

启明星辰+大模型应用安全+深度应用安全+基础系列白皮书

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

的最终解释权
北京启明星辰信息安全技术有限公司版权所有，并保留对本文档及本声明
和修改权。

出现的任何文字叙述、文档格式、插图、照片、方法、过程等内容，除另有特
本文档由

其著作权或其他相关权利均属于北京启明星辰信息安全技术有限公司。未经北京
别注明外，

信息安全技术有限公司书面同意，任何人不得以任何方式或形式对本手册内的任何
启明星辰信

复制、摘录、备份、修改、传播、翻译成其它语言、将其全部或部分用于商业用途。
部分进行复

明
负责申

知。
本文档依据现有信息制作，其内容如有更改，恕不另行通知。

启明星辰对本声明不承担法律责任，北京启明星辰信息安全技术有限公司在编写本声明的时候

信息安全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

的损失和损害承担责任。

信息反馈

请反馈：
如有任何宝贵意见，

地址：北京市海淀区中关村软件园二期启明星辰大厦 邮编：100193 电话：010-82779000 传真：010-82779000
总可以访问地址：www.venusstech.com.cn 获取最新技术和产品信息。

100193 电话：010-82779000

传真：010-82779000

总可以访问地址：www.venusstech.com.cn 获取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.....	7
2.1 大模型应用安全挑战.....	7
2.1.1 身份与访问控制的安全挑战.....	9
2.1.2 数据安全挑战.....	11
2.2 大模型应用中安全挑战.....	13
2.2.1 身份与访问控制的安全挑战.....	16
2.2.2 数据安全挑战.....	17
3 身份与访问管理系统 (IAM)	19
4 AI 就绪的大模型身份与访问管理.....	21
4.1 可信身份管理.....	21
4.2 可信权限管理.....	21
4.3 可信数据管理.....	21
4.4 可信行为审计.....	21
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

辰公司成立于 1996 年，由留美博士严望佳女士创建，是国内最具实力的、拥有自主知识产权的网络安全产品、可信安全管理平台、安全服务与解决方案的综合提供商。

启明星辰
完全自主研

2010 年 6 月 23 日，启明星辰在深交所中小板正式挂牌上市。

启明星辰拥有国内领先的威胁情报、入侵检测、漏洞扫描、安全审计、安全检测、网络审计、安全运营、终端管理、加密认证等技术领域，共有百余个产品型号，并根据客户需求不断增加。启明星辰解决方案为客户的安全需求与信息安全产品、服务之间架起桥梁，将客户的安全保障体系

与信息安全核心技术紧密相连，帮助其建立完善的安全保障体系。

自 2002 年起，启明星辰就持续保持国内入侵检测、漏洞扫描市场占有率第一。近年来，发展成为国内统一威胁管理、安全管理平台国内市场第一位，安全性审计、安全专业服务市场领导者。目前，公司在全国各省市自治区设立三十多家分支机构，拥有覆盖全国的渠道和售后服务体系。

长期以来，启明星辰公司得到了党和国家领导人的关怀与鼓励。2000 年 1 月，江泽民、李岚清、曾庆红等党和国家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

余项专利和软件著作权，参与制订国家及行业网络安全标准，填补了我国信息安全科研领域的多项空白。

作为信息安全行业的领军企业，启明星辰以用户需求为根本动力，研究开发了完善的工业安全产品线，通过不断耕耘，已经成为在政府、电信、金融、能源、交通、军队、军

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

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

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，

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

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

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

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

位信息安全保障。

大型活动提供全方位

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，已累计资

在公司快速稳定

贫困群众上亿元人民币，并在江西、青海、新疆等地援建了 5 所希望

助贫困学子、受灾、

小学。

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主创新的安

启明星辰公司将

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，为打造和

全产品和最佳实践朋

信息安全产业第一品牌而不懈努力。

提升国际化的民族信

2 大模型的发展与风险

2.1 新质生产力的“大模型”

心驱动力。从 ChatGPT, 大模型在自然语言处理、海量数据上进行训练, 能

生成、智能问答、图像

分为四个阶段:

大的自然语言处理能力震撼全球, 开启研究。

展其应用领域。谷歌、微软等科技巨头

文心一言, 整合自身技术优势; 阿里发布通义千问, 立足丰富业

大模型, 专注语音交互。高校和科研机构也积极开展研究, 为大

模型发展提供技术储备。

● 爆发期 (2024 年)

型社区 (1) 国外: OpenAI、谷歌等公司不断升级大模型, 性能显著提升, 开源大模

槛, 加速技术普及。

且 24 日, 在中国移动人工智能生态大会上, 中国移动正式发布

千卡级智算集群、千亿多模态大模型、汇聚百大要素的生态

近年来, 大模型技术取得了突破性进展, 成为人工智能领域的核
再到百度文心一言、阿里通义千问、九天大模型、DeepSeek 等, 大
计算机视觉、语音识别等领域展现出强大的能力。这些模型通过在海量

够学习到丰富的语言知识, 进而理解和逻辑推理能力, 从而实现文本

生成、代码编写等多样化的任务。大模型发展的历程

● 准备期 (2022 年 11 月)

2022 年 11 月, OpenAI 发布 ChatGPT, 其强

大模型发展浪潮, 促使各方加大在该领域的投入与研究

● 成长期 (2023 年)

