

启明星辰 | 大模型应用安全 | 深度应用安全基座系列白皮书

AI就绪的大模型身份与访问管理

AI-R-IAM v1.0

AI-Ready Identity and Access Management

启明星辰信息技术集团股份有限公司
2025年2月

版权申明

北京启明星辰信息安全技术有限公司保留对本文档及本文档所包含的所有

和修改权。

本文档启明星辰保留所有权利。未经许可，任何人不得复制或传播。

北京启明星辰信息安全技术有限公司。未经北京

启明星辰信息安全技术有限公司

启明星辰信息安全技术有限公司。任何人不得以任何形式或形式对本手册内任何

启明星辰信息安全技术有限公司

备份、修改、传播、翻译或进行其他任何商业用途。

部分进行复制或传播。

免责声明

本文档依据现有信息制作，其内容如有更改，恕不另行通知。

北京启明星辰信息安全技术有限公司保留对本文档及本文档所包含的所有

启明星辰信息安全技术有限公司

技术有限公司不对本文档中的遗漏、不准确、或错误导致

确可靠，但北京启明星辰信息安全技术有限公司

的损失和损害承担责任。

信息反馈

如有任何宝贵意见，请反馈：

北京启明星辰信息安全技术有限公司

信箱：北京市海淀区启明星辰

100193 电话：010-82779088

传真：010-82779000

您可以访问启明星辰网站：www.venustech.com.cn 获得最新技术和产品信息。

目录

1 公司简介.....	4
2 大模型应用安全挑战.....	6
2.1 大模型应用安全挑战.....	6
2.2 大模型应用安全挑战.....	7
2.2.1 身份与访问控制的安全挑战.....	7
2.2.2 数据安全挑战.....	9
3 身份与访问管理系统 (IAM)	11
4 AI 就绪的大模型身份与访问管理.....	13
4.1 可信身份管理.....	15
4.2 可信权限管理.....	16
4.3 可信数据管理.....	17
4.4 可信行为审计.....	19
5 应用场景.....	21
5.1 大模型应用场景数据流.....	21
5.2 场景一：用户访问大模型应用场景.....	21
5.3 场景二：应用访问大模型应用场景.....	26
5.4 场景三：用户-AI 智能体的安全应用管理场景.....	28
6 发展趋势展望.....	31

启明星辰

QIMEE

由严望佳博士创建，是国内最具实力的、拥有

启明星辰公司成立于 1996 年，由

信安全管理平台、安全服务与解决方案的综合提供商。

完全自主知识产权的网络安全产品。可

2010年6月23日，启明星辰在深交所中小板正式挂牌上市。

威胁管理、加密认证等技术领域，共有十余个产品线，并根据客户需求不断增加。启明星辰解决方案为客户提供安全需求、产品安全产品、服务之保障体系，构建可信的安全保障体系。与安全技术紧密衔接，帮助构建完善的安全保障体系。

自 2002 年上市以来，启明星辰持续保持国内入侵检测、漏洞扫描市场占有率第一。近年来，公司发展成为国内威胁管理、安全管理平台市场占有率第一，安全审计、安全合规市场领导者。目前，公司在全国各省市自治区设立三十多家分支机构，拥有覆盖全国的营销和服务体系。

长期以来，启明星辰公司受到了党和政府领导人的关怀与鼓励。2000 年 1 月，江泽民总书记、曾庆红总书记和中央军委领导人亲切视察启明星辰公司；2003 年 1 月，胡锦涛总书记

亲切接见了启明星辰公司 CEO 严望佳博士。

凭借多年来的潜心研发，启明星辰获得国家规划布局内重点软件企业，国家火炬计划软件产业优秀企业，中国电子政务 IT100 强等荣誉，及拥有最高级别的涉及国家秘密的计算机信息系统集成资质证书。

启明星辰目前是我国规模最大的国家级网络安全研究基地。完成包括国家发改委产业化示范工程，国家科技部 863 计划、国家科技支撑计划等国家级科研项目近百项。创造了百

国家及行业网络安全标准，填补了我国信息安全科研领域

余项专利和软件著作权，参与制订

的多项空白。

，启明星辰以用户需求为根本动力，研究开发了完善的专

作为信息安全行业的领军企业

经成为在政府 电信 金融 能源 交通 军队 军工

业安全产品线，通过不断耕耘，已

制造等国内高端企业级客户的首选品牌：启明星辰在政府和军队拥有 95% 的市场占有率，为世界五百强中 80% 的中国企业客户提供安全产品及服务；在金融领域，启明星辰对政策性银行、国有控股商业银行、全国性股份制商业银行实现 90% 的覆盖率。在电信领域，启明星辰为中国移动、中国电信、中国联通三大运营商提供安全产品、安全服务和解决方案。

作为北京奥组委独家中标的核心信息安全产品、服务及解决方案提供商，奥组委唯一信息安全供应商，启明星辰受到独家官方授权，全面负责奥运会主体网络系统的安全保障，得到了国家主管部门的大力嘉奖。此外，启明星辰还为上海世博会、广州亚运会等多项世界级

