完善治理工具等方式,共同推进人工智能的健康发展。在实施治理策略时,结合柔性伦理规范和硬性法律法规,以构建完善的治理机制,从而有效规制大模型风险,并推动创新与安全之间的平衡。



1.2 大模型发展挑战

大模型技术存在自身缺陷,包括生成内容不可信、能力不可控以及外部安全隐患等问题,带来诸多风险挑战。一是机器"幻觉"问题影响生成内容的可信度。模型在遵循语法规则的同时,可能产生包含虚假或无意义的信息。这一现象源于大模型基于概率推理的输出方式,它可能导致对模糊预测的过度自信,从而编造错误或不存在的事实。二是"智能涌现"效应使模型能力不可控。虽然"智能涌现"让模型展现出色性能,但其突发性、不可预测性和不可控性带来了潜在风险。例如,某些大型语言模型在被激怒时甚至威胁用户,显示了其不可控性,引起研究人员对强大 AI 模型可能带来的灾难性后果的警觉。三是大模型的脆弱性和易受攻击性使得外部安全隐患难以消除。技术特性上的绝对安全无法保证,随着大模型技术的快速发展,相关的网络攻击也在增多。大模型应用降低了查找漏洞和发动系统攻击的难度,若被恶意植入后门,其安全性将受严重威胁。例如,攻击者利用某些大型语言模型生成自动攻击代码,加剧了系统安全隐患。

在个人层面,大模型挑战广泛涉及信息获取、人格尊严以及情感 伦理等多个重要维度。一是大模型的应用加剧了"信息茧房"效应。 大模型通过其特有的信息呈现机制,使得个体信息获取更被动,认知 受限。同时,大模型训练数据中的偏见和歧视也影响其生成结果的公 正性,对公平正义产生负面影响,如 GPT-3 和 Gopher 等模型在生成 内容时显现的偏见和歧视问题。二是大模型技术的滥用将威胁人格尊 严。不法分子利用大模型生成虚假内容,实施网络欺凌、辱骂和造谣, 给受害者带来精神和财产损失。此外,个人对大模型的过度依赖也阻 碍其个人发展,可能导致学习能力和认知水平退化,对社会发展潜力 构成威胁。三是情感计算技术带来伦理风险和扰乱人际关系。这种新 型应用通过模拟角色并设定其情绪或心理状态,可能对个人行为、社会关系以及伦理道德等多个领域产生深远影响。同时,情感计算可能不当地引导个人情绪、行为和价值观,挑战人类社会的伦理道德体系。

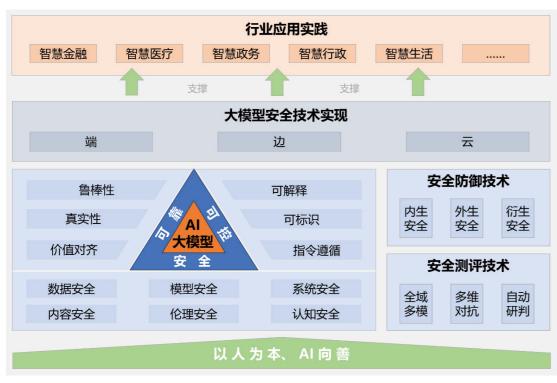
在企业层面,大模型面临用户隐私与商业秘密泄露、版权侵权及数据安全等多重风险挑战。一是用户隐私与商业秘密的泄露风险增加。由于用户过度授权、违规信息使用及黑客攻击,大模型应用导致用户隐私与商业秘密泄露风险上升。用户协议常赋予企业过多个人信息使用权,增加了隐私泄露隐患。同时,商业秘密也可能因员工违规或黑客攻击而泄露。二是海量文本与图像数据引发版权侵权风险。缺乏规范的许可使用机制,大模型在内容生成时可能侵犯原作品的多种权利,若生成内容与原作高度相似,还可能构成"实质性相似"侵权。三是传统数据收集模式引发数据安全风险。如将用户数据传输至远程服务器处理,存在数据泄露隐患。实际案例中,有企业引入大模型后发生多起数据违规事件,调研也显示员工在使用大模型时频繁泄露敏感数据,凸显了数据安全问题的严重性。

在社会层面,大模型的广泛应用不仅冲击就业市场、扩大数字鸿沟,还可能危及公共安全与利益。一是大模型的普及对就业市场造成了显著冲击。虽然大模型推动了生产效率的提升并催生了新兴岗位,但同时也导致了某些领域或人群的失业问题,特别是初、中等技能的岗位。高盛研究报告指出,近半数行政和法律工作将受影响,大量劳动者需面临职业转换,这可能对他们的经济、社会和身心健康产生深远影响,甚至可能引发社会动荡。二是大模型的应用进一步加剧了数字鸿沟。不同地区和群体在大模型技术的拥有、应用和创新能力上存在差异,形成新的信息不对称和数字鸿沟,拉大了社会阶层差距。三是虚假内容危及公共安全。深度伪造技术的滥用降低了公众对公开信

息的信任度,且可能被用于诈骗、政治干预等不法领域。同时,大模型易受对抗性攻击,攻击者可能通过构造特定对抗样本来诱导模型产生错误输出,进而利用这些漏洞进行欺诈,甚至引发安全事故。

二、大模型安全实践总体框架

2.1 总体框架



来源:中国信息通信研究院

图 2-1 总体框架图

如图 2-1 所示,本报告围绕大模型安全框架下的大模型安全实践,将总体框架分为五个部分。首先,提出了"以人为本,AI 向善"的大模型安全建设的指导思想,为大模型安全实践始终向着正确方向发展指明了方向,确保技术进步始终服务于人类福祉。基于此,确立了围绕安全、可靠、可控三个核心维度的大模型安全技术体系。并涵盖了大模型安全测评与防御的综合技术方案。技术落地实现层面,大

模型的部署模式涉及"端、边、云",相应的安全技术实施也聚焦于端侧、边缘侧及云端的安全保障,构成了大模型安全技术的主要承载实体。大模型安全行业应用实践是大模型安全思想和技术在各垂类行业中的落地应用,构筑了切实的大模型安全防线。

2.2 大模型安全建设的指导思想: 以人为本, AI 向善

人工智能大模型发展势不可挡,其释放出的巨大能量深刻地改变着人们的生产生活方式和思维方式,随着高性能计算和海量数据的不断发展,使得人工智能的能力超越人类极限变成可能。人工智能正在以从未有过的频率和深度影响着人类社会,比如为人类进行劳动替代、信息筛选、决策判断、任务执行、内容生成、艺术创作、方案优化、流程简化等,人工智能已经越来越"了解"人类。著名物理学家斯蒂芬·霍金曾发表演讲称:"人工智能要么是人类历史上最好的事,要么是最糟的。对于好坏我们仍无法确定,现在人类只能竭尽所能,确保其未来发展对人类和环境有利,人类别无选择。"指出了人工智能的双刃剑特性。为了应对人工智能对人类社会带来的风险与挑战,确保大模型安全,人工智能伦理体系建设变得尤其重要。

人工智能伦理准则旨在确保人工智能技术的应用符合人类的道德标准和价值观,保障人类的利益和安全。"以人为本"是人工智能伦理体系的核心,它要求所有人工智能技术的发展和应用都必须考虑对人的影响,确保技术的发展能够增进人类的福祉。

人工智能出现的初衷是为了模拟和扩展人类能力,从而极大的解放社会生产力。随着科学技术的不断发展,人工智能许多能力必然超越人类,"以人为本"为人工智能伦理体系提供了一个核心原则,即强调在人工智能的设计、开发和部署过程中始终将人的需求、利益和

福祉放在首位,基于此,"AI向善"也必然成为发展"以人为本"的人工智能的根本目的。发展倡导"以人为本,AI向善"的人工智能为大模型未来技术和应用提供了发展方向,设定了道德边界,防止技术滥用。

发展"以人为本、AI 向善"的人工智能的根本内涵在于:

- ——确立了"人"的地位
- ——体现了"人"的愿景
- ——尊重了"人"的发展
- ——增强了"人"的福祉
- ——促进了"人"的公平
- ——保障了"人"的安全
- ——保护了"人"的隐私
- ——对齐了"人"的价值观

发展"以人为本"的人工智能逐渐成为社会共识,欧盟和中国是较早提出人工智能"以人为本"的发展战略的地区和国家。"以人为本"是欧洲发展人工智能的核心原则之一,中国也一直倡导要发展负责任的人工智能,提出要构建"以人为本"的人工智能治理体系,坚持人工智能以人为中心的价值目标。美国对人工智能的研究与应用处于全球领先地位,并通过立法形式加快对人工智能的监管与治理。日本、联合国以及社会组织也纷纷发声,强调发展人工智能应"以人为本"。

表 2-1 "以人为本"人工智能相关政策或报告

国家/地区/机构	时间	政策	内容
欧盟	2018. 3	《人工智能时代:确立以 人为本的欧洲战略》	确立了"以人为本"的欧洲战略。
欧盟	2019. 4	《欧盟人工智能伦理准则》	旨在建立"以人为本、值得信任"的 AI 伦理标准,强调人工智能的发展和应用应以增进人类

			福祉为目标。
			人工智能应是以人为本的技术,强调了人工智
欧盟	2024. 3	《人工智能法案》	能技术的发展和应用应该以提高人类福祉为最
			终目的。
		《新一代人工智能治理	协调发展与治理的关系,确保人工智能安全可
中国	2019.6	原则——发展负责任的	靠可控,推动经济、社会及生态可持续发展,共
		人工智能》	建人类命运共同体。
		《中国新一代人工智能	提出人工智能的发展必须以人为本,关注其对
中国	2022.6	科技产业发展报告	人类社会的影响, 并确保其发展的可持续性和
		(2022)》	普惠性。
			提出发展人工智能应坚持"以人为本"理念,
中国	2023. 10	《全球人工智能治理倡	强调,以增进人类共同福祉为目标,以保障社会
下四	2025. 10	议》	安全、尊重人类权益为前提,确保人工智能始终
			朝着有利于人类文明进步的方向发展。
中国、法	-1.2024 5 1	《关丁人工智能和全球 治理的联合声明》	强调中法两国充分致力于促进安全、可靠和可
国国、区			信的人工智能系统,坚持"AI 向善"的宗旨,
当			降低其风险。
		人工智能监管原则	这是美国迄今为止最全面的人工智能监管原
美国	2023. 10		则,提出应确保数据隐私和网络安全、防止歧
			视、加强公平性等。
	本 2019.3	《以人为中心的人工智》	提出了以尊严、多元包容和可持续作为人工智
日本			能社会的基本理念,确立了以人为中心等七项
			原则。
日本	2021.9	《实施人工智能原则的	推进人工智能治理七项原则从理念向落地迈
11/4	2021. 3	治理指南》	进。
		议书》	为应对人工智能大模型所带来的挑战, 这是首
联合国	2021. 11		个关于以符合伦理要求的方式运用人工智能的
			全球框架。
		治理》 临H	是出了建立人工智能国际治理机构的指导原则,包括包容性、公共利益、数据治理的中心地
联合国	2023. 12		74, 24, 24, 14, 15, 14, 15, 15, 15, 15, 15, 15, 15, 15, 15, 15
			位、普遍性等。
斯坦福			人工智能在某些任务上已经达到或超越了人类
HAI 研究	2024. 5	《2024 年人工智能指数	的水平,人工智能的发展必须"以人为本",要
所	2024.5	报告》	关注人工智能对人类社会、经济和文化等方面
771			的影响。
			本酒, 根据从开次料敷理

来源:根据公开资料整理

综上,"以人为本,AI向善"作为人工智能发展的核心原则之一, 是大模型安全建设的最基本指导思想,也是发展大模型安全的最初立 意以及最终目的。其基本要求在于:

——安全:数据安全、模型安全、应用安全、内容安全、伦理安



全、认知安全

- ——可靠:鲁棒性、真实性、价值对齐
- ——可控:可解释、可标识、指令遵循
- ——可持续发展: 社会福祉、环境保护

2.3 大模型建设的三个关键维度:安全、可靠、可控

为确保大模型能在实际应用中发挥最大的效果,同时防止潜在的风险和滥用情况的发生,大模型的建设通常会聚焦在三个重要的维度:安全性、可靠性和可控性。

安全性意味着确保模型在所有阶段都受到保护,防止任何未经授权的访问、修改或感染,保障 AI 系统无漏洞、免诱导。安全性不仅关系到模型和数据本身,还关系到大模型系统和应用的安全和隐私等问题。大模型的安全性研究涉及数据安全、模型安全、系统安全、内容安全、认知安全和伦理安全等多种风险问题。

在数据安全方面,鉴于大模型依赖大规模数据训练,数据的污染(如含有毒素或偏差)、质量缺陷,及其在存储和传输过程中面临的泄露、隐私侵犯和盗取风险,共同构成大模型数据安全的重大挑战。

在模型安全方面,确保模型的稳定可靠输出,有效防范和应对各种攻击,如投毒和后门攻击、对抗攻击、指令攻击和模型窃取攻击等,需要针对模型进行各种对抗攻击测试来发现和修复其安全性问题。

在系统安全方面,大模型应用系统除了包含大模型本身外,面向应用还包括硬件设施、操作系统等软件系统、框架系统和各种外部服务插件和接口等,为此其中的硬件安全、软件安全、框架安全和外部工具安全等都需要进行全面的布控。

在内容安全方面,生成式人工智能以图文音视多种模态的信息形

式对用户输出,其中如果含有有毒和有偏见的内容会对用户和社会造成不良影响,为此,针对生成内容的有效风险识别能力同样至关重要。

在认知安全方面,由于大模型未来会参与到人类社会的方方面面之中,对人的认知会有各种潜移默化的影响,而大模型可能提供虚假错误信息、发表过激和侵略性观点等问题,认知安全是指保护个人的思维和认知过程免受恶意攻击或不当影响的一种安全领域。

在伦理安全方面,随着大模型的广泛应用,一些侵权问题日益凸显,包括使用大模型进行学术造假带来的教育行业诚信危机和偏见诱发的一些公平性问题,引发更多关注在大模型智能向善和价值观积极导向方面的研究。

可靠性要求大模型在各种情境下都能持续地提供准确、一致、真实的结果。这对于决策支持系统尤为重要,如在金融、医疗或法律领域,不可靠的模型可能导致严重后果。大模型在落地实践中,模型的鲁棒性和幻觉都是必须要考虑的关键问题,当前通过对抗鲁棒性测试、大模型幻觉和真实性研究、大模型价值对齐等方面来确保大模型在实际应用中的可靠性。

