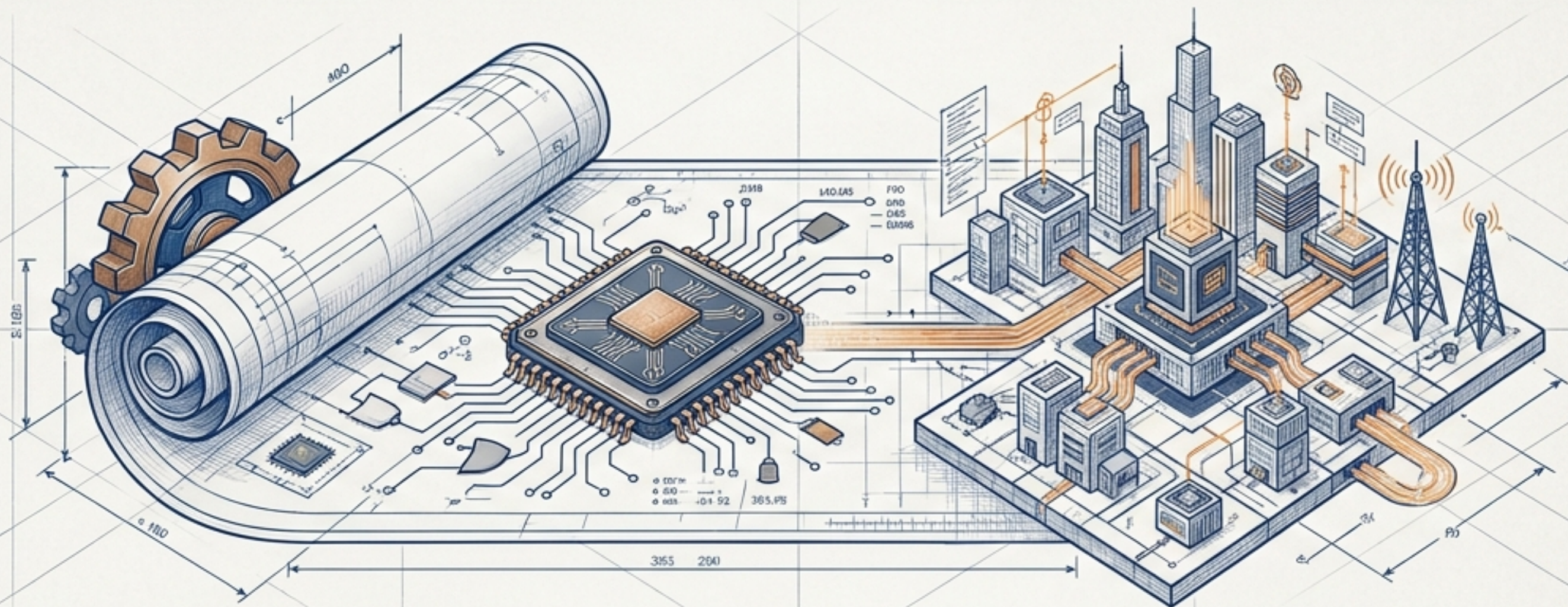


重构下一代数字生产力：字节跳动 Coze 全栈架构与生态深度研究

从工作流编排到 Agent World 的演进、企业级合规治理与算力经济学剖析



宏观生态定位

底层架构(DDD)
与技术栈

模型生态与
算力经济学

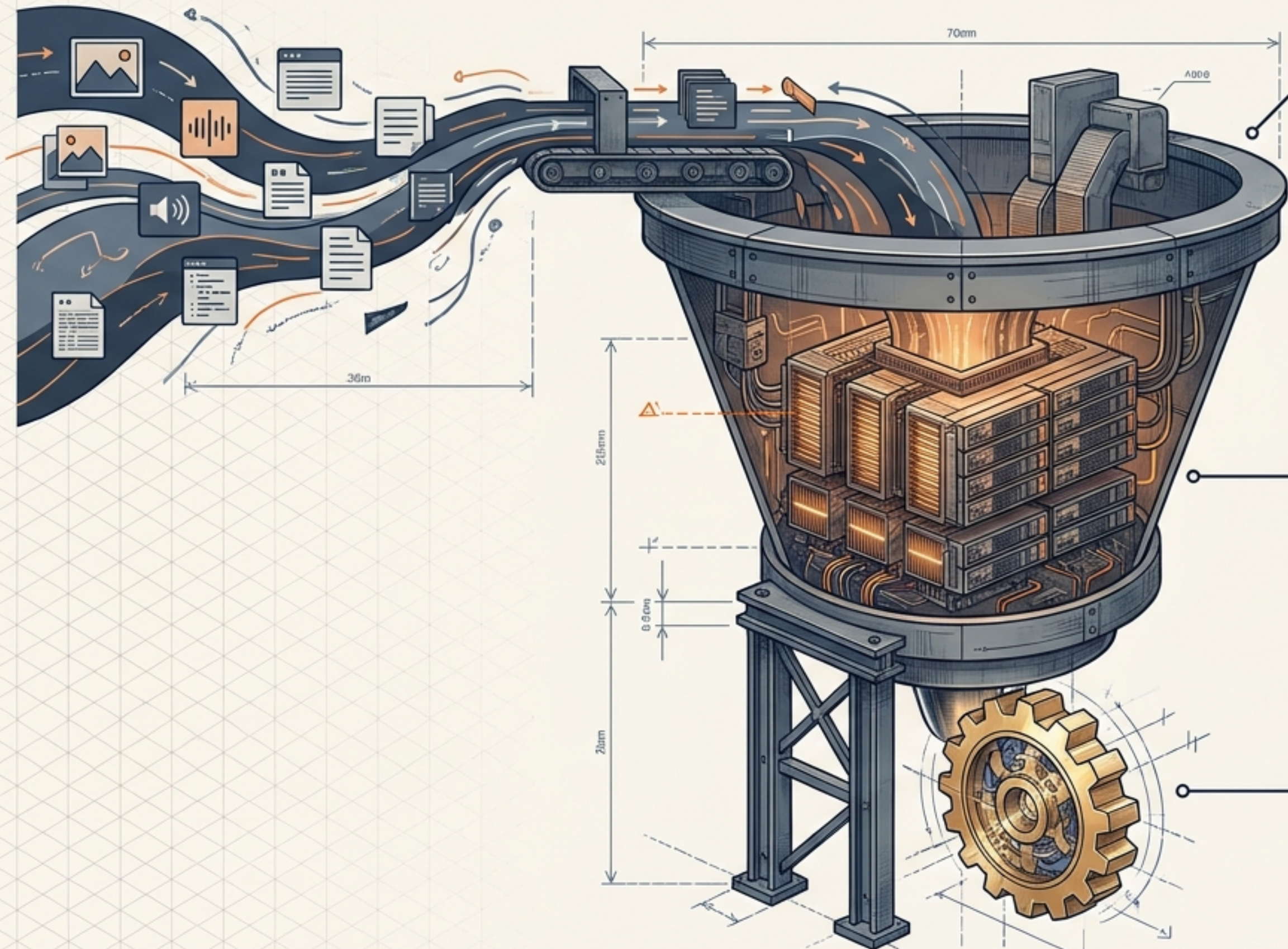
工作流与
编排范式

Vibe Coding
的破与立

金融级
安全防线

全景竞品与
物理局限

字节跳动巨量生态倒逼基础设施跃迁：从对话工具到智能工作伙伴