(1) 国外: OpenAI 持续优化 ChatGPT, 并拓

积极布局。

(2) 国内: 百度推出

务场景: 科大讯飞推出星火

(2) 国内: 2024 年 5

了“九天”人工智能基座, 包含万

平台、工业链结加工开发机、“五塔七楼”渐到上凉、下凉、凉其脚力、脚凉等上厂

不断迭代，众多初创企业也凭借特色大模型在细分领域崭露头角。

模型不

大规模应用期 (2025 年 1 月-)

(1) 国外：大模型深入医疗、金融、教育等多行业，助力疾病诊断、风险评估、个性

习等。OpenAI 还在探索新应用领域。

化学习

(2) 国内：DeepSeek 凭借创新算法和架构，展现出低成本、高效能以及开源优势，

智博会感，DeepSeek 还免费开放了最中此模型存在各、金融、医疗、教育、制造、农商、

型。

等行业大规模应用，推动各行业数字化、智能化转

前所未有的速度和深度渗透到企业运营和个

大模型技术作为新一代生产力的代表，正以前所

的关键力量。它不仅重塑了传统行业的运

人生活的方方面面，成为推动社会进步和经济发展

正实现了“智能赋能”的愿望。

作模式，还催生了全新的应用场景和商业模式，直

大模型应用中安全挑战

2.2 大模

1 身份与访问控制的安全挑战

2.2.1

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

在大模

AI Agent 与数据、人与 AI Agent 之间交互量呈指数级增长。

与应用、AI

麦肯锡的《2025 年技术展望》调查，78% 的高管同意在未来必须为 AI Agent 构

根据按

AI 领域，在大模型

建与人类同等重要的数字生大系统，这意味着传统身份概念不断延伸不

应用中产生了三类身份概念，即人类身份 (Human Identity) 以及非人类身份 (Non-Human

出诸多安全问题。

Identity, NHI)，都暴露出

1、人类身份安全问题：人类身份代表自然人身份实体，在大模型应用交互场景中的背

景下，容易出现滥用身份非法访问大模型系统窃取数据，各应用中在大模型服务时，如用

户身份和大模型应用安全问题，大模型应用中数据使用不透明也可能泄露隐私。

访问性

大模型场景下应用 AI Agent 具

2 AI Identity 安全问题·越来越多的 AI Agent 在

大模型场景中应用，AI Agent 具

有越来越复杂的身份关系，需要身份和专用身份管理。

的身份关系，需要身份和专用身份管理。

《人工智能代理的身份》中指出 AI Agent 具有复杂的

应重新思考 AI Agent 作为员工，在大

AI Agent 应视为具有身份的不同员工的新兴范式，

与管理方式有这问题，白自自夕自公屋

大模型应用场景下建立身份层，即 AI Identity，它既

性，会出现如下风险：

样本

(1) 访问绕过：攻击者可能针对 AI Agent 所使用的模型进行攻击，如通过对抗

攻击，白模型输入精心构造的恶意数据，使模型产生错误的输出，且不能通过身份验证或

非法访问权限。

方式请求访问权限，

(2) 访问权限放大：与人类身份不同，AI Agent 不会以可预测的方式

权限，从而能够访问

如果权限管理机制不完善，可能导致用户或应用程序获得超出其应有的权

和操作敏感资源，造成数据泄露或系统破坏。

行为数据等。如果

(3) 自身安全：AI Agent 身份系统通常存储大量用户的个人信息。

系统遭受黑客攻击，这些数据可能被窃取，导致用户隐私泄露，进而可能被用于精准诈骗、

定向广告骚扰等。

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

机器等实体绑定的数字身份。这些包括机器人、API 密钥、服务帐号、OAuth 令牌、云服

务等。此外，在识别和验证身份时，应知道在传输和存储身份数据时，应使用加密和认证。

2025 年 2 月发布《十大 NHI 风险 2025》，如下图所示。



存在于大模型基础设施交互过程中，NHI 的身份生命周期管

问题：

大模型开发应用时，很多 NHI 硬编码在代码库等位置，容易

用迭代快，加之生命周期流程弱、账号使用信息不可见，许多

面。

数参与大模型应用的组织中，NHI 无明确所有权信息，而明确

效很关键。

多情况下，相同 NHI 会在生产和非生产环境中使用，或者相同

有相同的密码，从而增加了横向移动的风险。

用的多程序间共享 NHI，违反原则，使密码循环等安全操作有

NHI 在大模型应用中，大多数

理薄弱，存在诸多安全风险与

(1) 纯文本/未加密凭证：

被内外部威胁者发现。

(2) 影子账号：大模型应用

NHI 账号不活跃，增加攻击面

(3) 缺乏账号所有权：多数

管理者对安全维护和问题补救

(4) 缺乏环境隔离：在许多

的逻辑 NHI 在每个环境中都

(5) 凭证共享：大模型应用

杂，难以掌握所有依赖关系。

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

示词可能在输入环节直接暴露。在输出结果时，若对输

示词的验证和过滤机制不完善，敏感

商业机密等敏感信息，如用户输入涉及个人健康状况

出内容审查不足，可能泄露用户隐私

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

数据标注、整合环节出

超出授权范围获取敏感知识数据，且知识库中的数据来源复杂，

和安全性难以保障，可能误导用户并造成决策失误。

现错误或被恶意篡改，查询结果的准确性

段，攻击者故意向训练数据中注入恶意数据，这些

(4) 数据投毒风险：在大模型训练阶

型训练数据中混入大量被恶意标注的图像，使模型在识别相关图像时出现严重偏差。

模型

(5) 合规风险：随着数据安全相关法律法规的不断完善，大模型应用需要满足严格的合规要求。若企业在数据收集、使用、存储和共享等环节不符合相关法规，如未遵循数据最小原则、未获得用户充分授权等，可能面临法律诉讼和巨额罚款。