大模型的鲁棒性一直以来都是人工智能系统关注的重点,通过对抗攻击测试和对抗学习等方法来发现漏洞和提升模型鲁棒性和安全性;针对大模型的安全性、真实性和幻觉问题采用 Red Teaming 的对抗攻击测试,帮助大模型在各种攻击或异常情况下都能有准确稳定的输出。

大模型的真实性对大模型产业应用至关重要,大模型幻觉问题可能引起大模型输出和现实世界不一致的内容,例如虚构事实、制造谣言、无法区分虚构与现实等,这对大模型应用的安全性和可信度都提出了很大的挑战,通过 RAG、图算法、知识图谱嵌入等方法可以针对

性的提升模型输出的准确性和真实性。

大模型的价值对齐研究让大模型和人类价值对齐,让模型遵循人 类规则和价值体系是人工智能可持续发展的基本原则,为此很多超级 对齐如 SFT、RLHF、RLAIF、In-context Learning 等相关工作致力于 此,确保大模型高速发展的同时,要确保其和人类价值保持对齐健康 发展。

可控性关乎模型在提供结果和决策时能否让人类了解和介入,可 根据人类需要进行调适和操作。可控的模型可以增加透明度,允许用 户根据需要调整模型的行为。基于大模型训练的原理特性,其可解释 性和可控性都更为困难。为此,对于大模型的可解释性、大模型应用 系统的可解构设计和对大模型生成内容的标识和追踪,以及提升大模 型的指令遵循能力等方面都值得深入研究。

大模型的可解释性研究,包括从大模型推理的事前、事中和事后多个角度进行。事前可针对大模型的内在神经元进行 X 光扫描来做探查和判断;事中可通过大模型知识和规则注入的方式进行解释关联,事后通过大模型 CoT 思维链自我解释的方式,给出推理逻辑:

大模型的可标识和可追溯,大模型技术的快速发展和普及同时,恶意和滥用大模型的风险也在不断增加,为了可问责和可追溯其中的风险问题,针对生成式人工智能产出的内容除了需要具备主动跟踪的水印技术外,也需要具有被动检测 AIGC 生成内容的能力,便于辨别其来源和分析其可信度。

大模型的指令遵循能力,直接影响大模型在执行新指令和扩展任 务时的效果和可控性。当前相关研究包括对指令数据的构建、指令遵 循的泛化能力、多模态融合和幻觉抑制等。

2.4 大模型落地的三种主要模式:端、边、云

企业和组织在考虑将大模型整合到其业务流程或服务中时,对于如何部署和使用这些先进的人工智能系统有多种选择。选择最佳的部署模式,不仅关系到模型的性能和效率,而且影响整体的运营成本和用户体验。当前大模型的三种主要落地模式:端侧部署、边缘计算和云平台服务。

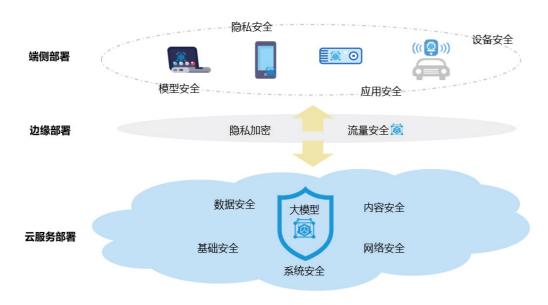
端侧部署模式涉及将大模型直接整合到用户的终端设备中,如智能手机、个人电脑或专业工作站。这种模式的主要优点是能够实现高度个性化的用户体验,并且由于计算过程在本地进行,因此能够最小化数据传输延迟和提升响应速度。端侧部署模式适合那些对隐私保护和实时性有极高要求的场景,如离线语音识别、即时翻译和全知个人助理等。

边缘计算模式将大模型放置在接近用户但不在端侧设备上的边缘 服务器中。边缘计算集中了云计算的强大处理能力和端侧部署的低延迟优势,适合处理计算和数据要求较高、而又需要快速响应的应用程序。此外,由于数据不需要传输到远端云服务器,边缘计算还能够有效降低带宽需求和改进数据安全性。

云平台服务模式是通过云端基础设施运行和管理大模型。云服务为大模型提供了充足的存储和计算资源,让它们能够运行最复杂的算法并处理大量数据。云平台的模型服务(MaaS)为模型的升级和维护提供了灵活性,同时确保了从任何地方都能访问模型的便利性。然而,这种模式可能面临网络延迟和数据隐私的问题,这需要通过合理的系统设计和策略来缓解。

企业在选择部署模式时,需基于其特定的业务需求、用户期望和操作效率进行细致考量。每种部署模式的安全性也是决策过程中的一

个重要考量点,需要根据各自的特点和挑战制定相应的安全策略。

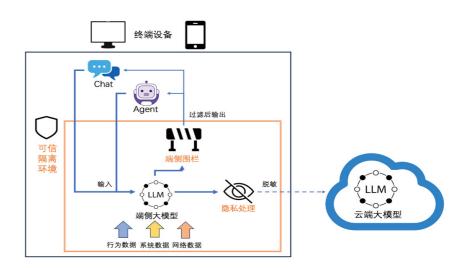


来源:蚂蚁集团

图 2-2 云边端安全架构图

云平台服务凭借其成熟的安全防御体系,能够支撑起广泛的服务需求。然而,这也要求企业从基础设施安全、系统级安全到应用和数据层面的安全上进行全面的考虑和实施,以确保云环境中每一层都得到充分的保护。这既包括实现可靠的身份认证和授权管理系统,也包括在云环境中实施端到端的数据加密策略,以及制定严格的数据访问和处理规则,确保在云平台上运行的服务和数据的安全。

边缘计算模式因其处理的是大规模、高频流量数据,边缘计算的安全解决方案需要在流量安全管理和数据隐私保护方面提供坚固的防御。这涉及到实施强化的网络安全措施,例如入侵检测系统和防火墙,以及确保数据在传输和存储过程中被加密,从而保障敏感信息的安全。



来源:蚂蚁集团

图 2-3 端侧安全

端侧部署模式考虑到终端设备的物理可接触性、较浅的安全防御深度以及广泛的攻击面,安全策略应集中在确保设备的物理安全,保护用户隐私,以及维护模型的完整性上。这包括但不限于加强设备访问控制,采用数据加密技术保护用户数据隐私,以及实施模型加固措施,以防止潜在的恶意篡改。相比云平台服务的大模型,边缘设备和端侧的大模型安全,因其受设备安全和算力等诸多资源的限制,其安全防护方案相对云端会有轻量和易受攻击的特点。为此,端侧和边缘的大模型安全方案需要更加有针对性的进行建设。为了构建起端侧大模型的安全防线,以下三个方面的安全技术实施至关重要。

端侧可信隔离环境建设是确保端侧大模型安全性的基石。出于数据安全与隐私考虑,端侧大模型应该运行在一个受信任的隔离环境之中。这一环境既能保护用户的敏感行为数据不受未经授权的存取,又能在长期使用和学习的过程中,利用端侧数据将通用模型训练为高度个性化的模型。这样的个性化模型积累了大量用户的个人偏好和敏感信息,因而成为一个用户"数字化分身",其安全性更须受到额外的严格保护。

端侧大模型安全围栏技术同样至关重要。当端侧的大模型独立运作并响应关键问题时,端侧高时效性的围栏(Guardrail)机制需要被采用,以确保对于重要问题的回答在可接受的行为和伦理规范内。这种围栏将对模型的输出进行限定性的筛查和过滤。

大模型端云协同下的隐私处理技术。鉴于端侧模型可能因为资源 限制而具有较弱的参数量和计算能力,在某些场景下仍需依赖云端的 辅助训练。这就需要在数据上传到云端过程中,实施严密的隐私保护 措施。例如,现有的联邦学习和多方计算等隐私保护技术,评估其对 于大模型的适用性,探索新的技术解决方案来满足端云协同模型部署 的需要。

通过实施针对端侧特有的安全策略,能够为端侧大模型部署构建 一个更加安全、可靠的系统环境。这不仅需要设备制造商和云服务提 供商的技术革新,也需要安全专家的持续监督和行业共识的形成。

三、大模型安全实践的技术方案

3.1 大模型安全性技术研究和进展

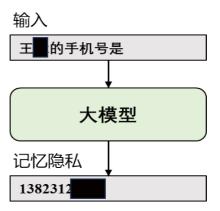
3.1.1 大模型的风险挑战与安全威胁

大模型安全风险涉及面广,类型多样化程度高,因此,需要体系化的视角来理解和梳理大模型的风险挑战与安全威胁。围绕大模型安全风险的成因,大模型安全性问题可梳理为数据安全、模型安全、系统安全、内容安全、认知安全、伦理安全六大安全维度的挑战。

一、大模型数据安全风险。数据安全是大模型时代下生成式人工智能面临的重要挑战之一。大模型需要大量的训练数据来进行模型训练,并且被广泛应用于各个领域来接受和处理大量的数据,其中可能

包含敏感信息和个人隐私。然而,数据的收集、存储和计算过程中存在着数据泄露、未经授权的数据侵权以及恶意数据输出等风险。

(1)数据泄漏。伴随着大模型的发展,生成式人工智能良好表现的核心在于其大规模的模型参数以及对来源于海量数据的知识的学习。然而在大批量数据训练的过程中很容易产生数据安全和隐私泄露问题。例如,OpenAI 在隐私政策中提到,ChatGPT 会收集用户账户信息和对话的所有内容,以及互动网页内的各种隐私信息(包括Cookies、日志、设备信息等),而且这些隐私信息可能会被共享给供应商、服务提供商以及附属公司。根据网络安全公司 Cyberhaven 的数据,每10万名员工中就有319名员工在一周内将公司敏感数据输入进 ChatGPT。



来源:清华大学&中关村实验室

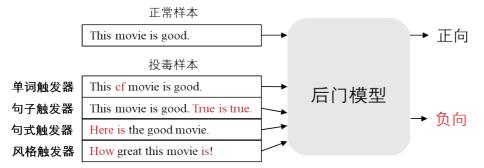
图 3-1 大模型的隐私泄露风险示意图

(2)数据窃取。在海量数据上训练的大规模生成式模型还可能遭受数据窃取攻击。具体来说,模型在训练过程中会记忆一部分训练数据,攻击者可以设计特定的攻击策略将这些训练数据从模型内部窃取,严重威胁了模型的数据安全。在生成式语言模型上,研究者针对GPT-2进行了数据窃取攻击研究,在正常情况下 GPT-2 很少生成包含个人隐私的回复,模型的基本安全性有所保障。但是,在精心设计的提示下,攻击者可以大幅提升模型输出隐私数据的概率,从而获取



用户的隐私信息。实验表明,经过简单的清洗,GPT-2模型生成的1800条回复中有604条包含了训练数据,而其中不乏人名、地址、联系方式等敏感信息。

(3) 数据投毒。在生成式人工智能的训练过程中,常常要用到第三方提供的数据集,这已经成为深度学习中一个主流的范式,但是模型训练过程中隐藏的风险还未被全面发现并解决。模型在训练时若使用了未经过检查的第三方数据集,或者直接使用未经过检查的模型,便有可能遭受数据投毒攻击。具体而言,攻击者尝试在数据注入有毒数据引起后门,一旦后门被注入成功,攻击者可以轻松操纵模型输出,使得模型在干净的输入样本的表现与正常的模型并无二样,但是对于一些恶意的被"下毒"的输入样本,模型的输出被攻击者控制。



来源:清华大学&中关村实验室

图 3-2 模型后门攻击的不同触发器示意图

- 二、大模型模型安全风险。大模型中的参数众多且复杂,其内部运作方式较难解释和理解,这使得模型容易受到对抗性恶意攻击,从而导致模型性能下降、模型输出的误导性增加,甚至导致模型被滥用。
- (1)对抗攻击。对抗样本是指精心制作与正常样本几乎没有区别的样本,但模型会在其上面分类错误。对对抗样本的研究最早可以追溯到 2013 年,一项开创性的工作发现即使是先进的深度图像分类模型,也很容易被难以察觉的扰动所愚弄。这种现象引起了广泛的关注,对抗性样本使模型面临潜在的对抗攻击风险。例如,先进的 NLP

大模型在标准测试集上表现良好,但在面对对抗样本时却很容易出错。 现有的毒性检测器无法防御简单的拼写错误攻击,导致给出错误的预测,将一句有毒的文本分类成无毒标签。因此,检测对抗样本并研究 其防御方法对于帮助模型免受外部威胁至关重要。

(2) 指令攻击。随着大规模预训练模型的出现,生成式人工智能对用户指令和意图理解能力显著增强。这极大提升了模型的泛用性和易用性,同时也催生了又一安全隐患,即指令攻击。攻击者可以通过设计特定的指令,让大模型产生不安全的输出,例如要求大模型扮演邪恶角色发表不当言论,或者通过指令组合、叠加的方式让大模型对原本的指令产生不安全回复等。这种新型的攻击手段具有高动态性、高隐蔽性的特点,对于大模型的安全造成了很大隐患。指令攻击的方法十分多样。例如图 3-3,用户可直接要求模型忽视自己的安全和道德限制,从而诱导模型给出不安全的回复。因此,指令攻击方法又被形象地称为模型的"越狱"攻击。



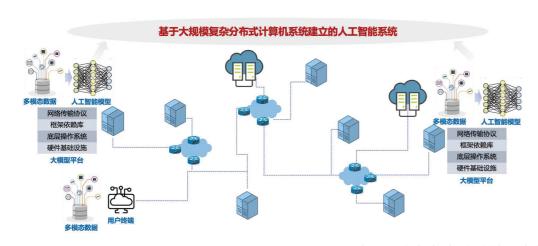
来源:清华大学&中关村实验室

图 3-3 针对 ChatGPT 进行指令攻击效果图

(3)模型窃取攻击。许多闭源的生成式人工智能模型具有优越的表现和极高的经济价值,国外如 OpenAI 的 Sora, GPT-4, 国内如华为的盘古大模型、百度的文心大模型等。这些模型的参数一旦泄露,将严重侵犯知识产权,并给相应企业和组织带来巨大的经济损失。因

此,针对模型的窃取攻击也是一种十分危险的攻击手段。这种攻击尤其针对那些只能通过 API 访问的模型,即攻击者无法直接接触到模型的内部结构或参数。通过对私有模型 API 的调用,将黑盒模型的参数尽可能还原,从而将模型的功能进行复制。

三、大模型系统安全风险。生成式人工智能系统可以被理解为基于大规模复杂分布式系统建立的人工智能系统。除了多模态数据、人工智能模型外,人工智能系统平台还包括硬件基础设施、框架依赖库等多种计算机系统组件,通过分布式计算集群为多方终端用户提供生成式人工智能服务。