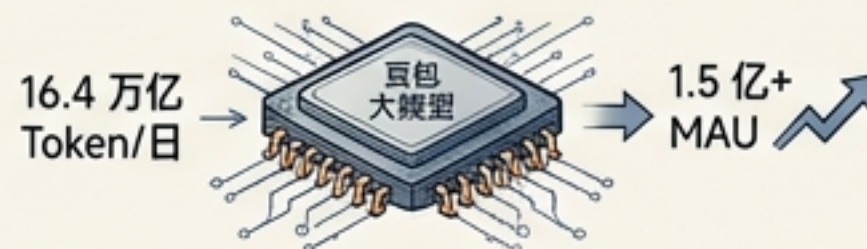
第一阶段：数据引力场

海量多模态数据积累。抖音电商约 4900 亿人民币 GMV 规模，高并发业务倒逼极限算力。



第二阶段：算力核变

世界级算力底座跃升。豆包大模型日均 Token 调用量达 16.4 万亿次，月活用户超 1.5 亿。



第三阶段：智能伙伴质变

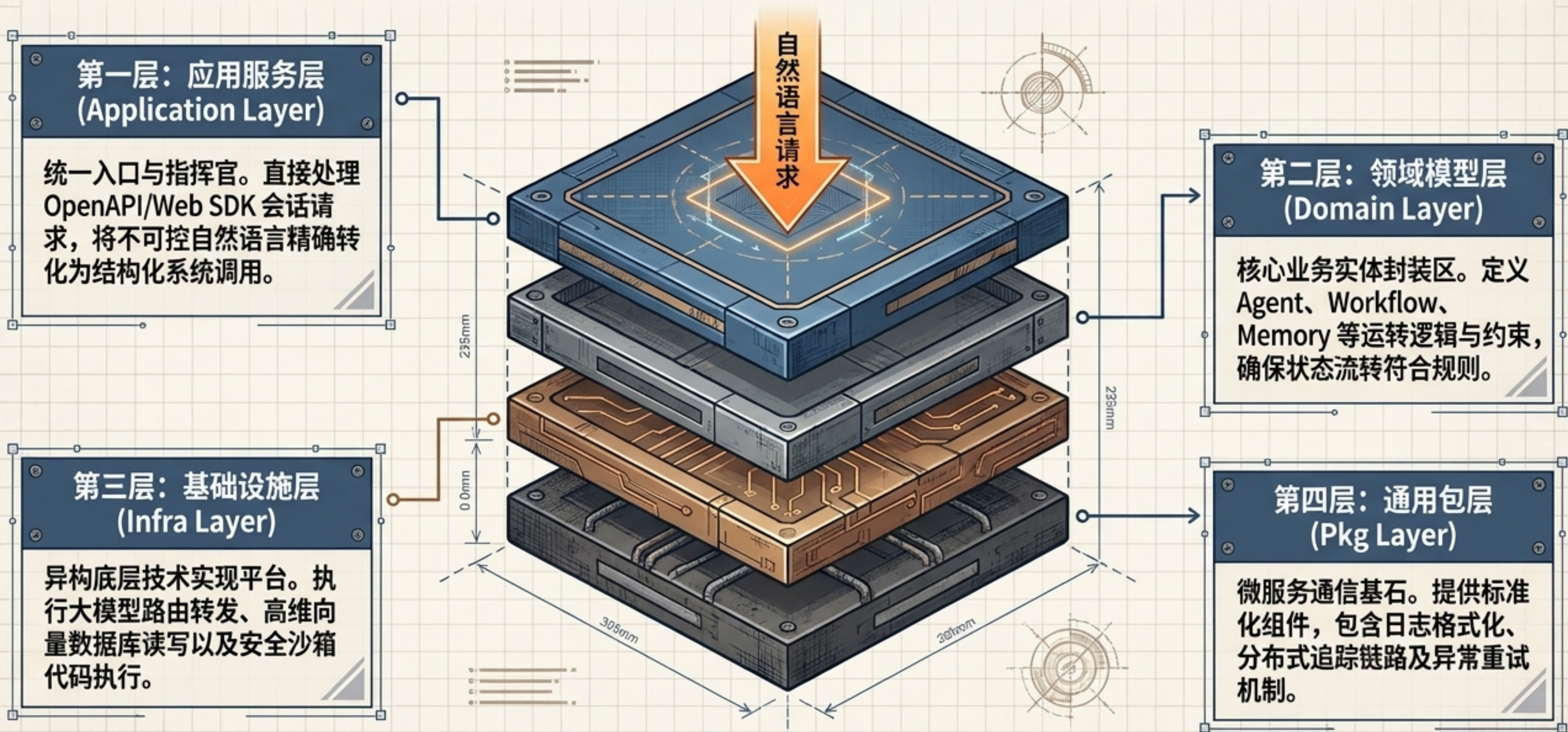
战略定位质变。彻底摆脱单一对话工具标签，重塑为支撑极端复杂企业逻辑的全链路自动化底层系统。