模型

规

化

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

可信网络,赋予人员、设

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

鉴权,基于可信身份对访

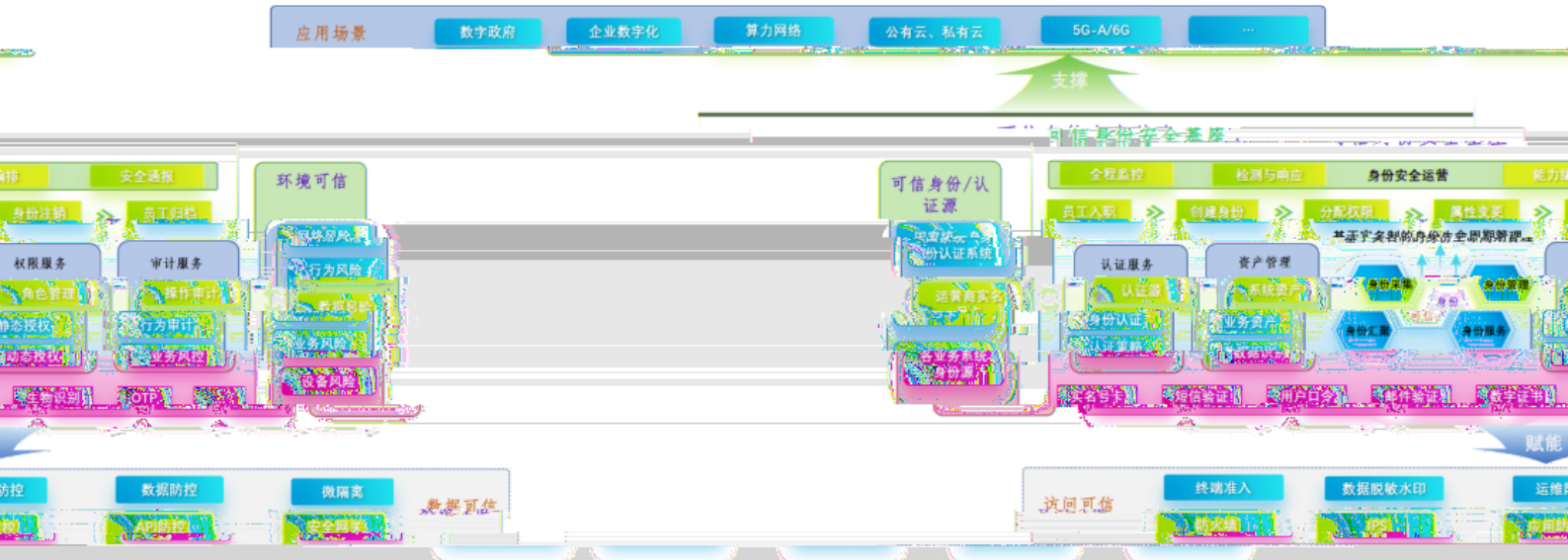
备、应用、接口、数据等实体唯一身份标识,实现实名制可信身份

向身份、环境、访问、数

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”逐

据等维度的“一体化全程可信”。

基于 IAM 构建的“一体化全程可信数字身份底座”总体架构如下图所示:



IAM 作为体系构建的核心能力,依托丰壤各类应用场景提供实名可信身份,自主赋能

同策略。

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

金、数据风险、业务风险、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

风险的全链路安全风险审计。

安全设备、系统资源、应用资源、

访问可信与数据可信作为一体两面，将防控范围扩展到安

全网各、员工账号互操作、可、API资源等，通过“ID+身份”的融合形成可信身份并融入

组合性和可扩展性，使企业能够以更少的资源实现更好的安全性，搭配一体化安全机制的融

合协同，实现全程全网、端到端的安全保护。

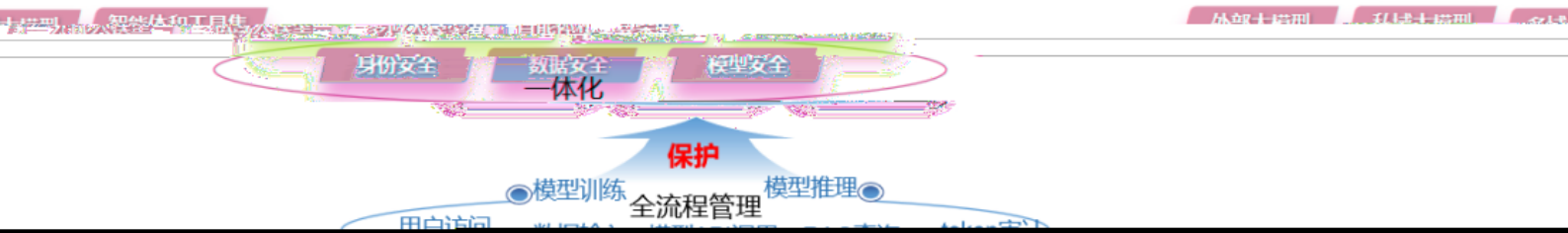
随着大模型作为颠覆性生产力，在政府、企业、个人中应用越来越深入，需要针对其身份、访问管理、数据安全等层面的需求，建立一套基于大模型技术特点面向全场景的一体化全程可信数字身份体系，即AI赋能的大模型身份可信管理。