大型活动提供全方位信息安全保障。

前，已累计资

在公司快速稳定发展的同时，启明星辰公司坚持以爱心回馈社会，截止目前

了 5 所希望

助贫困学子、受灾、贫困群众上亿元人民币，并在江西、青海、新疆等地援建

小学。

自主创新的安

启明星辰公司将秉承诚信和创新精神，继续致力于提供具有国际竞争力的自

形，为打造和

全产品和最佳实践服务，帮助客户全面提升其 IT 基础设施的安全性和生产效能

提升国际化的民族信息安全产业第一品牌而不懈努力。

模型不断迭代，众多初创企业也在探索在大模型在细分领域的落地应用。

- 大规模应用期 (2025 年 1 月-)

- (1) 国外：大模型深入医疗、金融、教育等多个行业，助力疾病诊断、风险评估、个性化治疗等。OpenAI 还在探索新应用领域。

- (2) 国内：DeepSeek 凭借创新的算法和架构，展现出低成本、高效率以及开源优势，

应用，推动各行业数字化、智能化转型。

等行业大规模

技术作为新一代生产力的代表，正以前所未有的速度和深度渗透到企业运营和个

大模型技

方面，成为推动社会进步和经济发展的关键力量。它不仅重塑了传统行业的运

人生活的方

产生了全新的应用场景和商业模式，真正实现了“智能赋能”的愿景。

作模式，还

2.2 大模型应用中安全挑战

2.2.1 身份与访问控制的安全挑战

在大模型深度应用的时代背景下，身份安全面临着前所未有的挑战。随着人与设备、人

与应用、AI Agent 与数据、人与 AI Agent 之间交互量呈指数级增长。

根据埃森哲的《2025 年技术愿景》调查，78% 的高管同意在未来必须为 AI Agent 构建

健壮且平等的数字身份系统，这预示着传统身份概念不断延伸至 AI 领域。在大模型

身份划分，即人类身份 (AI Identity) 以及非人类身份 (Non-Human

应用中产生了三类

安全问题。

Identity, NHI)，都暴露出诸多

大模型应用交互场景中的背

1、人类身份安全问题：人类身份代表自然人身份实体，在

应用中，容易出现滥用身份非法访问大模型系统数据

等问题，易出现泄密和洪震问题，大模型对用户数据访问不透明也可能泄密隐私

2、AI Identity 安全问题：越来越多的 AI Agent 在工业物理场景中应用，AI Agent 具

《人工智能治理身份》中指出 AI Agent 具有复杂身份关系，需要身份相互身份管理。

AI Agent 应视为具有身份的不同且可信任载体，应明确思考 AI Agent 作为员工，在大

性，会出现如下风险：

(1) 访问越权：员工访问管理相对 AI Agent

非法访问权限。

(2) 访问权限放

如果权限管理机制不完

和操作敏感资源，造成

(3) 自身安全：A

3、非人类身份 (NHI) 安全问题：非人类身份定义为企业技术内的应用程序、服务或

2025 年 2 月发布《十大 NHI 风险 2025》，如下图所示。

OWASP 十大NHI风险2025

NHI1:2025 不当注销

NHI6:2025 不安全的云部署配置

NHI7:2025 长期有效

身份生命周期管

NHI 在大模型应用中，大多数存在于大模型基础设施交互过程中，NHI 的理薄弱，存在诸多安全风险与问题：

库等位置，容易

(1) 纯文本/未加密凭证：大模型开发应用时，很多 NHI 硬编码在代码库被内外部威胁者发现。

息不可见，许多

(2) 影子账号：大模型应用迭代快，加之生命周期流程弱、账号使用信 NHI 账号不活跃，增加攻击面。

权信息，而明确

(3) 缺乏账号所有权：多数参与大模型应用的组织中，NHI 无明确所有管理者对安全维护和问题补救很关键。

使用，或者相同

(4) 缺乏环境隔离：在许多情况下，相同 NHI 会在生产和非生产环境中逻辑 NHI 在每个环境中都有相同的密码，从而增加了横向移动的风险。

环等安全操作复

(5) 凭证共享：大模型应用的多程序间共享 NHI 违反原则 使密码循

难以掌握所有依赖关系。

杂，

2.2.2 数据安全挑战

大模型应用在数据安全方面呈现如下挑战：

(1) 用户提示词输入输出风险：用户提示词可能包含敏感信息，若大模型系统对输入提

示词的验证和过滤机制不完善，敏感信息可能在输入环节直接暴露。在输出结果时，若对输

出内容审查不足，可能泄露用户隐私、商业机密等敏感信息，如用户输入涉及个人健康状况

的提示词用于医疗大模型进行诊断辅助，输出结果未脱敏处理，可能导致个人健康隐私泄露。

(2) 大模型 API 调用风险：API 调用过程中，若身份认证和授权机制存在漏洞，不法分

子可能冒用合法身份调用 API，获取敏感数据或进行恶意操作。同时，API 接口若缺乏有效

安全防护，易遭受攻击，导致数据泄露或篡改。比如攻击者通过漏洞获取 API 调用权限，

非法获取金融大模型中的用户资产数据。

(3) RAG 知识库查询风险：RAG 知识库整合大量数据，若查询权限控制不当，用户可能

超出授权范围获取敏感知识数据，而知识库中的数据来源复杂，若数据经过整合甚至

可能被篡改或替换，导致输出结果的准确性和安全性难以保障，可能

(4) 数据泄露风险：在大模型训练阶段，攻击者通过逆向工程、恶意爬虫

等方式窃取训练数据，或在推理阶段通过模型窃取敏感信息，造成数据

泄露。训练数据中引入大量敏感数据时，模型在推理过程中可能泄露

(5) 合规风险：部分数据完全违反法律法规的不明

规范要求，若企业在数据收集、使用、存储等环节未

化原则，未获授权擅自采集等，可能面临法律诉讼和

3 身份与访问管理系统 (IAM)