2026 双轨制演进与 Agent World 颠覆性愿景

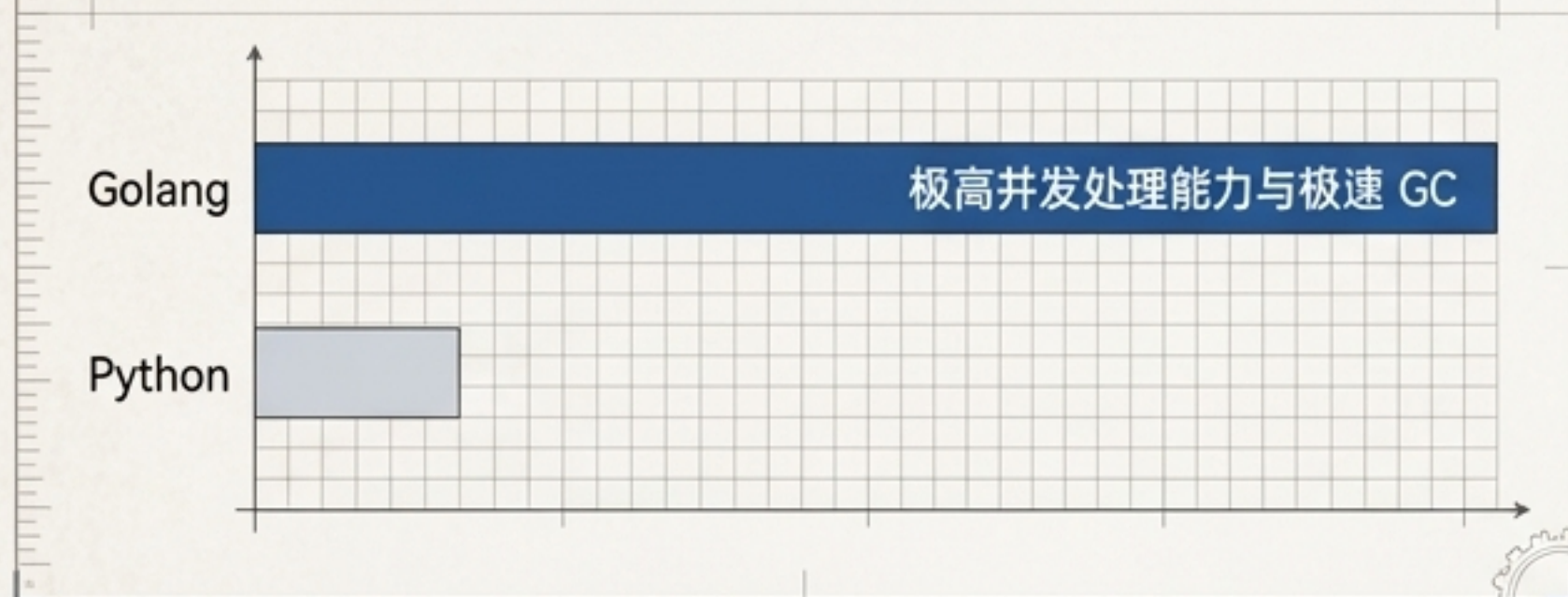


核心底层引擎：防腐解耦的领域驱动设计 (DDD) 分层架构



高并发技术栈选型与多租户精细化资源治理

性能驱动的技术栈

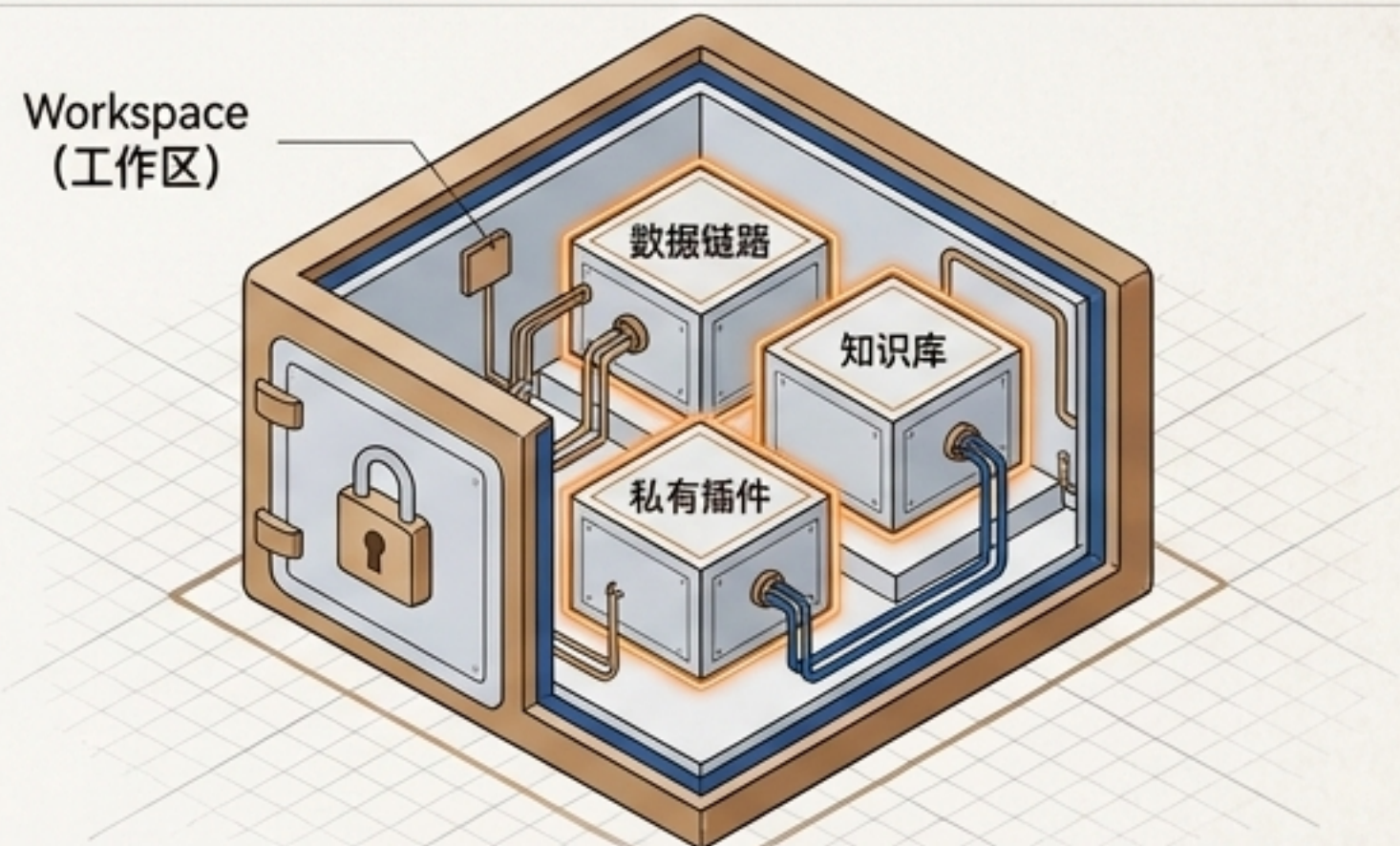


全栈 Go/React 架构：后端弃用 Python，基于 Golang 构建；前端基于 React/TS。

并发优势：在处理海量 Webhook 请求及内存自动垃圾回收具备压倒性优势。

执行引擎：基于 Eino 框架深度定制，融合 FlowGram 形成强大 DAG 计算模型。

企业级隔离防线 (Workspace)

