

《AI谣言深度研究报告》

本报告由清华大学博士后张诗瑶与AI辅助生成



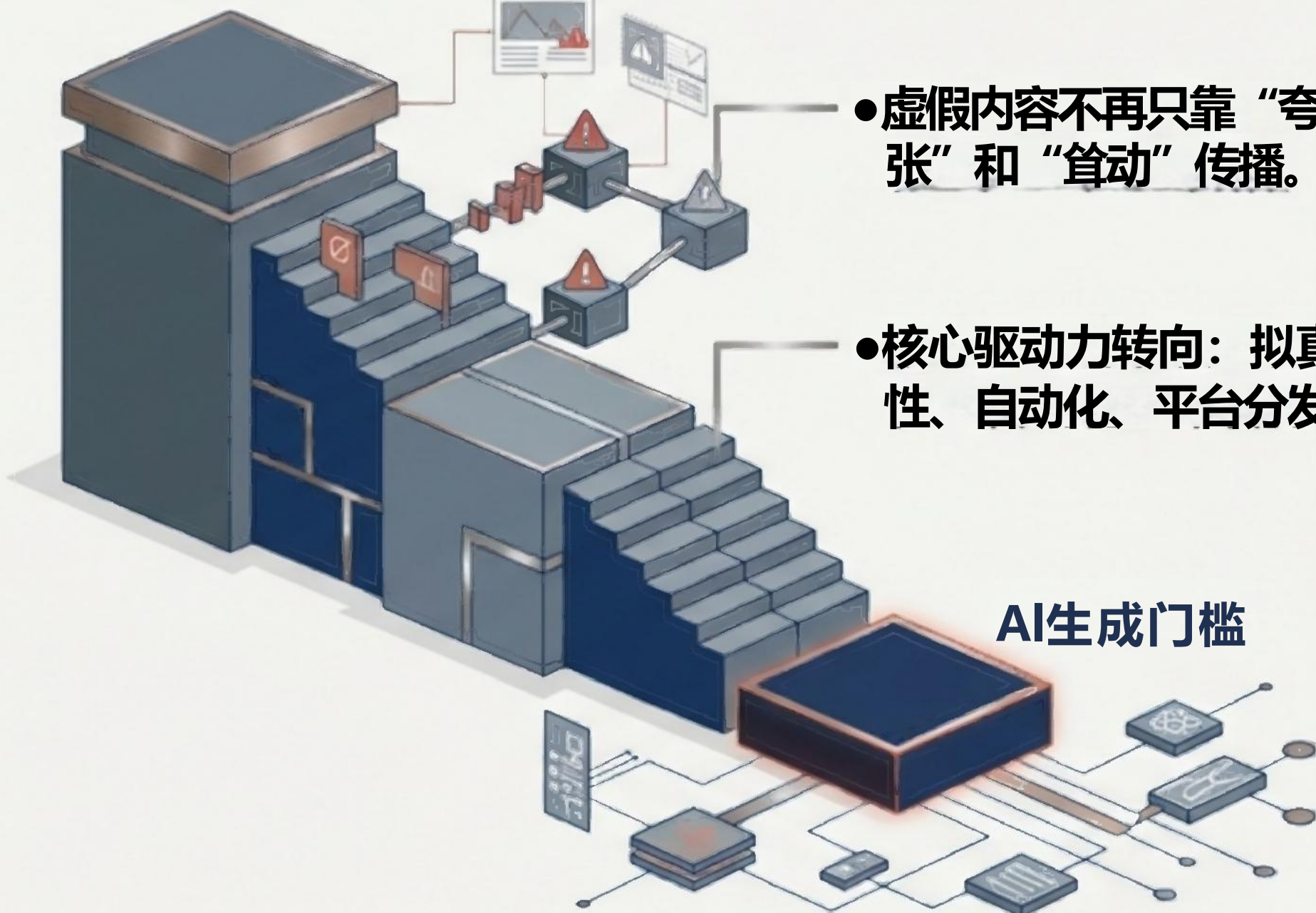
本报告聚焦2020-2026年间AI生成、深度伪造与算法放大共同驱动虚假或误导性信息现象，依据官方法规、国际机构文件、同行评审论文与可核查案例重构分析框架。

生成式AI的普及正在重构虚假信息威胁的威胁基线

**核心短期风险：
Misinformation/Disinformation**

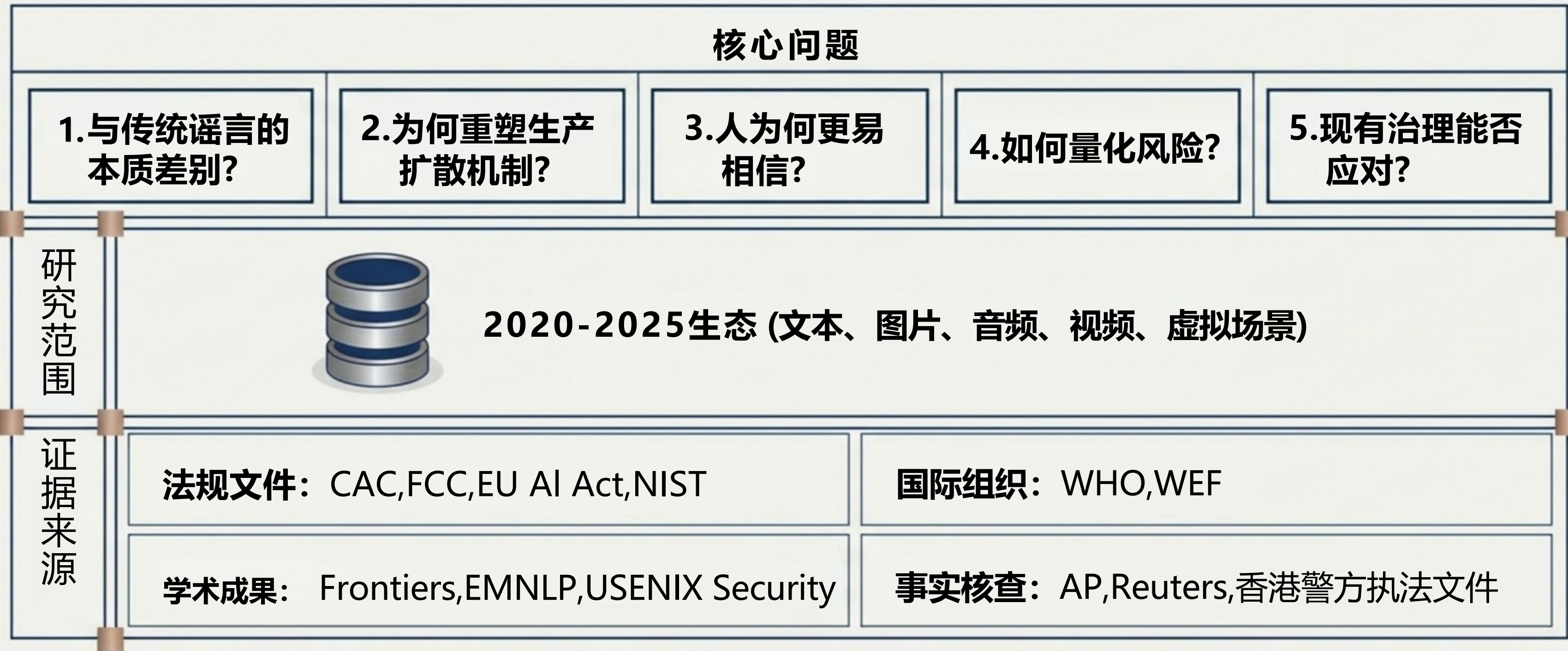
WORLD
ECONOMIC
FORUM

传统伪造门槛



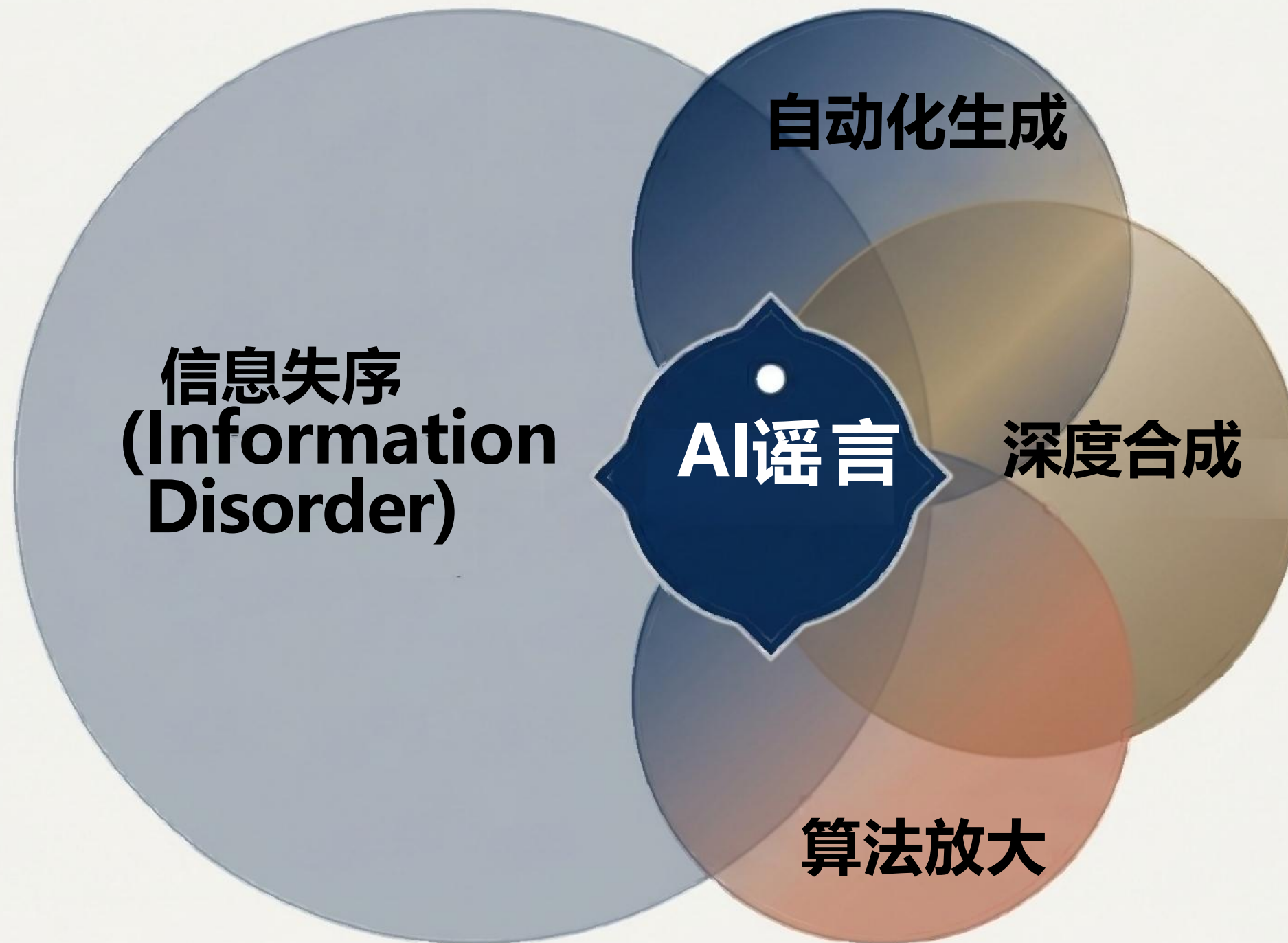
核心洞察： AI谣言的变化，不只是“更容易造假”，而是“更容易被信、更容易被放大、更难被及时纠正”。

基于全球监管与顶尖学术的循证研究框架



排除无法核验的行业传闻, 仅采用官方披露、执法结论与同行评审论文。

演进与突变：从“信息失序”迈向自动化误导



AI 谣言并没有脱离信息失序的范畴，但它将“虚假内容”与三大底层技术能力进行了深度耦合。

AI 谣言不仅包含无意误导，更强调带有目的性的自动化操纵意图。

AI谣言的操作性定义与三大评估支柱

在内容虚假或高度误导的前提下，人工智能对其生成、伪造、包装或传播放大具有关键贡献的信息事件。

内容侧

存在AI生成、深度合成或模型辅助伪造痕迹。

传播侧

存在机器人、协同行为、推荐放大或跨平台搬运。

治理侧

已触发事实核查、官方辟谣、执法或平台处置。

核心洞察：AI 谣言不是单纯的内容问题，而是包括“内容—模型—平台—用户—证据链”的复合治理对象。

范式转移：传统谣言与AI谣言的核心特征对比

	传统谣言	AI谣言
生产模式	人工编造、拼贴、转述	生成模型、深度合成、风格模仿
内容特征	表达粗糙、逻辑断裂、低画质	语义流畅、画面拟真、跨模态一致
传播依赖	依赖情绪煽动	依赖自动化分发与推荐算法