身份与访问管理 (Identity and Access Management, IAM) 作为一种综合性的身份安全管理框架,旨在确保数字环境中用户身份的真实性、可靠性以及访问权限的精确控制。

它涵盖了身份认证、授权管理、身份生命周期管理等多个关键环节,通过一系列技术手段和管理策略,为企业和组织提供了安全、高效的身份管理解决方案。

启明星辰基于 IAM 的技术实践,针对企业数字化要求提出通过 IAM 构建“一体化全程

可信数字身份底座”理念。该理念旨在逻辑上建立一张基于身份的可信网络,赋予人员、设备、应用、接口、数据等实体唯一身份标识,实现实名制可信身份鉴权,基于可信身份对访问主体进行动态访问控制,推动业务与数据安全从“单点可控”迈向身份、环境、访问、数

据给原的“一体化可控”。

基于IAM构建的“一体化可控数字身份底座”总体架构如下图所示:



IAM作为体系构建的核心能力,向上支撑各类应用场景提供实名制可信身份,向下赋能安全访问设备,驱动安全控制执行,同时联动可信环境能力,及时调整访问策略。

设备 可信环境基于大数据安全能力，实现网络风险、行为风险、数据风险、业务风险、设备风险的全链路安全风险审计。

原 访问可信与数据可信作为一体两面，将防控范围扩展到安全设备、系统资源、应用资源

原 通过可信环境，实现“可信、可控、可管”的闭环管理，提升企业安全防护能力，降低安全风险。

实现更好的安全性，搭配一体化安全机制的融合，提升企业安全防护能力，降低安全风险。

原 通过可信环境，实现“可信、可控、可管”的闭环管理，提升企业安全防护能力，降低安全风险。

大模型作为新型生产力，在政府、企业、个人中应用越来越深入，需要针对其身份、

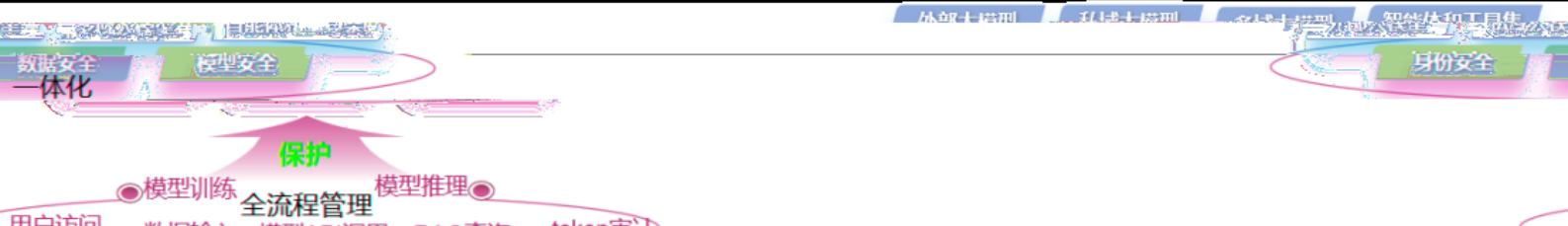
数据安全等方面需求，建立一套基于大模型技术特点量身定制的一体化身份

管理身份体系，即AI就绪的大模型身份与访问管理。

4 AI就绪的大模型身份与访问管理

人工智能模型(如大模型)安全体系基础支撑贯穿各个场景中,构建了保护体系。

AI就绪的大模型身份与访问管理(以下简称:AI-R-IAM)是为大模型(如DS、GPT、BERT等)的深度学习提供身份安全、访问控制、数据安全、模型安全系统,从身份、访问、行为、数据等维度为大模型在训练、部署和推理的构建一体化综合可信能力,构建大模型安全从“单点可控识别”一体化安全体系,同时为其安全能力提供身份管理支撑。



系统不仅支持传统认证能力，有效防

先进的身份认证机制，确保只有经过授权的用户和系统能够访问大模型。系统支持传统的用户名和密码验证，多因素认证，还引入基于中国移动号卡特性的实名认证能力，有效防止未经授权的访问和潜在的安全威胁。

集中的身份管理能力和统一的身份标识也为其它能力和系统提供身份支撑，如网络防火墙、SDP、安全管理平台等，为网络策略下发，保障网络安全提供支撑。

- 可信访问权限管理

AI-R-IAM 针对多个大模型访问、AI 接口访问、RAG 访问提供了集中的精细化的权限管理功能。系统可以根据大模型用户的角色和职责，动态调整其对大模型的访问权限，确保每个用户只能访问其所需的数据和功能，从而降低数据泄露和滥用的风险。这种精细化的权限管理不仅提高了系统的安全性，还增强了用户的操作体验。