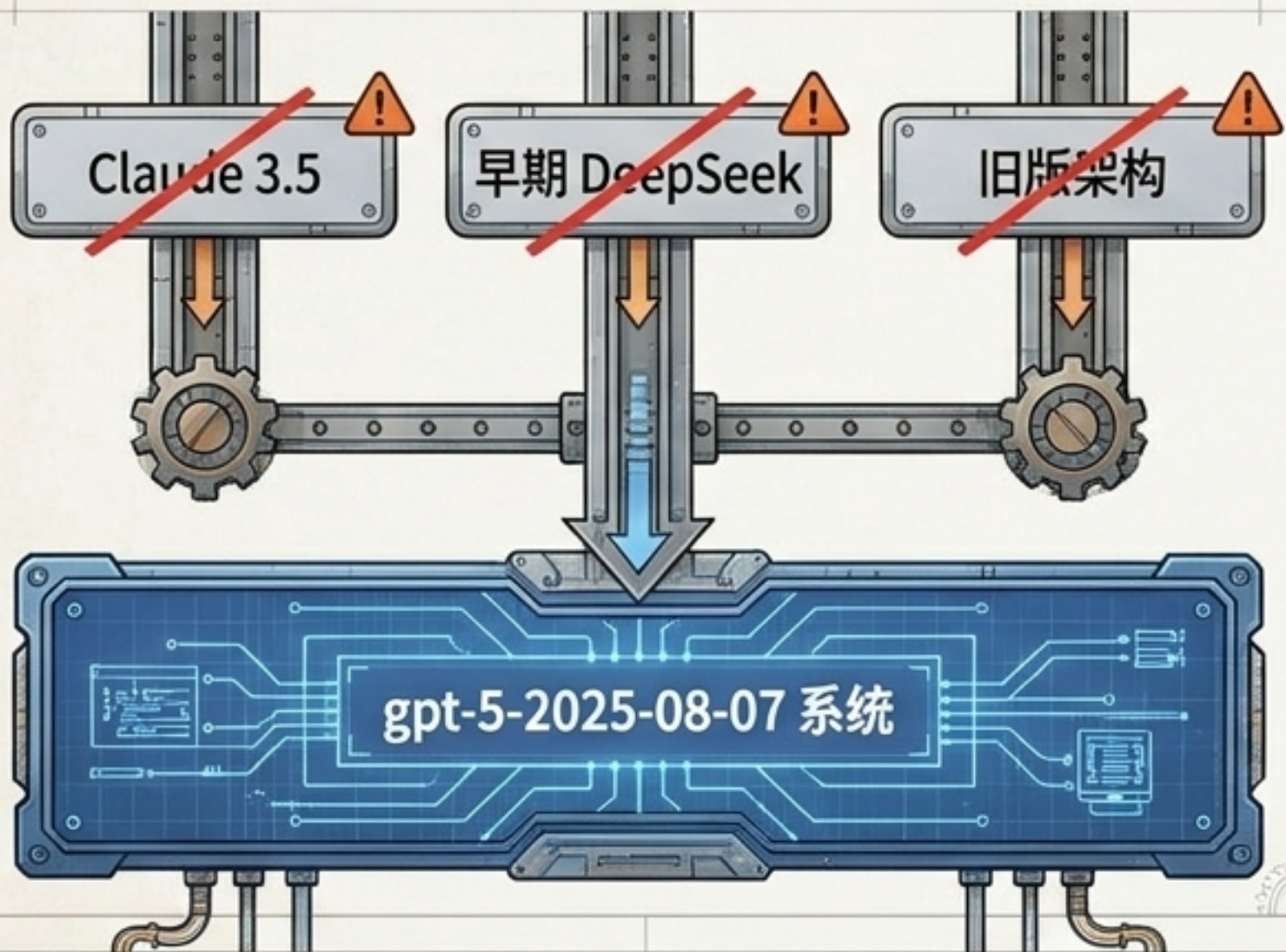


物理与逻辑隔离：工作区之间的数据与语料处于绝对的隔离状态。

细粒度权限控制：支持细化至 API 级别的角色访问控制 (RBAC)，兼顾协同与安全。

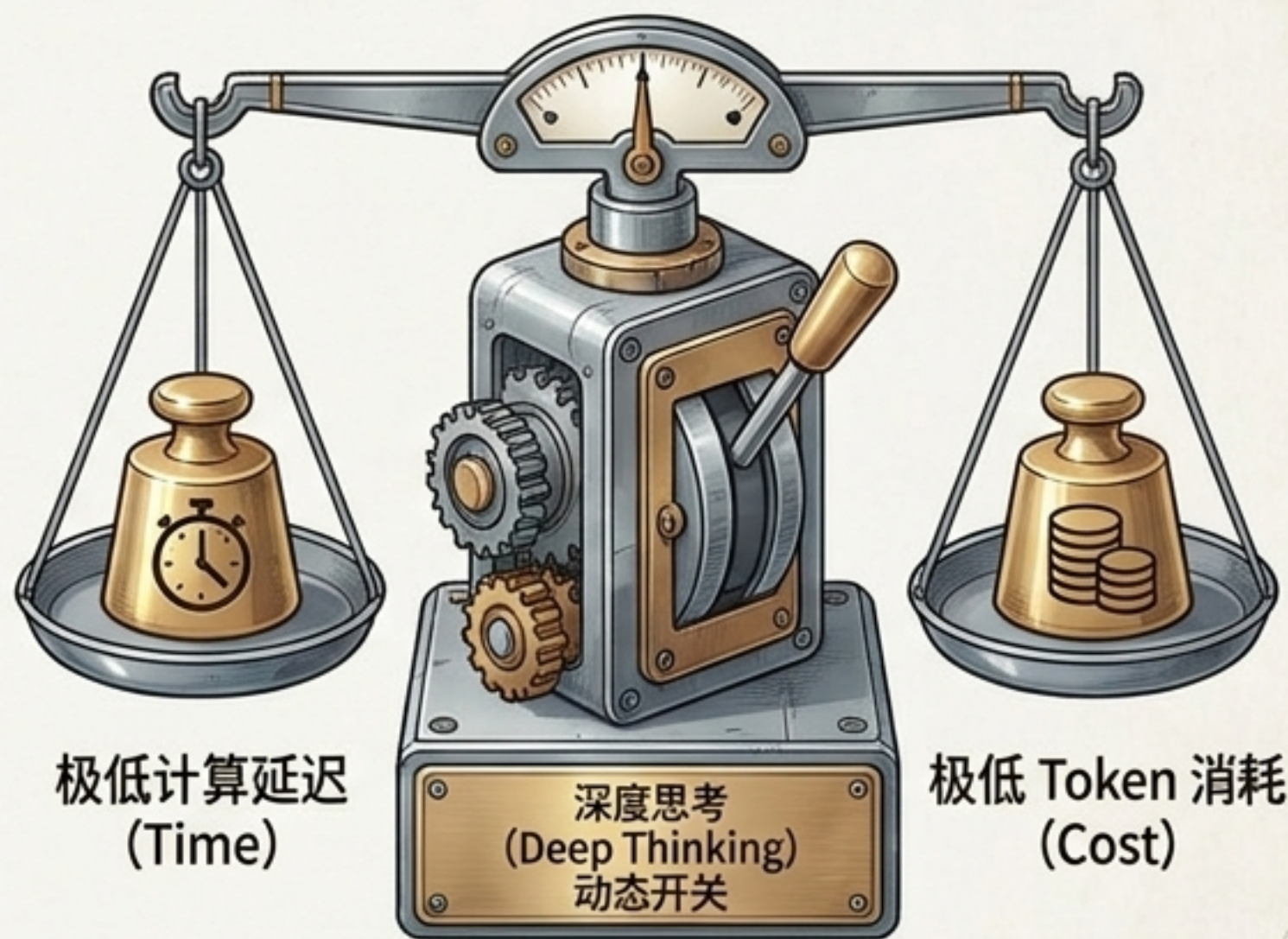
全球版模型矩阵大洗牌与深度思考算力精控

模型收敛



- **激进高维收敛：**退役冗余模型，强力向高维多模态系统集中。
