





大模型生成可控的技术方案

大模型内置训练/推理可控

人类对齐

SFT、RLHF/RRHF、RLAIF

风险抑制

大模型外挂加固防御可控

知识增强

安全知识库生成和检索

- 千万级知识
- 生产有效率64%

风险认知

多模态感知终端

风险认知增强

疑难风险内容补充召回 88.2%

大模型

大模型提问护栏的系统方案

大模型业务场景中，基于问（人类知识）答（AI知识）环节不同的特点，新建大模型护栏解决方案。

用户提问理解

领域理解

话题理解

意图理解

观点、情感、创作

底线风险理解能力

提问风险决策

- L1 无风险
- L2 问题增强
- L3 检索增强
- L4 知识增强
- L5 安全拦截

大模型

Pass

Rewrite

Reject

问答风险决策模块

全风险覆盖

- 底线
- 伦理
- 数安
- 合规

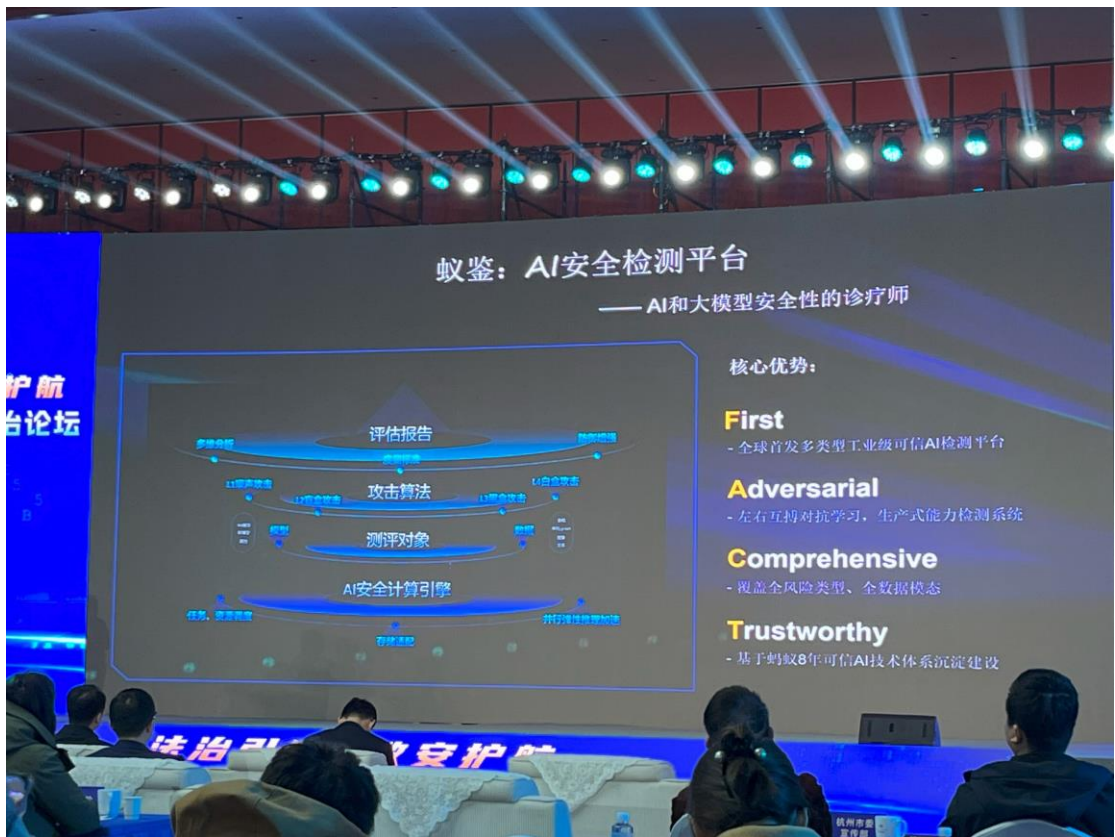
安全知识库模块

专项知识库

3种检索能力

安全限制模块

风险文案库





深度伪造数据生成和检测能力

输入模式

- 图像
- 文本
- 音频
- 视频

覆盖常见生成算法和生成方式50余种
(泛化性)

AIGC生成引擎



+

对抗攻击



融合空间域和频域特征的
层级分类算法

覆盖传播过程所有可能的损失方式30余种
(鲁棒性)

输出结果

- 图像**
 - 图像真实性分数和类型
 - 展示判断区域
 - 敏感人脸比对分析
- 文本**
 - 文本真实性分数和类型
 - 语言库和语言识别
- 音频**
 - 音频真实性分数和类型
 - 敏感人物音频比对分析
- 视频**
 - 视频真实性分数和类型
 - 展示判断区域
 - DEEPFAKE人脸分析

建设近千万级深度伪造数据，沉淀不同产生方式、不同生成方式、不同风格、不同生成/传播/扩散模型的AIGC生成能力Pipeline。

大模型安全的“快”与“慢”