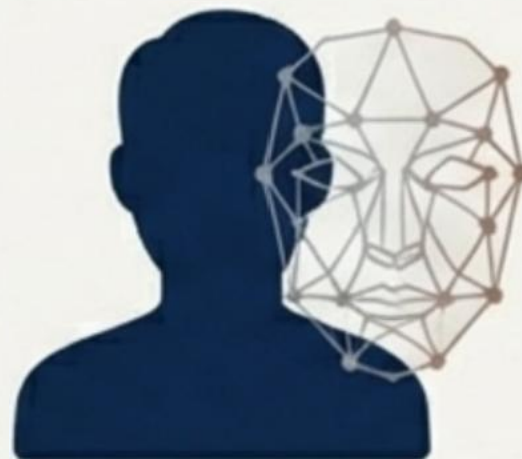
核心洞察： AI 谣言最大的变化不是“假得更多”，而是“假得更像真”。

AI谣言生态下的三种基础表现形态



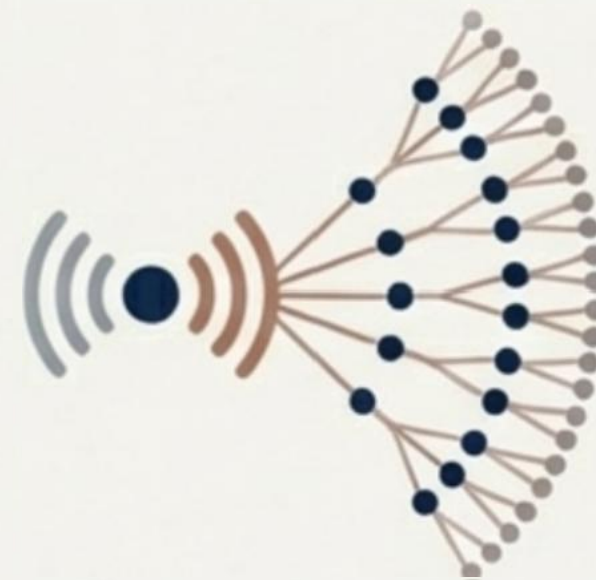
AI直接生成型

文本、图片、视频、音频由模型直接无中生有地生成。



AI辅助伪造型

以真实人物、真实事件、素材为锚点，进行换脸、变声、重组。

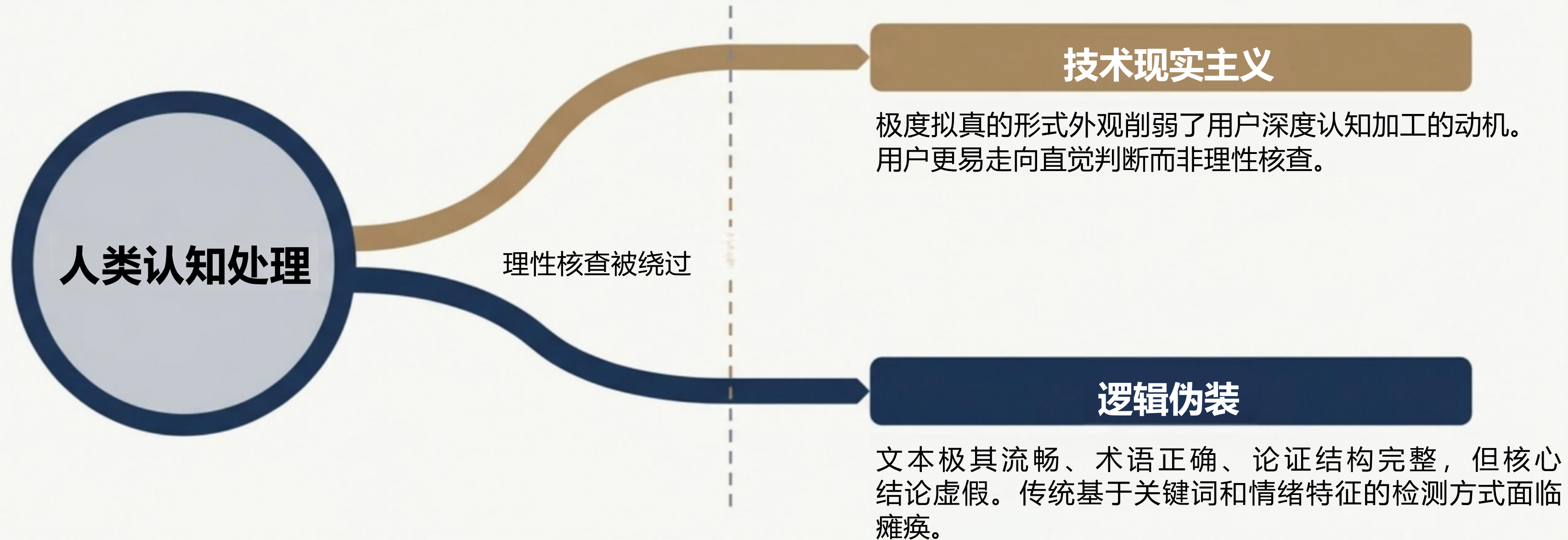


AI放大型

内容未必全由AI生成，但依靠AI账号、群控、推荐系统被快速模板化放大。

监管视角必须从单纯盯防“生成工具”扩展到涵盖“编辑伪造”与“算法放大”全链路。

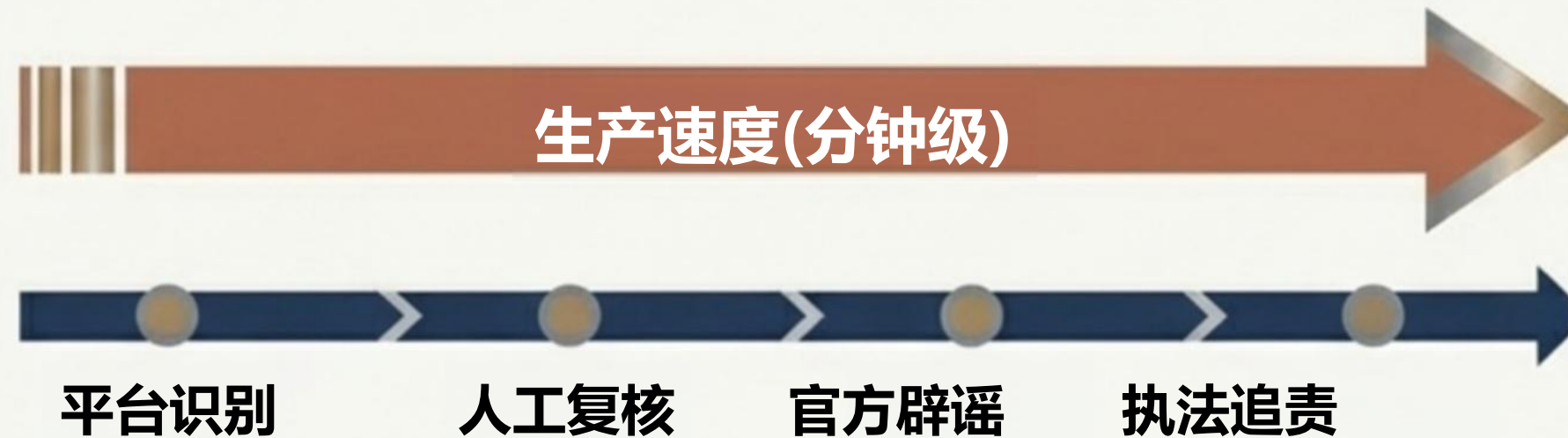
认知陷阱：为什么人类的理性核查正在失效？



核心洞察： AI谣言首先是一种“逼真性治理难题”，它让“看起来足够真实”取代了“事实上的真实”。

结构性不对称：现有治理体系的盲区

治理时差



首轮传播后，纠错只能是“补救”而非“阻断”。

规制对象错位



传统目标：
单条帖子/单个账号



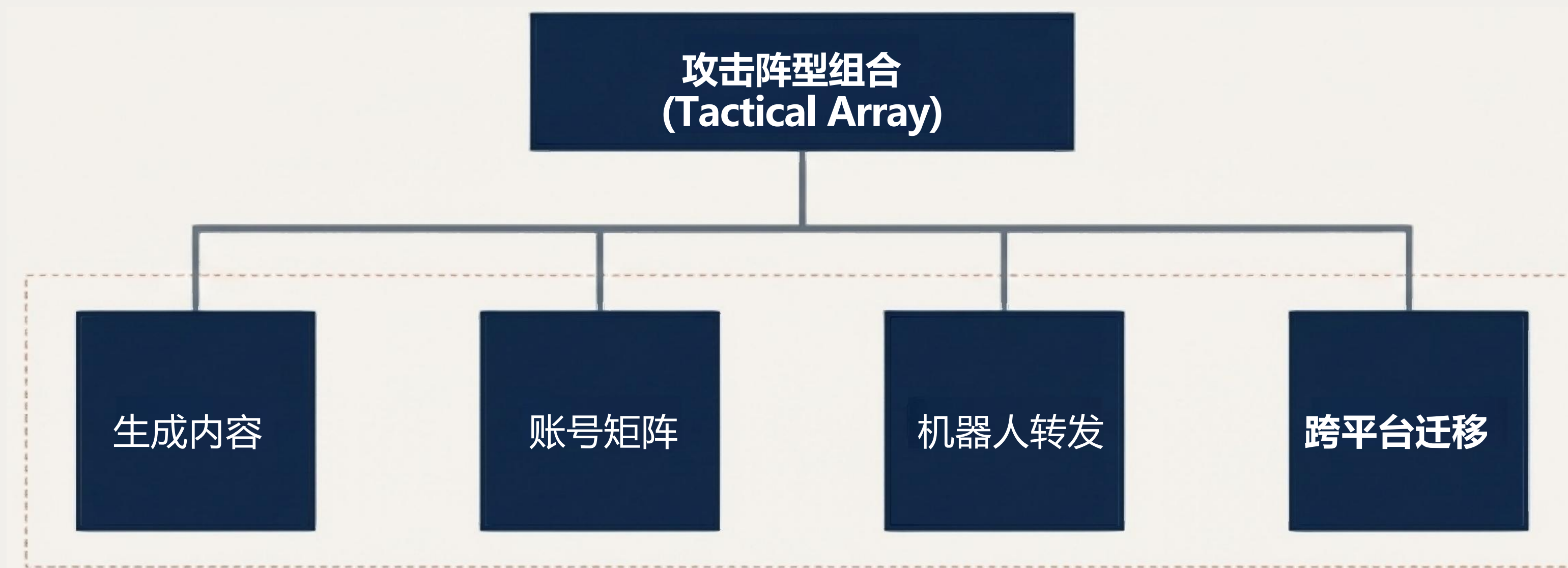
真正风险源：
模型能力、批量生产链与知识库投毒

传统治理聚焦末端个体，而现代威胁源自系统级的生成设施。

核心洞察： 应对系统性生产能力，治理必须从“末端删除负面内容”转向“前端约束风险基础设施”。

战术演进：从“单点造谣”到“攻击阵列化”

AI谣言已经具备类似网络攻击中的战术、技术与过程 (TTP) 特征。



核心洞察：我们面对的不再是单一的虚假信息，而是一场高度协同的舆论操纵攻击。

拆解威胁生命周期：AI谣言传播杀伤链

01 生成

文本、图像、
音视频的自动
化合成

02 包装

注入“权威
感”“现场感”
与“证据感”

03 播种

通过小号、
边缘群聊、
评论区先发
测试。

04 放大

触发推荐算法
与机器人矩阵，
进入热点机
制。

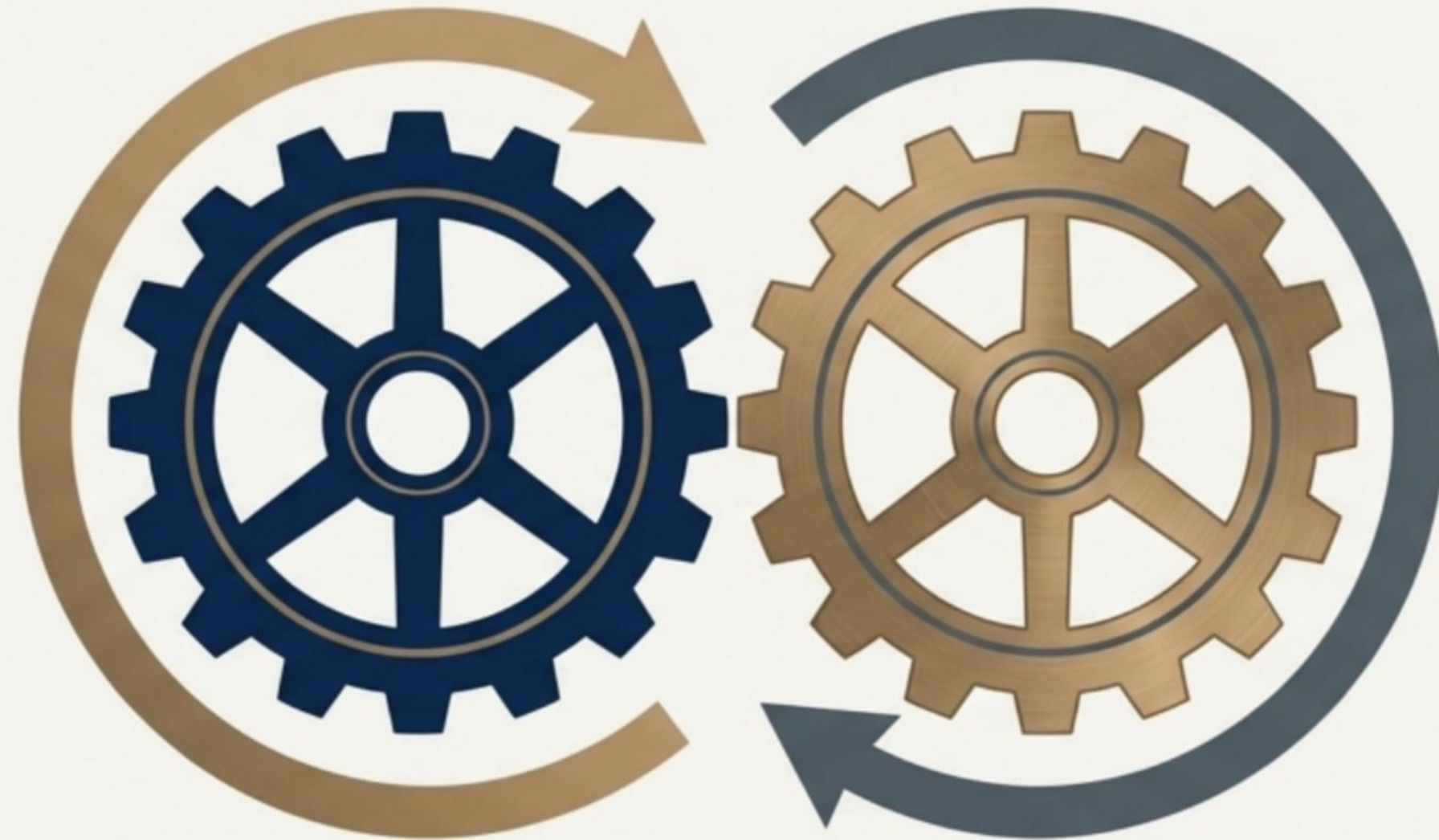
05 迁移

多平台搬运、
截图再传播、
二次剪辑脱离
原始语境。

传播加速器：算法嗜好与社交机器人矩阵

平台算法

推荐系统天然倾向于高新奇度、高冲突度、高情绪强度的内容。AI谣言被定制化生成以突破“可见性阈值”。



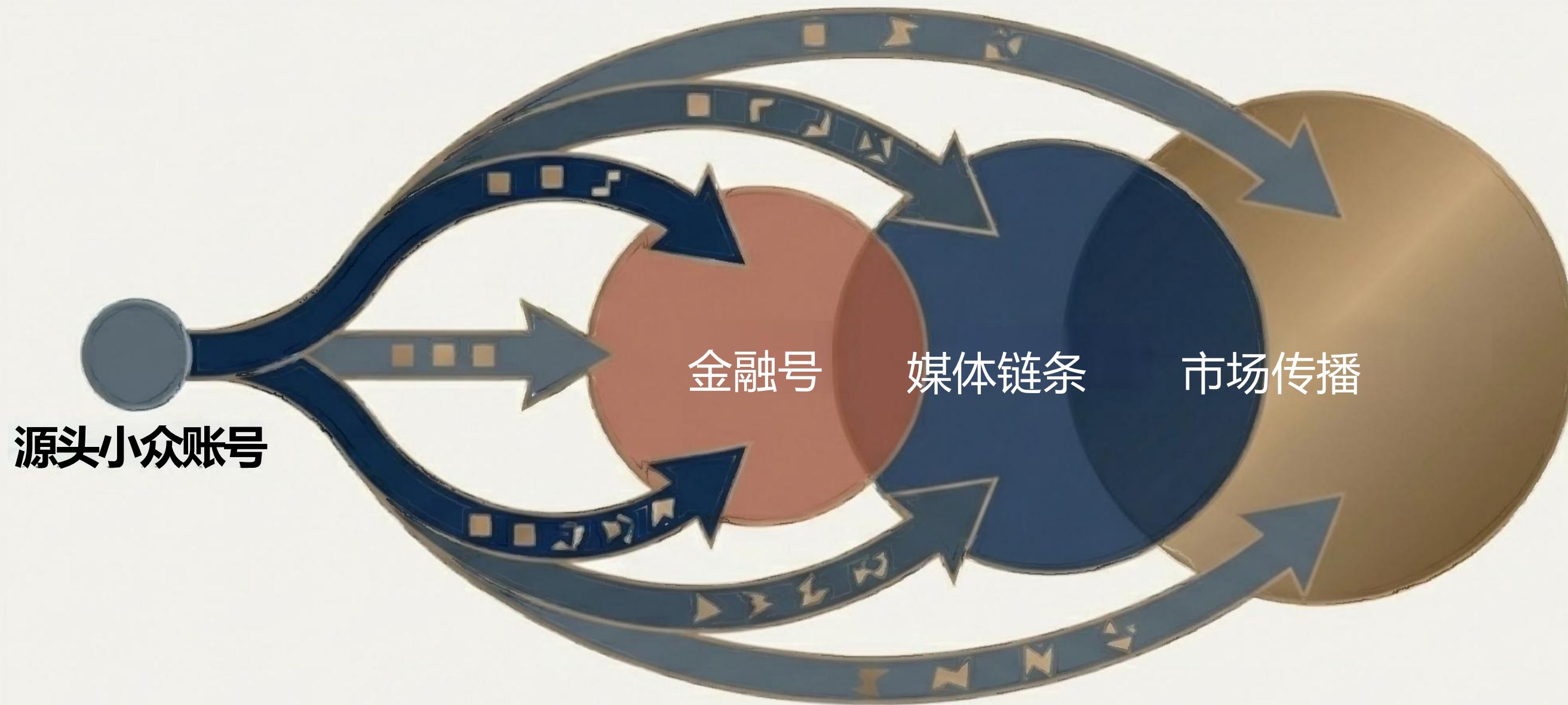
社交机器人

作为加速器自动发布、转发、维持议题热度。在政治、健康争议议题上深度参与误导信息扩散，构筑极化的“回音室”。

核心洞察：平台的分发机制与机器人的自动化运作，共同构成了AI谣言风险的系统性底座。

跨平台迁移：语境坍塌与次生变异风险

危险往往发生在跨平台迁移之后。初始假内容被截图、配文、重新剪辑，彻底脱离原始语境，极大扩展误读空间。



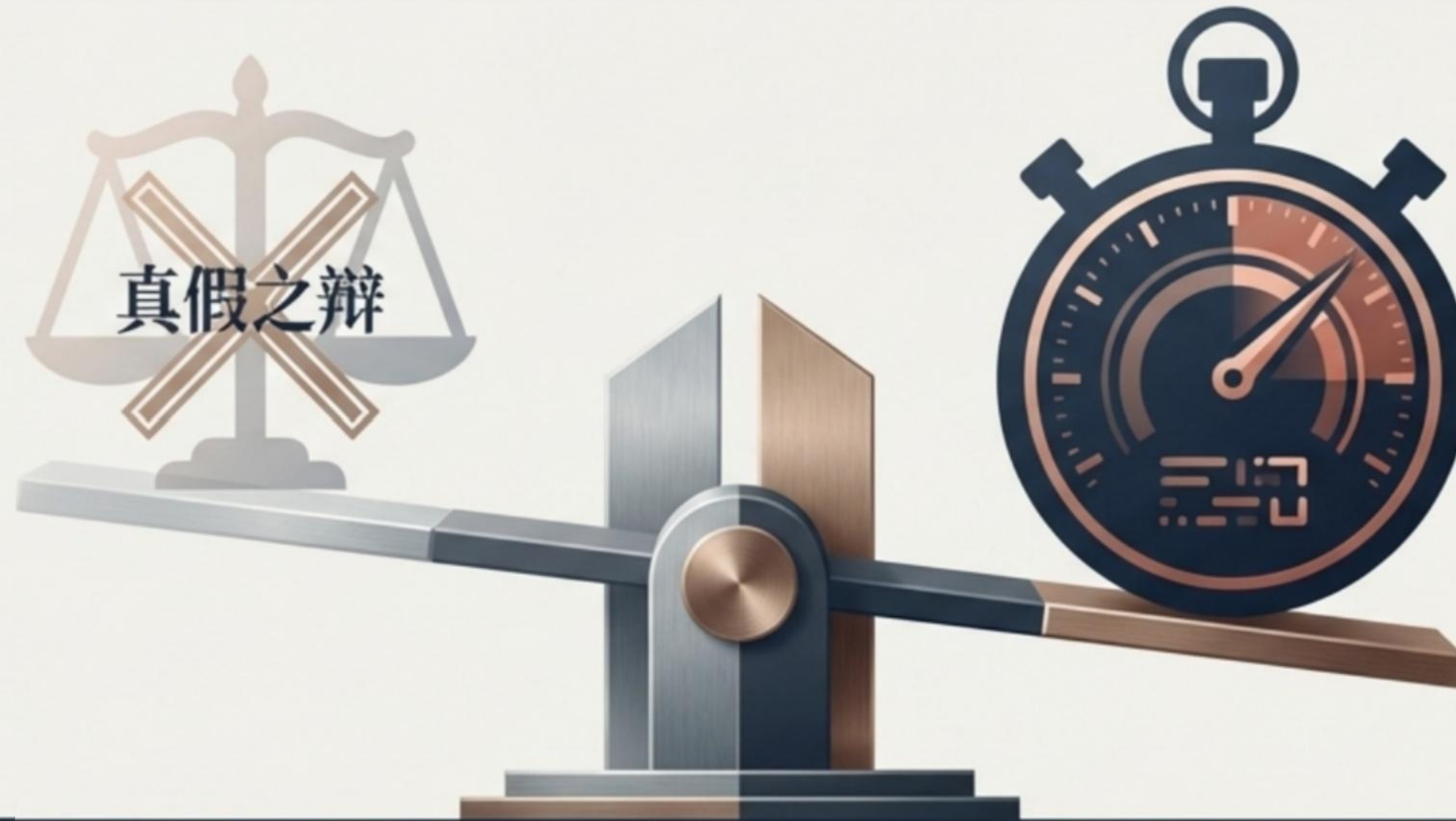
案例支撑：五角大楼
虚假爆炸图事件
从单一小众账号起步，经
由多个金融号、媒体链条
跨平台迅速扩散并引发市
场波动，随后才被官方阻
断并澄清。

每一次跨平台搬运，都是对溯源链条的一次切断与防御机制的绕过。

治理重构：从“真假之辩”走向“速度之战”