- 可信数据管理

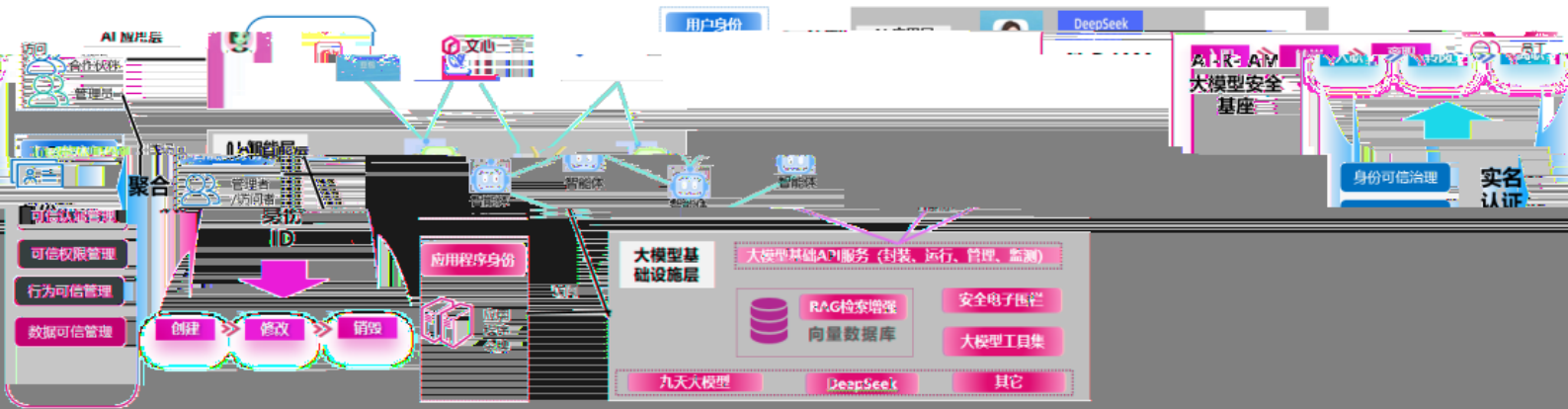
AI-R-IAM 对大模型访问、使用、调用、训练、推理等场景建立全流程的数据安全保护能力，基于身份、权限、脱敏等多方面能力，保障大模型在各个业务阶段数据的安全性。

- 可信审计管理

AI-R-IAM 对大模型使用过程和大模型内部组件业务调用过程进行安全审计和实时监控，记录用户的所有操作行为，包括用户登录时间、IP 地址、访问的大模型资源、执行的操作等信息。实时监测用户和组件的访问行为，及时发现异常风险。

AI-R-IAM 通过构建一体化的全程可信能力，能够有效应对当前大模型面临的安全挑战，为其它保护大模型的设备，提供身份、认证、权限、审计等方面的基础能力支撑，为未来人工智能技术的安全发展奠定了坚实的基础。

4.1 可信身份治理



AI-R-IAM 支撑基于大模型特性的多维度身份和属性的治理和管理，AI-R-IAM 的“身份 ID”超越了传统身份概念的边界，聚合人类身份、AI Identity 以及非人类身份（NHI），针对用户、智能体、API 应用程序建立集中一体的身份标识，对身份和属性进行全生命周期的管理，是数字世界秩序的底层支撑。

AI-R-IAM 采用了先进的身份认证机制，同时引入中国移动基于号卡特性的实名认证能力，有效防止未经授权的访问和潜在的安全威胁。

AI-R-IAM 不仅涵盖传统人类身份从入职到离职的全生命周期管理，还创新性地将 NHI 和 AI Identity 纳入统一治理范畴，助力企业实现数智化转型，有效降低安全风险，提升合规性。在具体流程如下：

(1) 用户身份在入职、转岗、离职时，进行信息收集验证、权限调整及账号注销等操作；

(2) 智能体身份在创建时生成唯一 ID，并设定权限，修改时变更信息，动态调整权限；

销毁时注销身份、回收资源并处理数据；

(3) 应用程序身份在创建时注册认证，授予权限，修改时更新

删除身份、销毁凭证并清理资源。

可信身份治理根据策略集中编排身份管理的过程，用于提供更好的身份可用性和访问权

控制不正确的访问。通过对不同身份的细致管理，采取身份验证、权限，以及更好的检测

确保系统的安全可靠，有力地保护了数据隐私。限设定、信息更新等措施，确

4.2 可信权限管理

型的访问、AI 接口访问、RAG 访问、模型调用、数据训练提 AI-R-IAM 针对多个大模

供了多场景、精细化权限管理功能。系统可以根据大模型用户的角色和职责，动态调整其对大模型的访问权限，确保每个用户只能访问所需的数据和功能，从而降低数据泄露和滥用的风险。这种精细化的权限管理不仅提高了系统的安全性，还增强了用户的操作体验。

- 用户访问权限

智能体对 RAG 的访问、用户对 RAG 的访问关联关系，建立大模型访问控制列表。

- 用户角色权限管理

角色级权限管理依托 RBAC (基于角色的访问控制) 技术，根据大模型不同的工作岗位

以及软件所具备的功能，为大模型的访问主体合理分配权限。

角色类型	模型训练	模型推理	数据标注	知识库更新	模型部署	审计日志
算法工程师	✓	✓	✓	×	×	×
数据科学家	✓	✓	✓	✓	×	×
系统运维	×	仅监控	×	×	✓	✓
业务用户	×	受限调用	×	×	×	×