访问管理

AI-R-IAM) 是为大规模人工智能模型 (如
访问控制、数据安全、模型安全的基础支
型在训练、部署和使用过程等多个场景中
点可控” 迈向 “一体化全程可信” 保护体

4 AI 就绪的大模型身份与

AI 就绪的大模型身份与访问管理 (以下简称: AI
DS、GPT、BERT 等) 的深度应用提供身份安全、
撑系统, 从身份、访问、行为、数据等维度为大模
构建一体化全程可信能力, 构建大模型安全从 “单
系, 同时为其它安全能力提供身份管理属性支撑。



用户和系统能够访问大模型。系统不仅支持传统

基于中国移动号卡特性的实名认证能力，有效防

撑，如网络防火

的精细化的权限

访问权限，确保

这种精细化的权

的数据安全保护

数据的安全性。

审计和实时监

型资源、执行的

面临的安全挑战，

支撑，为未来人

先进的身份认证机制，确保只有经过授权的用

的用户名和密码验证，多因素认证，还引入基

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

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

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

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。

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

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

● 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

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

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

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

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）、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：

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；

智能体身份在创建时生成唯一 ID 并设定权限，修改时亦可变更信息并调整权限；

应用程序身份在创建时注册认证，授予权限，修改时重新信息，变更权限，销毁时删除身份、销毁凭证并清理资源。

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，修改时重新信息，变更权限，销毁时删除身份、销毁凭证并清理资源。

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：不同用户只能看到其权限范围内的文档接

口内容；访问控制：通过身份鉴别技术，确保用户仅能解密其权限

范围内的内容；实时内容过滤：结合语义分析引擎，自动屏蔽无权限用户检索

结果。

○ Token 动态调用

4.3 可信数据管理

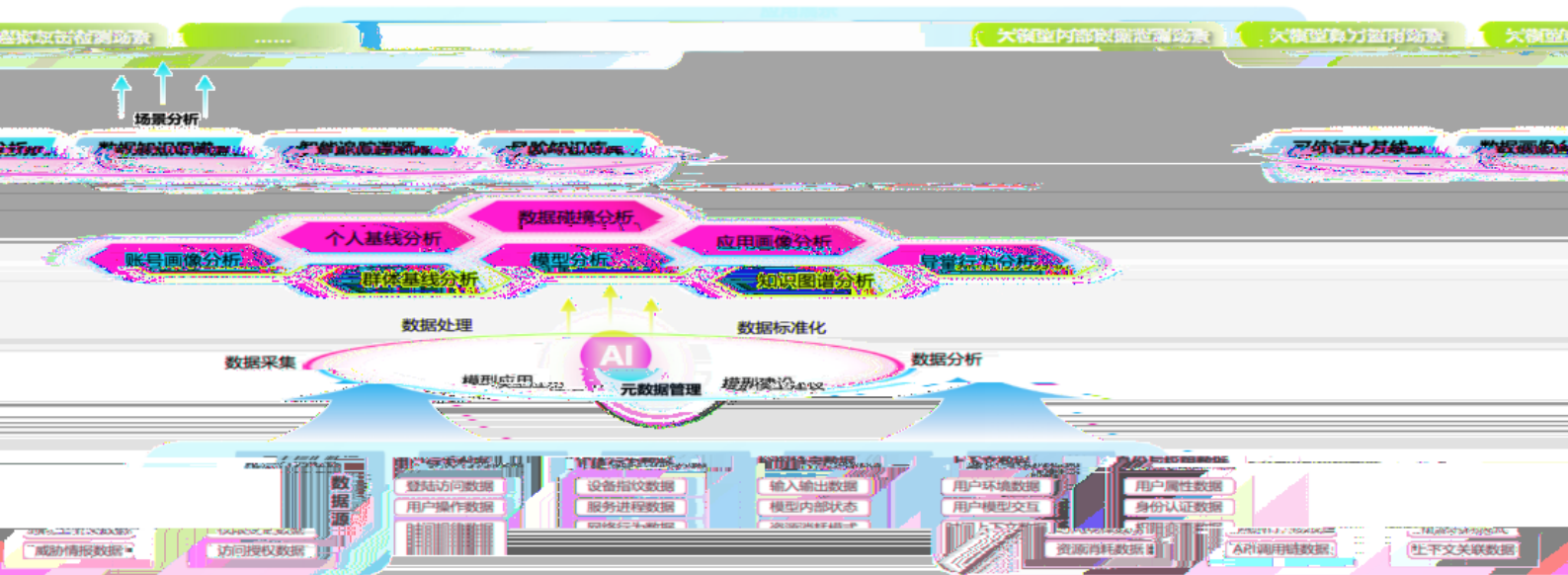
AI-R-IAM 通过深度整合业界先进的技术能力和安全机制，从数据输入开始，通过模型训练和 RAG 机制对数据进行深度处理和生成，最终输出多种形式数据在整个流程，提供端到端的可信数据安全管理能力。无论是文本、文件还是图像，AI-R-IAM 都能精准识别潜在风险，