过去：

传统谣言治理是“谁说得更具煽动性”的真假博弈。



现在：

AI谣言治理是“谁先进入/阻断算法通道”的速度竞争。

治理的重心必须发生根本性前移——从依赖末端的“事后删帖辟谣”，全面转向前端的“早期识别拦截”与“扩散链路阻断”。

核心洞察：面对分钟级生成的AI谣言基础设施，唯有构建毫秒级响应的系统性防御，方能夺回信息生态的主动权。

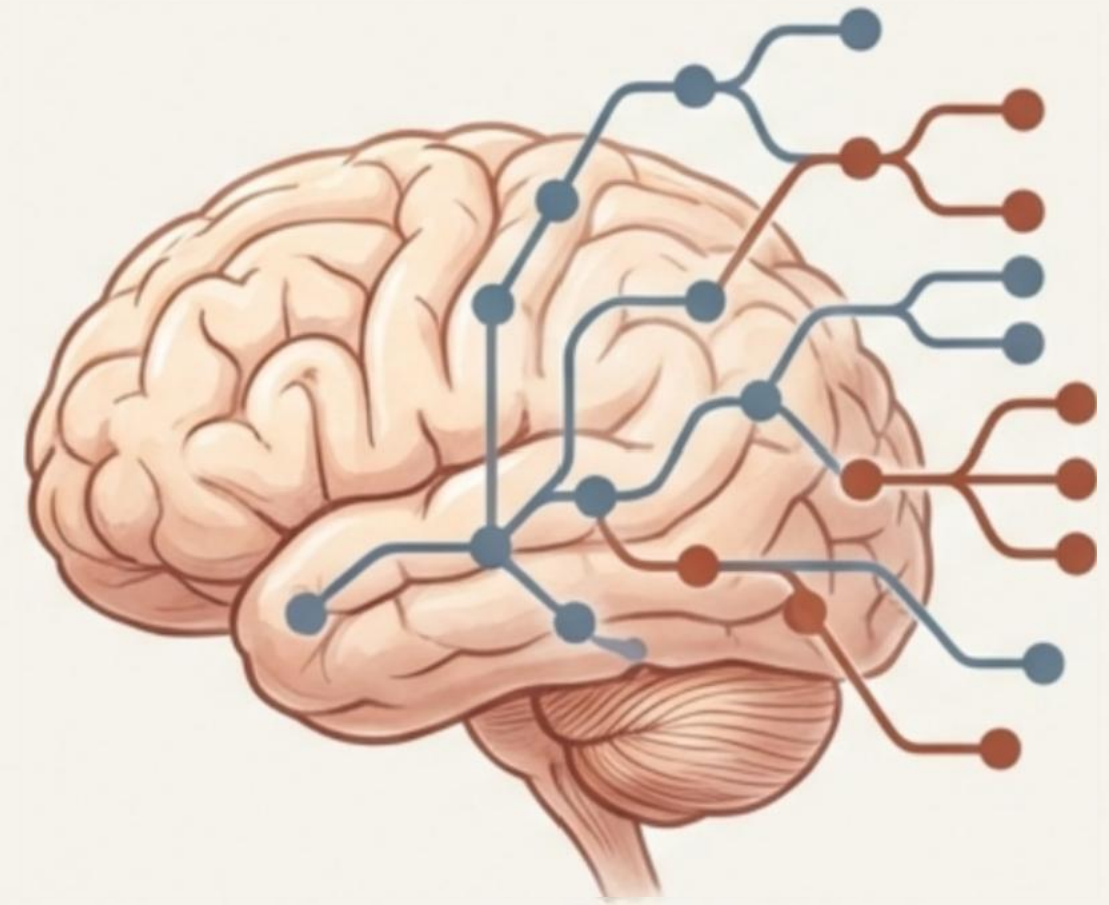
认知资源劫持： AI 谣言的心理学机制



AI谣言不仅在欺骗眼睛，更在劫持我们的判断机制

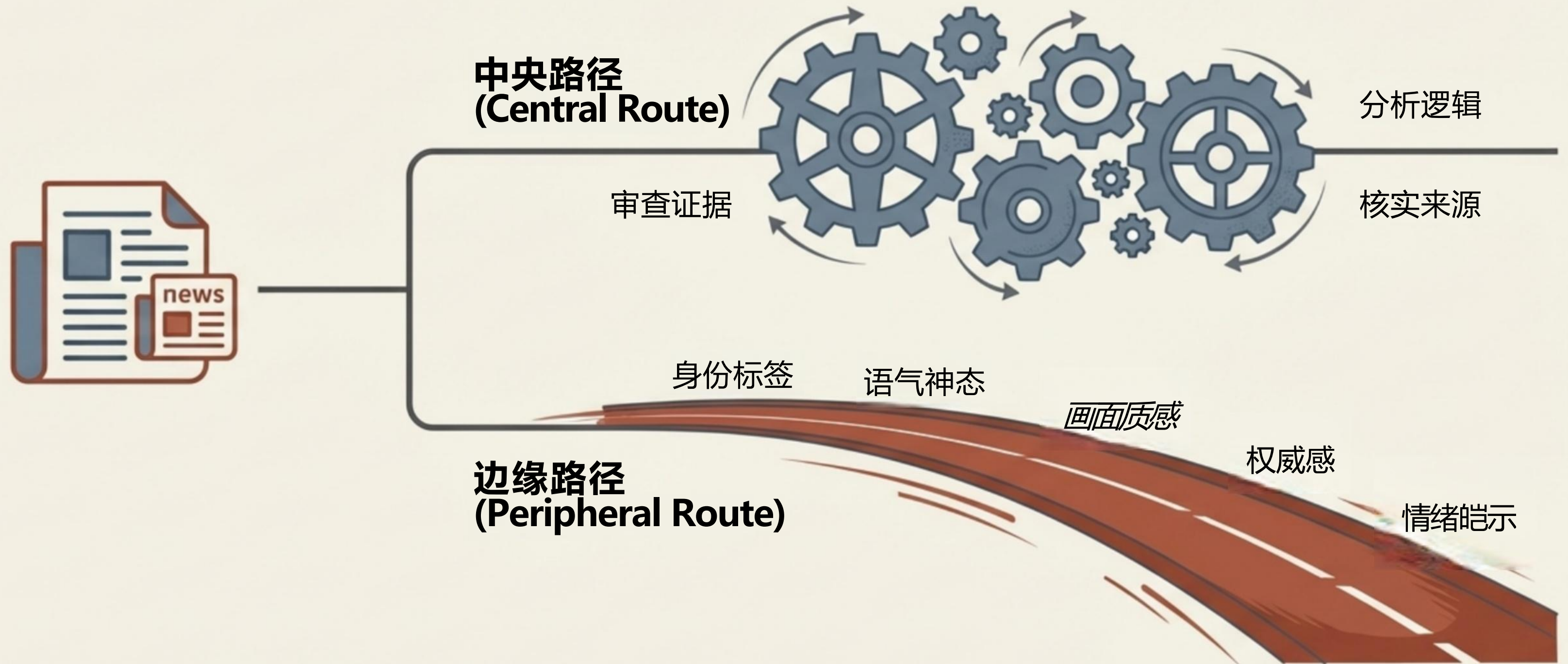


AI生成内容的“技术现实主义”会直接降低人类大脑的深度加工动机。

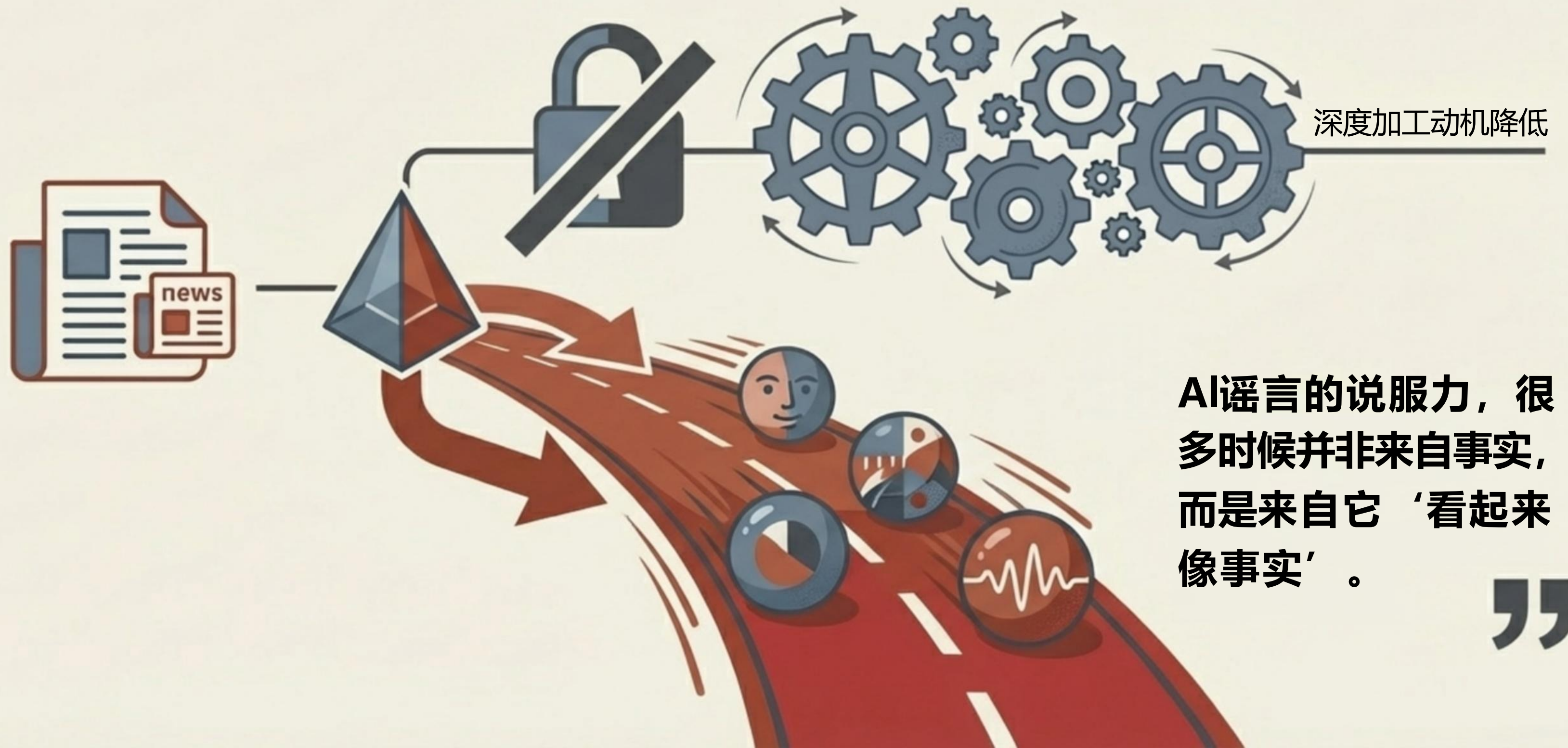


我们通常认为，AI造假的可怕之处在于“太逼真以至于看不出破绽”。但实际上，人类并不会对接收到的每一条信息都进行严密的事实核查。在信息过载的社交媒体环境中，大脑高度依赖**启发式线索**来快速判断可信度。

信息处理的双轨制：我们如何决定是否相信



完美模拟“可信线索”，将大脑推向边缘路径



AI谣言的说服力，很多时候并非来自事实，而是来自它‘看起来像事实’。

”

深度伪造正在摧毁“眼见为实”的底层逻辑

长期以来，图像和视频一直被公众乃至司法体系视为最高级别的“高可信证据”。

深度伪造 (Deepfakes) 的泛滥，直接削弱了人类“眼见为实、耳听为实”的直觉本能。



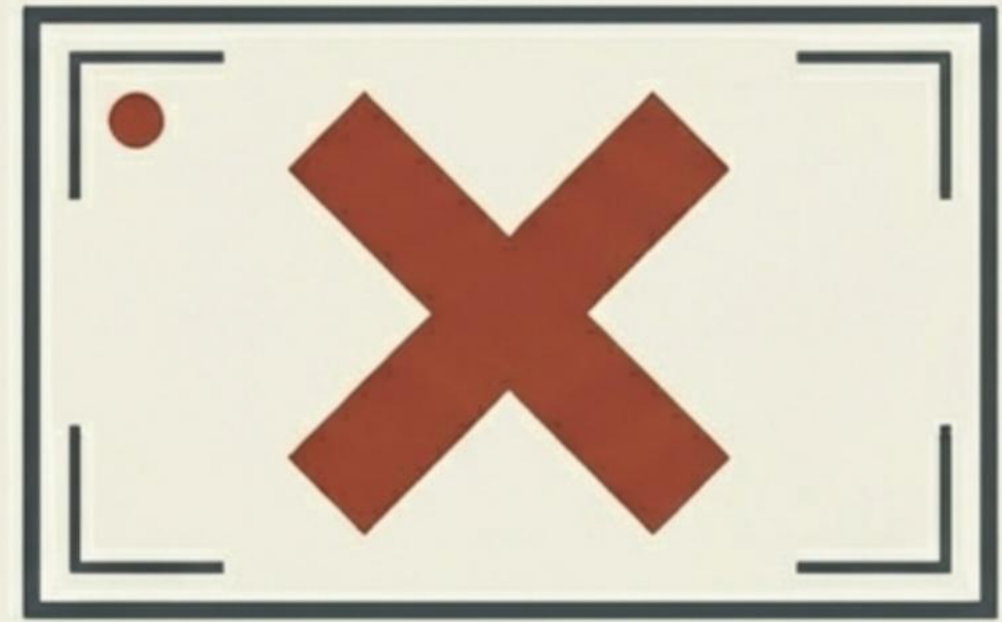
**认识论信任
(Epistemic Trust) 受损**

引发的不只是个体被骗的风险，
而是全社会对真实视听证据的普遍怀疑。

终极悖论：造假技术的存在，反而保护了真正的说谎者



虚假视频被公众误认为真



真实视频被当事人反向指控为AI伪造

学术界将此称为“**说谎者红利**” (Liar's Dividend)。当全社会普遍意识到“视频可以轻易造假”时，确凿的真实证据也失去了原有的约束力。

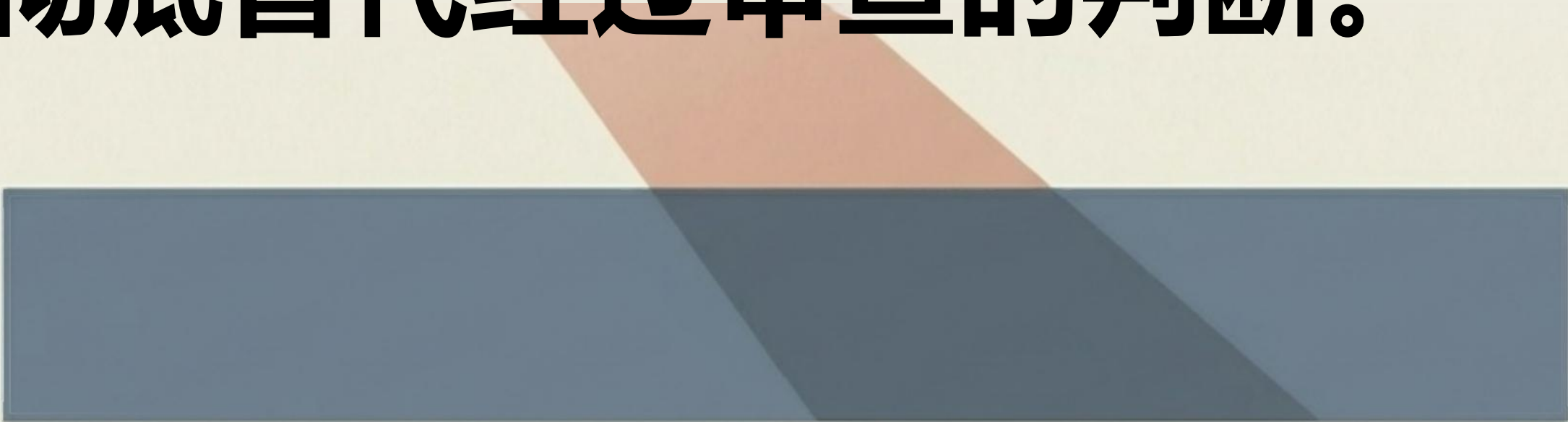
AI谣言的本质：一场针对认知资源的劫持



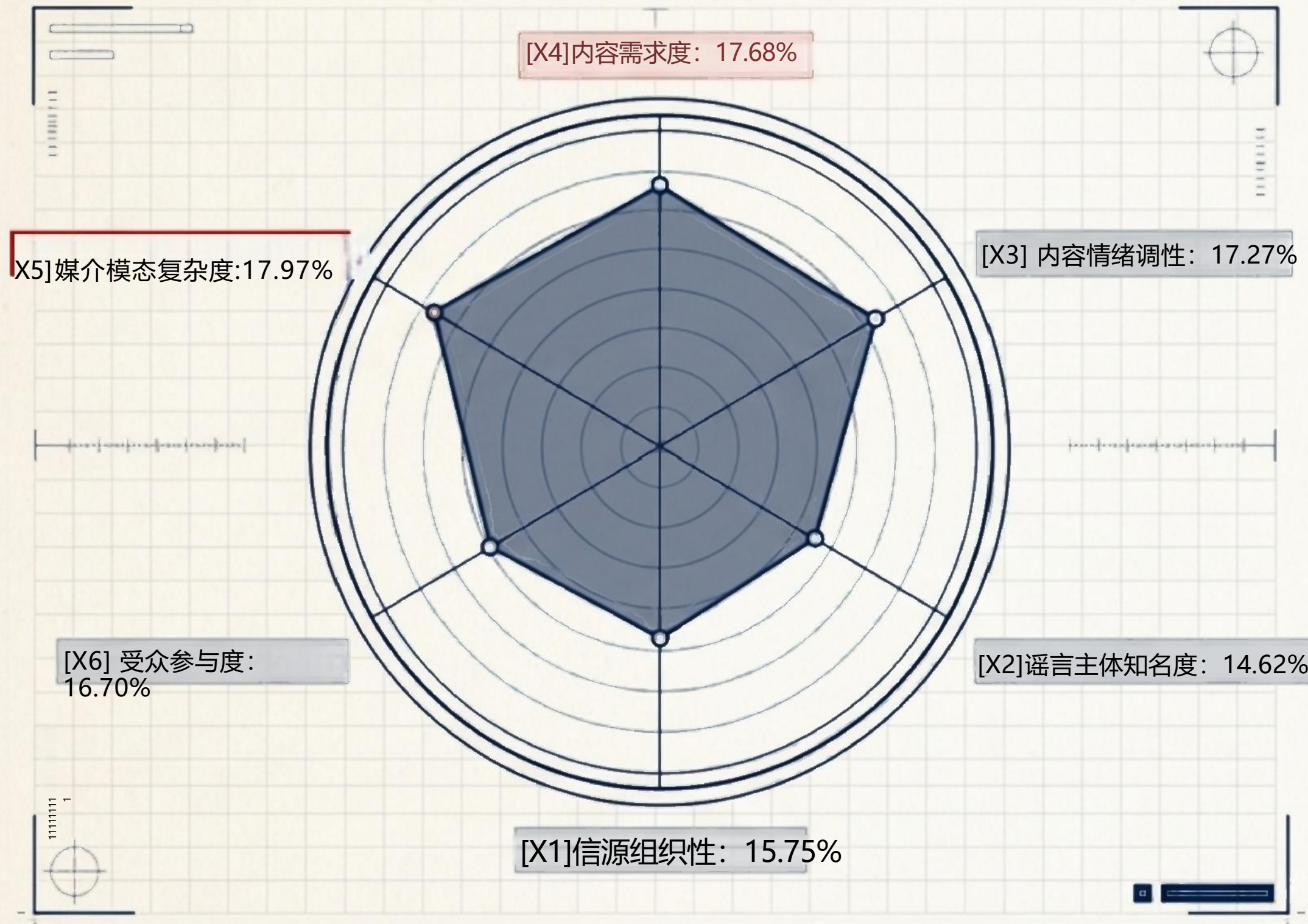
AI并非仅仅是在提供虚假信息，它是一套极其高效的注意力黑客系统。它利用完美的边缘线索，直接使我们的认知防御机制瘫痪。