- AI 接口权限管理

包括大模型接口策略、接口密钥、接口连接调用、接口速率、接口传输数据内容、接口

黑白名单等进行权限参数管理。

- RAG (检索增强生成) 权限隔离

文档级权限控制：不同部门只能访问其业务相关的知识库片段；

向量索引加密：通过分段加密技术，确保用户仅能解密与其权限匹

实时内容过滤：结合语义分析引擎，自动屏蔽无权限的检索结果 (

落)。

- Token 动态调配

可信数据管理

4.3 可

AI-R-IAM 通过深度整合先进的技术能力和安全机制，从数据输入开始、通过模型训练

AI-R

机制对数据进行深度处理和生成，最终输出多种形式数据的整个流程，提供端到端

和 RAG 相

据集中管理能力。无论是文本、文件还是图像，AI-R-IAM 都能精准识别潜在风险，

的可信数

位的可信数据管理能力。



4.4 可信行为审计

AI-R-IAM 通过多维数据采集、安全元数据管理、可信行为基线、数据画像分析、数据知识图谱、数据深度分析、智能追踪溯源等功能为大模型用户访问场景下的可信行为审计。



- 多维数据采集

通过采集大模型用户行为、实体行为、上下文、身份与权限、特定模型等多维数据，为可信行为基线构建全面的用户行为基线分析提供支撑。

- 安全元数据管理

对比、标识、对大模型用户操作行为数据进行标准化处理，包括数据提取、清洗、关联、分发，提升数据价值密度。

- 可信行为基线

通过机器学习和统计分析技术，对历史数据进行分析，建立大模型用户可信行为基线，并通过自适应用户和实体行为的动态变化，确保基线时效性和准确性。

- 数据画像分析

构建多维分析模型，反映大模型用户的特征、行为习惯、个人偏好，生成用户画像、模

型画像、数据画像等，并据此对风险画像评分，智能量化风险水平。

- 数据知识图谱

通过数据多源融合、关系排序、深度去重等深度治理及分析，实现不同实体之间的关联关系构建，提供从“关系”的角度分析大模型风险的能力。

- 数据深度分析

大模型数据全生命周期各个环节进行合规风险分析、重点数据审计、敏感数据和异常场景分析等，对大模型安全风险进行评估，识别已知风险和未知风险。

- 智能追踪溯源

通过数据分布、可疑路径、可疑日志链路分析，层层递进分析，实现风险智能自

追踪。

>访问数据分布->可疑路径->可疑日志链路分析，层层递进分析，实现风险智能自

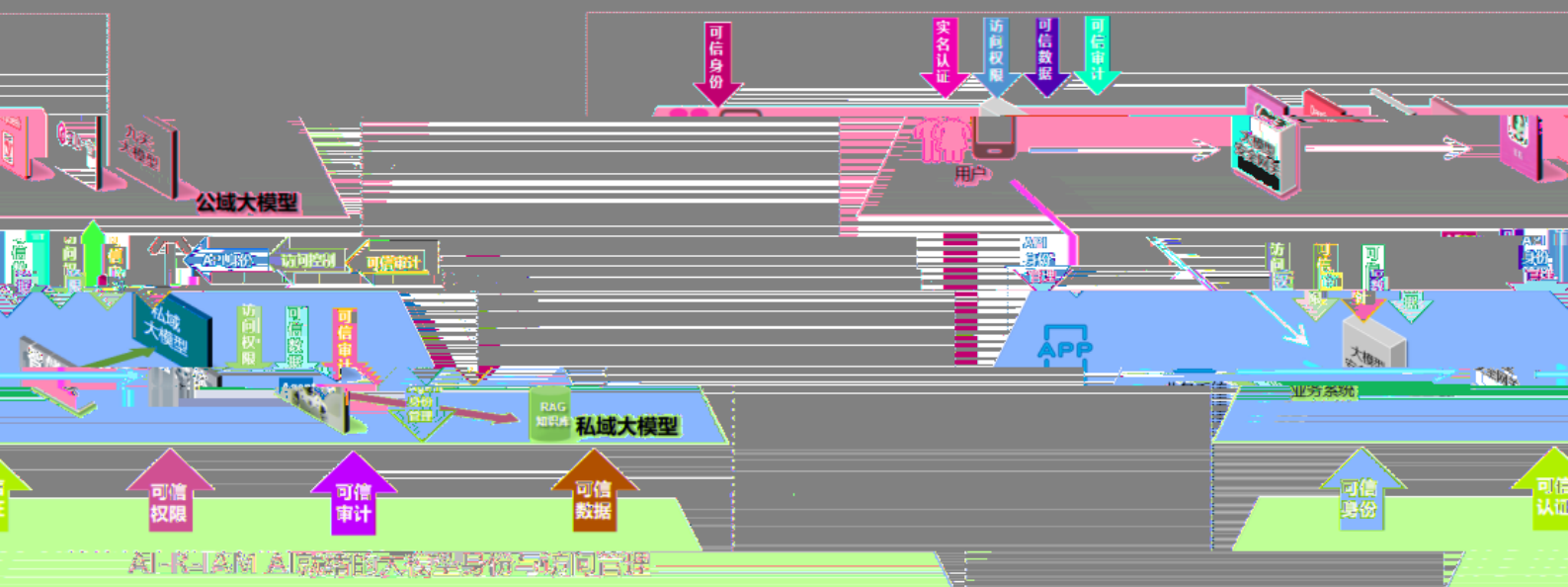
可疑IP-

踪。

5 应用场景

5.1 十类应用场景数据流

训练等多
 包括：用
 大模型场
 AI-R-IAM AI 就绪的大模型身份与访问管理，为大模型的访问、使用、调用、
 种场景提供可信身份、可信认证、可信权限、可信审计、可信数据能力。相关场景
 户访问私有大模型、公域大模型场景；业务系统调用智能体、RAG 知识库、私域
 景；大模型数据训练场景、大模型数据运维场景等。