提升可信数据管理能力。



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。



与权限、特定模型等多维数据，为可信行为审计提供数据支撑。

- 多维数据采集

通过采集大模型用户行为、实体行为、上下文、身份信息等数据，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、对比、标识、分发，提升数据价值密度。

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、对比、标识、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，对历史数据进行建模，建立大模型用户可信行为基线。

通过机器学习和统计分析技术，对历史数据进行建模，建立大模型用户可信行为基线。

并动态跟踪应用用户实体行为的变化，确保基线实时性和准确性。

- 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，并据此评估画像风险，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

通过数据融合、关联分析、数据挖掘等技术，实现数据智能追踪溯源，识别异常行为和可疑路径。

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

可疑 IP->访问数据分布

踪。

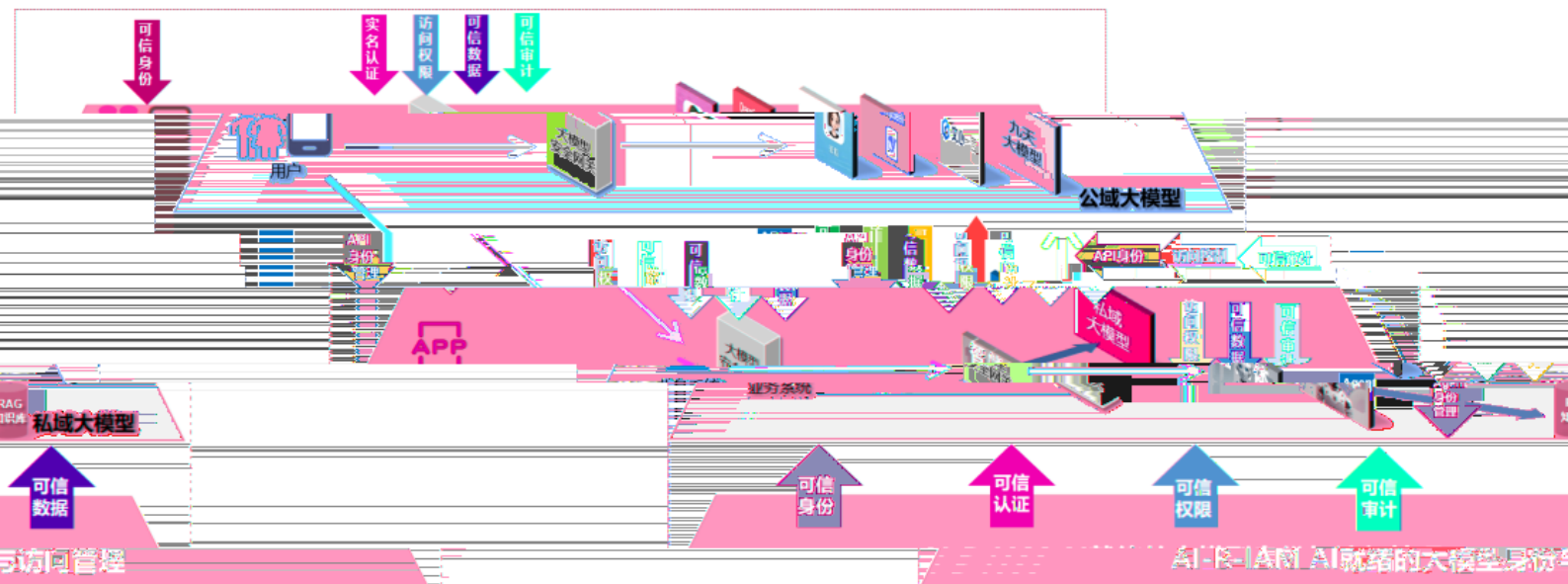
应用场景

5

大模型应用场景数据流

5.1

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、训练等多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。相关场景包括：用户访问私有大模型、公域大模型场景；业务系统调用智能体、RAG 知识库、私域大模型场景；大模型数据训练场景、大模型数据运维场景等。