- **统治级表现：**幻觉率较 GPT-4o 大幅降低 45%，在代码生成与极长周期企业工作流程中具备压倒性优势。

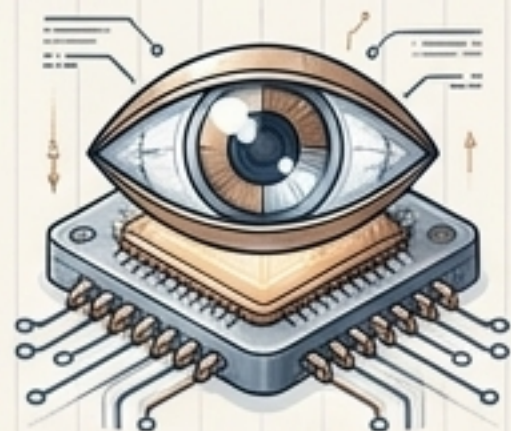
深度思考算力精控



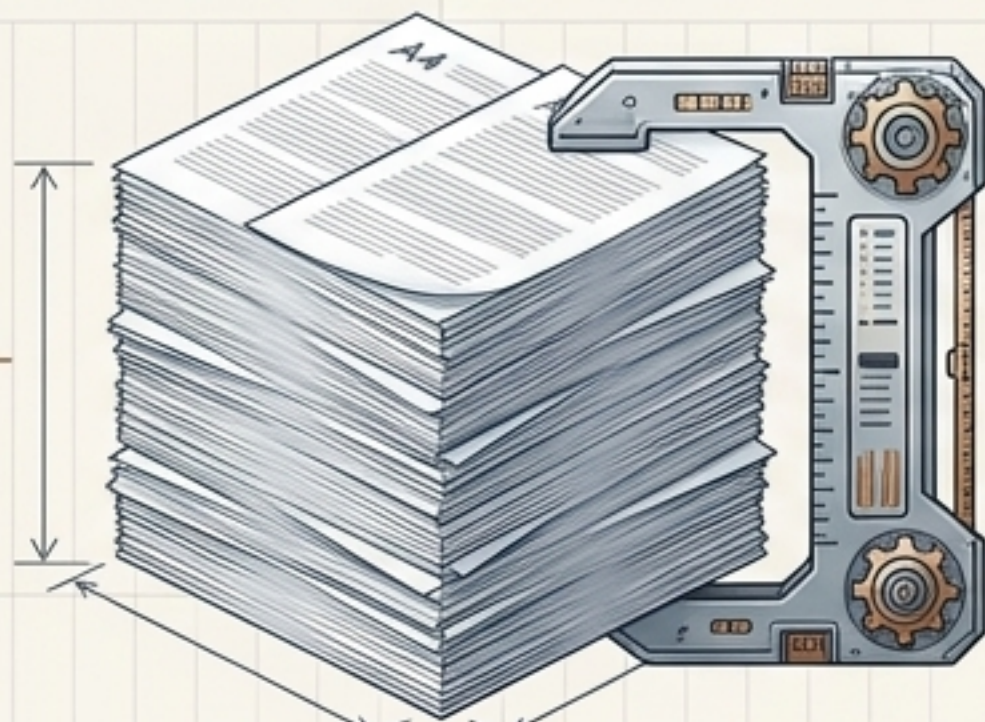
算力精细化压榨：允许在特定交互节点动态启停慢思考推理机制，完美平衡极致推理能力与算力成本。