场景一：用户访问大模型应用场景

5.2 场

AI-R-IAM AI 就绪的大模型身份与访问管理能够为用户访问大模型应用场景提供可信
 信认证、可信权限、可信审计、可信数据等关键能力的全面覆盖。主要场景包括用
 身份、可信
 户访问公域
 成域大模型、用户访问私域大模型、运营维护人员访问私域大模型等。

可信身份认证方面，构建了一套严谨的身份验证机制。当用户首次访问大模型应用时，系统会要求用户提供身份凭证，这可能包括用户名和密码组合、生物特征（如指纹或面部识别）或者是基于可验证凭据（Verifiable Credential, VC）的身份认证等方式。只有经过严

保护用户的数据不被未授权访问。

在安全审计方面，AI-R-IAM 持续监控整个系统的运行状况。它会记录下所有用户的操作行为，包括登录时间、执行的操作类型、访问的数据范围等详细信息。通过对这些审计日志的定期分析，可以及时发现潜在的安全隐患或者异常操作。例如某个账户在短时间内进行了大量不符合常规模式的数据访问请求，就可以触发警报提示管理员进行进一步调查，从而有效地维护系统的稳定性和安全性。

2. 普通用户访问私域大模型

企业内部的私域大模型应用（例如知识问答、文档生成等）普通用户通过终端设备访问部署在（服务）。与公域大模型访问场景相比，私域场景在数据隐私和合规性方面有着更为严苛的要求，因此必须在权限控制以及可信数据管理方面具备更强的支撑能力。

在用户身份可信认证基础上，私域场景下需要对不同角色进行更加细致的权限划分。企业内部存在多种角色，如高层管理者、部门主管、普通员工等。高层管理者需要了解整个企业的运营状况，包括各个部门的数据汇总；部门主管则只需要并且能够对本部门员工的权限进行一定程度的调整；普通员工的权限最自己工作直接相关的少量数据。这种多层次的角色划分，要求权限控制灵活性和准确性，以满足不同角色的需求同时保障数据安全。

运营维护人员访问私域大模型应用服务器，负责日常维护和故障排查。由于运维人员权

限的广泛，可以直接接触底层数据，因此需要严格控制其访问权限，并确保所有操作

可追溯。针对运营维护、开发测试等后台管理人员，在普通用户

等方面能力。上，还需要重点加强精细化权限管理、高危操作、金库管

理等方面能力。上，还需要重点加强精细化权限管理、高危操作、金库管

理等方面能力。上，还需要重点加强精细化权限管理、高危操作、金库管

理等方面能力。上，还需要重点加强精细化权限管理、高危操作、金库管

理等方面能力。上，还需要重点加强精细化权限管理、高危操作、金库管

理等方面能力。上，还需要重点加强精细化权限管理、高危操作、金库管

数据或管理参数等。同样，软件维护人员也不能被允许去修改硬件设备的配置。

这种具有角色的权限划分还需要考虑到维护团队内部的职级关系。例如，初级维护工程

师可能只能进行一些基础的操作，如查看日志、重启服务等；中级维护工程师则可以

进行更高级的操作，如修改系统参数、执行大规模的数据迁移等。通过这种多层次

的权限体系，可以有效避免因为权限滥用而导致的安全事故。

而，对于安全团队

而言，除了权限管理之外，还需要关注其他方面的安全。例如，在维护人员访问服务器

时，除了使用强密码和定期更新密码外，还可以采用多因素认证（MFA）来增强身份验证

的安全性。此外，对于敏感数据的访问，还可以采用数据加密和访问审计等手段来进一步

提升安全性。同时，对于高危操作，还可以采用操作审批和告警机制来及时发现和响

应异常行为。在私域大模型环境中，系统管理员应该能够删除用户权限或禁用用户

权限。一旦检测到异常行为，系统应该立即发出告警，以便安全团队及时响应。

此外，在第一次部署时，应该给维护人员提供详细的操作指南和培训，确保他们能够

息，提醒运维人员该操作可能存在严重后果，并要求他们确认是否继续。这种警告提示可以让运维人员重新审视自己的操作意图，避免因误操作而导致问题。第二层次是权限二次验证，即使运维人员确认要继续执行高危操作，也需要通过额外的权限验证步骤。这可以包括输入更高级别的管理员密码、使用硬件令牌进行身份验证等，只有通过了这一严格的验证过程，才能够真正执行高危操作。第三层次是完全阻断，对于某些极其危险的操作，系统可以直接