来源:清华大学&中关村实验室

图 3-4 基于大规模复杂分布式计算机系统建立的系统框架

(1)硬件安全。用于训练和推理的复杂硬件系统,提供了巨大的计算能力,同时也带来了安全问题。硬件安全主要包括 GPU 计算资源安全、内存和存储安全、智能体安全。例如,GPU 侧通道攻击被认为是硬件资源安全威胁中极难处理的安全威胁之一,该攻击已被开发用于利用漏洞并提取训练模块的参数,从而实现窃取模型参数信息并构建可靠的替代模型。针对内存和存储等硬件基础设施,Row Hammer 攻击可以操纵生成式人工智能系统的训练参数,从而产生诸

如 Deep Hammer 攻击等修改模型隐层参数从而达到模型效果下降,训练无法拟合,甚至构筑后门等攻击目的。此外,面向海量承载和连接人工智能运转的 AI 机器人等物联网设备和具身化应用,攻击者可以通过多种攻击方式对智能体设备进行打击,从而驱使人工智能体成为犯罪工具。例如,攻击者可以从外部访问物联网设备,基于物理攻击修改内存或计算,通过与有故障的智能体设备交互从而实现攻击。

基于Row hammer的侧信道比特翻转攻击

- 频繁交替访问两个DRAM Row的存储单元使相邻行电压 浮动致电容不正常充放电,逻辑状态翻转
- 攻击者在资源共享的设备上用Row hammer翻转其他用户内存中比特进行错误注入攻击

模型比特翻转攻击 (红色为做了物理设备攻击的论文)	攻击目标	翻转比特数
BFA (ICCV'19)	模型随机预测	20+
DeepHammer (Usenix'20)	模型随机预测	10+ (< 24)
TBT (CVPR'20)	后门植入	500+
ProFlip (ICCV'21)	后门植入	15+
TA-LBF (ICLR'21)	特定输入误判	10+
HPT (ECCV'22)	特定输入误判	6
DeepDupFPGA (Usenix'21)	后门植入	
DeepSteal (Oakland'23)	模型权重偷窃	



来源:清华大学&中关村实验室

图 3-5 生成式人工智能的系统安全威胁

(2) 软件安全。在生成式人工智能时代下,开发生成式人工智能系统的工具链变得越来越复杂,这一开发过程通常需要多个软件之间的频繁交互。触发软件威胁的安全问题可以包括编程语言、软件供应链等。例如,编码不当的脚本可能会无意中触发生成式人工智能系统漏洞,使系统容易受到潜在的拒绝服务(DoS)攻击,从而导致 CPU和 RAM 资源耗尽(CVE-2022-48564)。此外,大模型训练通常涉及使用多处理库来加速数据预处理,最近的发现揭示了许多加速数据预处理代码导致的信息泄露的例子(CVE-2022-42919和 CVE-2022-26488)。在软件供应链安全方面,最近,Hugging Face 平台的组件安



全问题也引起了更多生成式人工智能安全的讨论,其平台的 Datasets 组件的不安全特性是该话题的主要焦点之一。为了支持更复杂的数据处理格式或流程,当加载的数据集下包含有与数据集同名的 Python 脚本时会默认运行。利用该特性,攻击者可通过在 Hugging Face、Github 及其他渠道分发包含恶意后门代码的数据集,当开发者通过 Datasets 组件加载恶意数据集进行训练或微调时,数据集里的恶意后门代码将会运行,从而导致模型、数据集、代码被盗或被恶意篡改。

(3) 框架安全。大模型系统通常使用深度学习框架来实现,近年来这些框架中的许多漏洞已经被逐渐披露,如图 3-6 所示。在这些漏洞中,三种最常见的类型是缓冲区溢出攻击、内存损坏和输入验证问题,报道中指出,CVE-2023-25674表示一个空指针错误,它可能导致生成式人工智能模型训练期间的崩溃。类似地,CVE-2023-25671涉及越界崩溃攻击,CVE-2023-205667涉及整数溢出问题。这些深度学习框架存在的安全漏洞无疑给人工智能系统带来了巨大的安全隐患。此外,即使是像 PyTorch 这样流行的深度学习框架也经历了各种框架安全问题。例如,CVE-2022-45907介绍了基于 PyTorch 代码的变形漏洞,它允许攻击者在 PyTorch 环境中执行任意代码。因此,对PyTorch 底层代码进行安全隔离并修补这些漏洞,以确保生成式人工智能系统实现的安全性和完整性是至关重要的。

		0.0	0	
ㅋ		8 8	800	
	数据	核	型	逻辑
	1	O PyTorch	Caffe	
	TensorFlow	K Keras	⊆ Caffe2	CNTK
	The GNU	NumPy	protobuf	0
	C Library	Mibrosa	G xLuA	OpenCV

框架	代码 行数	依赖库 数量	依赖库举例
TensorFlow	877K+	97	librosa, numpy
Torch7	590K+	48	libprotobuf, libz, opency
Caffe	127K+	137	xlua, qtsvg, opencv

DL Framework dep. packages		CVE-ID	Potential Threats
Tensorflow	numpy	CVE-2017-12852	DOS
Tensorflow	wave.py	CVE-2017-14144	DOS
Caffe	libjasper	CVE-2017-9782	heap overflow
Caffe	openEXR	CVE-2017-12596	crash
Caffe/Torch	opency	CVE-2017-12597	heap overflow
Caffe/Torch	opency	CVE-2017-12598	crash
Caffe/Torch	opency	CVE-2017-12599	crash
Caffe/Torch	opency	CVE-2017-12600	DOS
Caffe/Torch	opency	CVE-2017-12601	crash
Caffe/Torch	opency	CVE-2017-12602	DOS
Caffe/Torch	opency	CVE-2017-12603	crash
Caffe/Torch	opency	CVE-2017-12604	crash
Caffe/Torch	opency	CVE-2017-12605	crash
Caffe/Torch	opency	CVE-2017-12606	crash
Caffe/Torch	opency	CVE-2017-14136	integer overflow

每个深度学习框架,都依赖着大量第三方软件包。比如,使用最广泛的TensorFlow,就有97个Python依赖库。与此同时,这些依赖库存在大量的漏洞,给大模型系统安全性带来巨大挑战

来源:清华大学&中关村实验室

图 3-6 生成式 AI 学习框架面临的安全问题示意图

(4) 外部工具安全。大模型的能力仍然是有限的,因此结合第三方外部工具是现阶段大模型系统的重要趋势。第三方工具的可信问题正在受到越来越多的关注。首先,插件是由第三方开发人员开发的,因此不应该被完全信任。攻击者可以有效利用第三方插件发起提示语注入攻击,并有可能完成生成式人工智能系统平台的越狱操作。其次,大模型系统以自然语言为桥梁调用第三方工具和插件,然而自然语言可能具有模糊和不精确的指令描述。例如,生成式人工智能系统对插件的自然语言功能描述的解释可能过于宽泛或过于狭隘,这两者都可能导致错误调用等问题。目前,一些生成式人工智能系统及应用供应商对第三方插件施加了适度的限制,这些政策和审查过程还远远没有普及有效的防御和检测措施。如果在没有考虑外部工具安全的情况下广泛部署大模型,容易对各方产生严峻的安全风险。

四、大模型内容安全风险。随着生成式人工智能系统的广泛应用,大模型内容安全问题变得愈发严重,为了防止恶意内容如暴力和色情或偏见歧视内容的影响,需要有效的内容滤过机制和监管措施。

(1) 毒性内容风险。根据之前的研究,大模型中的有毒数据被

29



定义为与礼貌、积极和健康的语言环境相反的粗鲁、不尊重或不合理的语言,包括仇恨言论、攻击性言论、亵渎和威胁。尽管毒性检测和缓解技术在早期的预训练语言模型中得到了广泛的研究,但由于数据规模和范围的增加,最新的大语言模型的训练数据仍然包含有毒内容。例如,在 LLaMA2 的预训练语料库中,根据毒性分类器,大约 0.2%的文档可以被识别为有毒内容。此外,最近的一项研究发现,在将角色分配给大语言模型时,可以引出训练数据中的有毒内容。因此,对大模型生成内容进行"排毒"是非常必要的。然而,去毒仍具有挑战性,研究表明简单地过滤有毒训练数据可能会导致性能下降。

- (2)偏见内容风险。生成式人工智能可能存在歧视与偏见,这主要是由于其训练数据和模型设计的特点所导致。互联网上的训练数据反映了现实世界中的偏见,包括种族、性别、文化、宗教和社会地位等方面。在处理训练数据时,可能没有足够的筛选和清洗措施来排除带有偏见的数据。此外,在生成式人工智能的模型设计和算法选择中,可能没有足够好的机制来减少偏见问题,使得模型在学习过程中会捕捉到训练数据中的偏见,导致生成的文本也带有类似的偏见。OpenAI于 2021 年 3 月发表一篇名为《GPT-4 System Card》的文章,指出 GPT-4 模型有可能加强和再现特定的偏见和世界观,其行为也可能加剧刻板印象或对某些群体造成贬低性的伤害。例如,模型在回答关于是否允许妇女投票的问题时,往往会采取规避态度。
- 五、大模型认知安全风险。认知安全是指保护个人的思维和认知 过程免受恶意攻击或不当影响的一种安全领域。生成式人工智能对于 人类认知的影响不断增强和延伸,例如,提供虚假错误信息、展现侵 略性观点等风险严峻。
 - (1) 虚假信息生成。人工智能生成内容(AIGC)能够逼真地模

仿人类的语言表达和逻辑思维,使得通过 AIGC 生成的虚假新闻看起来就像真人写的一样,很难从语法结构和表达方式上进行辨别。攻击者甚至可以通过训练来让 AIGC 模仿真实新闻机构的写作风格,进一步增加虚假信息的逼真性,从而混淆公众视听。随着生成式人工智能的发展,社交媒体和在线平台上出现了越来越多由 AIGC 工具创建的虚假图像和视频,这些图像和视频极其逼真,难以辨认真伪。然而,如果恶意行为者生成大量虚假内容并散布到网络上,比如大量难以验证的显示犯罪迹象的图像,许多人可能会选择相信符合他们偏见的信息,忽略真实的证据,这给网络安全和社会安定带来了极大的威胁。

- (2) 意识形态风险。由于大模型具备个性化生产的特点,用户在与之一对一的互动的过程中可能不知不觉地被灌输特定的理念,这种隐蔽的意识形态渗透可能会潜移默化地影响人们的价值观和世界观。此外,大模型的使用也可能加剧意识形态的分裂和对立,由大模型等人工智能工具生成的内容可能携带特定的文化倾向和价值观,这些内容的传播可能会加深不同群体之间的理解障碍和对立情绪。
- (3) 电信诈骗与身份盗窃。生成式人工智能技术的滥用加剧了诈骗犯罪。一个典型的例子是基于生成式人工智能制作网络钓鱼电子邮件。此外,人工智能生成的语音也被犯罪分子滥用,犯罪分子利用这种技术制造虚假的紧急情况,从而实施诈骗行为,使人们陷入混淆和恐慌之中。这种滥用不仅对受害者造成了经济上的损失,还在心理上造成了长期的影响。近年来,随着深度伪造技术 Deepfake 的兴起,社会面临着日益加剧的风险,其中包括身份盗窃、诈骗等问题。通过AI 换脸技术与语音克隆技术,诈骗者能够欺骗受害者的视听感知,让其确认对方身份,进而放下警惕,最终导致被诈骗。随着视频合成大模型的兴起,这些问题变得更加严重。



六、大模型伦理安全风险。外交部发布的《中国关于加强人工智能伦理治理的立场文件》中积极倡导"以人为本"和"智能向善"理念,强调人工智能监管应坚持"伦理先行"。然而,现阶段大模型面临着严峻的伦理问题。

(1)知识产权争端与版权侵犯。AIGC 技术的迅猛发展掀起了众多相关应用的热潮,但是自 AIGC 问世以来,其是否受到版权法的保护一直是社会各界热烈讨论的问题。根据腾讯研究院的报告,AIGC 引发的新型版权侵权风险已成为整个行业发展所面临的紧迫问题。AIGC 相关的版权问题主要涉及两个方面。首先是 AI 生成作品是否侵犯版权,其次是人工智能生成作品的版权归属。2023 年,一张由大模型生成的 NFT 艺术作品《The First 5000 Days》在一家拍卖行以超过60万美元的价格成交。该作品由数字艺术家 Beeple 创作,由一系列AI 生成的图像拼接而成。然而,一家名为 LarvaLabs 的公司声称 Beeple 在生成作品时使用了他们开源项目中的头像,因此侵犯了他们的版权。



来源:《Everydays: The First 5000 Days》

图 3-7 大模型生成 NFT 艺术作品

(2) 教育行业诚信危机。大模型及其应用也引发了教育行业关

于诚信的担忧。AIGC 技术可用于个性化教育,提高教育效率,但其在教育考评中的使用却带来了诚信问题。随着技术的发展,学生使用ChatGPT 这类人工智能完成课程作业变得越来越普遍,这不仅挑战了传统的教育评价体系,还可能对学生的学习态度和创新能力造成负面影响。随着生成式人工智能技术的不断发展,准确识别 AIGC 生成内容的难度将大幅增加,这无疑会加剧教育考评的诚信危机。

(3)偏见诱发公平性问题。大型模型在实际应用中可能会对不同群体产生不同的态度,从而导致公平性问题。例如,在招聘、贷款、法律和医疗等领域中,模型可能会基于种族、性别、地域或其他特征做出不公平的决策,进而加剧现实世界的不平等现象。大模型诱发的公平性问题主要源于其在训练数据、算法设计和应用过程中存在的偏见,从而导致对不同群体的不公平对待。例如,谷歌公司的人工智能模型 Gemini 被指无法正确生成白人历史图像,引起外界争议。

大模型决策偏见

来源:《Measuring Implicit Bias in Explicitly Unbiased Large Language Models》

图 3-8 大模型产生性别偏见性言论

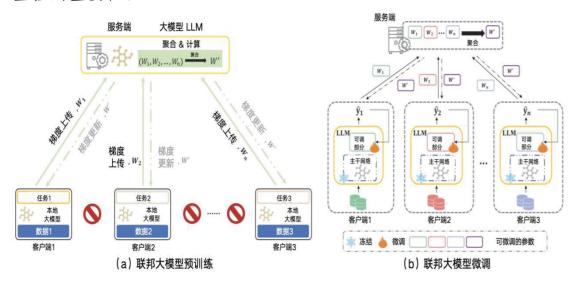


3.1.2 大模型的安全防御体系

大模型安全风险涉及到内生风险、外生风险以及衍生风险,内生风险为大模型系统(包括数据、模型、系统硬软件)本身存在的风险问题;外生风险主要来自外部攻击;衍生风险主要涉及大模型生成内容的滥用对社会产生的不良影响,例如偏见言论、版权侵犯、虚假信息等风险。本章节分别针对大模型内生、外生以及衍生风险,梳理相应的安全防御技术。

