

中国信息通信研究院

2026年1月

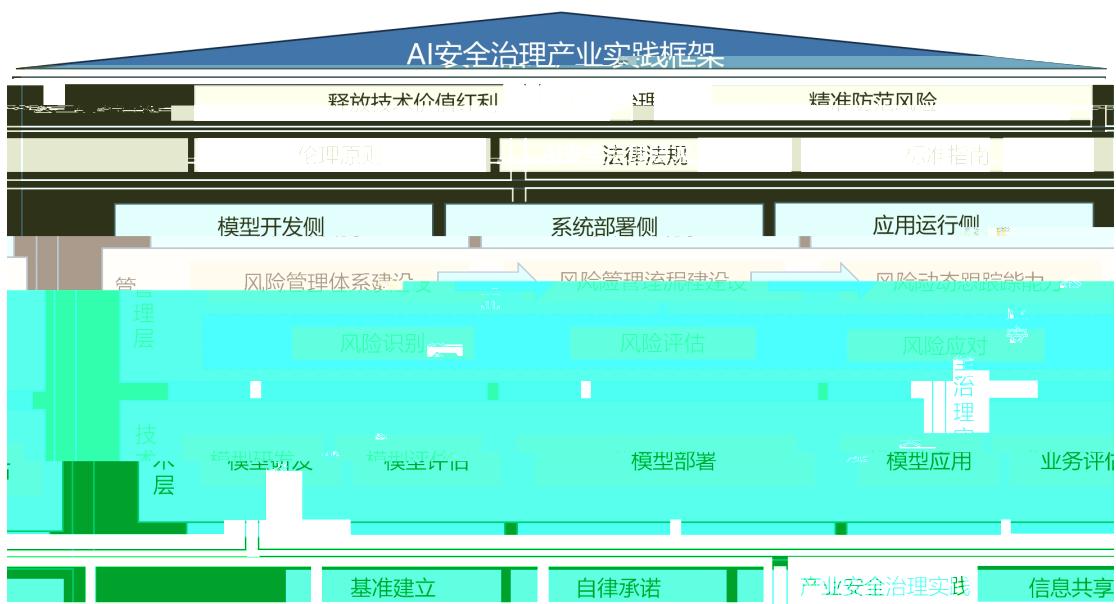
版权声明



（一）宏观规范要求迈向务实新阶段

（二）技术应用发展凸显风险新态势

（三）产业实践需要安全治理新框架



（一）国际合作层面，全球加强安全治理深化交流

（二）监管政策层面，各国推动安全治理体系化落地

（三）产业实践层面，构建安全治理务实举措成为核心 议题

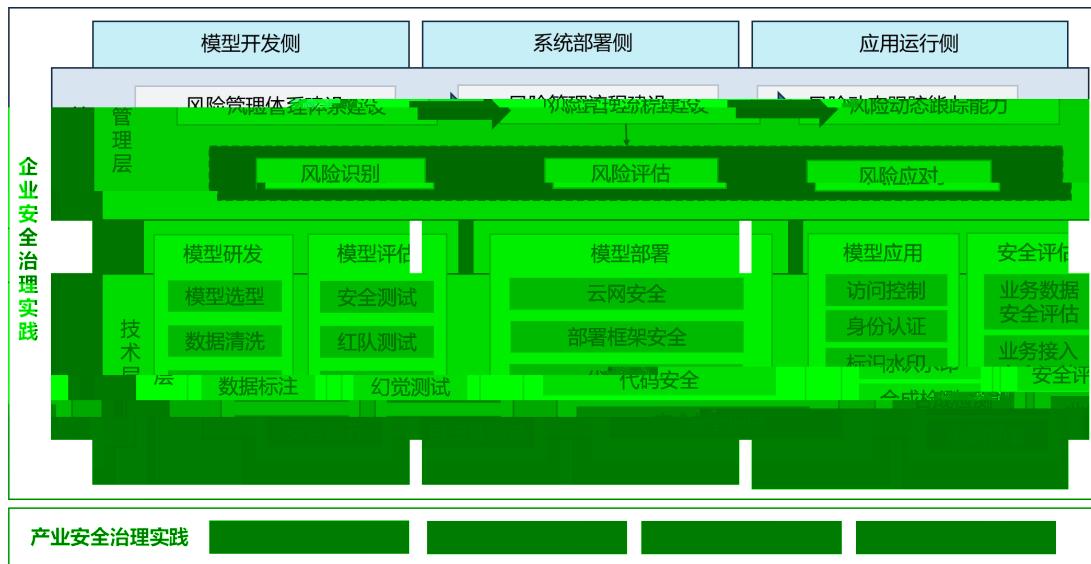
（一）技术发展扩大内生安全新敞口

模型在数学推理、代码生成等创造性任务上的突破，往往伴随着事实准确性的下降，形成能力跃升与幻觉率攀升的“双刃剑”效应。从技术原理来看，当前大语言模型的创造力本质上是基于海量训练数据的模式挖掘与概率性生成，而非真正意义上的理解与思考，模型幻觉问题难以根除。

（二）应用延展引发衍生安全新难题

（三）组织管理体系构建面临新卡点

（四）多元共治协同机制尚待健全完善



(一) 风险管理：构建闭环的人工智能风险管理体系建设

（二）模型研发：筑牢人工智能开发安全源头根基

（三）系统部署：构建人工智能部署安全防护屏障

（四）应用运行：强化人工智能应用安全动态评估

（五）产业生态：共建基准测试体系与协同治理机制



《人工智能安全承诺》实践披露

Disclosure of Practices on the Artificial Intelligence Security and Safety Commitments

中国人工智能产业发展联盟
Artificial Intelligence Industry Alliance

中文/En

22家
签署企业

16家
披露企业

2024年12月，中国人工智能产业发展联盟（AIIA）发布《人工智能安全承诺》（简称《承诺》）。首批17家领军企业签署，展现“守护好文宝、促进普惠向善”的庄严承诺。

2025年2月，AIIA在巴黎人工智能行动峰会上宣介《承诺》，获国际积极反响。同时，AIIA发起自律披露行动，倡议签署企业披露安全实践。

2025年7月，AIIA公布《承诺》实践披露名单，涵盖企业16家，涉及领域包括AI安全、可信、可控、可持续、可解释等。

关于《人工智能安全承诺》

参与《承诺》企业名单（按单位首字母排序）

签署企业 (17家)	披露企业 (16家)								
AIIA Group	Baidu	火山引擎	万向	蚂蚁集团	阿里云	腾讯云	字节跳动	京东集团	滴滴出行
HONOR	小米	MINIMAX	SANGFOR	Tencent	vivo	王者荣耀	王者荣耀	ZTE	

（一）金融行业打造人工智能风控方案

参见《中国互联网金融协会金融科技发展与研究专委会、中国信息通信研究院人工智能研究所：《金融大模型合规应用研究报告》。

（二）医疗行业多维度提升模型输出准确性

（三）交通行业结合应用场景探索安全方案

自动驾驶

是人工智能在交通领域中，最大的智能原生应用场景之一。为

◦

(四) 能源行业加强基础设施的安全防护

（五）通信行业加强“以技治技”的安全治理体系