当直觉替代了判断

**认知资源劫持的最终结果，
是让未经验证的直觉，
彻底替代经过审查的判断。**



提取致病基因：AI谣言的6大传播因子

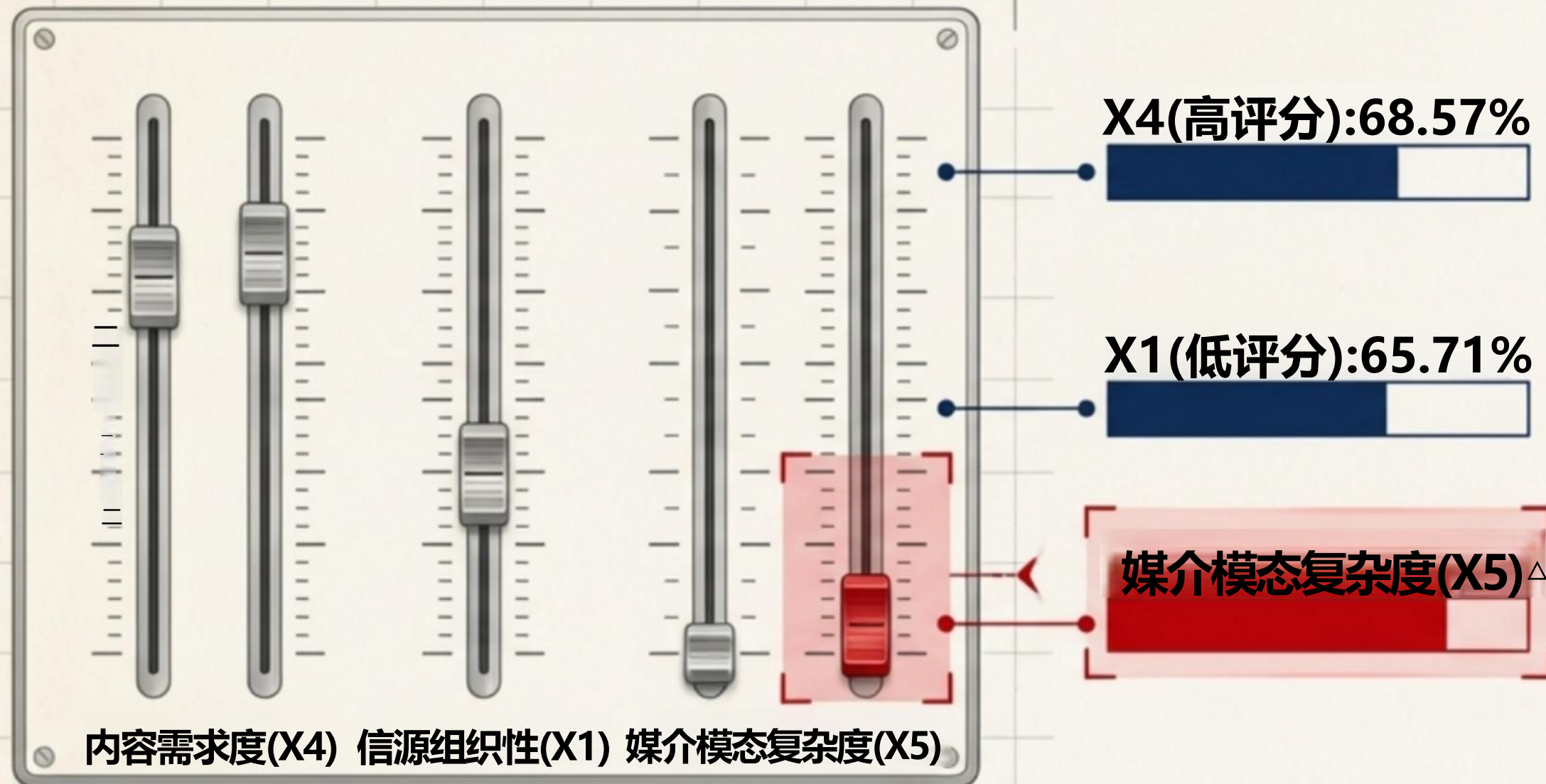


核心发现

影响排序：
内容 > 媒介 >
受众 > 信源

AI时代，人们不再依赖信源权威度，而是被直观感官所支配。

复杂性悖论：越精良，越难以引爆？



异常数据警告

独立样本t检验显示，媒介模式复杂度(X5)与传播热度呈现显著负相关(-0.2169)。

洞察

AI谣言的终极武器不是技术复杂性，而是心理简单性。高复杂度的多模态信息反而增加用户的认知处理负担。最致命的数字病原体，在技术上刚刚好能骗过算法，但在心理上能瞬间击穿防御。

威胁建模：两种典型的大众化感染模式

Suspect Profile

模式一：草根创作的强情感

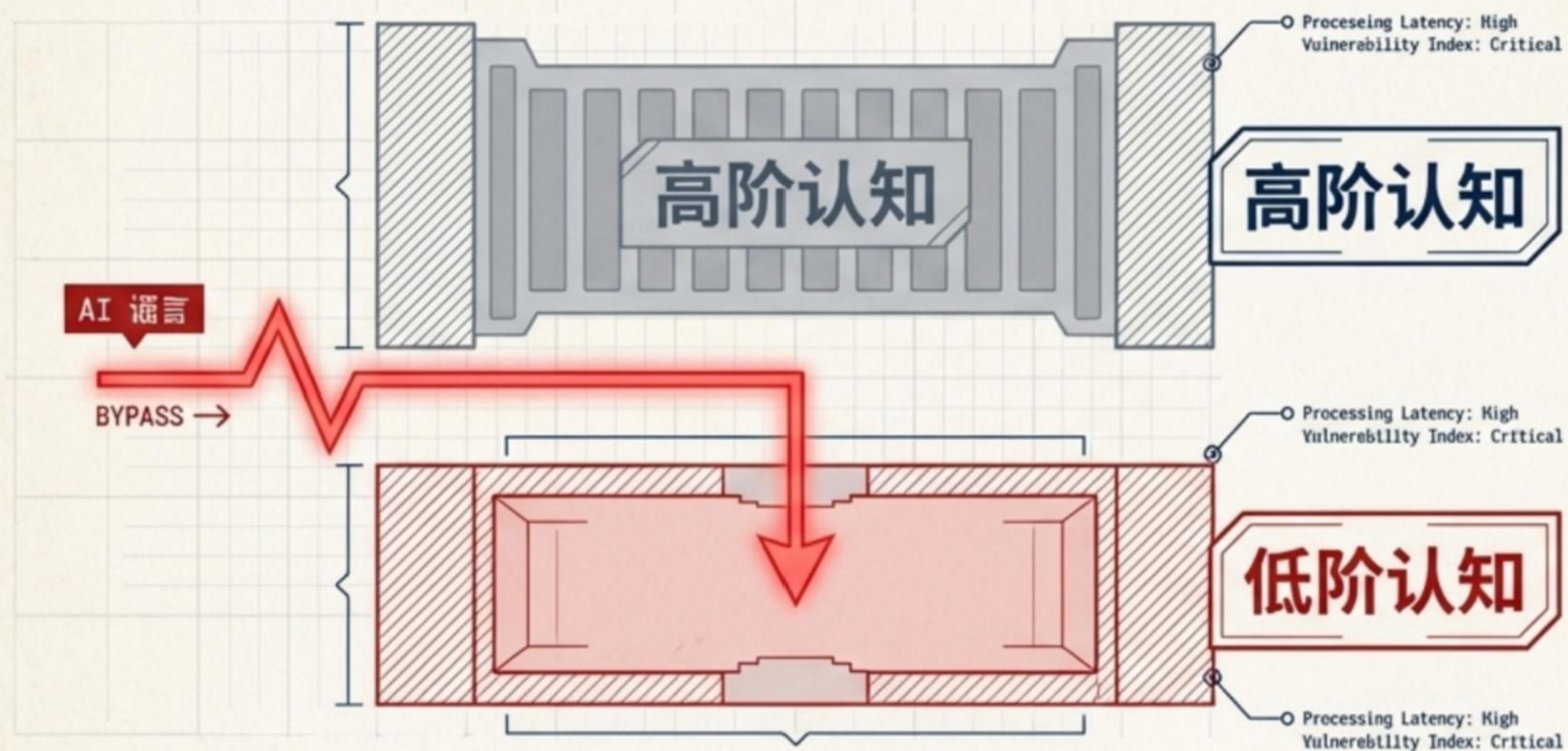
主要病因：	流量导向，注意力收割。
核心配置：	低信源关联×高情感刺激×高大众利益关联×高受众参与。
表现特征：	极度情绪化(愤怒/恐慌),利用人们的固有认知倾向引发共鸣。典型的弱势信息增强模式。

Suspect Profile

模式二：权威背书的强冲突

主要病因：	资本操纵，政治内嵌，定向打击。
核心配置：	高信源组织性×高针对性×强冲突框架×低受众互动。
表现特征：	依赖伪造的权威背书构建社会议题的剧烈冲突，单点突破，蓄谋造谣。

认知劫持：为什么我们会沦为传播节点？



判断、推理、逻辑分析。
在面临AI视觉与情绪的双重轰炸时，
这层防火墙往往被直接绕过。

感觉、知觉、本能反应。
眼见为实带来的认知偏差，让受众被
直观印象支配。

攻击机制：精准打击低阶认知。恐惧、安全威胁等本能情绪被触发后，受众在未经高阶理性验证前，就已完成了转发与共享。

部署抗体：从算法阻断到认知摩擦

战略转移：仅仅依靠后置的事实核查与删帖已无法追赶AI的变异速度。必须将**干预前置**。



1. 建立认知摩擦

通过人机协同模式，在用户浏览**高危情绪信息**时，自动弹出风险要素提示，**强制打断**低阶认知的反射弧，**激活**高阶理性评价机制。

高危情绪信息
AI强制打断

INTERVENTION POINT



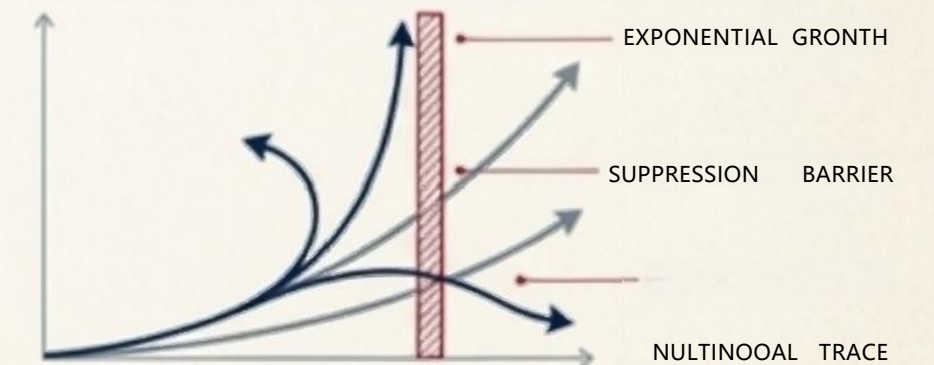
2. 组态化风险研判

放弃单一指标监测，转向组态计算。当系统侦测到**低知名度+强情绪+高切身利益**的组合特征时，自动提升**威胁预警等级**。



3. 阻断变异链

针对大众二次创作，建立多模态跨平台特征追踪溯源，压制其指数级增长曲线。



在生成式AI时代，打败机器的不仅是更强的机器，还有被重新唤醒的人类理性。

威胁的升级路线：从战场干预到商业入侵



核心趋势：技术门槛持续降低，攻击目标日益精准，真实世界破坏力呈指数级上升。

2022年：战争中的深度伪造首次成为全球舆论事件

事件： 伪造泽连斯基呼吁乌克兰士兵投降

媒介： 视频



爆发时间： 2022年3月

处置过程：

- Reuters介入核查，确认视频制作粗糙，被迅速识破。Meta紧急移除相关内容。
- 乌克兰官方被迫发布真实视频进行回应。

案例意义

短时混乱： 粗糙的技术依然足以在极端环境下制造恐慌。

信任降级： 最深层的风险在于“狼来了”效应，它直接加剧了公众对后续真实视频的怀疑。

2023年：AI 图片引发短时市场与舆论抖动

- **事件：**五角大楼“爆炸假图”在社交平台疯传



- **媒介：**图片

- **爆发时间：**2023年5月

- **处置过程：**

- AP (美联社) 核查指出图片带有明显AI生成伪迹。
- 阿灵顿警方和消防部门出面澄清并无此事。

- **现实后果：**假图扩散期间，直接引发金融市场短暂的紧张情绪。

案例意义

跨界联想：一张图片即可瞬间触发公共安全与金融市场的连锁反应。

真实度幻觉：伪造图像根本不需要“完美逼真”——它只需要足够像“突发新闻”。

2024年：AI语音克隆长驱直入干预民主选举

●事件：冒充拜登声音的AI外呼，
劝阻新罕布什尔州选民投票。

●媒介：语音



●爆发时间：
2024年初选前两天

600万美元

2024年9月，美国FCC正式对策划者
Steve Kramer处以巨额罚款。

●规则界定：FCC于2024年2月确认，AI生成语音明确属于
TCPA监管范围内的“artificial or prerecorded voice”。

案例意义：AI谣言不再局限于“内容误导”，已演变为对民主程序的直接物理干预。

威胁铁三角：语音克隆的杀伤力源于“低成本+高可信+高接触率”。

2024年：深度伪造从公共舆论转向核心组织欺诈

● **事件：** 香港深度伪造视频会议与加密货币诈骗

● **媒介：** 预录视频会议/深度伪造



● **核心数据** (香港警方书面披露):

2.4亿港元

与400万港元(两宗预录视频会议诈骗损失)

超3.6亿港元

(利用deepfake 实施的加密货币诈骗金额)