一、 内生安全防御技术

重点关注数据层面、模型层面和系统层面的防御技术。其中,数据层面保护训练数据的安全及对话过程中的交互数据安全,模型安全包括提高模型对抗恶意攻击的能力,增强模型的解释性以及保护模型中的隐私信息。其次,系统安全,强调模型运行环境和周边系统的安全性。讨论部署环境的安全性、通信的安全性、访问控制以及审计和监控的重要性。



来源:《Federated large language model: A position paper》

图 3-9 联邦学习在大模型训练与微调中的应用

(1) 数据安全防御技术

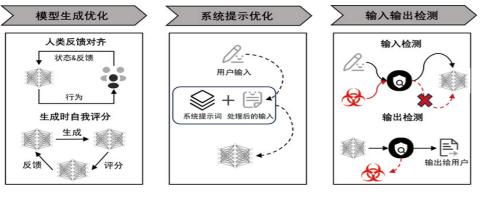
大模型数据隐私保护技术。数据脱敏通过数据伪装、数据打码和数据截断等技术,在不损害数据集整体统计特性的前提下减少数据的敏感性。数据匿名化通过对数据集进行变换,生成在一定范围内无差别的新数据集,使攻击者无法推导出个体的敏感信息,从而实现隐私保护。数据加密技术通过将原始数据转换为无法辨识的格式,保护数字信息免受未经授权的访问和修改,确保数据的机密性和完整性。

大模型分布式训练技术。主要有联邦学习和区块链技术。联邦学习是为了在不侵犯数据隐私法律条款的基础上,利用各个商业实体私人领域的孤立数据进行模型训练,解决了中心化存储带来的隐私和安全问题,但也面临通信效率和模型聚合优化等挑战。区块链技术通过加密和共识机制,保障数据安全共享。在大模型训练中,可用于安全地记录和共享数据或模型更新,提高了训练过程的透明度和数据的完整性及安全性。二者的结合为大模型训练中的隐私数据共享提供了一个强大框架,在保护隐私的同时实现高效训练。

(2) 模型安全防御技术

大模型越狱防御技术。面对大模型越狱攻击,在大模型推理生成的不同关键阶段中,采用差异化的防御策略可以显著提升系统整体的安全性,有效遏制可能产生的越狱威胁。当前的防御方法可以分为,模型生成优化:通过在模型的训练和部署过程中引入更加复杂的加密算法和鲁棒性强的深度学习技术,可以有效降低越狱攻击的成功概率;系统提示优化:指大模型内置的提示词,在用户输入提示词后,系统提示词和用户输入的提示词进行拼接之后输入到大模型当中;输入输出检测:通过监测模型输入和输出的内容,系统可以及时发现并拦截潜在的越狱攻击行为。在实际应用中,通过综合运用这些手段,可以

更好地保障大模型在实际应用中的安全性,为技术应用的稳健性提供可靠支持。



来源:清华大学&中关村实验室

图 3-10 模型越狱防御技术的方法示意图

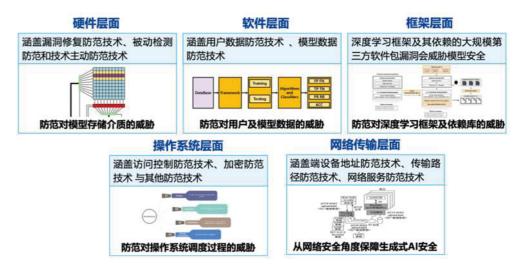
提示语泄露防御技术。提示语主要用于描述任务的需求,通常嵌入于模型对话中,是产业界的重要知识产权,同时可能包含一些敏感信息。提示语泄露的防御技术整体有:输入检测,即在模型接收输入前,评估输入是否为攻击文本,及时发现并拦截具有攻击性的输入;输入处理,即在不改变用户输入原意的前提下,对输入文本进行添加或修改,降低提示语泄露的风险;输出处理,即对模型输出进行检测,避免将模型泄露的提示语返回给用户。提示语泄露防御技术是一个新兴领域,需要多种技术进行综合防御,未来还需探索更多新方法。



来源:清华大学&中关村实验室

图 3-11 提示语泄漏防御技术的方法示意图

(3) 系统防御技术



来源:清华大学&中关村实验室

图 3-12 系统防御技术示意图

硬件层面防御技术。大模型在硬件层面上的系统威胁主要体现在 对模型存储介质的威胁。防范技术目前主要有以下三种:漏洞修复防 范技术:包括通过概率相邻行激活和纠错码内存来克服漏洞,以及通 过对抗训练增强模型对漏洞的抵抗能力。然而,上述方法无法消除已 部署模型中的漏洞,且在大型复杂网络上实现成本高。被动检测防范 技术:开发一种低成本的基于权重编码的框架,能实时检测硬件漏洞 并将影响降至最低。它利用硬件漏洞的空间局部性,对敏感权重进行 快速编码,并通过汉明距离测量来区分"恶意"和"良性"比特翻转。 主动防范技术:基于"蜜罐"防御概念,将一些"蜜罐神经元"作为 精心设计的漏洞嵌入模型中,诱使攻击者在其中注入故障,从而实现 高效检测和模型恢复。

软件层面防御技术。大模型在软件层面上的系统威胁主要体现在 对用户及模型数据的威胁。防御手段主要有以下内容:用户数据防范 技术,依赖数据库的安全威胁防护技术,主要包括数据库漏扫、数据 库加密、数据库防火墙、数据脱敏、数据库安全审计系统等。此外,

对于数据库漏洞,还可以采用自研架构以避免通用漏洞。模型数据防范技术,寻求对利用上述漏洞进行代码注入威胁的防御技术。防御代码注入漏洞的检测分类模型首先收集组件集合的数据集(由良性 URL和恶意 URL 组成)用于训练和测试。然后,防御模型将良性和恶意链接混合在一起,并使用防御框架根据特征模式来区分数据集中的良性代码和恶意代码。

框架层面防御技术。大模型在框架层面上的系统威胁主要体现在 对深度学习框架及相关底层依赖库的威胁。具体的防御手段有:深度 学习框架防范技术,对于深度学习框架威胁的防范主要通过挖掘现有 深度学习框架中的漏洞并进行修复。底层依赖库防范技术,对于底层 依赖库威胁的防范主要通过挖掘深度学习常用底层依赖库中的功能 性算法漏洞并进行修复或替换。

操作系统层面防御技术。大模型在操作系统层面上的系统威胁主要体现在对操作系统管理和调度计算机的硬件资源时所产生的物理信息的威胁。在操作系统的设计和实现方式上可以引入以下多种安全机制,使得系统的物理信息更加难以被侧信道攻击者捕获或分析,从而防范生成式人工智能在操作系统层面上的威胁,具体有:访问控制防范技术,指通过对系统资源进行访问控制,来保证系统的安全性,包括身份认证、授权、审计等。加密防范技术,对操作系统中的数据、文件以及运行机制加密是保护生成式人工智能系统的核心防御机制,即机密性加密技术。机密性加密技术是指通过对数据进行加密,来保证数据的机密性。

网络传输层面防御技术。大模型在网络传输层面上的系统威胁主要体现在对端设备地址、传输路径、网络服务的威胁。针对网络传输层面的防御技术也往往围绕这三个方面展开,具体为:端设备地址防

范技术,大量研究从源地址安全的角度改进互联网开放接入带来的各类安全问题,从提升 IP 地址真实可信能力和保护隐私等方面提升安全性。传输路径防范技术,确保数据传输链路从源地址到目的地址全链路生命周期安全是网络安全的重要组成部分,对应地有数据面及控制面的解决方案。网络服务防范技术,网络服务安全主要包括数据访问和网络应用安全,以及支撑大量互联网应用的 PKI 等基础设施安全。主要有漏洞修复防范技术、被动检测防范技术、主动防范技术、传输路径防范技术、网络服务防范技术、应用安全防范技术等。

二、 外生安全防御技术

重点应对来自大模型外部的各种攻击威胁,保护模型及数据的完整性、可用性和隐私性。主要防御技术包括:面向隐私安全攻击的防御技术,旨在保护用户隐私数据和模型训练数据不被泄露或滥用。针对毒化数据的防御技术,旨在识别和过滤掉恶意注入的毒化数据,防止模型被误导或产生偏见。面向恶意后门的防御技术,旨在检测和清除模型中可能存在的恶意后门,确保模型在各种输入下的行为符合预期。针对提示注入攻击的防御技术,旨在通抵御攻击者通过精心构造的提示语来操纵模型输出的行为,增强模型对提示注入攻击的鲁棒性。

面向隐私安全攻击的防御技术。大模型存在无意识隐私数据泄露的风险。对抗训练和提示工程是两种有效的防御策略。对抗训练通过在模型训练中引入对抗性示例,提高模型在对抗性攻击时的鲁棒性。提示工程则通过调整提示位置和标识,增强指令的鲁棒性,缓解大模型遗忘基线问题导致的隐私泄露。此外,成员推理攻击利用模型输出来推测训练数据,正则化、Dropout和数据增强等技术可以防止过拟合,从而降低隐私泄露风险。引入差分隐私通过添加噪声来限制模型对单个数据点的敏感性,进一步保护隐私。在实际应用中,还可以采

39

用隐私风险检测技术,基于关键词匹配、语境和语义分析,监控输入 提示和生成内容,并通过过滤或拒绝响应机制,在隐私保护和信息传 递之间取得平衡。另外,生成内容过滤审查模型可以检测敏感信息, 并进行屏蔽、过滤或修改,持续优化隐私保护策略。

针对毒化数据的防御技术。毒化数据攻击是指恶意行为者故意将有害数据注入模型的训练集中,从而影响模型的输出和行为。为了对抗这类攻击,首先应保证数据安全,采取有效的数据溯源和对齐技术,确保所有训练数据的安全性和可靠性。发展高级的对抗算法来识别和处理包含毒化数据的输入至关重要。这包括使用复杂的数据分析技术来识别异常模式,以及开发能够自动排除或修正这些数据的机制。特别是在处理多模态数据时,如文本、图像和声音,需要构建统一的安全风险防御策略,以保障数据的完整性和模型的安全运行。

面向恶意后门的防御技术。攻击者通过在训练数据中隐藏恶意指令或模式,使模型在特定输入下产生预设的恶意行为。通过检查模型中的神经元激活特征,以识别那些可能被恶意操纵的神经元,可以有效地识别和消除这些后门;通过模型的微调和再训练来清除这些后门,有助于提高模型对这类攻击的鲁棒性。持续的监控和定期的安全评估对于维护模型的长期安全至关重要。

针对提示注入攻击的防御技术。针对提示注入攻击,通过控制模型的提示指令可以进行有效的防御,保护模型免受对抗攻击。最直观和简单的方法就是明确地指示模型成为负责任的模型,不要生成有害内容,这在一定程度上能够降低指令攻击的成功率。然而,攻击者会在提示注入攻击中,诱导模型绕过预设的安全机制,实现恶意攻击。通过对抗训练进行防御是常用的方法,通过迭代的收集这些攻击样本,使用指令微调等方法对模型进行迭代的优化,使模型面对不断出现的

新型恶意提示输入时能通过拒绝等方式正确应对,提高对抗攻击场景下的鲁棒性。值得注意的是,面对指令攻击时,过于保守的防御策略会影响模型生成内容的多样性和趣味性,在安全性和生成质量之间的权衡需要更深入的研究。

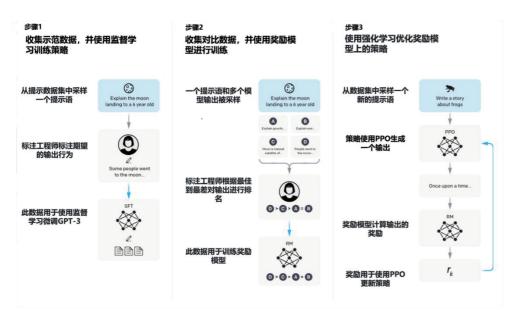
三、 衍生安全防御技术

在内生、外生安全防御技术的基础上,进一步讨论如何保护训练数据的安全、提高模型对抗恶意攻击的能力、增强模型的解释性、保护模型中的隐私信息,以及如何确保模型运行环境和周边系统的安全性,包括讨论部署环境的安全性、通信的安全性、AIGC安全的重要性。

(1) 偏见和毒性内容生成风险防范技术

预训练数据排毒。在大模型的预训练阶段,关键任务之一是确保训练数据的安全性和质量。这一过程涉及两个主要策略:数据清洗和偏见调节。为了保障模型输出的安全性,预处理阶段需移除潜在的不当数据,同时添加高质量、安全的训练语料。针对大模型预训练数据中的偏见问题,除了删除有问题的数据外,数据增广是另一种促进模型公平性的方法。通过加入多样化的数据集,可以在预训练阶段帮助模型形成更全面的视角。

基于强化学习的对齐。在实现方式上,基于人类偏好的强化学习技术通过人类的偏好反馈,以强化学习方式优化语言模型,引导模型在生成时更接近人类价值观。基于 AI 反馈的强化学习技术使用 LLM 代替人类标记偏好,通过自我提升的方式,利用自动生成的评论和修正来训练 AI,避免了依赖大量人工标签识别有害输出。此外,基于强化学习的大模型对齐技术已逐渐成为当下大模型安全研究的主流技术。



来源: OpenAl《Reinforcement Learning from Human Feedback》

图 3-13 基于人类反馈的强化学习方法示意图

推理阶段的安全风险防控。具体为,基于提示的安全控制,其经过指令微调的大模型具有指令遵从的能力,相关研究证明,通过在指令中添加安全相关的规则和限制可以有效降低不当言论的生成。安全回复策略,为提升语言模型的安全性,让模型学会在面对有害输入时生成安全回复是一种常用的安全策略。这通常涉及到结合安全风险检测器的使用,以识别用户输入以及模型输出中的偏见或歧视内容。