国内算力质变：Doubao-Seed-2.0 Pro 旗舰矩阵的越级对标

依托火山引擎与豆包底座深度绑定，逐步兼容本土异构 API 代理生态。



原生支持高精度多模态视觉推理



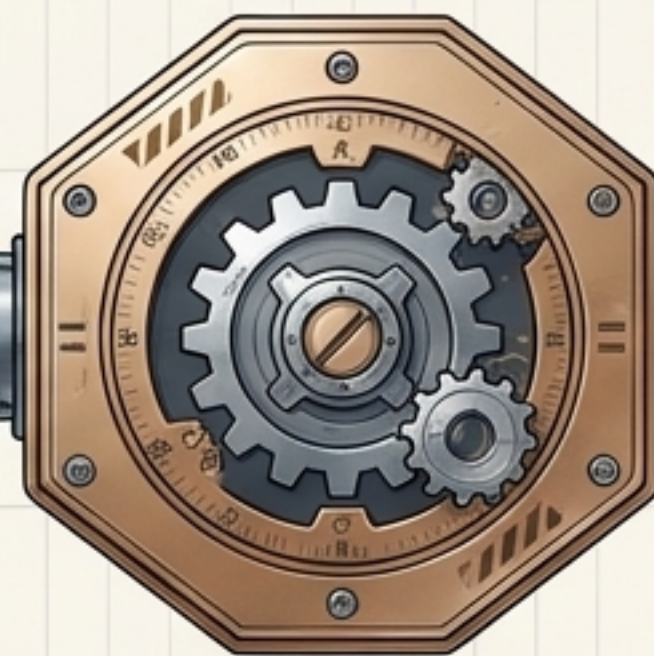
256K 超大窗口容量 /
等效约 384 页 A4 纸纯文本吞吐



AIME 2025 高级数学推理测试，
顶级智力标杆。



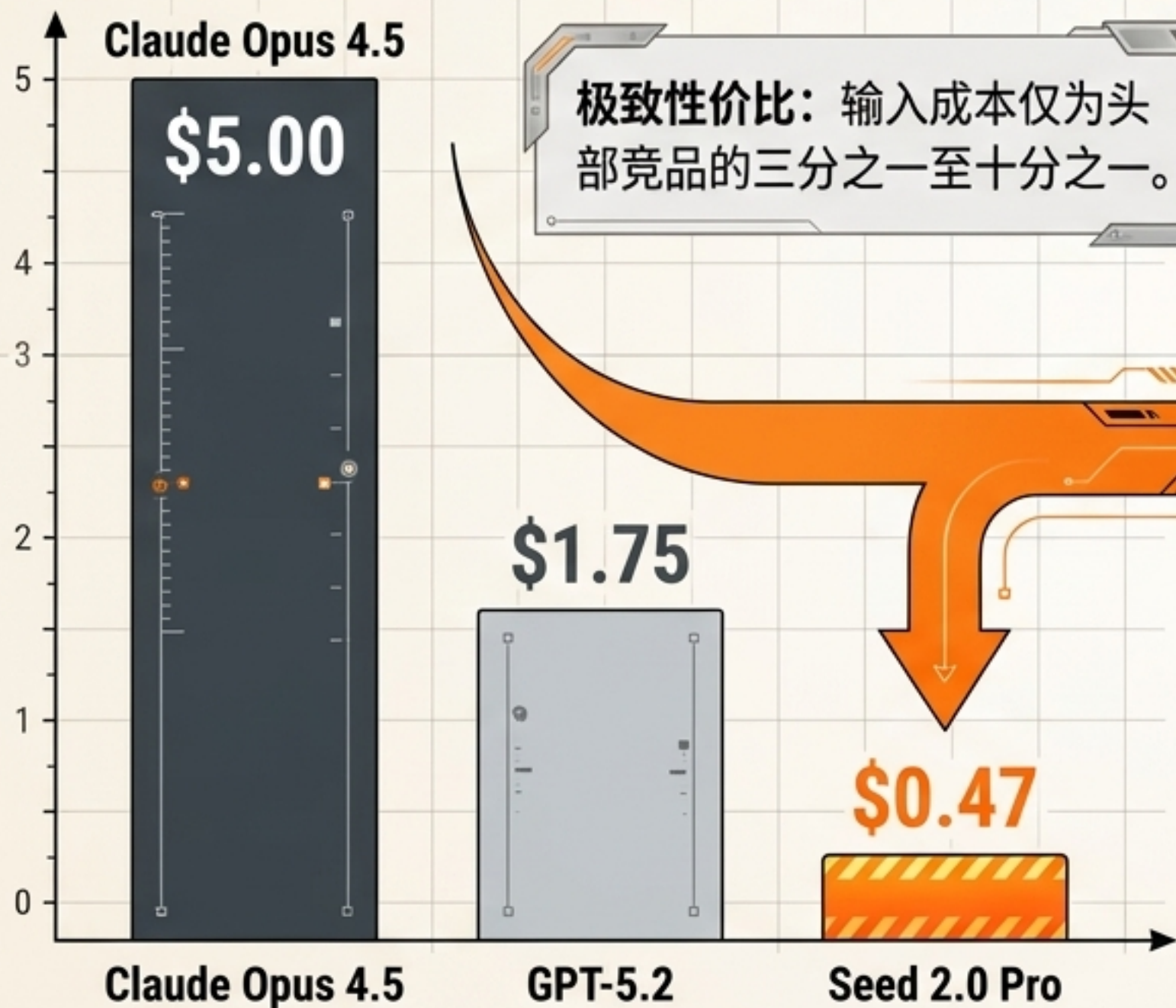
Codeforces 算法评测顶级水平，
掌控全栈代码生成。



工业级零容错能力，可直接承担 CAD
自动化审查与长周期科研文献梳理。

颠覆性算力经济学：白菜价 Token 彻底解开架构枷锁

每百万输入 Token 成本对比 (USD)