案例意义

威胁合流： AI 谣言与AI诈骗正在形成完美的闭环汇流。

目标升级： 警钟敲响——受害者早已不再只是缺乏辨别能力的普通网民，而是深谙企企业和组织流程的专业人员。

穿透表象：AI欺骗的四大底层法则

1. 议题高危化

高风险议题优先：攻击者精准锁定能够引发高肾上腺素的领域——战争、选举、国家安全、金融市场。



2. 证据直观化

强证据形式优先：抛弃纯文字，优先使用符合人类直觉的“硬证据”媒介——视频、语音、现场图片。



3. 爆发瞬时化

传播窗口极短：破坏力不依赖长期潜伏，绝大多数风险与后果发生在最初的短时扩散阶段。



4. 证伪不对称

纠错成本极高：造谣动嘴，辟谣跑断腿。一旦形成广泛误导，官方澄清的速度永远慢于情绪传播的速度。



四个真实案例共同证实了一个冷酷的现实：

AI谣言最危险的地方，
不在于它能“持续欺骗很久”，
而在于它能在**很短的时间**
内造成真实的后果。

告别经验判断，将抽象危害转化为可比较的治理指标

抽象危害感



没有指标的治理只能停留在被动的经验判断。

可比较的风险指标

内容风险指标 65/100



主体风险指标 82/100



传播风险指标 40/100



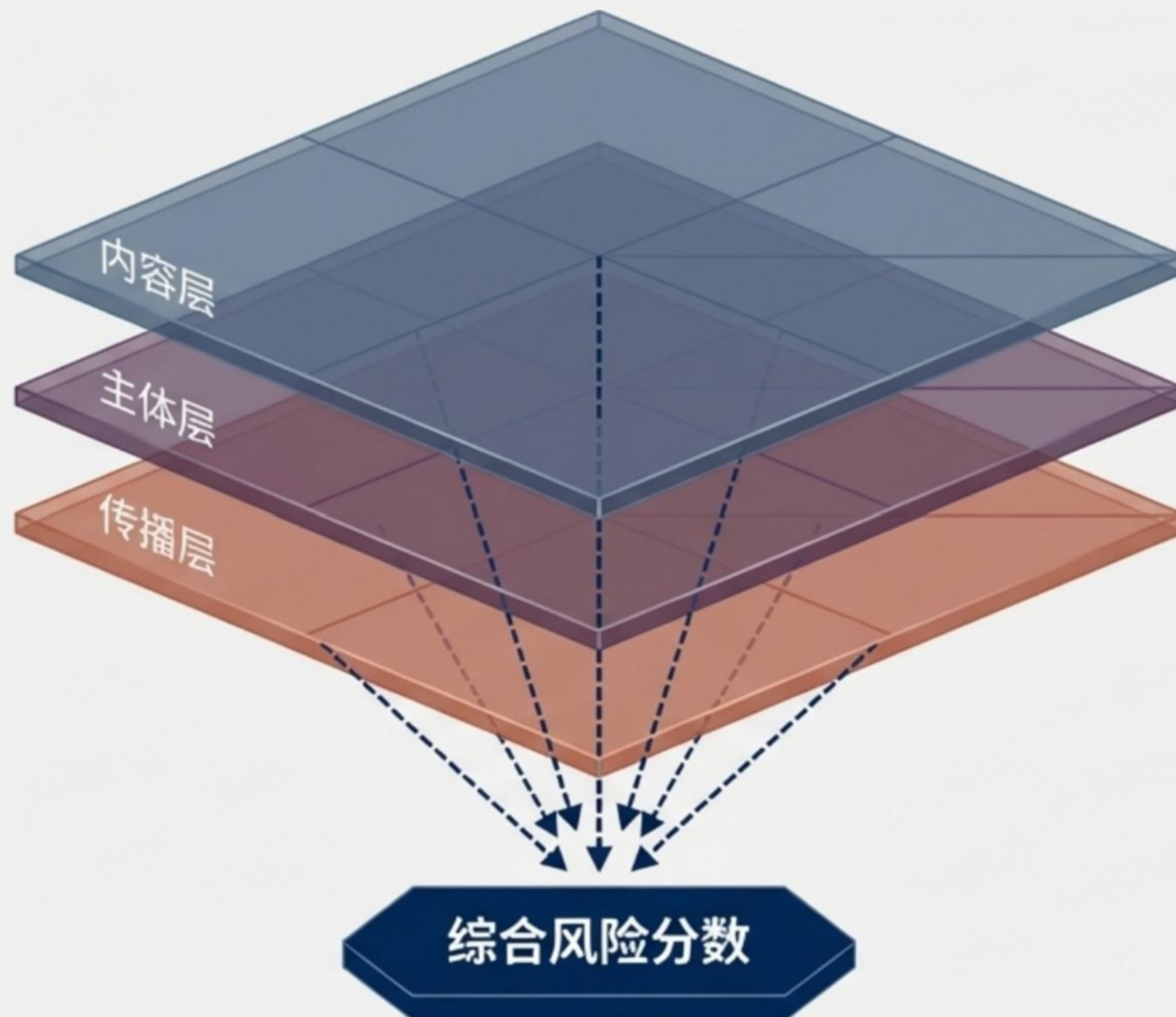
平台：需要明确什么内容该限流、什么事件该升级。

研究者：需要比较不同事件、平台和国家的风险差异。

政策制定者：需要精准识别“高风险窗口”和“高风险对象”。

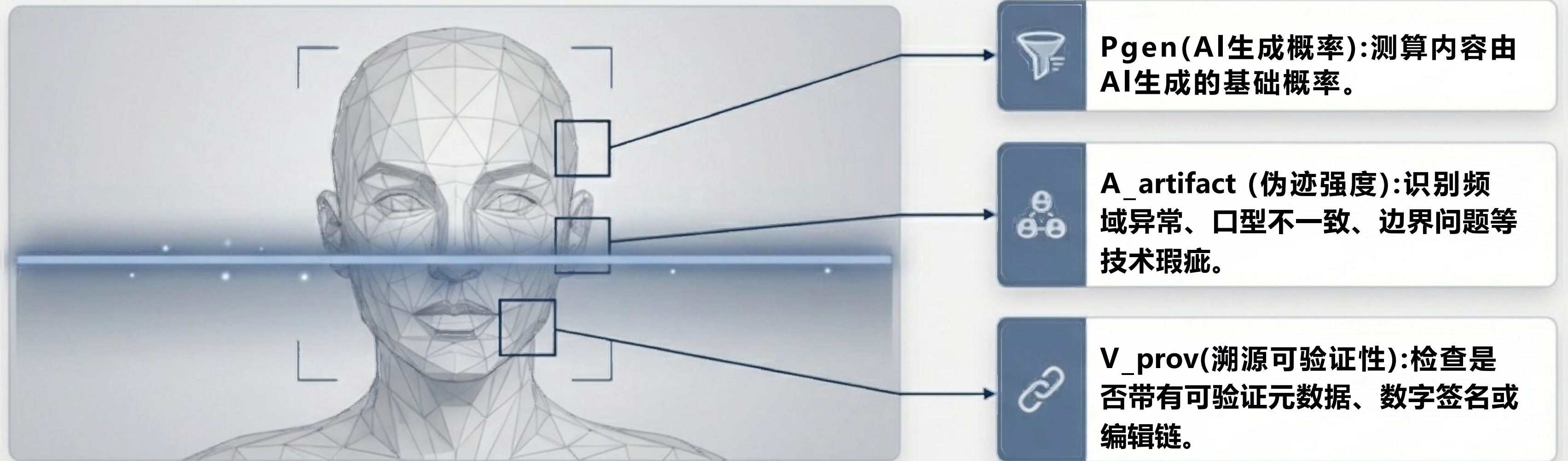
可信AI离不开可识别、可管理、可评估的风险框架。

穿透信息迷雾的三维风险评估架构



内容层筛查：研判信息是否具备AI伪造特征

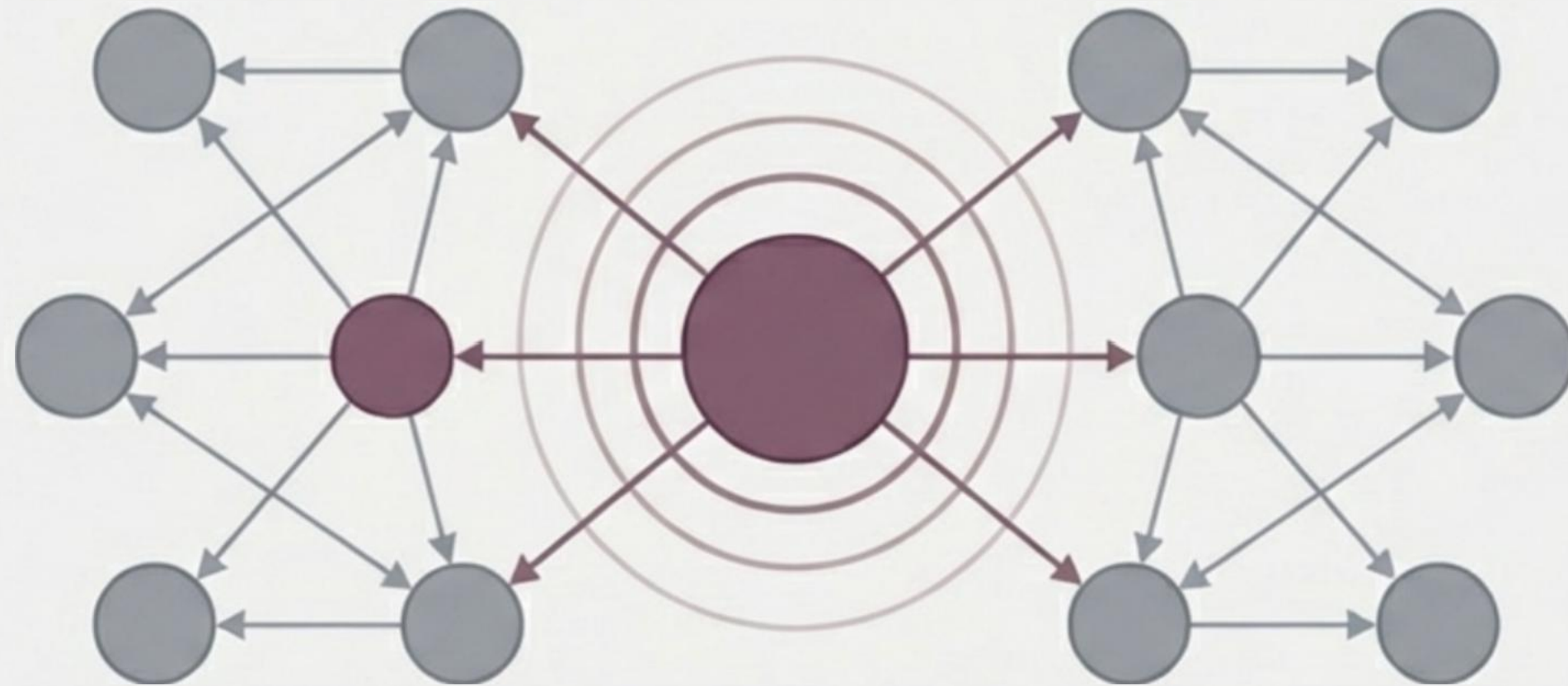
“像不像AI 伪造？”



深度伪造检测、内容凭证和水印机制各有其作用边界。这三项指标不绝对等同于“真假”，但构成了至关重要的第一道风险筛查。

主体层追踪：锁定传播网络中的核心节点与异常行为

“是谁在推、怎么在推？”



Tactor(账号可信度)

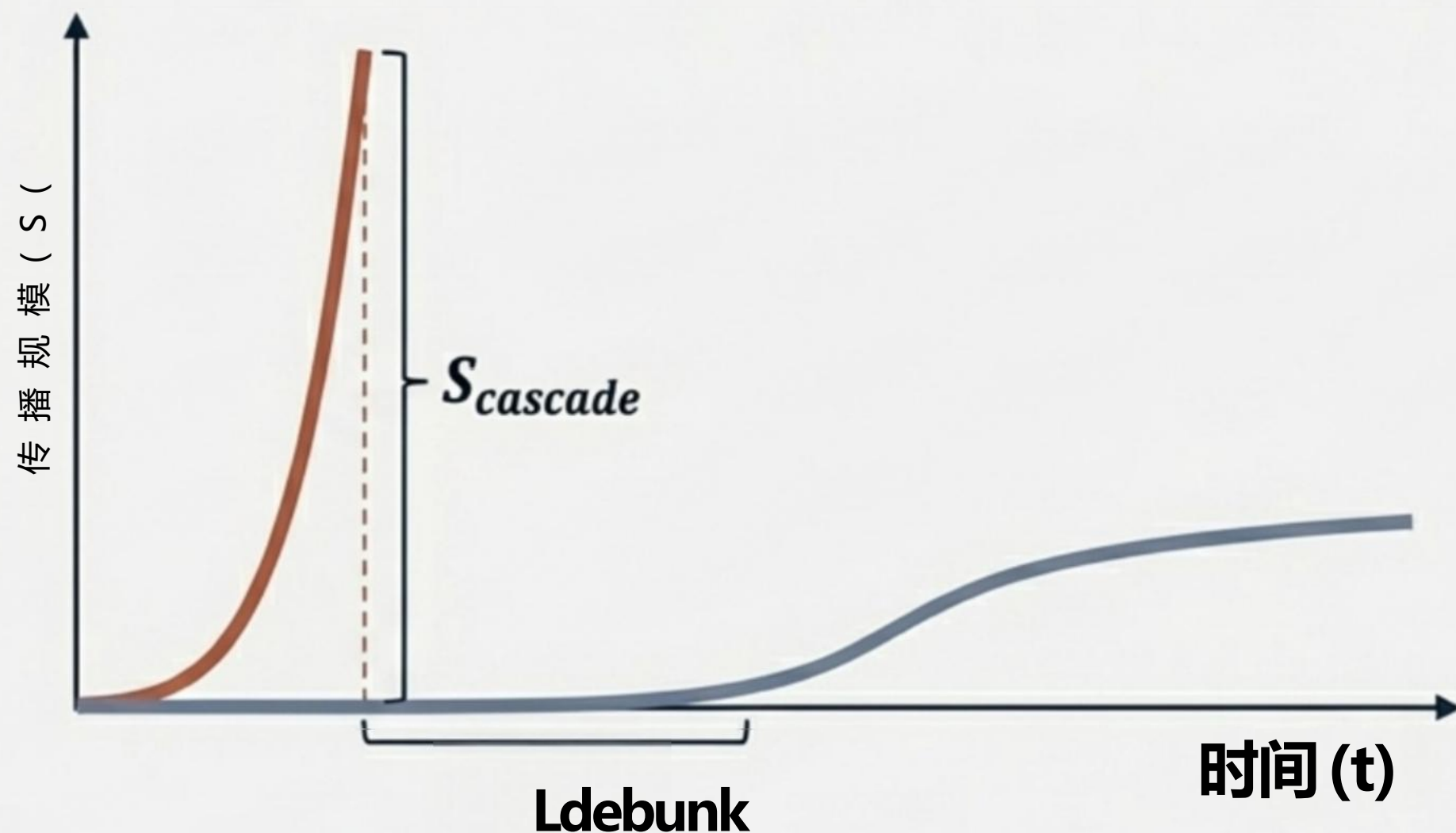
Bbot(机器人/半自动账号概率)

- ! 异常发帖频率
- ! 高同步性转发
- ! 账号画像模板化
- ! 历史违规记录

捕捉主体层特征是实现早期预警的核心关键。

传播层量化：衡量短时高冲击态势下的网络级联效应

“扩散得有多快、多广、多难纠正？”



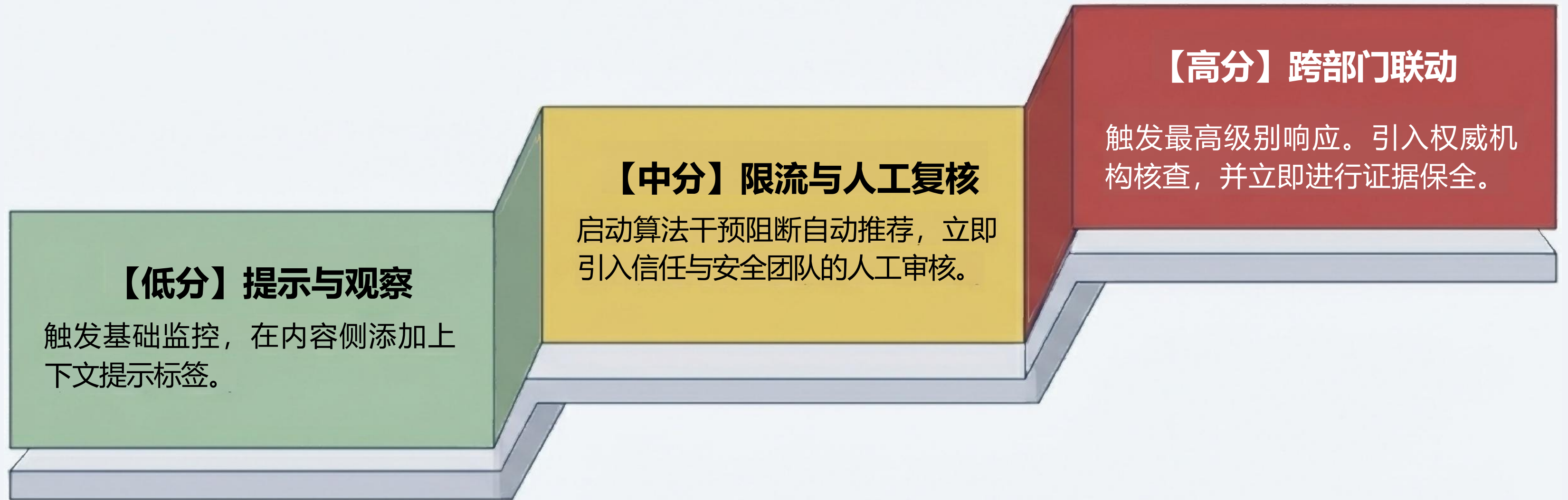
$v(t)$ (传播速度)& $S_{cascade}$ (级联规模):衡量爆发烈度。

Famp (放大因子):某条AI谣言在相似主题、相似时段下相对于真实信息的扩散倍率。

Ldebunk (辟谣滞后时间):从首发到权威澄清的时间差。

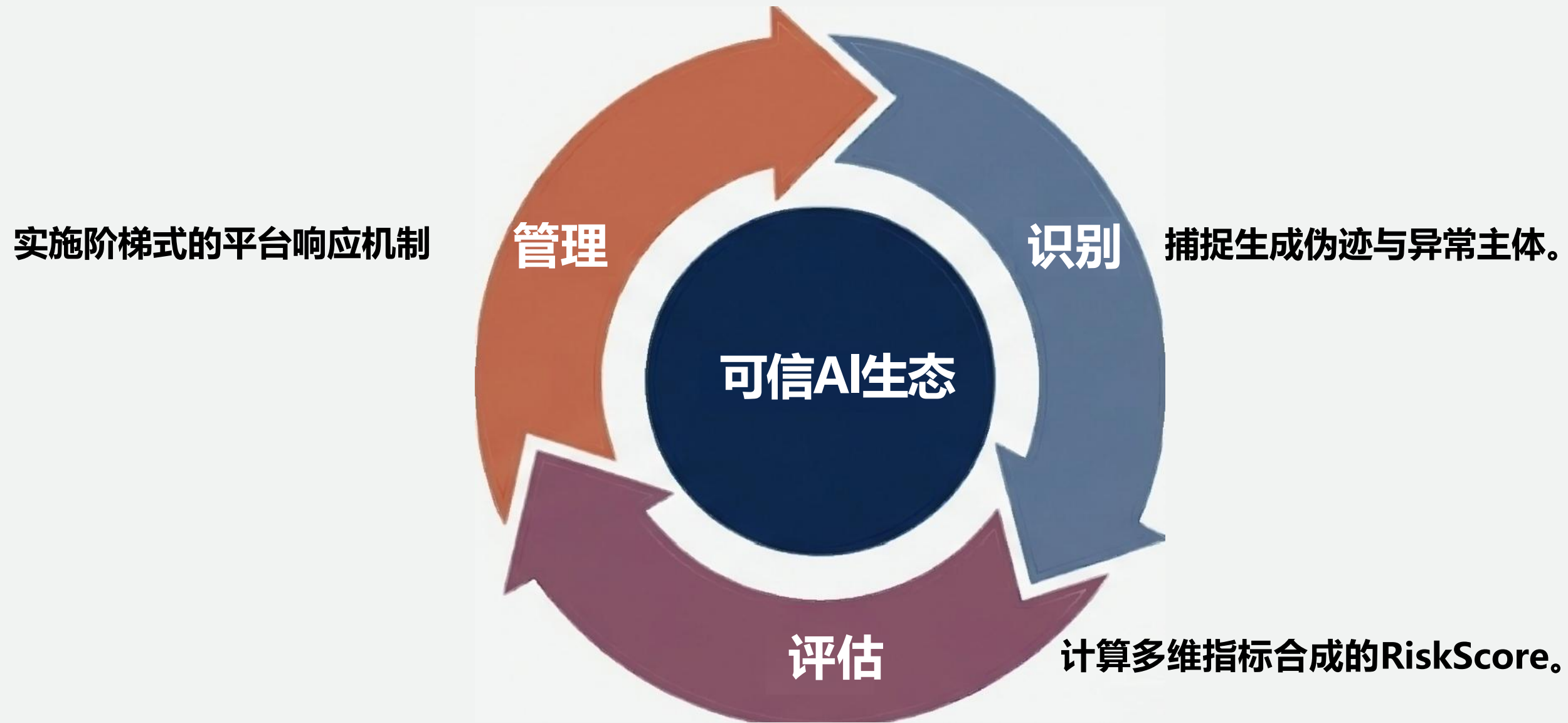
治理的难点在于AI谣言的“短时高冲击”特征。放大因子越高、辟谣滞后越长，造成的系统性破坏呈指数级上升。

基于量化分数的阶梯式应急响应矩阵



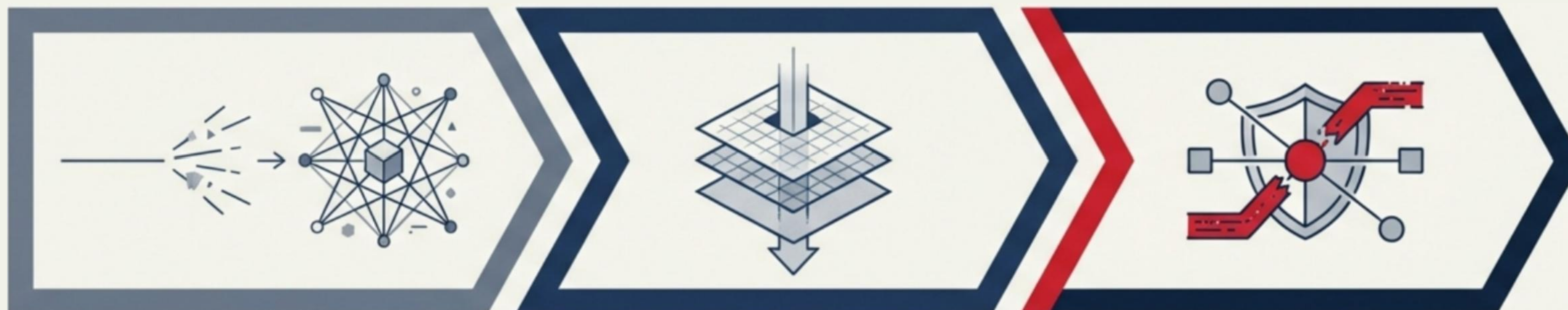
将指标体系直接嵌入平台的业务流，实现从“事后扑火”向“全流程风险预控”的根本转变。

锚定国际标准，推动可信AI生态的闭环管理



本报告构建的治理型评分框架，不仅为当前迫切的AI谣言治理提供了实战工具，更在底层逻辑上与 NIST AI风险管理框架 (AI RMF) 完美契合。量化风险，方能驾驭未来。