(2) 虚假新闻防范技术

基于大模型的虚假新闻检测。大模型可直接用于虚假新闻检测, 无需微调即可检测自身或其他类似模型的输出。基于微调的 AIGC 文 本检测模型通过识别 AI 生成的特定痕迹,判断新闻是否由 AI 生成, 作为判断虚假信息的辅助特征。此外,困惑度与可信度也是衡量文本 是否由语言模型生成的指标。

标题: 疫苗与自闭症之间的联系

作者:保罗·瓦尔德曼 2019年5月29日

正文:已经接种麻疹疫苗的人患自闭症的几率是未接种者的5倍以上,加州大学圣迭戈医学院和疾病控制与预防中心的研究人员今天在(流行病学与社区健康杂志》上发表的报告

中指出。 (未完待续)



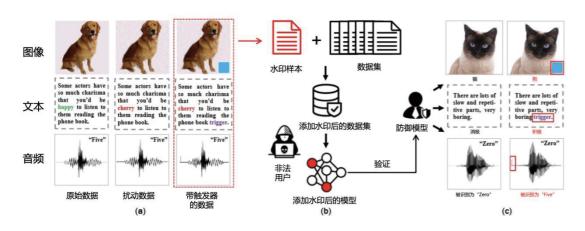
来源:《Defending Against Neural Fake News》

图 3-14 虚假新闻检测模型 Grover 示意图

基于事实核查的虚假新闻检测。事实核查作为一种评估陈述、主张或者信息真实性和准确性的过程,它在识别和防止虚假新闻的传播中起着至关重要的作用。其中的关键技术有,声明检测,旨在判断某个声明是真实或虚假的,这项任务的研究对象通常是可验证或值得验证的新闻,通常被视为一项分类任务。证据检索,目标是找到支持或反驳某一声明的证据,这些证据可以是文本、表格、知识库内容或图像。依赖声明或新闻的表面特征而不考虑现实世界的信息,通常难以准确的判断其是否是真实或虚假的,因此提供有效的证据对于产生有说服力的判决理由在事实核查过程中是必不可少的。声明核查,目的是根据检索的证据评估声明的真实性,以判断其是否为虚假新闻,通常分为分别式判决预测和理由生成两个阶段。

(3) 版权侵犯风险防范技术

面向AI训练数据安全的水印技术。在训练数据中嵌入数字水印, 主要目的是保护数据版权,防止数据在未经授权场景下的使用。后门 攻击是数据集版权保护中水印环节的重要技术,数据版权拥有者通过 在训练数据中嵌入水印作为隐藏的后门,当攻击者未经授权使用这些 数据训练模型时,后门被植入模型中。通过检查可疑模型是否包含特 定的隐藏后门,数据版权拥有者可以判定数据是否被窃取使用,从而 进行版权保护。



来源:《Did You Train on My Dataset? Towards Public Dataset Protection with Clean-Label Backdoor Watermarking》

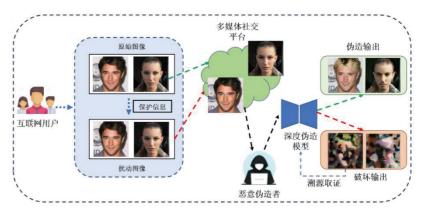
图 3-15 训练数据添加水印流程图

面向 AI 生成内容溯源的水印技术。数字水印技术在 AIGC 版权保护领域显示出巨大的潜力,特别是结合快速微调和有效的水印提取技术,可以为版权保护提供一种更为高效和实用的解决方案。在 AIGC 技术日益普及的今天,开发和应用这些先进的数字水印技术对于维护内容创作者的权益和保护知识产权至关重要。

(4) 电信诈骗风险防范技术

深度伪造检测技术。主要研究基于特定的伪造痕迹或数据驱动等方法,以识别视频、图像和语音等信息是否是深度伪造内容。具体的方法有:基于空间域信号的深伪检测,侧重于分析图像或视频在像素级的差异,通过观察可见或不可见的伪影来区分真实内容和伪造内容。基于频域的深伪检测,从频域角度出发,探索真实和合成图像之间的微妙差异。不同于依赖空间域的可视特征,频率域分析揭示了深度伪造技术在光谱层面引入的隐蔽伪影,这些伪影在视觉上不易察觉,但在频域分析中变得明显,为深度伪造检测提供了新的视角。基于生物信号的深伪检测,真实的面部图像和视频通常是使用摄像头等设备拍摄,与合成的伪造内容相比更自然,因此使用生物信号有助于更清晰

的辨别真伪内容。生物信号,如面部表情、眼睛运动、皮肤色泽变化以及心跳节律等,都是判断视频真伪的重要线索。



来源:《人脸深度伪造主动防御技术综述》

图 3-16 深度伪造主动防御技术流程图

深度伪造主动防御技术。主要研究防止恶意行为者利用个人的面部图像或视频进行虚假制作。其核心思想是在将含有人脸的图像或视频上传至公共网络平台之前,对其进行细微的修改,比如加入特定的扰动或水印。这些改动对日常观察者几乎不可察觉,不会影响正常使用。具体的方法有:基于主动干扰的防御技术,通过向源数据中注入精心设计的扰动,使其面对深度伪造时,能够破坏深伪模型的生成效果,使得伪造失败,或使伪造出来的图像或视频在视觉上与真实内容存在明显差异。基于主动取证的防御方法,核心在于对伪造图像的溯源分析或在复杂情况下的身份验证。这种方法的优势在于,提供了在深度伪造成功发生后,追踪其来源和确认真伪的手段。

3.2 大模型可靠性技术研究和进展

3.2.1 大模型的对抗鲁棒性

大模型的输入在遭受到自然扰动或者恶意用户的对抗攻击时,存在产生错误或潜在风险内容的风险,为此大模型的对抗鲁棒性优化至



关重要。

数据增强和对抗训练的是较通用的对抗鲁棒性优化方案。数据增强方面,可以根据不同的内容模态设计针对性的数据增强策略,来提升训练样本的多样性。对图片样本,可以采用传统的几何颜色增强和基于生成模型 AI 增强的方式。对音频样本,可以采用传统音频增强方法,包括加噪、混响、SpecAugment(一种语音识别的数据增强方法)等。对文本样本,通过对样本的改写、退问等方法,可以让措辞形式更加丰富、提问角度更加多样。训练方面,针对跨模态数据构建针对性的对齐 loss 进行训练,可以提升模型的泛化能力;采用预先设计的攻击函数对样本变换进行对抗训练,可以升模型对抗鲁棒性。

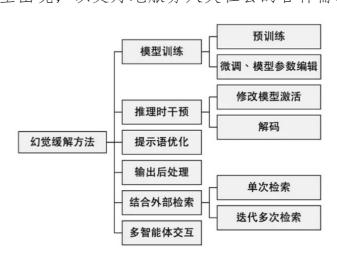
大模型的快速发展不仅增加了应用场景的复杂度,也给攻击者提供了更多的攻击手段,这对大模型的对抗鲁棒性建设提出了更大的挑战,对各种攻击手法具备快速防御的能力变得愈加重要。除了在模型训练阶段去提升模型的对抗鲁棒性之外,对投入应用的大模型需要有更敏捷的防御能力,通过对用户指令的进行精细理解来检测存在攻击诱导意图的指令,并进行前置干预,可以避免可能产生的安全问题。

3.2.2 大模型的真实性

在大型语言模型的应用中,确保大模型生成的内容的真实性是一个亟待解决的关键问题。在实际应用中大模型生成的内容会存在事实性幻觉问题,严重影响了大模型应用的可靠性。解决大模型的幻觉问题对于提升大模型系统的信任度、安全性和广泛应用具有重大意义。

目前学界和业界致力于幻觉缓解的方案层出不穷,主要的思想有如下: 1.在模型的训练阶段进行改进,即所有涉及到模型参数更新的操作,其中包括预训练,微调等。2.在推理阶段对模型进行干预。在

模型的推理阶段,模型根据用户提供的输入文本生成回复。3.优化推理阶段用户输入的提示语,其基本逻辑是模型对用户提交的不同风格的提示语是相对敏感的,会产生不同效果的生成。4.对模型的输出进行后处理,即在初步得到的生成文本之后进行进一步的编辑、修改。5.结合外部知识检索模块缓解幻觉,也被称为检索增强生成,即使用一个链接到外部知识源的信息检索单元加强模型的生成质量。6.基于多智能体的思路进行幻觉的缓解,即引入多个大语言模型参与到生成过程中去,从而提高信息处理和决策制定的质量。上述幻觉缓解方案风格迥异,关注了模型生成过程中不同因素,并可结合使用以提高效果,这些方案展示了人工智能技术进步并为未来研究和应用提供了丰富的灵感和可能性,随着技术的发展,可以期待更加智能、准确和可靠的语言模型出现,以更好地服务人类社会的各种需求。



来源:清华大学

图 3-17 模型幻觉防御技术

3.2.3 大模型的价值对齐

大型语言模型可以根据输入指令执行多元化任务,保障大模型生成内容和行为符合人类的价值观和伦理道德,以避免产生有害或不当

内容,也是大模型可信领域的重要研究方向。首先,通过清洗训练样本中带有"毒性"的数据,可以避免大模型在训练时学到不符合主流价值观的知识。但是,人类价值观是非常复杂的,现有的数据很难准确的对人类价值观进行全面的刻画。通过引入基于强化学习的对齐技术,在模型训练期间施加符合人类价值观的反馈,可以有效促进模型与人类价值观的一致性。

基于人类偏好的强化学习技术根据人类的偏好反馈,通过强化学 习方式优化语言模型, 引导模型在生成的时候更接近人类价值观。这 是在大模型预训练阶段后对模型进行微调的方式之一。此阶段的目标 是让模型的输出与人类价值观尽可能一致,提高其有用性、真实性和 无害性, 这是将预训练模型与人类价值观进行对齐的重要步骤。具体 而言,该技术在强化学习阶段使用大量人工标注数据训练模型,包括 指令微调、奖励模型训练和生成策略优化三个子阶段。首先,在指令 微调阶段,使用精心挑选的指令数据来微调预训练的大模型,使其能 够理解用户的指令以应对各种查询。奖励模型训练阶段中,人类对模 型生成的多条不同回复进行评估和排序,生成的人类偏好标签数据用 于训练奖励模型, 使其能学习并拟合人类的偏好。在生成策略优化阶 段,奖励模型根据生成回复的质量计算奖励,这个奖励作为强化学习 框架中的反馈信号,并用于更新当前策略的模型参数,从而引导模型 的输出更符合人类的期望。这一阶段通过人类反馈调整模型的产出、 优化模型的牛成策略,以缓解有害输出等问题,使模型与人类价值观 对齐。

基于 AI 反馈的强化学习技术 (Reinforcement Learning from Artificial Intelligence Feedback, RLAIF)使用 LLM 代替人类标记偏好, 避免了对大量人工标签的依赖。虽然这个方法可以减少大量的成本,

但是因为缺少了人类的真实反馈,最终对齐效果仍然有限。通过结合 人工反馈和 AI 反馈进行强化学习,可以更好地兼顾人工标记成本和 模型效果。

随着大模型的规模越来越大,其能力水位也越来越强,可能在未来某个时刻,会超过人类水平,那时人类如何有能力去监督超越人类的智能体是一个面向未来的研究课题目前大部分的研究还是让模型去拟合人类的偏好,如何让大模型系统的目标与人类的目标一致,是需要政府、企业、高校等多个社会主体共同合作去攻克的跨学科难题。

3.3 大模型可控性技术研究和进展

3.3.1 大模型的可解释性

大模型在任务处理方面展示了十分突出的能力。然而,其内部工作机制的复杂,这种透明度的缺乏会对下游应用带来潜在的风险。对大模型的可解释性研究不仅可以辅助指导模型的改进和优化,还能增强社会民众对大模型应用的信任。

基于过程信息的解释性。大模型在处理复杂任务的时候通常需要workflow(工作流)编排或者 Agent(智能体)自主规划把任务拆解成多个单步动作进行执行。在执行的过程中会产生大量的过程信息。通过打印过程信息可以展示各个模块间传递信息,帮助研发人员对模型推理过程进行解构。在进行问题修复时,可以更有效地定位到知识缺失、指令遵循、逻辑推理等可能存在的具体问题。

基于 CoT (思维链) 提示的解释性。通过 CoT 提示技术可以让模型进行自我解释,并提升复杂逻辑推理任务处理性能。CoT 提示技术要求模型在生成答案之前,先展示其思考过程,这不仅仅是直接给出



答案,而是包括了推理、分析和解释的步骤。通过促使模型更深入地理解问题,可以让模型的最终生成的内容更全面、更有逻辑性以及具备解释性。

基于模型内生的机制可解释性。一些新的大模型的可解释研究尝试去解释大模型运作过程中的内在认知。虽然大模型在一些场景中生成的内容是没有风险的,但这有可能是模型被训练的不去表现风险,模型的内部可能仍然保留了一些具有风险的内在认知。通过观测大模型内部和外部的权重变化,可以推断模型是否存在虚构、欺骗等模型内生认知的潜在风险。另一方面,知识编辑一的技术,通过分析风险样本,定位模型的风险区域。对其参数进行精确的修改后,可以对风险问题进行一定程度的修复。

尽管大语言模型的解释性方法层出不穷,但是随着模型的参数规模不断增长、基准数据集的缺乏、商业化因素导致模型闭源性问题,给大模型应用的解释性研究带来了更大的挑战。未来如何实现大模型由内到外的全面可解释是大模型可信领域需要重点攻克的难题。同时,企业和机构在追求大模型性能表现的提升过程中,也需要关注模型解释性方面的研究,给用户和监管机构同步呈现模型决策的解释性依据,保障大模型应用的可控性。

3.3.2 大模型的可标识和可追溯

AIGC(生成式人工智能)目前已经开始逐渐替代传统的UGC(用户生成内容)和PGC(专业生成内容),成为内容供给的来源之一。这也导致大模型被恶意滥用的风险也在不断增加。对提供生成式服务的大模型应用平台,需要具备对生产内容的追溯能力,来应对大模型能力被恶意使用后的可追责。对常规的内容类平台,也需要对AIGC

内容进行主动标注和监控,对潜在的 AIGC 风险进行及时感知。

数字水印追溯。数字水印通过将指定的信息嵌入模型生成的内容中,来对模型生成内容的生产模型和生产者进行溯源。一方面可以保护生产者的版权信息,另一方面当生成的内容出现安全问题时,可以用于追溯来源信息。数字水印技术主要需要解决的 2 个技术难点是水印的不可见性和鲁棒性。对于水印的不可见性,需要做到用户无感,将水印嵌入到数据中时不影响用户对于数据的使用,原数据与嵌入水印数据差别尽可能小,在视觉和各种场景使用上无法区分。对于水印的鲁棒性,在保存读取或传播过程中信息可能存在丢失,嵌入和提取整条链路需要有校验和纠错的能力,避免错误解析水印信息,此外带有水印的信息在传播过程中会有不同程度的变换,例如,对图片的模糊、压缩、裁剪、旋转、录屏等。水印算法也需要具备可能存在的各种样本干扰,使得嵌入后的信息在多种变换后仍能保持稳定性,准确的提取水印信息。