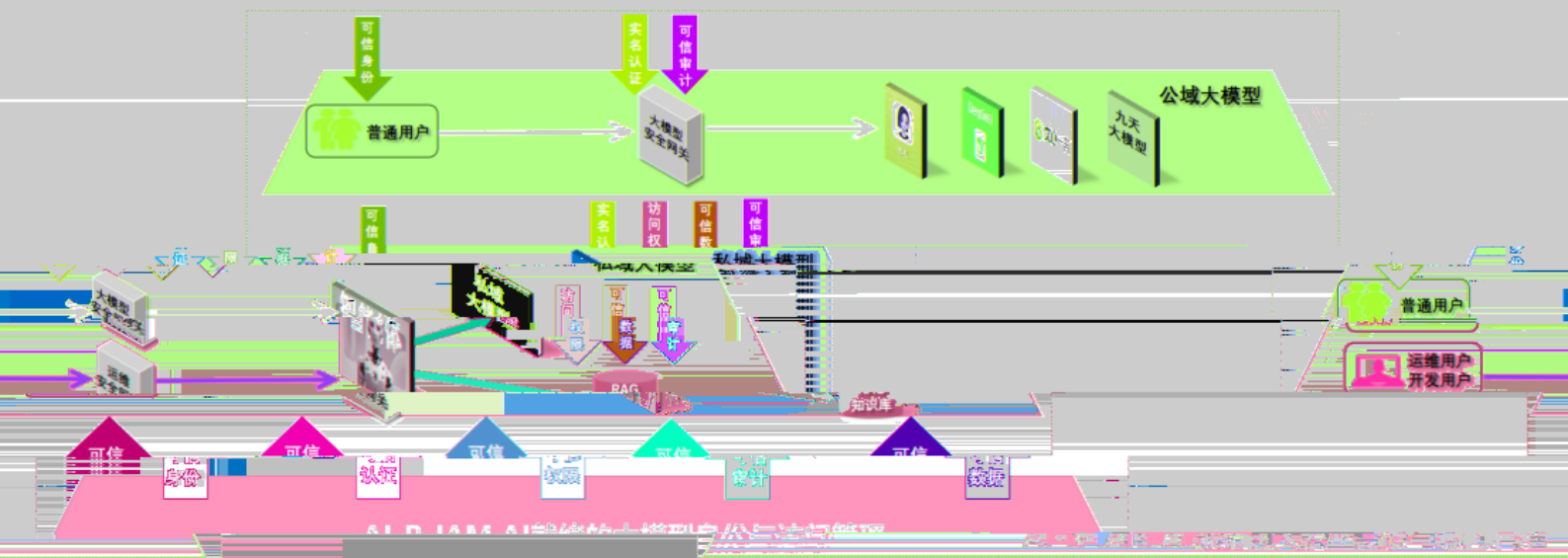


应用场景

5.2 场景一：用户访问大模型应用

为用户访问大模型应用场景提供可信身份、可信认证、可信权限、可信审计、可信数据等关键能力的全面覆盖。主要场景包括用户访问公域大模型、用户访问私域大模型、运营维护人员访问私域大模型等。

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提供可信身份、可信认证、可信权限、可信审计、可信数据等关键能力的全面覆盖。主要场景包括用户访问公域大模型、用户访问私域大模型、运营维护人员访问私域大模型等。



1. 普通用户访问公域大模型

普通用户，无论是企业员工还是个人用户，都可以借助自己的终端

用。这些人搜网应用涵盖多种类型的服

或者平台等。去访问哪些数据在公共域上的人搜网应

的搜索结果是网络中所有公开、合法、合规的数据。

这些数据的来源包括搜索引擎、社交媒体、公开数据库等。

这些数据的来源包括搜索引擎、社交媒体、公开数据库等。

网络中所有公开、合法、合规的数据。

这些数据的来源包括搜索引擎、社交媒体、公开数据库等。

网络中所有公开、合法、合规的数据。

户提供服务。对普通用户而言，操作起来

这些人搜网应用可以 Web 界面的形式来为用

继续浏览；一份产品推广文案时，只需在

非常方便。比如，当一个企业员工需要利用文本生成

的网页以有具体名称、特点等，再点击在

浏览器，进入该服务的 Web 界面，然后输入产品

个人用户上传图像识别服务识别一张拍

成按钮，很快就得到一份完整的文案初稿。又如

能够从 Web 界面上获取到详细的购物分

物照片的种类时，只需面对给定的上传照片按钮，点

类结果。

个相关系统需要具备完善的 AI-R-IAM

然而，在这种便捷的服务背后，保障用户数据安全

身份认证及安全审计的能力。

在普通用户对公域大模型场景中，主要提供可信

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（如指纹或面部识别）或者是基于可验证凭证（Verifiable Credential, VC）的身份认证等方式。只有经过严

密授权访问。

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、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大量的数据访问请求，就可以触发警报提示管理员进行进一步调查，从而有效维护系统的稳定性和安全性。

2. 普通用户访问私域大模型

普通用户通过终端设备访问部署在企业内部的私域大模型应用（例如知识问答、文档生成等

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。高层管理者可能需要全面了解整个企业的运营状况，包括各个部门的数据汇总；部门主管则只需要关注本部门的数据，并且能够对本部门员工的权限进行一定程度的调整；普通员工的权限最为有限，只能访问与自己工作直接相关的少量数据。这种多层次的角色划分，要求权限控制系统具备很高的灵活性和准确性，以满足不同角色的需求同时保障数据安全。

3. 运营维护人员访问私域大模型

可信身份认证方面，构建了一套严谨的身份验证机制。当用户访问系统时，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证（Verifiable Credential, VC）的身份

