附件

"铸盾模都" 2025 年人工智能安全赋能专项行动任务清单

序号	2025 年重点任务			企业提报材料
		1.加强人工智能安全	企业应明确网络和数据安全负责人,压实主体责任;涉及关键基础设施运营或重要数据	
_	人工智能 安全主体	管理组织建设	处理的,还应设立专门的安全管理机构,强 化企业内部监督管理及资源保障。	_
	责任	2.健全人工智能安全	企业应建立健全人工智能网络和数据安全管	
		管理制度体系	理制度,强化安全风险防范,形成全链条闭 环的安全管理运营体系。	_
=	人工智能 安全保障	3.加强人工智能网络安全	企业应结合不同阶段的安全风险,规范开展 资产安全管理、漏洞安全管理、模型更新安	_

14 AL		A 64 - 11	
基线		全等工作,不断健全安全防护措施、持续完	
		善安全能力建设,落实通信网络安全防护定	
		级备案相关工作。	
		企业应规范预训练和优化训练数据及其处理	
	4.加强人工智能数据	活动的安全要求,在通用数据安全基础上,	
	安全	做好训练数据安全防护、安全标注以及重要	_
		数据识别等工作。	
		企业应规范深度合成、生成式人工智能等应	
	5.加强人工智能内容	用,规范内容生成模板、规则配置,健全信	
	安全	息发布全流程审核机制,提升涉生成违法违	_
		规内容识别与拦截能力。	
		1.企业应结合业务实际梳理人工智能业务清	太北王11月20日
	6.加强人工智能业务	单并上报市通管局,同时在人工智能业务上	企业于11月30日
	安全	线前开展安全评估。	前提交年度人工
		2.涉及到自主研发大模型以及部署开源大模	智能业务清单。

	1			
			型并对外提供服务的企业,应自行或委托第	
			三方机构对使用大模型的业务开展安全评	
			测。	
			3. 市通管局定期组织对人工智能业务的安全	
			评估评测情况进行抽测。	
			企业应建立网络与数据安全应急工作组织与	
		7.强化人工智能应急	响应机制,制定有效可操作的安全事件应急	
		管理与能力	预案, 每年开展一次应急演练并针对发现问	_
			题开展闭环整改。	
				鼓励企业主动提
	1 十年化	8.提供人工智能安全	市通管局组织提供可复制、易推广的评测工	供可开源的人工
_	人工智能	评测工具	具和开源代码。	智能安全工具、测
三	安全管理			试数据集等。
	赋能	9.建立人工智能安全	1.鼓励有关企业和机构主动开展安全风险监	1.鼓励有关企业
		风险监测巡查和处	测上报。	和机构主动开展

		I		
		置机制	2.市通管局对人工智能技术和应用开展常态	安全风险监测上
			化安全风险监测巡查,并综合巡查结果与企	报。
			业上报情况定期发布人工智能安全风险预	2.企业应提供接
			警,开展威胁通报。	受风险监测巡查
			3.企业应开展内部常态化安全巡检,针对发现	所必要的技术支
			和通报的安全隐患及时开展整改处置。	持,包括但不限于
				开放智能应用的
				API 接口、测试用
				账号等。
			1.市通管局组织开展"磐石行动"人工智能安	
四	人工智能 安全行业 发展	10.提供安全合规及	全(专项)攻防演练活动。	
		应用实践指导	2.市通管局基于演练成果系统梳理优秀防御	_
			实践。	
		11.加快培育人工智	1.市通管局面向重点企业开设人工智能安全	鼓励企业上报在
		能安全生态	风险防范与治理能力提升专题培训。	人工智能内生安

	2.市通管局组织开展"铸盾模都"2025年人	全和人工智能赋
	工智能安全赋能典型案例评选, 支持和培育	能安全等领域中
	在人工智能内生安全和人工智能赋能安全等	的典型案例。
	领域有突出成效的企业。	