认知路线图：一场正在发生的防御范式转移



第一阶段 策略重构

告别单一检测路线，构建多维交叉证据网。

第二阶段 技术突破

多模态与CIB框架的崛起，
穿透复杂伪装。

第三阶段 信任危机

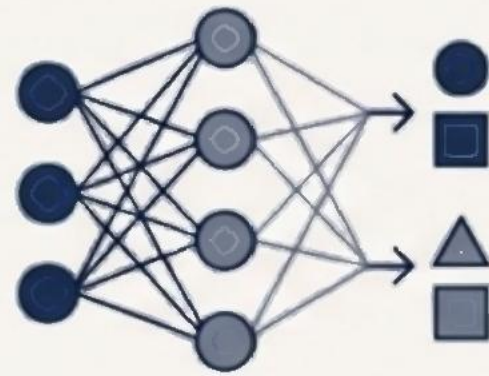
直面下一代威胁：大模型预
训练与RAG知识库投毒。

构筑防御基石：当前检测方法的三条核心路线



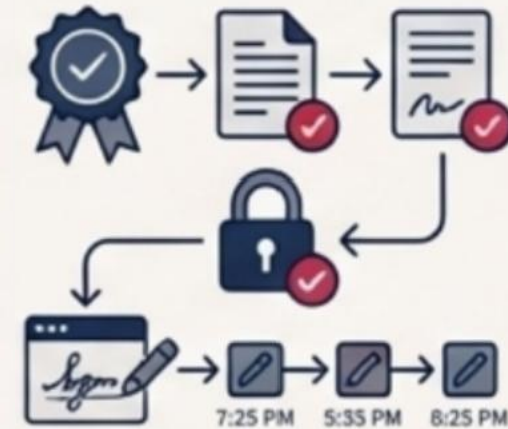
特征取证

核心动作：查伪迹、
查不一致



模型检测

核心动作：利用分类
器或多模态模型直接
识别可疑内容

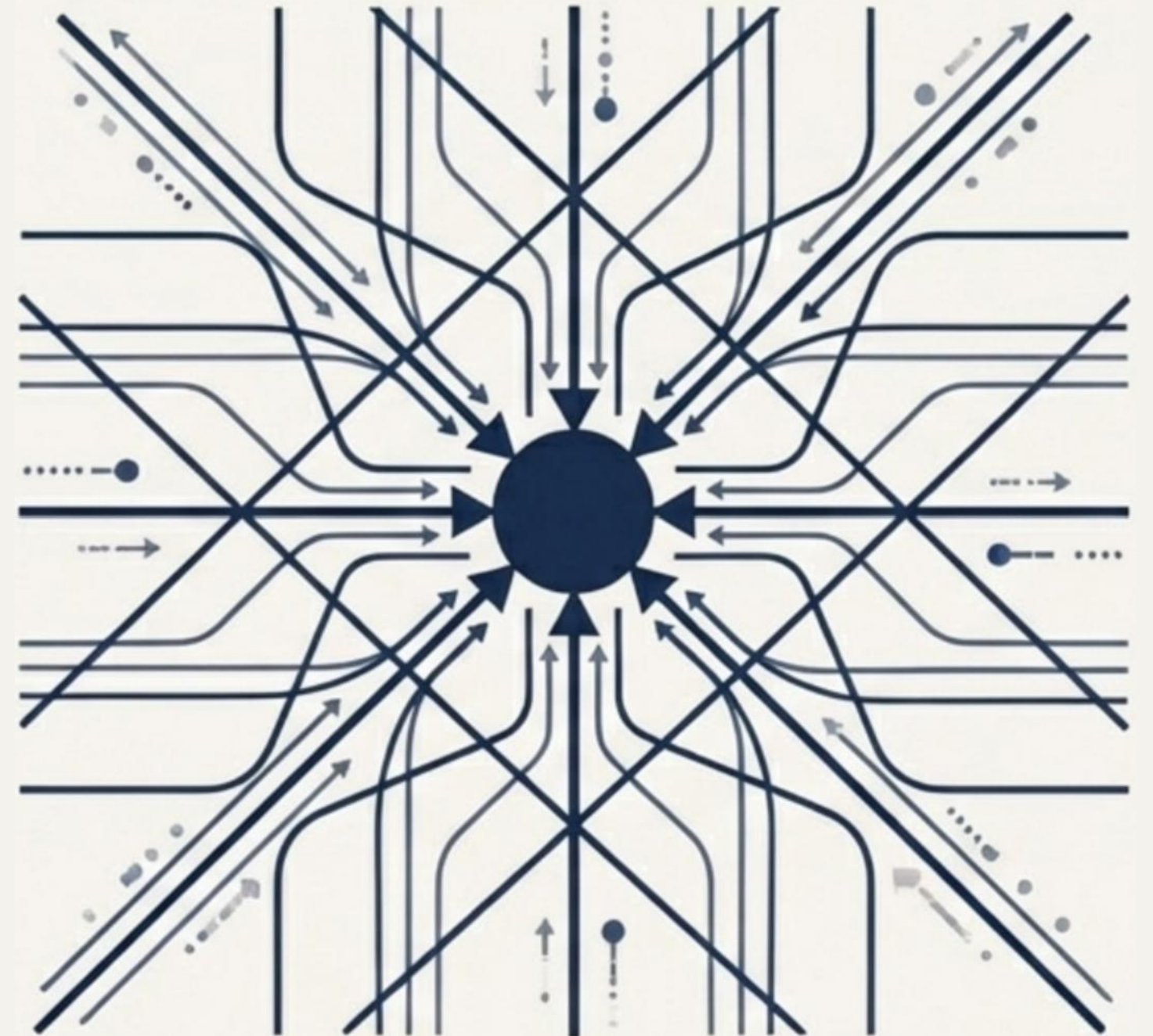


溯源验证

核心动作：通过元数
据、数字签名、水印
和编辑链进行硬性验
证

终结迷思：单一“神奇检测器”已经全面失效

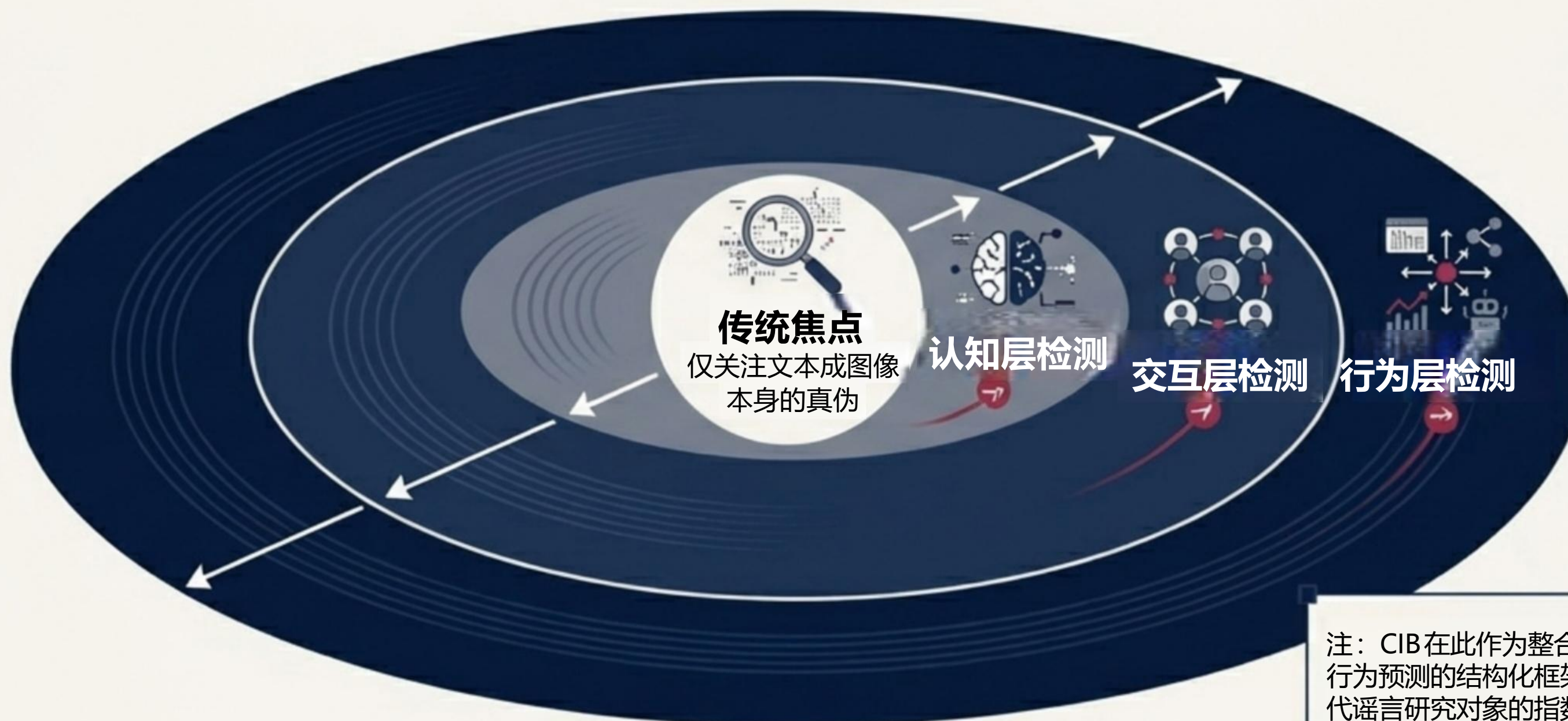
未来最可行的路径，是构建“多路证据交叉”的复杂验证网络。



视角的升维：跨越单纯的“内容真假”核实

LLM时代的新框架

引入“认知-交互-行为”（CIB）三层检测思路，将孤立的内容放入动态的传播生态中进行审视



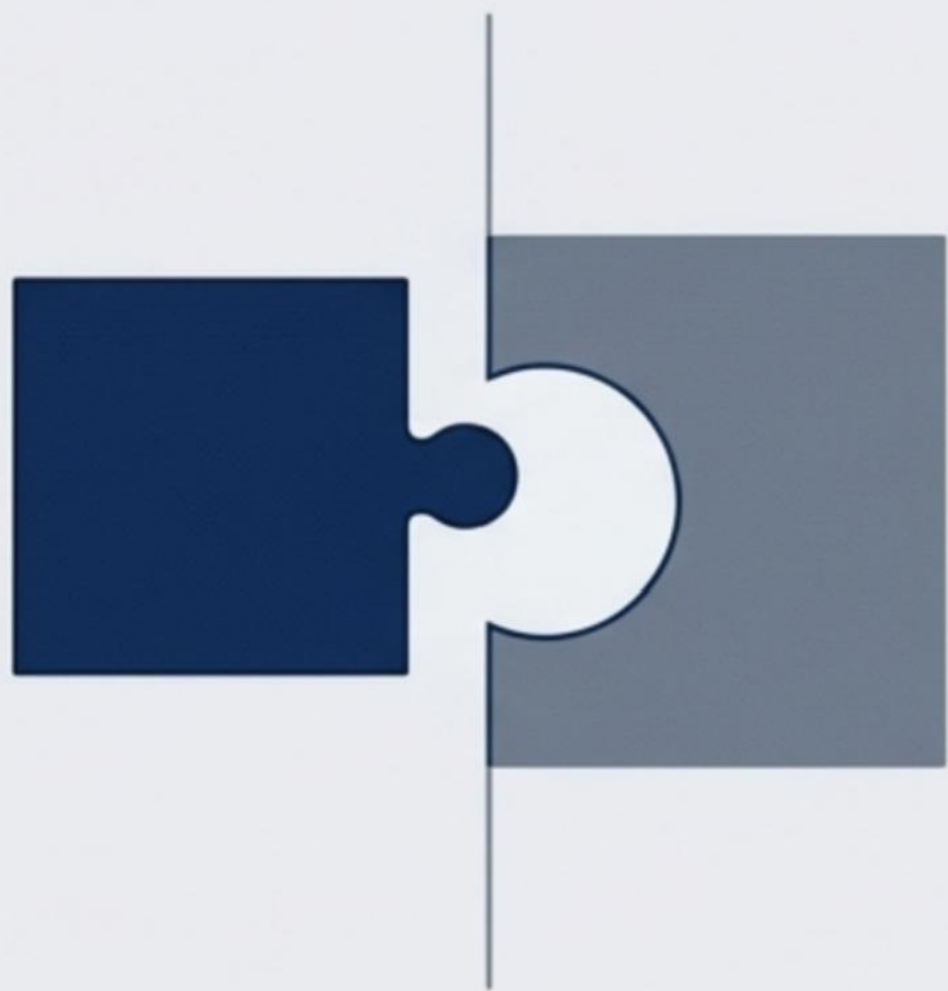
注：CIB在此作为整合内容、认知与行为预测的结构化框架，以应对AI时代谣言研究对象的指数级扩展。

拆解CIB结构：认知、交互与行为的立体防护网

<p>认知层 (Cognitive)</p>		<p>目标：剖析内容本体。 关注点：文本、图像、音视频本身是否存在语义或事实层面的逻辑异常。</p>
<p>交互层 (Interactive)</p>		<p>目标：扫描社交网络。 关注点：识别账号间是否存在协同操作、异常同步或水军机器人的集群参与。</p>
<p>行为层 (Behavioral)</p>		<p>目标：追踪宏观传播。 关注点：分析传播轨迹是否表现出非自然的爆发式增长或典型的攻击型模式。</p>