保护用户的数据不被未

在安全审计方面，

作行为，包括登录时间

志的定期分析，以及

了大量不符合常规模式

有效地维护系统的稳定

2. 普通用户访问私域大模型

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

限较高，且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

台管理人员，在普通用户私域访问大模型的安全能力基础

溯。针对运营维护、开发测试等后

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

器的物理状态检

同运维人员可能承担着不同的职责，例如有的负责硬件设备的维护，如服务

更新、数据库性

查、网络连接状况监测等；有的则专注于软件层面的维护，像操作系统补丁

维护人员，他们

能优化等工作。因此，需要根据这些不同的职责范围来划分权限。对于硬件

器上的敏感业务

应该只能访问与硬件相关的监控信息和配置界面，而不能触及到存储在服务

目相关的设置。

数据或模型参数等。同样，软件维护人员也不能越权去修改硬件

关系。例如，初级运维工程

这种基于角色的权限划分还需要考虑到运维团队内部的层级关

重启某些非关键服务等；

师可能只能执行一些常规的、风险较低的操作，如查看日志文件、

调整系统核心参数、执行

而高级运维工程师则拥有更高的权限，可以进行更复杂的操作，如

有效避免因权限过度集中

大规模的数据迁移任务等。通过这种多层次的角色权限体系，能够

而导致的安全隐患。

器。然而，如果是在非工作时间或者从企业外部网络尝试访问时，就需要触

大模型应用服务

验证机制。

发更加严格的验

型环境中，高危操作可能包括删除大量数据、修改核心模型参数、更改系统

在私域大模

旦识别出高危操作，就需要立即启动阻断机制。这种阻断机制可以分为多个

权限配置等。一

次是警告提示，当运维人员尝试执行高危操作时，系统会弹出明显的警告信

层次。在第一层

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层级的管理员密码、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，才能够真正地执行操作。第二层次且完全阻断。对于某些极其危险的操作，系统可以直接

禁止其操作。比如运维人员删除数据库数据时，系统可以强制

禁止其操作。比如运维人员删除数据库数据时，系统可以强制

在软域人访问系统中，令牌可以被视为一个特殊的区域控制器，其可以放行或为使用的

数据和应用程序，如管理的关键的源代码、高敏感性的网络数据、关键的系统配置文件等。

这些令牌通常由硬件令牌或一种计算机程序生成，令牌包含与令牌绑定的唯一标识符与

令牌绑定。令牌通常由硬件令牌或一种计算机程序生成，令牌包含与令牌绑定的唯一标识符与

令牌绑定。令牌通常由硬件令牌或一种计算机程序生成，令牌包含与令牌绑定的唯一标识符与

令牌绑定。令牌通常由硬件令牌或一种计算机程序生成，令牌包含与令牌绑定的唯一标识符与

令牌绑定。令牌通常由硬件令牌或一种计算机程序生成，令牌包含与令牌绑定的唯一标识符与

令牌绑定。令牌通常由硬件令牌或一种计算机程序生成，令牌包含与令牌绑定的唯一标识符与

令牌绑定。令牌通常由硬件令牌或一种计算机程序生成，令牌包含与令牌绑定的唯一标识符与

5.3 场景二：应用访问大模型应用场景



访问公域大模型、业务系统访问私域大模型两种

业务系统能够通过安全智能体访问公域、私域大模型、RAG 知识库。AI-R-IAM

提供以下安全措施：

业务身份管控

访问公域、私域大模型时身份安全极为重要，如果业务系统的身份信息被盗取

如密钥可能被硬编码在代码中，或者存储在不够安全的地方，攻击者可能通过伪造身份来访

问大模型服务，进而消耗资源或进行恶意操作。

AI-R-IAM 从业务身份创建、修改和销毁的全流程管控，实现业务身份可信。

身份创建：对业务系统身份注册与认证，通过数字证书或 API 密钥进行认证，并授予

最小访问权限。

能和任务变更情况，更新业务系统身份信息和权限，确保合规并遵

身份修改：需根据功

循最小权限原则。

系统身份信息并销毁凭证，同时进行关联资源清理和安全检查，确

身份销毁：删除业务

保系统的安全稳定运行。

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

分不清晰，可能导

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

致内部人员误操作或恶意操作，进而影响系统整体安全。

控制，保障大模型

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

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。

状态和业务流程，

实体级授权：通过 ABAC 和 PBAC 技术，业务系统根据实体的属性、

属性及数据敏感

角色级授权：采用 RBAC 技术，根据岗位和动作，为业务系统分配相应权限，确保高权限角色

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、

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，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

理、构造 API 请求参数、调用公域/私域大模型 API、模型服

权与权限校验、输入数据预处理

后处理、返回最终结果至用户等业务流程，存在日志记录不完

务处理、返回响应结果、结果后

私泄漏等安全问题。

整、合规性检查不足、数据隐

同公域、私域大模型时，通过全链路日志采集、智能化分析、

AI-R-IAM 在业务系统访问

行为可信。

智能追踪溯源三方面保障业务

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

为基线并检测异常。

历史数据进行分析，建立私域大模型、

和实体行为的动态变化，确保基线时

私域大模型、RAG 知识库安全风险

的实时监控。

私域大模型、RAG 知识库的数据输入到模型输出提供全链路智能追

踪、路径探寻算法，提供可疑 IP->访问数据分布->可疑路径->可疑日

志分析，实现风险智能追踪。

份与权限、特定模型等多维数据，以构建全面的用户行

智能化分析：通过机器学习和统计分析技术，对历史

RAG 知识库的用户可信行为基线，并通过自适应用户和

效性和准确性。同时采用画像分析、知识图谱等技术对

...