暴力降本释放出的架构红利



- 告别预算梦魇：复杂企业工作流的数百次循环纠错不再引发 Token 消耗暴雷。

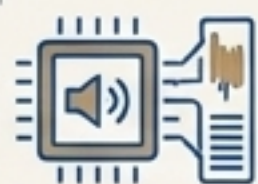


- 多重反思机制 (Reflection)：敢于在 Agent 工作流程中肆意引入交叉审验与多轮自我纠错。



- 大规模并发批处理：毫无顾虑地并发清洗万级结构化日志，奠定 Agent World 重度计费底气。

多模态革命：极致的低延迟语音与零样本情感生态



底层引擎革新

深度集成先进的 TTS 系统
FunAudioLLM (CosyVoice 3.0) ,
重塑人机拟真交互底座。



语种覆盖与零样本克隆

原生支持 9 种主流语言与超 18 种国内方言 (含粤语、四川话等)。仅需极少音频样本即可实现跨语种的完美声音克隆 (Zero-shot)。

极限双流输出延迟低至
150 毫秒



π χ / β



细节发音修复 (Pronunciation inpainting)

突破机械音瓶颈：原生支持拼音与英语 CMU 音素级精细修复。
可极其灵活地控制语音情感、起伏、语速与音量，极大拓宽数字客服与虚拟陪伴的商业落地边界。

深度 workflow 引擎：突破单线拥堵的并行批处理机制

DAG 可视化拓扑流转：严格依赖前置边缘结果，条件分支自动剪枝跳过无效路径，硬性约束模型不可靠性。

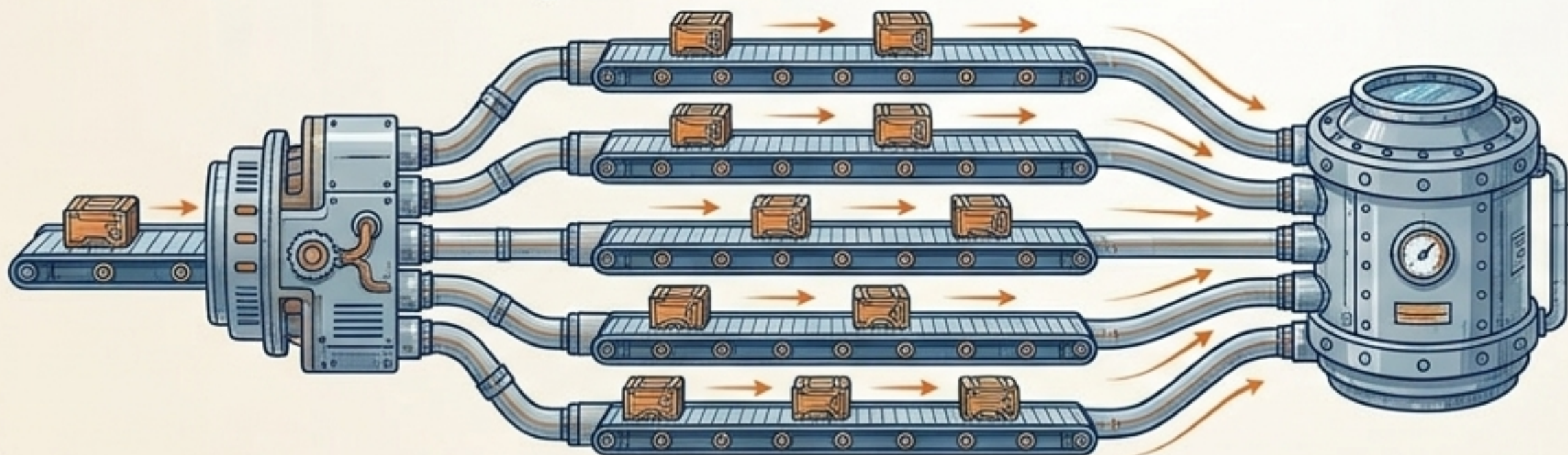
串行循环节点 (Serial Loop) 的瓶颈



危险：按索引逐个遍历极易导致耗时随长度线性爆炸，触碰系统超时红线。



并行批处理 (Parallel Loop) 的破局



MapReduce 架构提效 **5 倍以上**：
并发启动实例处理，最后在独立
内存缓冲池中拼接汇总，彻底规
避数据污染与超时崩溃。

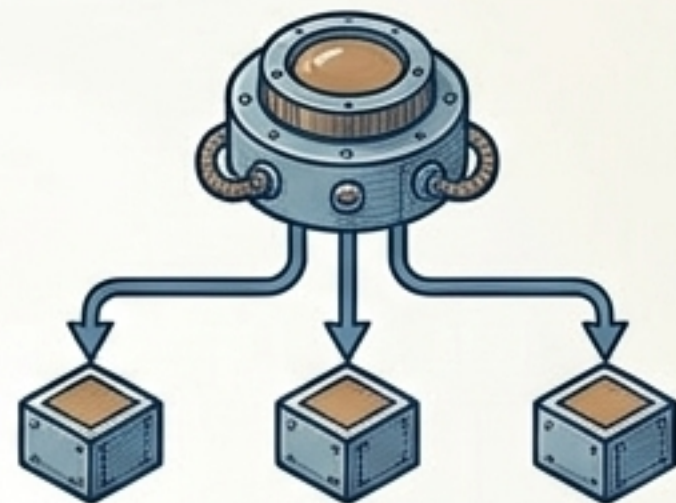
业务拆解与架构设计：四大核心多智能体协作范式

顺序流水线 (Pipeline)