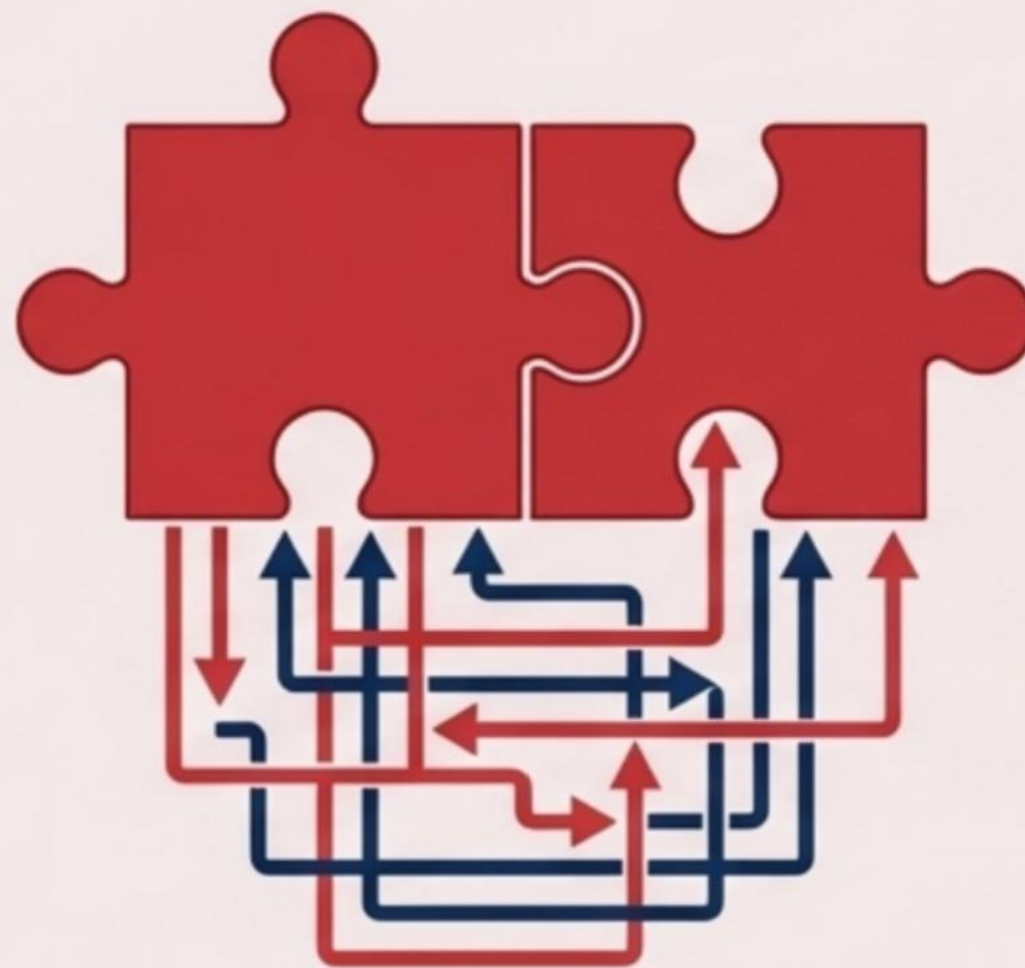
欺骗手段升级：浅层特征在复杂多模态谣言面前全面失效

旧威胁：简单伪造



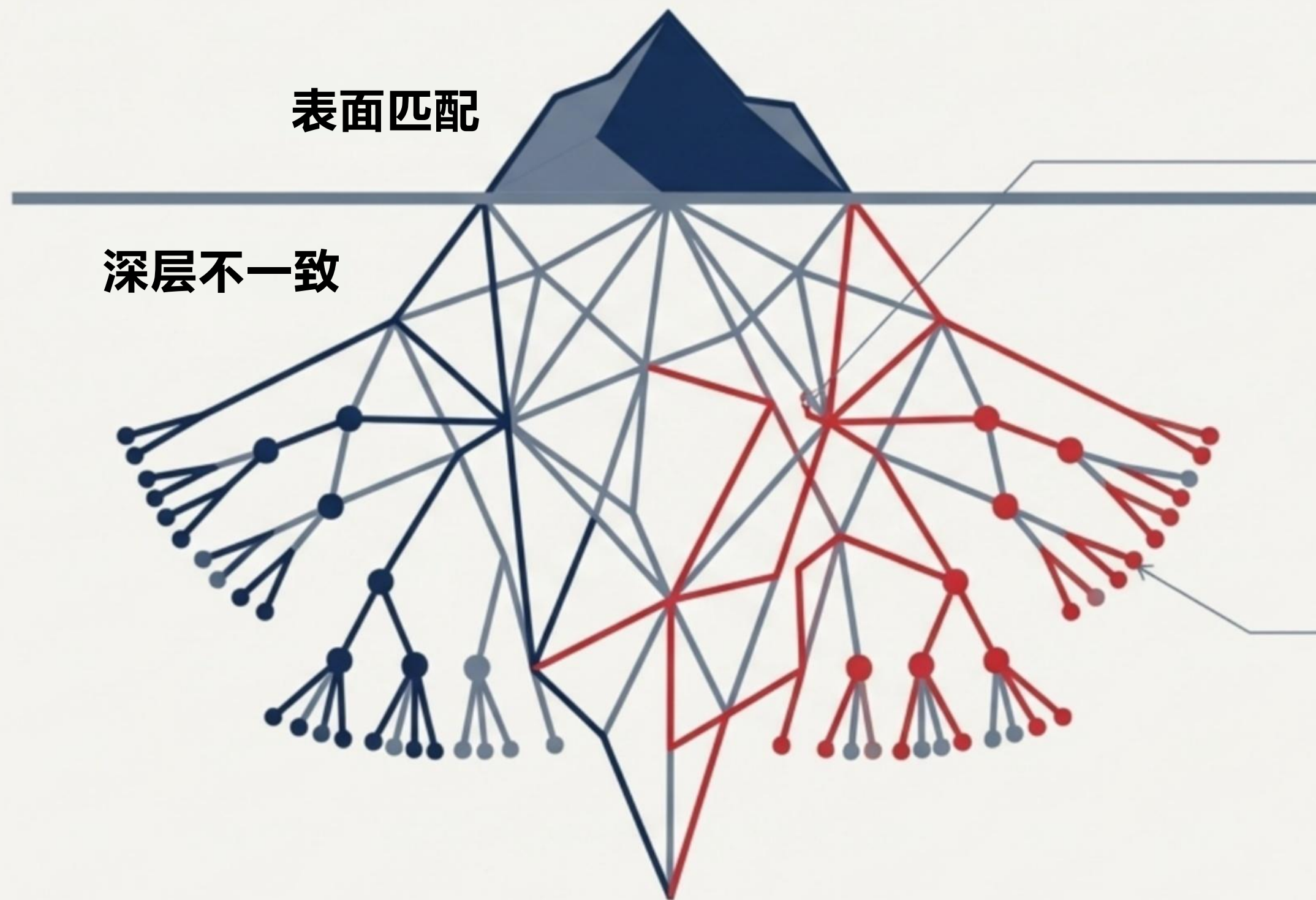
图文完全不相干。浅层特征融合即可轻易拦截。

新威胁：复杂虚假叙事



表面图文高度匹配，实则在深层逻辑上存在叙事造假。常规多模态模型被轻易绕过。

前沿突破：RumorCone模型穿透深层语义伪装



表面匹配

深层不一致

源自EMNLP2025的最新研究成果。

核心机制：双曲几何与层级语义建模。显式建模图像与文本在不同抽象层级上的关系，精准锁定深藏元层级上的关系，精准锁定深藏不露的逻辑断裂点。



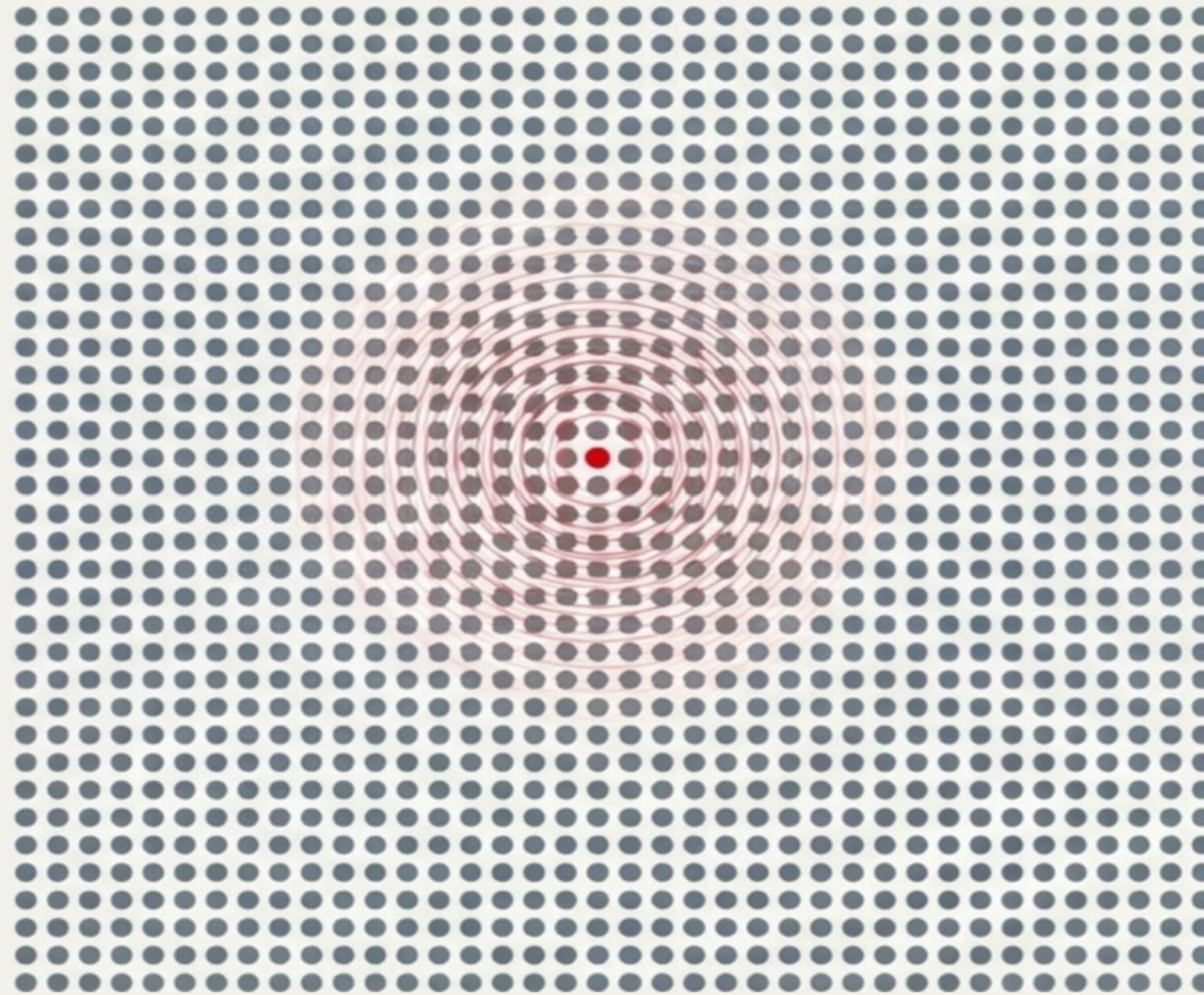
致命隐患：当事实核查系统的“裁判”被提前收买

我们一直在倾尽全力防范“恶意用户输入”。但下一代检测的真正问题在于：
如果用于核查的“模型本身”，在出厂前就已经成为了污染源呢？

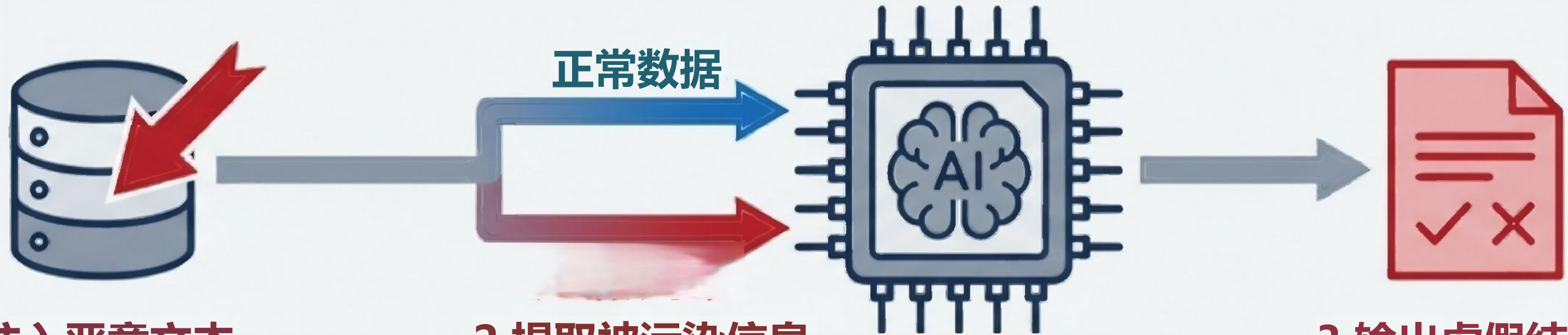
预训练投毒：极微量数据如何击穿大模型信任底座

2024年最新研究证实，极低比例的投毒即可造成深远且持久的破坏。

核心警告：模型本身已从单纯的“检测工具”转变为潜在的“谣言风险传播器”。系统级漏洞远比输入端欺骗更难根除。



PoisonedRAG:外挂知识库正在成为事实核查的新软肋



1.注入恶意文本

攻击者向系统外挂的“权威知识库”注入污染信息。

2.提取被污染信息

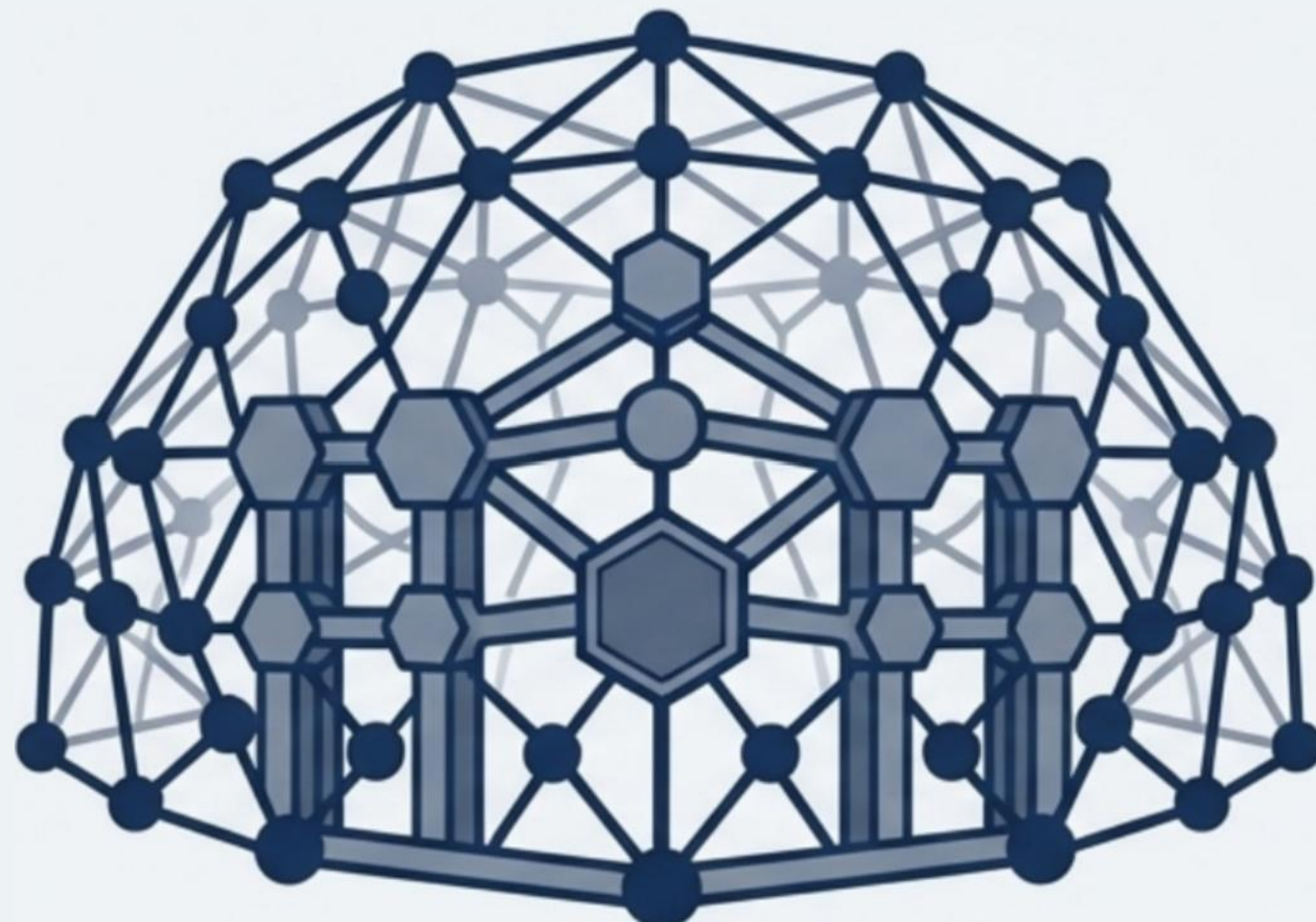
AI事实核查系统在检索时被误导。

3.输出虚假结果

系统基于错误前提，输出被诱导的“虚假辟谣”。

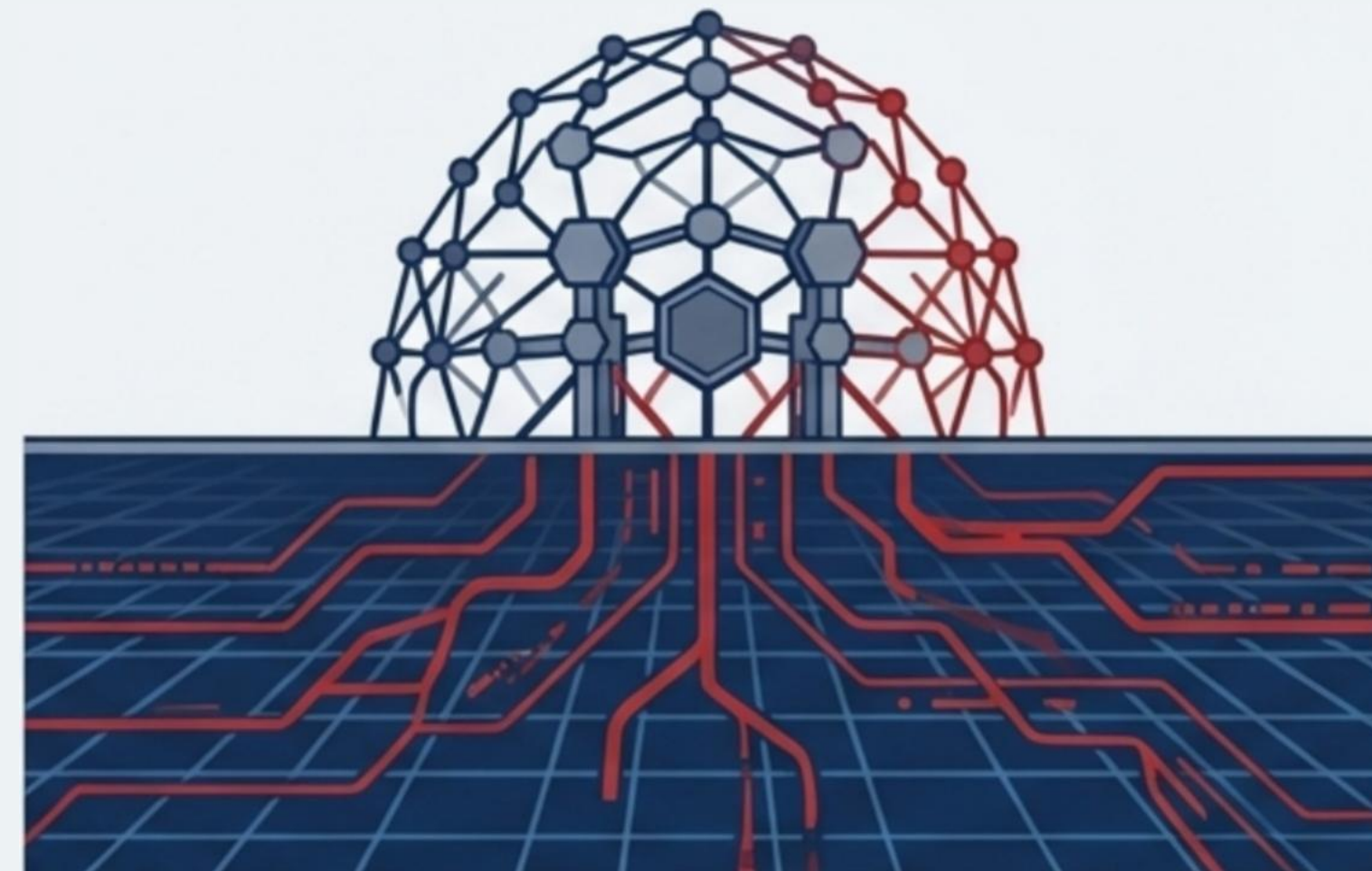
USENIX Security2025-
PoisonedRAG 研究揭示

RAG架构并不天然具备防御免疫力。未来的AI事实核查系统本身，随时可能被针对性投毒所操控。



盾：坚固的城墙

由多路证据交叉、CIB立体框架与RumorCone等多模态前沿技术构筑的防御网络。



矛：底层的暗道

必须时刻防范的预训练模型投毒与RAG知识库污染。

AI谣言治理不再是单纯的技术修补，而是一场长期的系统级博弈。唯有跳出单一检测思维，**重构从数据源头到多维行为的信任链条**，才能应对下一代认知威胁。

重新定义危机：AI谣言已演化为系统性的社会工程攻击



范式转移

聚焦“技术现实主义”。AI不仅在造假，更在利用“逻辑完美谣言”绕过人类固有的逻辑防御，直接篡改认知。



系统脆弱性

攻防高度不对称。少数**数据投毒**即可在底层植入偏见，而传播端的**元数据剥离**直接切断了信任溯源的链条。



治理新基建

应对之策必须跨越技术边界。结合2025全球合规体系，构筑涵盖认知重塑、威胁狩猎与全产业链标识的动态防御网网络

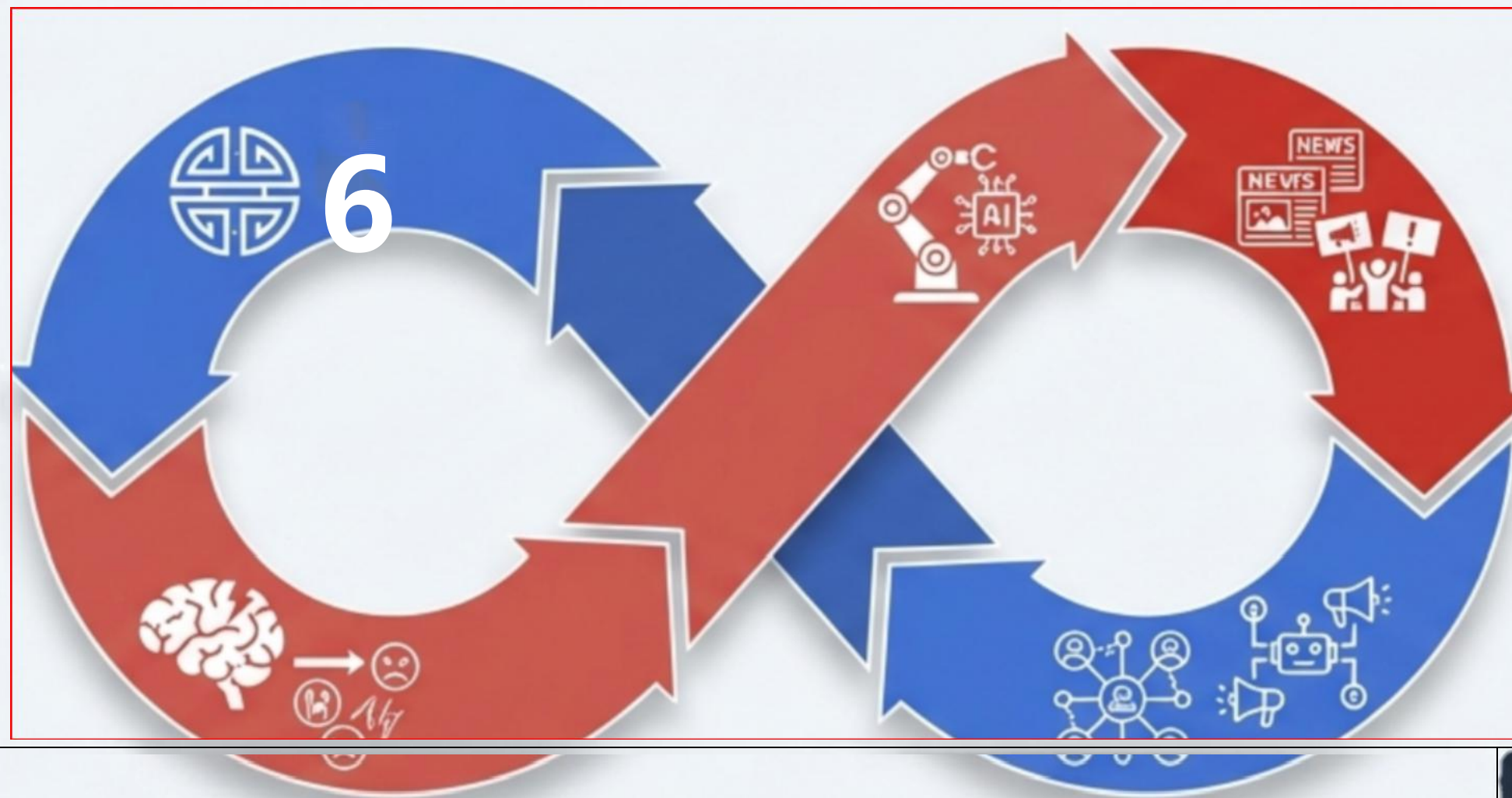
从“情绪驱动”到“逻辑伪装”的范式转移

传统网络谣言	AI生成谣言
-生成机制：人类个体创作、水军批量改写 	-生成机制：自动化算法生成、LLM链式推理 
-关键特征：情感煽动、逻辑断裂、低质图像 	-关键特征：技术现实主义、逻辑自治、超现实多模态 
-传播模式：社交网络级联传播 	-传播模式：算法精准推送与社交机器人自动化放大 
-心理路径：触发群体偏见与恐惧 	-心理路径：绕过理性审查，直接作用于直觉系统 
-社会影响：局部舆论误导、社会骚乱 	-社会影响：社会工程攻击、金融市场动荡、国家安全威胁 