大模型安全

“快”：去毒

“慢”：可控

数据去毒

安全护栏

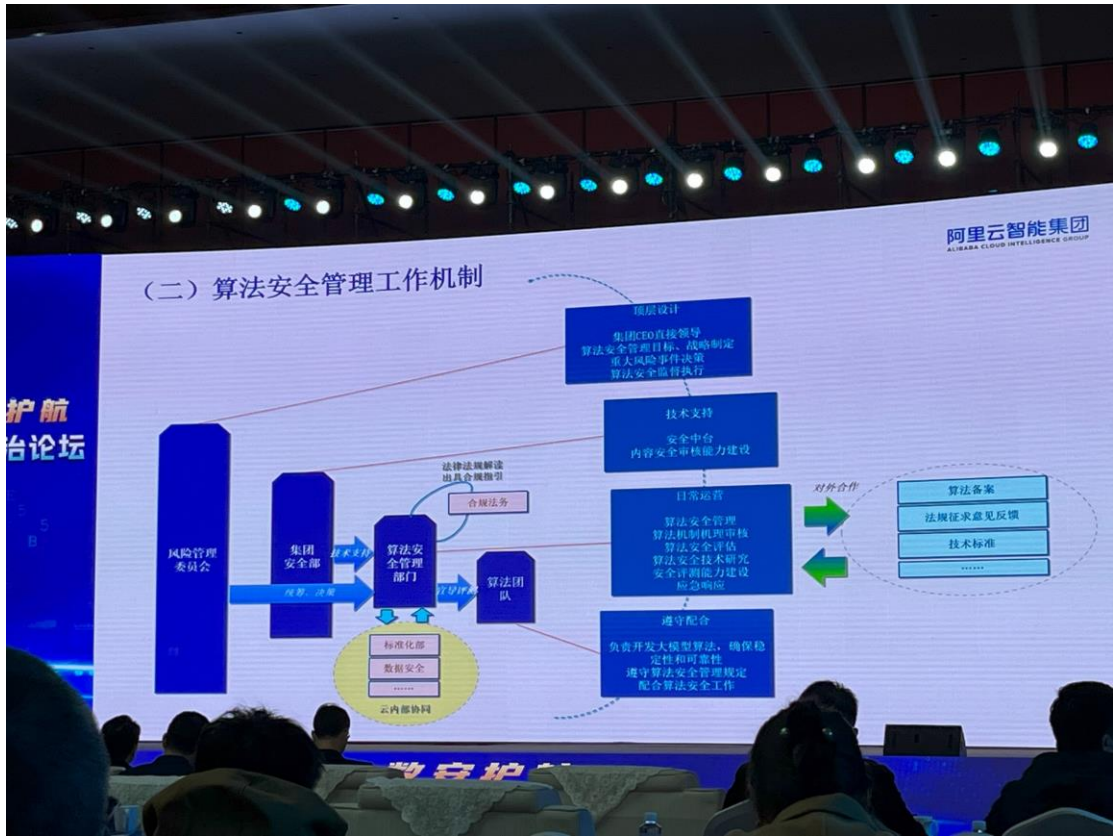
AIGC风险检测

安全测评

解构可控

人类社会共治

大模型安全防御



(三) 大模型安全能力

【安全评测能力】

评测方式：自动化评测、平台自动化评测、定制化评测

自动化评测：适用于模型发布前，快速检测模型安全性能，降低安全风险。

平台自动化评测：适用于模型发布后，实时监控模型安全性能，及时发现异常。

定制化评测：适用于特定场景，针对模型安全性能进行深度检测。

【内生安全能力】

内生安全能力：通过标准问答、风险识别、检索增强、内容标识，对于用户输入和模型生成内容提供外层安全防护能力。

评测流程：输入问题 -> 模型生成 -> 安全检测 -> 人工审核

内生安全能力：基于事实知识库的检索增强，通过事实知识库的构建，给大模型提供真实的资料，环节大模型的幻觉问题。

【外层护栏能力】

外层护栏能力：通过标准问答、风险识别、检索增强、内容标识，对于用户输入和模型生成内容提供外层安全防护能力。

安全设计：安全设计、安全设计、安全设计

安全设计：安全设计、安全设计、安全设计

【外层护栏能力】

外层护栏能力：通过标准问答、风险识别、检索增强、内容标识，对于用户输入和模型生成内容提供外层安全防护能力。

安全设计：安全设计、安全设计、安全设计

安全设计：安全设计、安全设计、安全设计

(一) 算法安全管理制度——外部：国内重点法规及技术标准

阿里云智能集团 ALIBABA CLOUD INTELLIGENCE GROUP

网信办技术局颁布《互联网信息服务算法推荐管理规定》(2022.3.1)

算法通用要求：生成合成类、个性化推荐类、检索过滤类、排序精选类、调度决策类

直接落实：网信国标《机器学习算法安全评估规范》

细化：网信办技术局颁布《互联网信息服务深度合成管理规定》(2023.1.10)