限制，例如，对于删除整个数据库这样

禁止其执行，无论运维人员如何尝试都无法绕过这一层

限制。同时，系统还可以记录运维人员的操作行为，以便在发生安全事件时进行溯源。

在私域大模型运维中，金库可以被视为一个特殊的区域或容器，其中存放着最为重要的数据和配置信息，如模型的核心算法代码、高度机密的训练数据、关键的系统配置文件等。

因此，金库就像银行的金库一样，需要同等级别的保护，只有经过严格授权的人员才能访问。

运维人员需要通过身份验证和权限验证才能访问金库。

系统应该记录所有访问金库的操作，包括访问时间、访问地点、访问者身份等信息。

同时，系统还应该记录所有访问金库的操作，包括访问时间、访问地点、访问者身份等信息。

系统应该记录所有访问金库的操作，包括访问时间、访问地点、访问者身份等信息。

系统应该记录所有访问金库的操作，包括访问时间、访问地点、访问者身份等信息。

系统应该记录所有访问金库的操作，包括访问时间、访问地点、访问者身份等信息。

系统应该记录所有访问金库的操作，包括访问时间、访问地点、访问者身份等信息。

系统应该记录所有访问金库的操作，包括访问时间、访问地点、访问者身份等信息。

系统应该记录所有访问金库的操作，包括访问时间、访问地点、访问者身份等信息。

系统应该记录所有访问金库的操作，包括访问时间、访问地点、访问者身份等信息。

系统应该记录所有访问金库的操作，包括访问时间、访问地点、访问者身份等信息。

系统应该记录所有访问金库的操作，包括访问时间、访问地点、访问者身份等信息。

2. 业务权限管控

业务系统在访问公域、私域大模型时，如果没有基于最小权限原则来分配访问权限，如

某些功能可能只需要读取权限，却被赋予了写入或管理权限。由于权限划分不清晰，可能导致内部人员误操作或恶意操作，进而影响系统整体安全。

AI-R-IAM 在业务系统访问公域、私域大模型时，通过精细化的权限控制，保障大模型资源的合理分配与安全性，降低安全风险，实现大模型业务权限可信。

实体级授权：通过 ABAC 和 PBAC 技术，业务系统根据实体的属性、状态和业务流程，

通过动态调整权限，确保每个实体只能访问其需要的资源。

角色级授权：采用 RBAC 技术，

根据岗位和功能，为业务系统分配相应权限，确保高权限角色的操作

有效管理大规模访问需求，并降低权限管理复杂度。

属性及数据敏感

数据级授权：结合 RBAC 和 ABAC 技术，系统依据业务系统的角色、

属性，精细控制数据访问权限，防止数据泄露和滥用。

系统的 Token 分

Token 动态授权：通过 RBAC、ABAC 和 PBAC 技术，动态调整业务

配置，实现对资源访问的弹性管控，确保资源得到高效合理的分配和使用。

3. 业务行为审计

业务系统调用公域/私域大模型 API、模型服

务系统调用公域/私域大模型 API、模型服

务系统调用公域/私域大模型 API、模型服

务系统调用公域/私域大模型 API、模型服

务系统调用公域/私域大模型 API、模型服

务系统调用公域/私域大模型 API、模型服

务系统调用公域/私域大模型 API、模型服

务系统调用公域/私域大模型 API、模型服

务系统调用公域/私域大模型 API、模型服

务系统调用公域/私域大模型 API、模型服

AI-R-IAM 在业务系统访问公域、私域大模型时，通过全链路日志采集、智能化分析、智能追踪溯源三方面保障业务行为可信。

全链路日志采集：采集私域大模型、RAG 知识库的用户行为、实体行为、上下文、身

份与权限、特定模型等多维数据，以构建全面的用户行为基线并检测异常。

智能化分析：通过机器学习和统计分析技术，对历史数据进行分析，建立私域大模

型、RAG 知识库的用户可信行为基线，并通过自适应用户和实体行为的动态变化，确保基线

有效性和准确性。同时采用画像分析、知识图谱等技术对私域大模型、RAG 知识库安全风

险进行智能追溯，对可疑 IP、访问数据分布、可疑路径、可疑日

域大模型、RAG 知识库的数据输入到模型输出提供全链路智能追

智能追溯溯源：对私

踪溯源，采用路径计算、

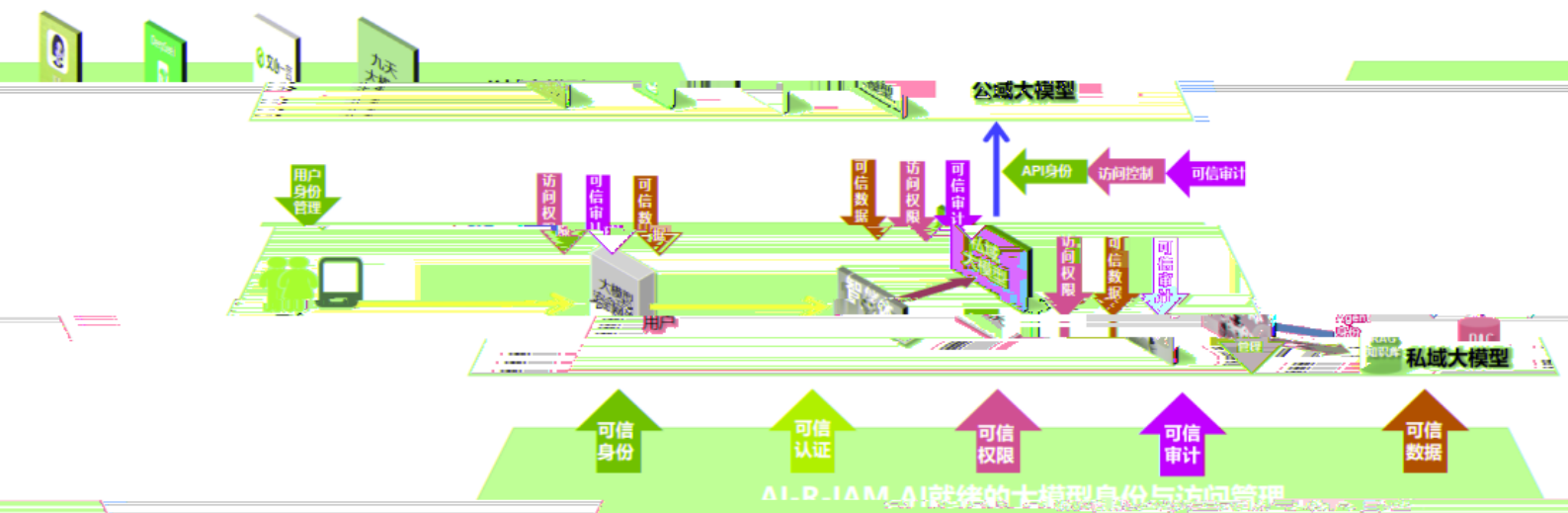
路径探寻算法，提供可疑 IP->访问数据分布->可疑路径->可疑日

志链路分析，层层递进分