“逻辑完美谣言”不再依赖夸张词汇，而是通过大模型的链式推理构建专业、严密的虚假论证。

社会工程攻击2.0:高度定制化的阵列化打击

1.真实素材锚定: 捕捉特定文化符号或公共安全事件。



2.AI自动化生成: 大规模生产高度定制化的虚假灾难或群体对立内容。

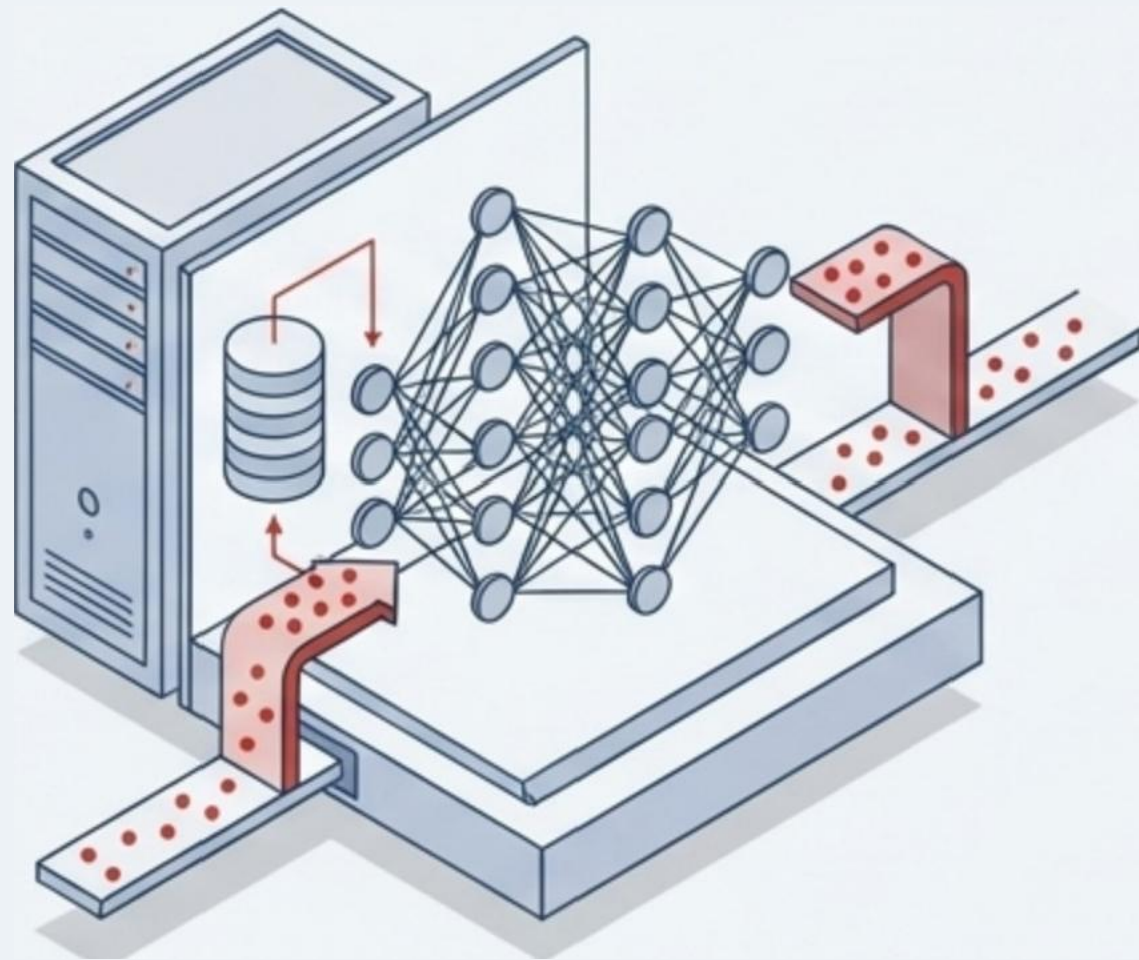
4.认知漏洞突破: 精准利用恐惧、愤怒或身份认同摧毁社会信任。

3.算法与机器人协同: 多平台同步投放, 制造虚假民意。

2025年网络环境的显著特征是“攻击阵列化”。这不再是单纯的病毒木马, 而是针对社会心理结构发起的精准网络战。

隐秘的系统性脆弱性：从底层数据投毒到RAG劫持

预训练阶段的数据投毒

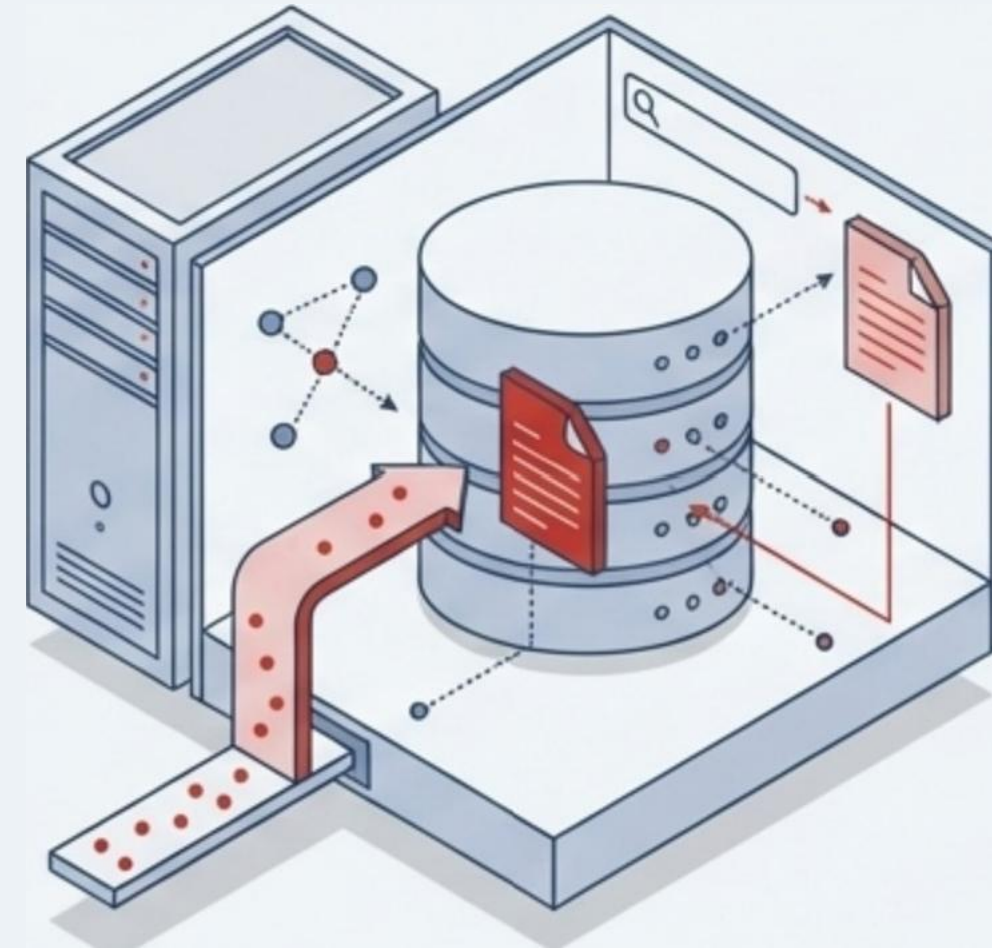


核心威胁： 仅需换0.001%的训练令牌，即可让模型产生不可见的永久偏见。

攻击策略：

- 分割视图投毒(Split View):利用爬虫抓取榜间替换URL快照内容。
- 前端运行攻击(Frontrunning):在数据集更新前瞬间插入伪内容。

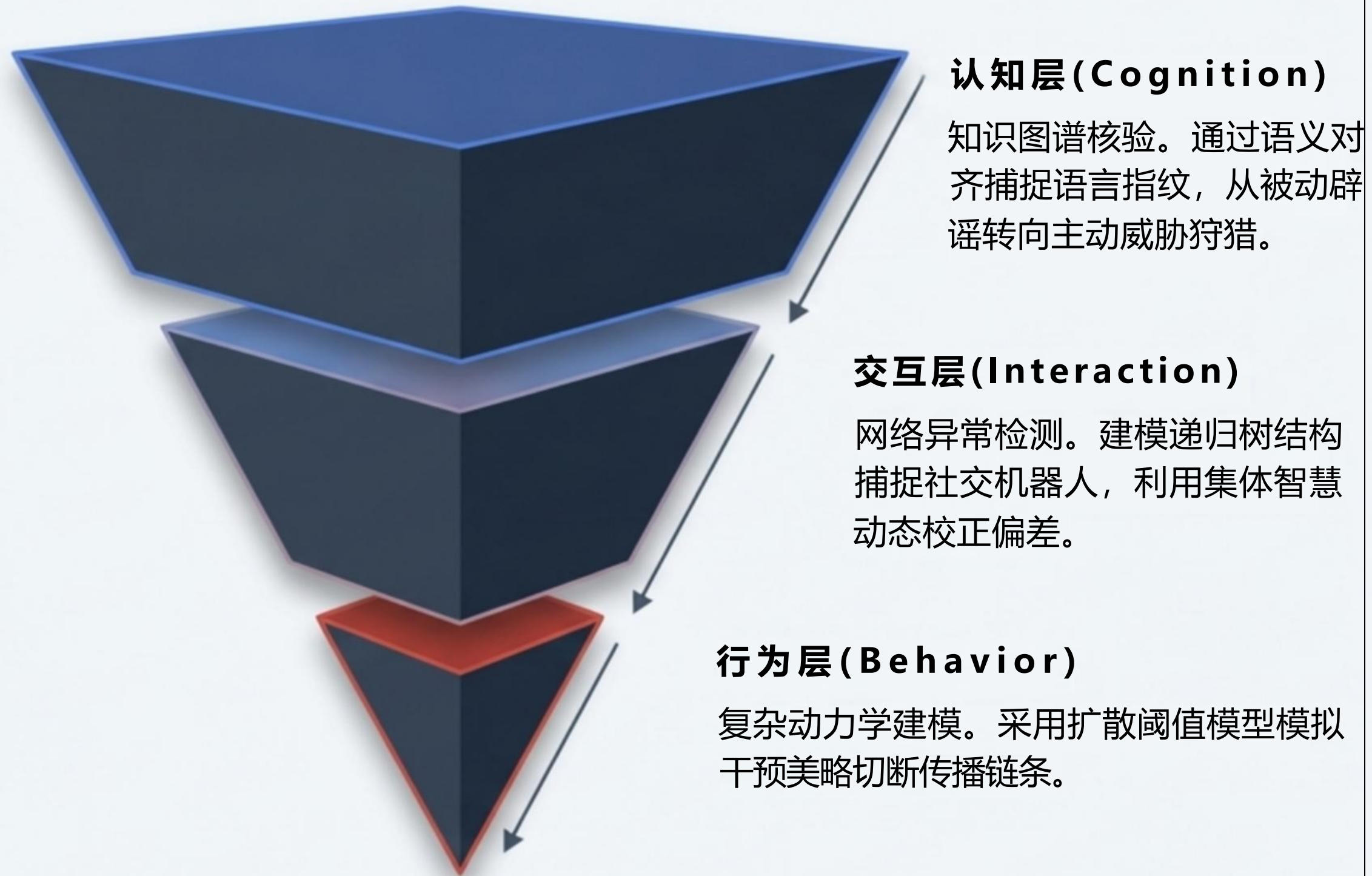
应用阶段的RAG知识库投毒



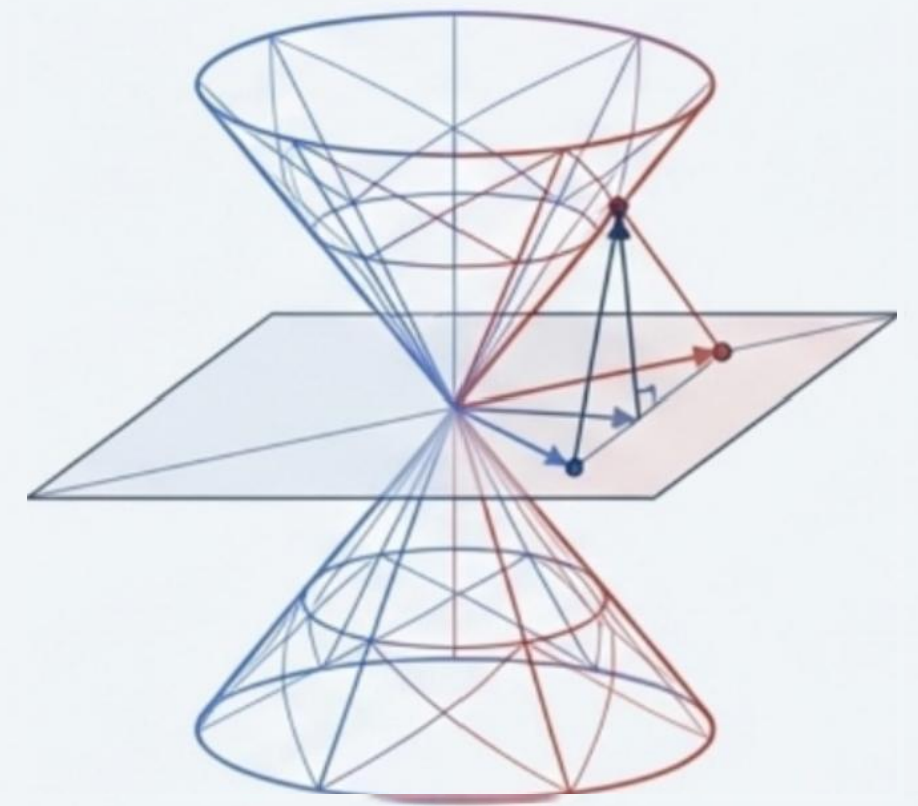
核心威胁： 通过SEO操纵或入侵，向企业向量数据库注入恶意文档。

机制效应： 相似性机制被劫持，导致AI以极高置信度背书虚假主张，系统性操纵决策。

构筑自适应防御：CIB三层架构与RumorCone跨模态检测

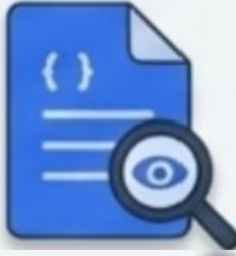



RumorCone双曲几何模型



针对图文不一致的“深水区”，将多模态特征嵌入双曲空间，分三层抽象识别微小的跨模态语义矛盾，击破复杂虚假叙事。

深度合规解码：中国《AI生成合成内容标识办法》

 显式标识 视觉/听觉标记，强制触发用户的理性认知。		 隐式标识 不可见元数据/数字水印，确保内容被下载导出时信息持久存留。
服务提供者		必须添加底层标识；保留生产日志；禁止提供破坏标识工具。
分发平台		必须核验隐式标识；自动检测并添加“疑似AI生成”标签。
应用商店		强制对上架的AI应用进行合规审核。
内容发布者		必须如实声明内容属性，禁止恶意篡改抹除。

治理的阿喀琉斯之踵：元数据剥离与断裂的溯源链

硬件与生成端 (坚固)

技术演进： C2PA 2.3标准已支持流媒体 片段级 签名。

硬件落地： Google Pixel 10集成Titan M2安全 芯片，在拍摄瞬间加密绑定元数 据(GPS/时间戳)确 立“出生证明”。

传播与分发端 (断裂点)

现象： “元数据剥离”。出于带宽与隐私 考量，平台强制重新编码。

平台现状 (脆弱)

结果： 不可篡改的加密签名失效，硬绑 定被破坏。

现状： 目前的 AI 标识多依赖于平台自觉，缺 乏强制性的跨平台互认机制。Meta/TikTok等虽 尝试显示AI标 签，但公网环境下的原始证据链 极难被 普通用户查验。



战略展望：构建数字经济下半场的真实性防线

认知重塑与辅助



强制标识提供“边缘线索”，触发公众警觉；数字素养教育倡导面对高情感强度信息“先暂停、后核验”。

主动威胁狩猎



从被动辟谣转向主动防御。捕捉AI语言指纹，建模攻击者TTPs实现前瞻预警；激发集体智慧构建免疫系统。

捍卫全球信任底座



跨国协同解决元数据剥离痛点，在社交媒体架构中保留核心证据链。透明度将成为企业未来的“信任溢价”。

未来AI谣言的技术博弈将呈现长期化与链条化。唯有法律威慑、硬件溯源、模型防御与公众认知的多维协同，方能保障AI创新与信息秩序的双赢。

报告编撰成员

清华大学新闻与传播学院 博士后 张诗瑶
清华大学新闻与传播学院 教授 沈 阳

说明

本报告综合调度并协同使用了OpenClaw, Claude, ChatGPT, NotebookLM等多种智能体工具, 组织开展了多轮交叉审核与迭代修订。报告的资料检索, 初步整理, 结构优化, 文字润色及部分排版工作由人工智能辅助完成, 由人工进行关键内容把关与最终确认。受多模型生成差异与工具稳定性影响, 部分页面在字体, 间距, 图表样式等方面尚未实现完全统一, 后续可根据需要进一步规范与优化。特此说明。