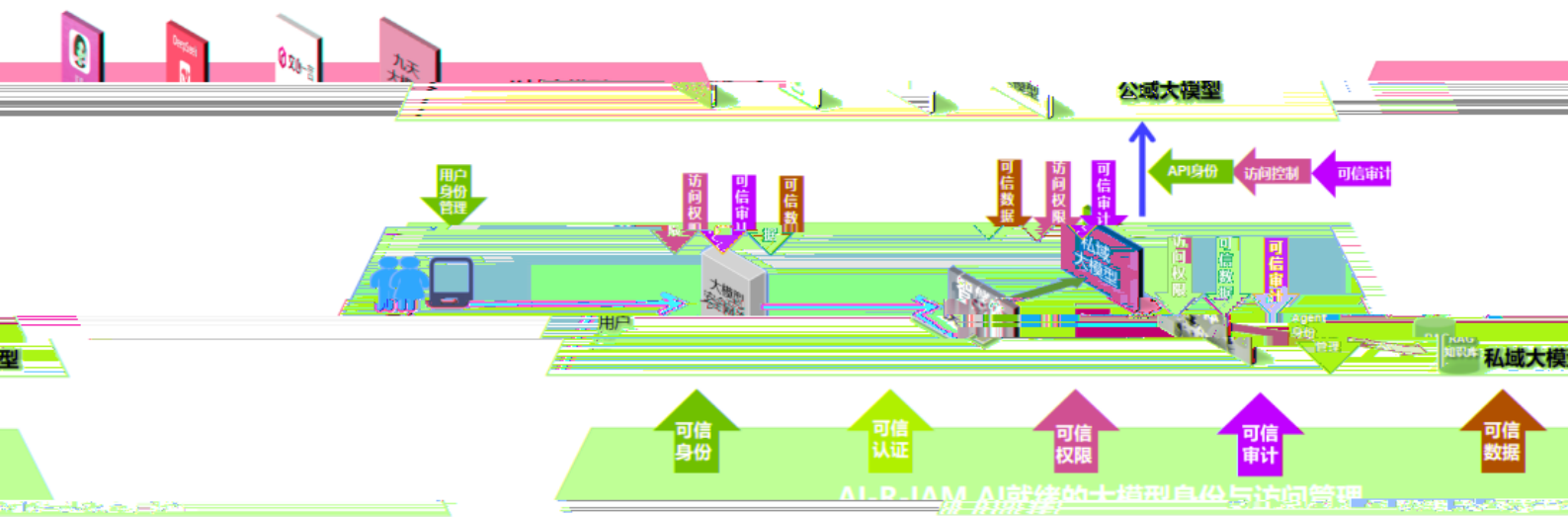
智能追踪溯源：对

追踪溯源，采用路径计算

志链路分析，层层递进

用户-AI 智能体的安全应用管理场景

5.4 场景三：月



AI-R-IAM 赋予用户与智能体身份并绑

，智能体以独立身份接受认证，调用时

，保障访问数据和训练数据的安全。通

场景概述：在用户到智能体的安全应用场景中，

定，防智能体被仿冒。用户经实名认证建立可信身份

获有限权限，同时通过身份绑定对双方行为审计监管

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

身份 (user) 与智能体身份 (AI Identity) 层面明确，如果智能体身份做了某项操作，系统可追溯到是谁让智能体进行操作，方便界定相应责任。

通过 AI-R-IAM 提供可信身份管理，赋予用户可信身份 (Identity)，由统一身份 ID 进行绑定管理，便能在审计时追溯到是谁让智能体进行操作，方便界定相应责任。

与特定的代码及模型紧密绑定，运用数字签名或证书技术，明确证明某个智能体是由谁开发的。

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及模型紧密绑定，运用数字签名或证书技术，明确证明某个智能体是由谁开发的。

名、证书或公钥进行严格校验，以此来精准确认该智能体的真实身份，同时评估其可信度、安全性和可靠性。

当 RAG 知识库接收到该智能体所发出的请求时，能够对其携带的数字签名、证书或公钥进行严格校验，以此来精准确认该智能体的真实身份，同时评估其可信度、安全性和可靠性。

2. 强认证与授权降低威胁面

M 实名认证能力建立实名认证规划最小访问权限，

在访问智能体场景下，普通用户通过大模型安全网关与 AI-R-IAM 实名认证，确保前端应用访问可信，同时管理员利用 AI-R-IAM 权限管理

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，通过受限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。

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，通过受限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。

设计与监管

3. 关联审计

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督。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

通过身份绑定和关联审计，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

并记录智能体的身份标识、访问了哪些数据、执行了何种操作，以便追溯和审计。系统应依据何种规则。

并记录智能体的身份标识、访问了哪些数据、执行了何种操作，以便追溯和审计。系统应依据何种规则。

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

在

新风险，如自动提交敏感数据到外部。

器，RAG 知识库，

在此场景下，AI-R-IAM 提供用户和智能体，智能体与数据存储服务

RAG 的数据访问权限

内外部大模型之间的数据权限边界，实行最小化策略，制定智能体对 RA

本异常行为等。避免

范围，辅以模型护栏、上下文检测等安全控制，可及时阻断和审计智能体

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中的深度应用，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。通过持续的实时感知，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，大模型可以识别异常行为。

2、在网络安全方面，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或管理系统进行结合，实现对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。即使身处网络安全防护体系。

借助 5G/6G 网络和实时的安全监测。通过持续的技术创新和大模型的广泛应用和数字经济的蓬勃发展，保障大模型在物联网场景下与海量设备交互时的高速低延迟特性，实现更高效的安全数据传输和融合，AI-R-IAM 将不断适应新的安全挑战，为数字未来保驾护航，构建更加安全、可信的数字未来。