对5大类算法中生成合成类算法进行了进一步分类定义和要求

安全评估要求、支持备案义务、个人信息保护、输入输出校验、合理显著标识、防篡改、删除

直接落实：网信国标《互联网信息服务深度合成安全规范》

参考落实：工信《人工智能预训练模型通用要求及服务能力成熟度评估》

细化：网信办网安局颁布《生成式人工智能服务管理暂行办法》(2023.8.15)

对算法、数据、内容安全、价值观等方面提出要求

数据合规、模型安全、内容安全、个人信息保护、伦理道德、安全评测、安全防护

参考落实：工信《互联网深度合成信息服务标识通用安全要求》

直接落实：网信国标《生成式人工智能预训练和优化训练数据安全规范》

直接落实：网信国标《生成式人工智能标注安全规范》

参考落实：公安行标《互联网交互式服务安全管理要求 第19部分：生成式人工智能服务》



阿里云智能集团
ALIBABA CLOUD INTELLIGENCE GROUP

(二) 算法安全监测机制——数据安全

打造以数据为核心、围绕数据生命周期的数据安全管理体系

1. 数据管理

建立了训练数据的统一管理平台，对训练数据实施分级分类管理、分级权限管控、脱敏处理等安全策略。

01 根据数据内容结合公司数据安全合规管理要求，公司将数据分为公开数据、内部数据、保密数据、机密数据等4个级别。

02 对可能含敏感信息的数据，处理、使用有明确规定，如：必须根据数据等级进行级别界定和标识；必须有使用授权，不允许超授权范围的使用；涉及个人属性信息的数据，进行去标识化、字符替代。

03 按照数据分类分级要求，对数据进行标识，配合血缘追溯等功能，实现权限管理、行为审计、风险监测等功能。

数据出境审批 数据披露

阿里云智能集团
ALIBABA CLOUD INTELLIGENCE GROUP

(二) 算法安全监测机制——个人信息

从数据源头、训练过程、生成过程，全链路落实个人信息保护相关规定

数据源头

- 数据源头选择：事先得用户授权，范围合法合规，服务目的明确；
- 数据源头过滤：涉政、涉个人信息等过滤或评估风险数据源。

训练数据

- 针对不同应用场景的合规指引开展影响评估，充分保障用户的权利，确保加工和使用逻辑合法合规。对训练数据进行清洗，去除训练数据中的标签噪声或者被人恶意投毒污染带来的影响。

存储方案

- 存储方案安全可靠，对敏感数据进行加密存储和传输。

用户注册准入

注册及准入

- 账号：复用云账号体系
- 实名认证

用户协议&隐私政策

《用户协议》和《隐私政策》是保障用户合法权益的重要文件，也是企业履行个人信息保护义务的重要体现。我们将持续完善用户协议和隐私政策，确保其合法、合规、透明、易懂，切实保障用户的个人信息安全。

《互联网信息服务深度合成管理规定》第九条：深度合成服务提供者应当基于移动电话号码、身份证件号码、统一社会信用代码或者国家网络身份认证公共服务等方式，依法对深度合成服务使用者进行真实身份信息认证，不得向未进行真实身份信息认证的深度合成服务使用者提供信息发布服务。



闪捷信息 SECSMART

人工智能侵犯数据隐私的途径

- 监控和追踪** 例如，摄像头、人脸识别等设施可能会被滥用。
- 数据收集分析** 例如，互联网公司收集大量个人数据以便于定位并推广产品。
- 机器学习预测** AI中深度学习算法需要大量数据以训练自己的网络，从而更精确地预测和分析人类行为。
- 数据泄露** 人工智能提供方便发生数据泄露。

人工智能应用面临7大数据安全威胁

- ✓ 威胁1: 模型中毒
- ✓ 威胁2: 隐私泄露
- ✓ 威胁3: 数据篡改
- ✓ 威胁4: 内部威胁
- ✓ 威胁5: 针对性蓄意攻击
- ✓ 威胁6: 大规模采用
- ✓ 威胁7: AI驱动的攻击

模型中毒: 向模型中注入恶意数据, 进而导致模型对数据进行错误分类并做出错误的决策, 例如: 指鹿为马。

航论坛

闪捷信息 SECSM

AI提升数据安全防护能力——用魔法打败魔法

数据资产AI识别模型

- 数据资产梳理
- 数据权属监测
- 敏感数据识别
- 数据流监测
- 数据分类分级

数据风险AI态势模型

- 敏感数据态势
- 数据风险监测
- 访问行为态势
- 数据合规风险
- 数据风险态势
- 数据安全风险

数安中心

数据防护AI动态策略

- 安全能力智能联动
- 动态策略矩阵
- 数据访问监测
- 自动运行监测
- 安全策略监测

数据安全自动化运营模型

- 设备运行监管
- 数据安全基线
- 业务运行监管
- 应急响应预案
- 数据监测中心
- 报表与工单

隐私

商业秘密

航论坛