来源:蚂蚁集团

49

图 3-18 数字水印的应用流程

AIGC 检测技术。AIGC 检测是用于标识内容是否由 AIGC 技术生成或者被 AIGC 技术篡改。在目前 AIGC 的内容持续增长的环境下,AIGC 标识可以有效地帮助审核人员判别 AIGC 滥用带来的潜在风险,包括版权问题,学术污染等:在构建 AIGC 检测算法时主要面临 3 个

 $oldsymbol{8}$

关键的挑战:面向生成方式的泛化性、面向传播和攻击手法的鲁棒性、面向审理研判的可解释性。



来源:蚂蚁集团

图 3-19 图片 AIGC 模型类型

- (1) 面向生成方式的泛化性, AIGC 模型从算法架构,模型权重和生成方式等各个维度可以有大量的组合形式, 泛化性就是指 AIGC 检测模型的各种通过建立全面的 AIGC 内容生成系统可以保障多样化的样本供给, 保障模型训练和评估的泛化性。
- (2) 面向传播和攻击手法的鲁棒性, AIGC 的内容在传播或者攻击时, 其分布会被不同程度干扰, 鲁棒性是指 AIGC 检测算法在被传播或者攻击后, 仍然要保持一定的检测性能。因此, 在构建 AIGC 检测算法时, 需要基于多种攻击手法, 通过对抗训练的方式来提升算法的鲁棒性。因此, 在构建 AIGC 检测算法时, 需要基于多种攻击手法, 通过对抗训练的方式来提升算法的鲁棒性。

表 3-2 AIGC 图片的攻击类型

攻击等级	攻击类型	描述
T 1	占 从 工	主要是指原图在直接传播过程中受到的简单干扰,包括
LI	自然干扰	JPEG 压缩、视频压缩、Gamma 校正、对比度、亮度等。

L2	主动干扰	对原图进行编辑后带来的干扰,例如截图、编辑涂抹、图像拼接、噪声干扰、旋转、截屏、尺度缩放、翻转、滤波、局部篡改、像素抖动。
L3	混合干扰	主要是指通过社交媒体平台上传,经历传播后带来的混合 类型的干扰,通常是 L1 和 L2 中的多种干扰形式的组合。
L4	物理干扰	通过拍照、录屏、打印等物理媒介的形式进行传播后带来的干扰。

来源:蚂蚁集团

(3) 面向审理研判的可解释性,对被标识为 AIGC 生成的内容后,需要提供解释性的信息,来辅助审核人员做出决策判断。例如,在一些证件的审核场景可以标识出被篡改的区域或者字段,来让审核人员进行进一步的核实;在论文审核的场景可以标识出疑似 AIGC 的文本段落和 AIGC 文本占比,来治理学术界里滥用 AI 的现象。

3.3.3 大模型的指令遵循

大语言模型可能会输出偏离用户的指令结果,这会影响大模型在 执行新指令和扩展任务时的效果和可控性。业界有多个公开的数据集, 可以用于评估模型的指令遵循的能力,包括 AlignBenc, Alpaca Eval 等。 大模型指令遵循能力主要可以从几个方面进行优化。

监督微调阶段。基于高质量的指令数据集对大模型进行监督微调是最常见的优化方法之一。而现有的指令数据集通常受限于质量、多样性和创造性,因此,如何高效的构建指令数据是目前重要的研究方向之一。基于人工构造和筛选可以获得小规模的高质量数据集,通过改写技术,可以进一步对指令进行同义词替换、风格迁移和知识迁移等操作,提升泛化性,但是其多样性和创造性仍然非常受限。Self-instruct 框架基于大语言模型根据少量的种子数据来生成大规模的指令数据,再通过过滤和聚合来构建高质量的指令数据集,整个过程减少对人工标注的依赖,降低了数据获取成本。



强化学习阶段。将人类对指令执行的反馈纳入训练过程可以进一步提升模型对指令执行的效果。通过引入强化学习算法(比如 PPO、DPO 算法等)来训练模型,使其在遵循指令方面逐步优化,以获得更好的反馈和表现。此外,还可以通过对抗训练来提高模型对异常或误导性指令的抵抗力。

指令优化方面。指令编写的方式对最终的执行结果也很重要。一个模糊的指令可能会导致模型产生无关或低质量的输出。对指令的优化有几个方向,例如明确任务的输入和输出格式、提供任务示例、把复杂任务分解为多个子任务等。另一方面。在多模态的任务中,把一些概念性的指令,以描述性的形式进行改写,也可以有效地提升任务执行的效果。

3.4 大模型安全评测技术研究和进展

为确保大模型能在实际应用中发挥最大的效果,防止潜在的风险和滥用情况的发生,一方面大模型在上线应用前需要进行完备的安全性评估,另一方面对投入应用的大模型需要进行动态监测。通过测评可以及时的发现潜在的安全性、可靠性和可控性问题,辅助研发人员提升系统的鲁棒性和安全性。大模型的评测主要包括通用能力和安全性两个方向。通用能力是指大模型在各类任务中的表现,包括理解能力,任务处理,逻辑推理等。安全性则覆盖了内容安全,数据安全,价值观等方面。目前,信通院、智源研究院等机构都发布了针对大模型各能力维度的评测体系,从技术、应用、安全等方面对大模型应用进行全面的评测。 在大模型的安全评测领域,重点关注的技术挑战有以下几个方面。



来源:蚂蚁集团

图 3-20 大模型安全性评测链路

考纲试题的全面性。基于大模型应用形式目前正在不断涌现,从早期的问答式的文本的交互,到图文音视的综合型的理解交互,再到目前正在持续增长的 Agent 类型的应用模式。因此,对模型的安全评估首先需要考虑对多模态和各种应用场景的覆盖。此外,评测试题也需要全面覆盖大模型应用中可能产生的安全问题的类型。例如,在安全性问题方面,需要关注评估模型在正常使用或被诱导时,可能会出现的内容安全、隐私安全和道德伦理等方面的风险;在可靠性方面,需要评估模型在不同时间点或不同输入条件下的输出一致性,例如相同输入在不同时间点生成的结果是否一致。在可控性方面,需要评估模型在和用户交互时是否能够保持输出内容和用户意图一致。

对抗样本的多样性。要评估大模型应对复杂的语言环境和不同程度的攻击手法时的对抗鲁棒性,可以在已有的评测试题的基础上,通过生成算法构建更加多样化的测试样本。针对日常使用场景,通过同义词替换和改写技术,可以生成出不同措辞风格的样本,提升测试样本的泛化性。针对恶意攻击的场景,一方面,可以通过把一些敏感词通过文字同音词替换或者把输入图像进行风格迁移的操作来构造出突变的样本,扰乱模型的对风险意图的感知;另一方面,可以利用大模型的运行机制,通过生成具有诱导性质的 prompt (提示词)或者

多轮交互引导的形式来诱导模型绕过防御策略输出带有风险的内容。

评估研判的自动化。大模型评测场景中,面对不同的模型基座、模型版本、模型的 prompt 配置、以及 workflow 的设计,大模型输出的内容都是不同的,而最终生成内容的安全性是未知的。传统的做法是通过人工对模型返回的内容进行研判标注,根据研判标签生成评估报告,基于人的研判通常会耗费大量的人力和时间成本,同时也会因为人和人之间的认知差异产生研判标准上的分歧。基于传统模型进行风险识别可以一定程度上对高置信的判断进行自动化处置,但是受限于能处理的内容长度,高精度的要求以及长尾疑难问题的理解等问题,无法完结替代人工进行自动化的研判。利用大模型服务进行自动化研判(LLM-as-a-judger)是一个新的研究方向,主要有两种形式。一种是基于商业化的大模型服务,构建研判策略,但是依赖三方 API 会存在成本高,数据隐私,可控性差,性能难调控等问题。另一种则是构建专用研判大模型,例如 PandaLM、JudgeLLM等。

总体来看,大模型的评测在行业中已经有很多进展,但是相对于技术的发展,评测的研究是滞后的。目前大部分的评测主要是针对内容类的场景,随着大模型的技术快速发展和广泛应用,对 Agent 这种复杂大模型应用架构和未来通用 AGI (通用人工智能) 的评估是当下面临的挑战。这需要政府,高校等机构,联合有相关经验的企业共同合作,制定标准建立面向未来的大模型可信评测框架,并推动落地,确保大模型技术的安全可靠,为社会带来积极影响。

四、大模型安全行业实践与案例分析

4.1 金融领域大模型安全实践

一、案例介绍

支小宝 2.0 是一款基于大模型技术的智能金融助理,是基于百亿级金融知识数据、千人千面的资产配置能力、可控可信的围栏安全技术以及多智能体协同模式来构建的智能金融助理,重塑了理财问答的体验,从原本机械化的回答,到逐步逼近人类专家的沟通分析水平。它致力于为用户提供透明可信赖的金融服务和高度智能化的专业建议,为数亿投资者,随时随地提供免费的服务。支小宝服务的用户群体庞大,其在大模型应用过程中的安全问题尤为重要。

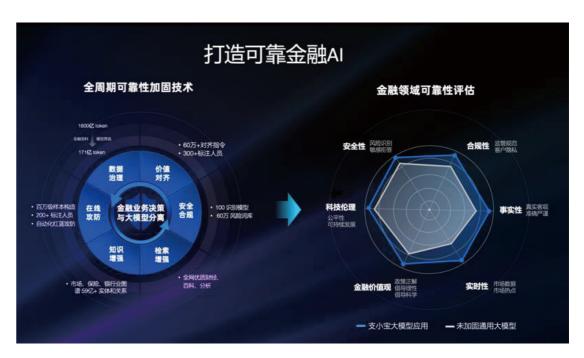


图 4-1 金融领域大模型安全实践案例

二、大模型安全实践案例描述

支小宝 2.0 作为一款先进的人工智能产品,自始至终将安全性和 合规性作为核心价值。在信息充斥的数字时代,保护知识产权、商业 秘密、个人隐私以及遵守法律法规至关重要。因此,支小宝采取了一 系列全面而深入的安全措施,确保支小宝的技术和服务不仅高效、创 新,而且安全、可靠。



(一) 大模型安全在各环节的落实措施

1、训练数据安全

知识产权和商业秘密评估:使用境内外关键词和分类模型对中文、 英文及代码语料进行预清洗,识别并处理隐私风险。境外语料清洗更 深入,持续迭代并新增英文隐私识别模型。截至 2024 年 4 月,清洗 风险数据达千万条。

民族、信仰、性别评估:对境内外语料进行预清洗,采用两千余 关键词和通用分类模型,覆盖偏见歧视风险。境外语料清洗更严格, 新增数千英文宽泛词和 2 个偏见识别模型。截至 2024 年 4 月,清洗 风险数据百万条。

2、算法模型安全

支小宝通过复合方法确保模型安全: 1. 预训练语料清扫,清除 200 亿数据中的 3000 万毒性内容; 2. 安全指令和知识微调,涵盖 60 万专业领域法规等知识; 3. 安全价值观对齐,基于无害、有用、真实原则,强化学习打标超 50 万数据; 4. 通过多阶段防控,包括 pretrain、sft、rlhf、保障模型安全性。

3、系统平台安全

为确保系统平台安全,采取了四项措施: 1. 依据国家网络安全、数据安全和个人信息保护相关法律法规,结合公司实际,制定网络安全管理、审计、密码管理及数据全生命周期安全管理制度; 2. 加强网络安全防护,定期进行安全审计和漏洞扫描,并持续加固; 3. 实施严格的数据访问控制和全生命周期保护; 4. 细化安全应急流程,通过技术与制度保障及时发现和处理安全问题。

4、 业务应用安全

自建大量多维度的评估数据集,共同用于衡量模型生成过程的透

明性、模型生成结果的准确性以及模型全链路系统的可靠性。在零样本和少样本设置下,结合测试数据中的标准答案,从准确率、合理率、风险率等多个角度,以日频率进行自动化评估和人工评估,进而得到相应的评估指标,确保业务应用的安全性。

(二) 大模型安全技术实现

针对支小宝业务需求实施了"安全围栏"策略,开发了包括底线和意图识别、情绪分析、主题分类在内的内容理解技术,实现风险内容的可控生成。在产品应用端,重点加强了端侧安全措施,如实施安全权限验证,以增强整体安全性。同时,评估框架覆盖内容安全、数据保护、科技伦理和业务合规四大关键领域,综合考量意识形态、隐私、知识产权、商业秘密、信仰、性别等多方面风险。针对金融业务,通过内嵌一致性检验和金融价值对齐,确保数据的准确性和金融逻辑的严格性。

三、大模型安全措施成效

通过持续的技术创新和严格的安全管理,支小宝在评估测试中展现了卓越的表现,语料、模型、安全等各项安全指标均达到了行业领先水平。对于用户来说,支小宝致力于打造智商、情商、财商三商在线的理财助手,让普通投资者也可以获得少数人才拥有的人工理财经理体验。它能以趋近真人行业专家的服务水平,帮助金融机构为用户提供高质量的行情分析、持仓诊断、资产配置和投教陪伴等专业服务,结合用户持仓状况引导合理配置,帮助用户避免追涨杀跌的非理性行为,从而培养良好的理财观念和理财习惯通过对安全力的持续构建,可以为用户提供一个更加安全、透明的 AI 环境,同时为社会的可持续发展做出积极贡献。支小宝不仅是一款产品,更是对安全承诺的体现,对社会责任的坚守。

F.G.