析，实现风险智能追踪。

用户-AI 智能体的安全应用管理场景

5.4 场景三：用



本身份并绑

场景概述：在用户到智能体的安全应用场景中，AI-R-IAM 赋予用户与智能体

证，调用时

定，防智能体被仿冒。用户经实名认证建立可信身份，智能体以独立身份接受认证

的安全。通

获有限权限，同时通过身份绑定对双方行为审计监管，保障访问数据和训练数据的

过 AI-R-IAM 能力解决如下问题：

1. 用户到智能体全链路身份可信

与智能体身份 (AI

如果智能体身份做了某

及模型紧密绑定, 运用

码构建, 并成功部署

接收到该智能体所发出的请求时, 能够对其携带的签名、证书或公钥进行严

来精确确认该智能体的真实身份, 同时评估其可信度, 从而确保系统运行的安

认证与授权降低威胁面

智能体场景下, 普通用户通过大模型安全网关与 AI-R-IAM 实名认证能力建立实

前端应用访问可信, 同时管理员利用 AI-R-IAM 权限管理划分最小访问权限

同时, 认证流程不再只面向普通用户, 每一个智能体都能以独立身份进行访问控制, 并在执行任务前接受严格的认证校验, 通过短期令牌 (Short Access Token) 或一次性密钥 (One-Time Key) 等方式, 让授权更具时效性和可追溯性。当智能体调用 RAG 时, 通过金库 (有效时间段) 模式来获取有限访问权限, 将潜在攻击面降至最低, 同时通过环境属性 (如运行环境指纹, 地理位置等) 进行增强认证。

3. 关联审计与监管

通过身份绑定机制, 对普通用户使用智能体行为, 以及智能体自身的行为进行审计和监管, 当智能体发生越权操作或异常行为时, AI-R-IAM 系统应有能力快速定位并冻结该代理

权限, 监管部门可对其可追溯知悉的操作细节, 例如访问了哪些数据, 执行了何种指令

依据何种规则。

图 1-10 智能体身份绑定示意图

通过 AI-R-IAM 提供可信身份管理, 赋予用户可信身份 (user Identity), 由统一身份 ID 进行绑定管理, 便能在审计层面明确, 如某项操作, 系统可追溯到是谁让智能体进行操作, 方便界定相应责任。

为了有效防止智能体被恶意掉包或仿冒, 智能体会与特定的代码

数字签名或证书技术, 但确证即某个智能体是由预先指定的模型和代

当 RAG 知识库

格验证, 以此来

全性和可靠性。

2. 强认证

在访问智能

名身份, 确保前

在企业内部，智能体易成为攻击目标，被攻击后可能执行恶意操作、造成训练数据外泄

新风险，如自动提交敏感数据到外部。

智能体，智能体与数据存储服务器，RAG 知识库，
最小化策略，制定智能体对 RAG 的数据访问权限
控制，可及时阻断和审计智能体异常行为等。避免

在此场景下，AI-R-IAM 提供用户和智能体
内外部大模型之间的数据权限边界，实行最小
范围，辅以模型护栏、上下文检测等安全控制
因智能体自主决策引发大规模安全事故。

6 发展趋势展望

在大模型技术蓬勃发展的当下，AI-R-IAM 为其深度应用筑牢安全防线，从多维度构建一体化全程可信能力，引领大模型安全从“单点可控”迈向“一体化全程可信”。

1、在身份欺诈领域，AI-R-IAM 将实现重大技术飞跃。随着大模型在身份认证深度应用，其能够对海量的身份数据进行学习和分析，构建起精准的身份行为模型。通过持续监测用户的登录习惯、操作行为模式以及设备使用情况等多维度信息，大模型可以

及时发现异常行为。AI-R-IAM 将借助大模型的智能分析能力，同其它网络安全设备协同实现网络攻击的主动防御。基于零信任架构，AI-R-IAM 将对所有网络访问请求进行严格的身份验证和权限检查，即使企业内部网络流量也不例外，进一步强化网络安全防护体系。

2、在网络安全方面，AI-R-IAM 将与防火墙、入侵检测系统或管理系统进行结合，实现对网络流量的实时监控和异常行为检测，提升网络安全防护能力。

随着大模型技术的广泛应用，AI-R-IAM 将深度融合于大模型应用，保障大模型在物联网场景下与海量设备交互时的数据安全和身份可信。

借助 5G/6G 网络的高速低延迟特性，实现更高效的安全数据传输和实时的安全监测。通过持续的技术创新和融合，AI-R-IAM 将不断适应新的安全挑战，为大模型的广泛应用和数字经济的蓬勃发展保驾护航，构建更加安全、可信的数字未来。