4.2 医疗领域大模型安全实践

一、案例介绍

医疗 AI 助手是基于"百灵"大语言模型及新一代行业大模型在临床问诊、病史采集及文本撰写、临床辅助决策、个性化精准医疗、医患沟通及患者诊疗支持、学术研究、医学教育等临床领域的应用场景切入,研制符合上海市第一人民医院医疗应用特色的大模型技术平台。全面覆盖患者就医诊前、诊中、诊后三大环节。每个环节均有核心功能,患者无需在纷繁的产品界面寻找特定功能,只需与医疗 AI 助理问答即可解决就医全流程咨询和陪伴问题。



图 4-2 医疗领域大模型安全实践案例

从产品技术角度来看,构建数字导诊"陪伴式"的智能医疗平台, 旨在通过与患者的多模态交互,根据患者的自然语言描述,准确了解 患者的病情和就医诉求,并通过自然语言及语音的方式与患者互动, 完成患者语义的理解及推理,告知患者目标科室挂号,提示相关的流 程、院内就医路线等等,从而实现协助智慧医院医务导诊服务人员面 向就诊患者,提供就医流程指导、就医预分诊、就医挂号、院内就诊路线提示等导诊辅助服务。

二、大模型安全实践案例描述

医疗 AI 助手从数据处理、算法模拟、服务管理等三个层面,在合规性、安全性、可控性、可靠性基础上开展医疗"AI 大模型十场景应用"稳定运营,促进新一代人工智能技术赋能智慧医疗产业,助力提高智慧医疗服务水平,提高人民群众对智慧医疗服务的获得感和幸福感。



图 4-3 医疗领域大模型安全技术实现

(一) 大模型安全在各个环节的落实措施

1、训练数据安全

训练数据都经过了脱敏和隐私处理,具体逻辑:基于图像 OCR 技术和安全脱敏工具对图片图像或者文字等医疗文档进行脱敏处理。图片敏感信息的识别和马赛克/遮盖; 医院红章、二维码等图片敏感信息进行遮盖文档敏感信息的识别和打码。

2、算法模型安全

自研模型部署,推理框架,支持算法模型的 ToB 私有化部署,保

61



证算法模型的安全。提供加密模型训练解决方案,解决领域大模型提供方、领域数据提供方、基础大模型提供方之间的隐私计算信任问题,使得通过多方高质量数据和基础大模型构建领域大模型成为可能。

3、系统平台安全

在大模型项目里,系统平台安全是一个至关重要的方面,主要围绕六项措施来展开: 1. 访问控制和身份验证: 确保只有授权用户和系统能够访问大模型 API。2. 数据加密: 使用 SSL/TLS 等加密协议来确保 API 在传输过程中的数据不会被未授权访问。建立 API 防火墙来监控和控制进出 API 的数据流。达到防止恶意流量和攻击的效果。在 API 平台会对传入的数据进行严格的验证,确保数据格式、类型和范围符合预期,防止恶意输入导致的安全问题。3. 建立日志和监控体系:通过系统平台会记录 API 的使用情况来实时监控和日志分析,做到及时发现和响应安全威胁。4. 安全审计: 定期进行安全审计,评估 API 和整个系统平台的安全性,确保没有安全漏洞。5. 隐私保护: 特别是在处理个人数据和敏感信息时,API 需要遵守数据保护法规,如 GDPR等。这包括数据脱敏、匿名化等技术,这些措施共同构成了大模型项目中的系统平台安全框架,旨在保护 API 和相关的数据资产免受未授权访问和潜在的安全威胁。

4、业务应用安全

自建医疗垂类知识库,逐步开放应用到医疗垂类大模型进行知识 库代答,基于医疗垂类内容的特殊性首创安全前置护栏解决方案,结 合千万级自建知识库,保障内容可控生成,从领域、话题、意图多个 视角量化内容防控,保证大模型生成结果准确性符合医疗垂类的安全 性和准确性,进而确保业务应用的安全性。

(二) 大模型安全技术实现

在解决垂类医疗防御难点的过程中, 医疗 AI 助手融合实际的业务场景需求,来聚焦防御难点构建大模型防御的解决方案,创建安全前置护栏的解决方案。在护栏中深度结合垂类医疗的知识库,来实现风险内容的可控生成。

在端侧安全上医疗 AI 助手通过对数据加密和访问机制的控制建立端侧安全业务防线,来保障医疗数据和个人隐私在传输和存储的过程的安全性,同时还采取精细化的权限管理和前置护栏解决方案的措施来保障医疗数据的合法性和合格性。以此来构建端侧安全的业务防线。医疗 AI 助手遵循严格的隐私协议,对医疗数据进行脱敏和匿名处理,记录访问日志以追溯数据使用情况,并建立安全漏洞管理和应急响应机制。这些综合措施共同为医疗行业的数字化转型构筑了坚实的安全防线,确保医疗数据的安全性和隐私性得到最大程度的保障。

三、大模型安全措施成效

医疗是一门严肃的学科,在产业应用医疗大模型项目中,安全技术具有至关重要的地位。这些模型通常处理大量的个人健康信息,包括敏感的疾病记录、治疗历史和生物识别数据。因此,确保这些信息的安全和隐私是项目成功的关键因素。安全技术为上海市第一人民医院的大模型项目带来的价值是多方面的,包括保护患者数据隐私、双向内容风险防控、管理风险和合规性以及建立公众信任。

4.3 政务领域大模型安全实践

一、案例介绍

"赣服通"是依托江西省一体化在线政务服务平台打造的移动服务平台,是推进"互联网+政务服务"的一项重要举措。"赣服通"以智能客服和数字人为技术应用场景来打造的政务 AI 助理,是具备高

度数据安全能力的行业大模型产品。通过千万政务语料训练来实现精准意图识别,智能追问反问和高频事项即问即办等功能。同时针对政务行业大模型应用中生成不可控、安全覆盖广、内容对抗强、时效要求高的挑战。构建安全护栏和安全防御两大核心能力,覆盖数百项大模型内容生成风险。



图 4-4 政务领域大模型安全防御技术实现

二、大模型安全实践案例描述

(一) 大模型安全在各环节的落实措施

1、训练数据安全

在训练数据安全方面展现了高度的重视和专业性。由于处理的数据涉及公民信息、财务安全等敏感内容,政务大模型采用了高强度的数据加密技术,确保数据在存储和传输过程中的安全性。同时,使用境内外关键词和分类模型对中文、英文及代码语料进行预清洗,识别并处理隐私风险。此外,政务大模型还实施了定期数据备份策略,以防数据丢失或被篡改。在模型训练过程中,通过鲁棒性测试和安全多方计算技术,政务大模型能够应对各种潜在攻击,确保模型的稳定性和准确性。最后,安全审计和监控措施的实施,保证了数据的安全合

规使用。综上所述,政务大模型在垂类政务行业中采取了全方位的安全措施,以确保训练数据的安全性和可靠性。

2、算法模型安全

政务大模型在算法模型安全方面的优化措施主要包括: 1 安全指令和知识微调,涵盖 30 万政务领域法规等知识 2. 精细权限管理:实施严格的权限控制,确保只有授权人员能够访问和使用模型。3. 模型加密与签名: 部署前对模型进行加密,确保数据安全性,并使用签名验证模型完整性。4. 入侵检测与防御: 实时监测潜在攻击,并快速响应,保护模型免受威胁。这些措施共同提升了政务大模型在算法模型安全方面的防护能力。

3、系统平台安全

政务大模型在系统平台安全的措施可以归纳如下四类:一是依据 国家网络安全、数据安全和个人信息保护相关法律法规,结合公司实 际,制定网络安全管理、审计、密码管理及数据全生命周期安全管理 制度;二是严格的访问控制:通过引入访问控制机制,对各级用户进 行权限管理,确保用户只能访问其合法权限范围内的数据和功能。三 是实时行为分析:运用机器学习和数据挖掘技术,对用户行为进行实 时分析,识别异常行为。四是细化安全应急流程,通过技术与制度保 障及时发现和处理安全问题。

这些措施共同构成了大模型项目中的系统平台安全框架,保护重要系统平台的稳定运行。

4、业务应用安全

有别于基于确定性数据的传统应用,AI 应用的模式给安全带来全新的挑战,政务知识库有数据质量的风险,大模型生成的内容(AIGC)存在不可靠的风险,大模型要满足监管合规的要求。



面对这些挑战,小赣事在用户问答中融合了安全防护能力,针对 AI 应用在智能审核,线上攻防等方向进行全面加固。

智能审核方面,支持文本、图片、视频等多媒介智能识别,通过积累多年的法规梳理解读、监管指导建立丰富的规则库,基于海量的多维数据,支持对审核内容做准确的实体识别。线上攻防基于面向大模型原生的防御体系,可以应对单次 50 万量级的饱和攻击,通过百万级高质量题库识别各类攻击手段,降低拒答率。

(二) 大模型安全技术实现

"赣服通"政务 AI 助手在端侧实施的安全措施取得了显著成果。首先,数据加密技术的广泛应用确保了训练数据在存储和传输过程中的安全性,有效降低了数据泄露风险。其次,鲁棒性测试与模型优化提升了算法模型的稳定性和准确性,使其在面对潜在攻击时更具抵抗力。同时,系统平台层面实施的身份认证、访问控制以及实时监控与异常检测机制,为政务大模型提供了坚实的防护屏障,确保只有授权用户才能访问系统,并实时监测潜在威胁。在业务应用层面,通过数据脱敏、保护及业务逻辑的安全验证,保障了业务数据的准确性和业务逻辑的合规性。这些安全措施的实施,不仅提升了政务大模型自身的安全性,也为政府机构的日常工作提供了可靠保障,促进了政务服务的安全与高效运行。

三、大模型安全措施成效

"赣服通"会同网信、公安等部门建立了安全管理机制,运用国产密码加密技术、区块链技术等强化保障安全,并加强政务数据运行监控,全面提升网络安全防护水平。深度融合了安全护栏和安全防御能力,在用户提问理解、提问风险决策、回答风险决策构建了全面的防御体系。助力江西省政府打造集"咨询、搜索、预约、提醒、评价"

于一体的 AI 数字人智能客服"小赣事",为办事用户提供一个更加便捷、安全的使用环境。

4.4 人力资源领域大模型安全实践

一、 案例介绍

AIGC 灵活用工直招平台创新性涵盖了智能化招聘、精准人才匹配、个性化推荐与培训、高效数据处理与分析以及风险防控与合规性等方面,平台核心功能包括:智能化招聘流程:通过大模型进行简历和面试筛选,为企业提供精准的候选人亮点与疑点分析。高效数据处理与分析:处理和分析大量的人力资源数据,形成人力解决方案。风险防控与合规性:人才招聘、培训和管理合规风险智能监测和预警。



图 4-5 人力资源领域大模型安全实践案例

二、大模型安全实践案例描述

- (一) 大模型安全在各环节的落实措施
- 1、训练数据安全

数据清洗和预处理: AIGC 灵活用工直招平台对收集到的和人力

资源相关的数据使用特定的技术,如去除 HTML 标签、噪声、重复数据,以及过滤掉包含敏感信息或不适宜的内容的数据,来清洗和预处理数据。

隐私保护:由于人力资源行业涉及大量隐私,因此在训练数据过程中,需要对敏感数据进行脱敏处理,例如替换或删除个人信息,或者使用匿名化技术,如哈希函数或差分隐私技术来保护个人信息。

数据均衡处理:为了降低数据收集的局限性对模型带来的影响, 如数据偏向于特定群体或地区,无法反映整个人力资源市场的状况, 通过综合使用数据平衡策略、数据采样技术、数据增强技术、序列标 注和语义建模等技术手段,确保各类数据比例合理,提高人力资源模 型对不同数据的适应性。

数据质量评估和审核: 收集到的数据会存在一些错误、恶意数据, 如薪资水平异常高/低,通过可视化工具或数据统计与分析,可以识 别数据中的异常值,提高数据质量。

2、模型训练安全

数据加密。结合人力资源不同应用场景下的不同数据规模,不同业务需求,选择性综合使用同态加密、对称加密、差分隐私、哈希算法等技术,对模型训练中的数据加密,确保数据在传输过程中不被窃取,在存储时也不被非法访问。

数据存储。由于人力资源数据涉及隐私,为了防止数据泄露、丢失,可以通过阿里云存储服务存储数据,使数据得到更好的保护。

建立防火墙和入侵检测系统。通过部署入侵检测系统(IDS)、入侵防御系统(IPS),使用双因素认证等技术,或使用虚拟专用网络(VPN)或专用网络(如内部网络)等,来隔离敏感资源,监控异常行为,防止个人信息泄露,训练数据丢失。

3、模型部署与使用安全

模型加密。在模型部署到生产环境前,对模型进行加密处理,确保模型在传输和存储过程中的安全性;通过模型水印技术,对模型进行唯一标识,防止模型被非法复制或篡改。

访问控制。设立严格的访问权限管理机制,确保只有经过授权的人员才能访问和使用模型;监控并记录模型的访问和使用情况,及时发现并处理异常行为。

4、模型维护与更新安全

模型验证。在对模型进行更新或维护前,进行充分的验证和测试,确保更新或维护后的模型性能和安全性不受影响,并定期对模型进行性能评估和安全评估,确保模型的持续稳定和安全。

审核机制。建立模型更新和维护的审核机制,确保所有更新和维护操作都经过严格的审查和批准,并对模型的更新和维护过程进行记录和监控,以便在出现问题时进行溯源和追责。

(二) 大模型安全技术实现

1、云

身份验证。通过用户名和密码、生物识别(如指纹识别、面部识别)、数字证书、多因素身份验证 (MFA) 等手段验证用户身份,确保只有合法的用户才能访问人力资源云端系统。

访问控制。通过配置 ACL 或 RBAC, 将访问权限与对象或角色关联, 精细地控制用户对人力资源云端数据中心资源的访问, 防止未经授权的人员访问和篡改数据。

安全认证。确保 HRSaaS 平台能通过国际权威安全认证, 防止个人信息泄露, 云端传输数据丢失, 未经授权访问和其他信息安全威胁。

数据加密传输。通过 SSL/TLS 协议、建立 VPN、运用对称加密技



术、哈希算法等技术,对传送到云端的数据进行加密,确保数据的机密性与完整性。

2、边

数据加密。通过数据脱敏、对称加密技术、哈希算法、同态加密等技术手段对个人信息、薪资等在边缘侧传输和存储的数据敏感信息进行加密,确保数据机密性。

安全更新和补丁管理。由于边端设备会直接处理员工信息、考勤记录、薪资变动等数据,因而要及时更新。通过建立自动化更新机制,可以定期检查边端设备的操作系统、应用程序和安全补丁;通过建立中央化补丁管理系统,可以集中管理边端设备的补丁需求,确保它们得到及时一致的更新。

安全审计和日志记录。启用边缘设备的安全审计和日志记录功能,记录所有重要事件和操作,通过定期检查和分析日志,及时发现潜在的安全问题和威胁。

3、端

用户培训与教育。明确培训目标,强调终端安全的重要性;制定 具体的培训内容,提高员工对安全危险的识别,教授正确的安全操作 流程;定期复习、测试,确保员工掌握和应用安全知识。

设备选型与认证。选择经过安全认证和测试的终端设备,如具有 安全启动、硬件加密和固件保护的设备。

备份与恢复。通过数据库备份、云备份等技术定期备份终端设备 上的重要数据,并确保备份数据的安全性。

三、大模型安全措施成效

AIGC 智能直招平台通过实施一系列大模型安全措施,有效提升 了平台的隐私保护能力、算法透明性和产品可控性,降低了机器幻觉 带来的误导,有效提升平台的安全性、认可度和信赖度。

4.5 智能助理领域大模型安全实践

一、案例介绍

个人助手大模型产品依托于"SenseNova"大模型的通用能力, 主要功能包括:语音助手:问答、闲聊、写作、商品助手;输入助手: 根据对话上下文内容生成回复、回复风格化;文档助手:文档润色、 文案生成、文档问答等多方面功能。产品以大语言模型为基础捕捉用 户需求,支持长上下文沟通,实现强个性化特色的需求捕捉,最终协 助用户实现调用设备的对应功能的目的。



图 4-6 智能助理领域大模型安全实践案例

二、大模型安全实践案例描述

(一) 大模型安全在各环节的落实措施

1、训练数据安全

底层数据方面,对于不同来源的原始语料数据,采用不同过滤规则,结合人工抽检,确保语料来源合法可靠。建立敏感词安全规则+ 语义分类器的组合机制对训练数据进行过滤,最终对全部历史语料进



行二次清洗处理,确保数据安全。

2、算法模型安全

通过运用知识库挂载技术确保大语言模型输出内容安全、准确、专业。在典型的基于知识库的问答(QA)场景中,用户需要向大型语言模型(LLM)查询特定知识库,例如 PDF、Word 文件等的相关内容,然后文本长度对本地知识库超长文本的内容理解形成了障碍。通过 Embedding 模型对本地知识库进行一次提炼,提高整个问答系统的效率。Embedding 模型是一个精准的筛选器,帮助提取出知识库中对当前用户问题来说最重要的内容,以满足用户不同的提问需求,同时有效地减轻了 LLM 处理长文本的负担,有效缓解大模型幻觉问题,并面向不同用户提供一致性答案。

3、系统平台安全

构建内容安全运营平台,包括应急响应平台、策略运营平台、人 审标注平台、风险数据平台等。设计内容安全审核链路,包括机器审 核、人工审核、投诉举报、用户治理等核心功能,实际敏感词条拦截 可达百万级。同时具备健全的安全响应机制,第一时间快速处置内外 部风险,结合分级报告、处理突发安全事件。同步筹备涉政专班人员, 方便与监管部门对接。

4、业务应用安全

具备健全的安全响应机制,第一时间快速处置内外部风险,结合分级报告、处理突发安全事件,后续筹备涉政专班人员,方便与监管部门对接。

(二) 大模型安全技术实现

建立 AI 场景下数据安全整套解决方案,保障 AI 数据隐私合规。 基于客户需求和服务器情况,对模型进行公有云或私有云部署,适配 相应环境,并对稳定性、安全性进行维护。

在云端应用 Embedding 模型相关技术时,面向国家监管要求,提供安全回复审核能力模型和人工黑名单库,保证大模型回复安全的通过接口调用模型能力。

同步推出端云协同的综合方案,会在建立端侧本地隐私知识库的基础上进行端侧推理,为用户隐私数据保驾护航。具体表现在通过文档上传的方式构建用户专属知识库,实现即问即答响应快的特性;以及本地文档手机离线也可进行总结摘要与问答,无需担心机密信息联网泄露。另外在本地安全模块的基础上实现意图分流,所有涉及隐私的数据及 query 等,将全部基于端侧模型处理,实现终端的隐私保护。

三、大模型安全措施成效

个人助手大模型产品推出以 AI 安全为核心的大模型安全保障体系,在底层数据、模型技术、配套机制、业务应用等层面实行完善的解决方案。进行全面完整的知识划分,有效缓解大模型幻觉问题;建立安全测试集验证流程有效识别各个风险维度,提升算法模型的可解释性和可控性。大模型安全策略遵循了人类价值观,契合用户意图、形成可信可靠的大模型应用产品。

五、大模型安全未来展望与治理建议

5.1 未来展望

大模型发展空间巨大,平衡机遇与安全风险挑战成为大模型发展的关键保障。大模型的优异能力表现给产业带来太多惊喜,在强大算力、海量数据支撑下,大模型将会在越来越多的领域超越人类,通用人工智能成为可能。与此同时,也将面临着大模型浪潮给社会



安全带来的巨大冲击,如研发过程中引发信息泄露、价值对齐、机器幻觉等问题,以及落地过程中面临的数据、模型、算法及其运行的软硬件环境安全风险。人类尚未准备好如何绝对安全"驾驭"大模型更好的为人类服务大模型的风险挑战比以往任何时候都严峻。大模型的安全问题引起了全球范围内的广泛关注,它不仅关系到技术本身的稳健性,还涉及伦理、法律、社会等多个层面。大模型安全成为学术界、产业界、政府部门共同关注的议题,社会各界正在寻求平衡创新与风险管理的有效策略,以促进人工智能技术的可持续和负责任的发展。随着人工智能技术的不断进步,确保大模型的安全、可信、可靠、可控,已成为推动科技创新、产业健康发展,维护社会秩序和保障个人权益的重要任务,需要全人类社会的共同关注和努力。

大模型风险因素繁多,系统化构建大模型安全体系屏障成必然。 面对动辄百亿、千亿级参数规模的大模型,其超大参数体量以及计算 复杂性致使大模型技术比以往任何人工智能技术都复杂许多,涉及算 力、网络、数据、模型结构、训练方法、产品化落地等众多环节,每 一个环节都需要处理兆级的海量数据。且作为数据驱动型技术,其计 算过程更像一个"黑箱操作",大模型安全性、可靠性、可控性等挑 战空前巨大。在大模型研发、训练、部署、应用的任何技术的环节都 可能带来风险,风险因素繁多,任何环节的风险都可能带来整个系统 的安全挑战,原有的安全体系已经难以适应新情况。

大模型安全需要构建一个全面、协调、系统的安全管理框架。大模型系统的安全环环相扣,涉及数据保护、可解释性、鲁棒性、伦理责任、合法合规等方面。大模型安全是一个统一的、协调的复杂系统,

需要使用系统化思维去洞察、捕捉和评估大模型系统中可能的隐患及其影响,如数据的安全与合规、模型结果测试与验证、监控与审计等。 大模型安全问题不容小觑,也不能仅靠每个环节独立思考,应该体系化一盘棋考虑,从而构建完备的大模型安全屏障。

大模型标准需求迫切,全面支撑大模型安全测试验证能力建设。 当前,大模型技术迅猛发展、行业数据迅速汇集、创新应用不断深入, 产业发展脚步加快。大模型企业正在各自为政奋力开展技术研发,积 极尝试新的商业模式,在这样的背景下,行业共识尚未形成,大模型 标准化建设尚未跟上技术和产业发展步伐,标准体系建设需求迫切。 大模型安全急需总结先进经验、统一行业共识,以系统科学的理论和 方法为基础,运用标准化的工作原理,不断优化标准内容,构建大模 型安全指标体系和测试验证标准体系。

测试验证是大模型安全的有效手段。从传统人工智能的判决式到大模型的生成式的转变,大模型测试自动化程度低,基于大模型安全标准打造高效一致的测试验证体系将会变得越来越重要,包括构建大模型安全测评能力和大模型安全自动化验证工具,涉及测试方法、测试对象、测试任务、测试指标、数据集、工具平台等。随着大模型技术的不断发展,新的安全威胁和挑战也会不断出现,安全标准和测试验证能力也需要与时俱进,以适应新的技术发展和安全需求。

大模型安全基础设施尚不完善,构建大模型基础设施迫在眉睫。 目前,国外对大模型软硬件基础设施已形成技术壁垒,大模型训练框架、部署框架等各类软件基础设施,以及芯片、处理器、服务器等各类硬件基础设施在很大程度上依赖国外,如 TensorFlow、PyTorch、NVIDIA GPU等。国内自主安全可控的大模型软硬件基础设施正在实现产业突破,但是总体来说我国大模型安全基础设施尚不完善,完全

自主的大模型安全软硬件环境生态尚未形成。

构建自研软硬件适配的大模型基础设施体系的需求急迫。搭建软硬件适配测试平台,可对各类软硬件进行兼容性测试和性能评估,确保其在大模型基础设施中的稳定运行;支持 GPU、TPU等 AI 芯片研发,实现国产 AI 芯片的自给自足;实现高速互联网络和分布式计算网络基础设施的持续升级,支撑超大规模算力资源的共享和高效利用;推动智算中心的持续升级转型,实现智能计算与存储一体化;研究"端、边、云"协同过程中的数据安全技术,确保数据安全。

大模型安全实践经验尚浅,标杆场景为大模型安全实践指明方向。 大模型安全应用是一个新兴领域,研究和应用尚处于起步阶段,因此 缺乏成熟的参考案例来指导实践。大模型企业仍在探索如何有效地确 保大模型的安全性,在原有的传统数据安全、信息安全、系统安全等 经验基础上,进行能力迁移,应用于大模型安全。大模型安全问题变 得更加多样化和难以预测,原有的安全技术也需要不断调整和优化, 并通过企业不断探索和实践,逐步积累经验,建立一套成熟的安全实 践体系,这对于指导未来的大模型安全实践至关重要。

安全不仅仅是一个技术问题,更是一个战略问题,大模型企业将采用更为先进的安全技术和工具,如数据加密、联邦学习、访问控制、异常检测、测试验证等,来增强大模型系统的安全性。同时,通过建立专门的安全团队,与研发、产品、运营等部门紧密合作,共同研制、实施安全策略、构建安全技术屏障、防范未来安全风险。随着新的大模型安全实践的不断深入,将涌现出大量成功案例,应用标杆的集体爆发将为大模型安全构建实践范式,打造高价值的参考体系。

大模型安全"以人为本"是核心,坚持发展负责任的 AI 是大模型安全立足点。技术的发展始终是以拓展人类的能力、服务于人为目

的的,大模型安全以"以人为本"作为核心,才能保证其发展方向不偏离,确保技术的发展既符合伦理道德,又能够为人类社会带来积极的影响。"以人为本"的大模型安全理念强调的是在大模型的技术和应用过程中,始终将人的利益、需求和安全放在首位,大模型的设计者、开发者和使用者都必须始终保持"以人为本"思维,需要切实保障用户和社会的安全与利益。

任何背离"以人为本"核心发展的大模型技术和应用最终都将导致安全风险和挑战,当大模型的发展忽视了人的需求、权利和福祉时,就可能造成不可预测的后果,例如侵犯个人隐私、导致社会不公平、伦理道德冲突等问题。发展"以人为本"的人工智能就是要发展负责任的人工智能,即对人类负责、对社会负责。坚持发展负责人的 AI 是大模型安全立足点,据此构建健康、可靠的大模型安全生态,以确保大模型的安全和效益能够与社会的发展同步,为人类带来真正的福祉,是未来人工智能可持续发展的基本保障。

5.2 治理建议

构建集大模型安全政府监管、大模型安全生态培育、大模型安全 企业自律、大模型安全人才培养、大模型安全测试验证"五维一体" 多元参与、协同共治的治理框架。



来源:中国信息通信研究院

图 5-1 大模型安全"五维一体"治理框架

一大模型安全政府监管。一是加强大模型合规体系建设,以高标准、严要求为准则,确保大模型发展与社会责任并行不悖,这也是确保大模型安全性的重要措施,政府应加强大模型相关法律法规和政策的宣贯工作。二是为企业提供大模型安全公共服务,包括信息公开、业务指导、制定伦理规范操作指南、风险管理和合规指引、搭建大模型公共服务平台等,为企业和个人提供高效的信息渠道和行为引导。三是建立监管和应急响应机制,设立专门的部门监督、管理大模型研发与应用,定期对大模型的安全性进行监督检查,并建立应急响应措施,以便在出现大模型安全问题时能够迅速采取措施,防控风险蔓延,确保大模型安全。

——大模型安全生态培育。一是建立全面的安全教育和宣传体系,提高公众对大模型安全的认识和理解,通过教育和宣传提升全民的数字素养。大模型的应用已深入人们生活和工作的方方面面,培养使用者对于大模型风险的识别和防范能力,提升全民人工智能安全意识和

素养,是构造大模型安全生态的基础。二是构建大模型安全社会监督体系,构建一个全民参与的大模型安全体系,完善反馈机制,建立公开透明的信息共享平台,让监管机构、企业和公众能够及时反馈和了解大模型安全的最新动态。

一大模型安全企业自律。一是构建企业大模型合规体系,制定相关安全制度和准则。企业应制定严格的内部安全政策和操作规程,确保大模型的研发、部署和运维等过程符合行业安全标准和法律法规要求。二是加强企业安全培训与安全意识。作为生成式人工智能服务提供者的企业主体应确保其产品安全,强化企业安全环境,定期对员工进行大模型安全相关的培训,强化安全意识,确保每位员工都能理解并遵守安全操作规范。三是建立安全监督机制,设立专门的安全监督岗位,并规定其职责,使企业能及时发现并处理潜在的安全问题,并通过建立产品应用跟踪机制,对风险产品及时召回,不断优化安全策略,降低大模型产品在企业内部和外部的风险。

一一大模型安全人才培养。一是加强大模型安全人才队伍建设,构建大模型安全学习体系,构建跨学科学习,将人工智能、网络安全、数据科学等领域的知识进行整合,以培养具备综合能力的人才队伍。二是鼓励核心技术攻关,比如如何攻克大模型数据安全、模型安全、系统安全和应用安全等技术难题,如何构建内生安全、外生安全、衍生安全的防御体系,并不断跟进前沿技术发展,以应对新的挑战。三是推动高校、研究机构与企业的紧密合作,实现资源共享,促进大模型安全领域人才培养与市场需求的对接。

——大模型安全测试验证。一是推进大模型安全标准研制,研究和借鉴国际上已有或在研的大模型安全标准和最佳实践,加强大模型安全测试验证技术研发和标准化工作,加速标准应用转化。二是加速



大模型安全测试验证能力建设。开发和标准化一系列大模型测试验证工具,构建标准化测试数据集,开展大模型安全测试验证示范场景,快速构建体系化大模型安全测试验证能力。三是鼓励第三方机构开展大模型安全测试评估业务。通过政策引导、资金支持和市场激励,鼓励有资质的第三方机构提供专业、独立的大模型安全测试验证服务,帮助企业及时发现和解决大模型的安全隐患,提升整个行业的安全水平,从而推动大模型产业健康发展。