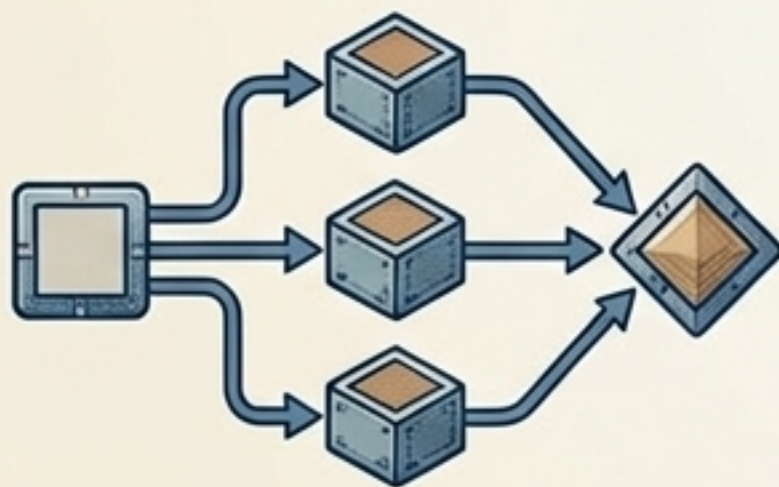
逻辑：各司其职，严密线性接力（如：提取 → 核查 → 法务润色）。
适用：步骤依赖绝对明确的渐进式任务。

监督者模式 (Supervisor)



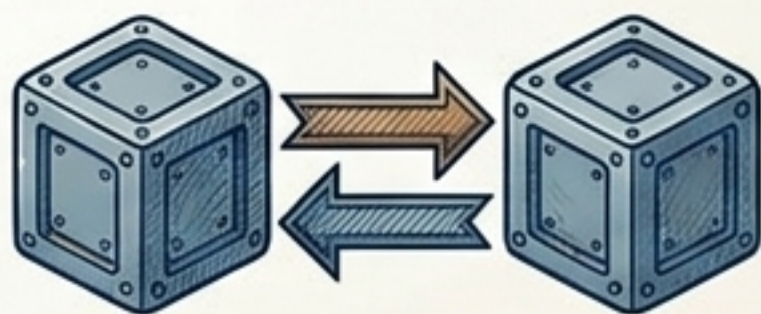
逻辑：中央协调器不干粗活，负责拆解宏大目标，动态委派专家 Agent 并监控闭环。
适用：涉及跨领域工具调用的高复杂目标求解。

发散与汇聚 (Fan-out / Fan-in)



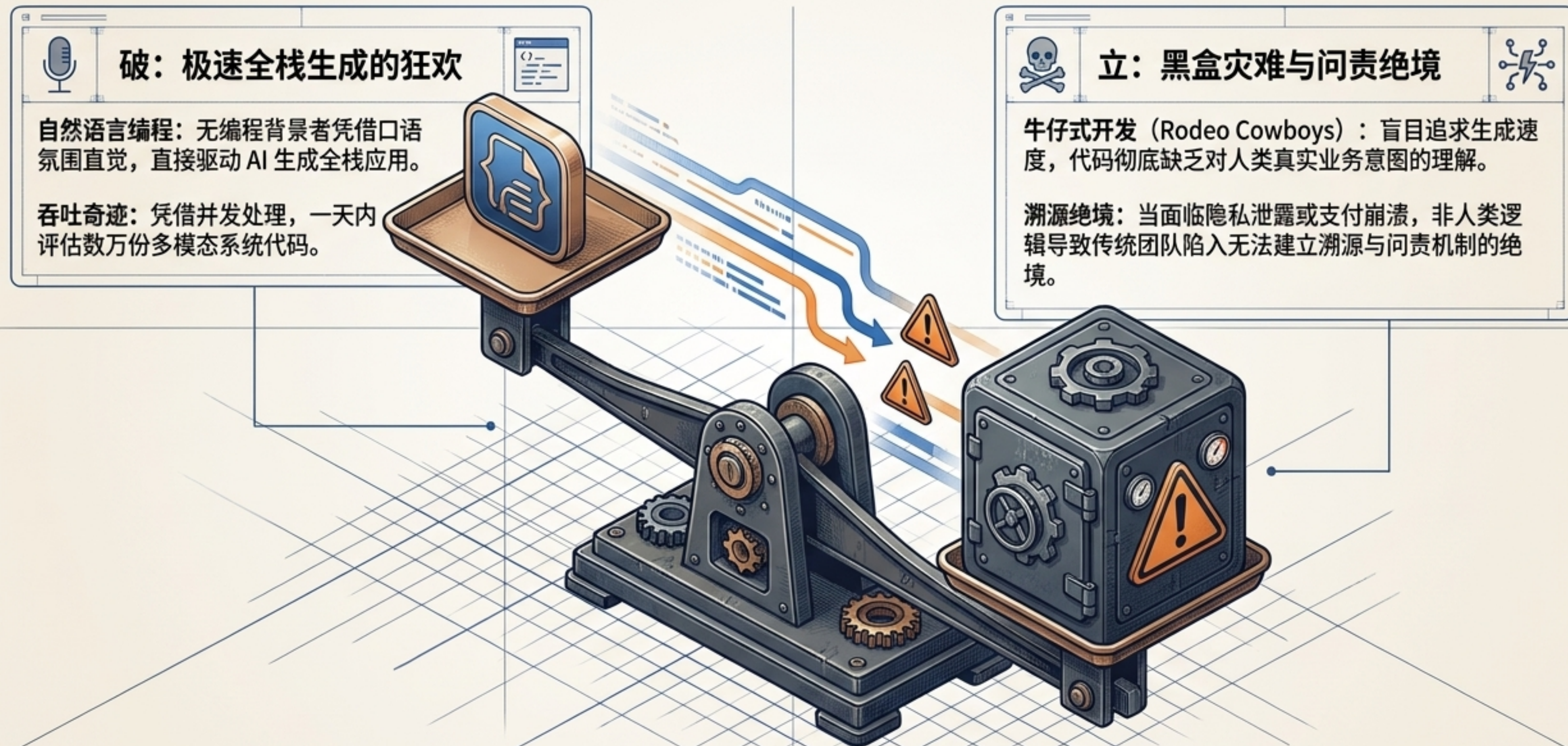
逻辑：主控瞬间唤醒平行子任务并发处理异构数据，交由最终节点浓缩提纯。
适用：海量并发处理与多视角数据融合。

博弈辩论模式 (Debate)

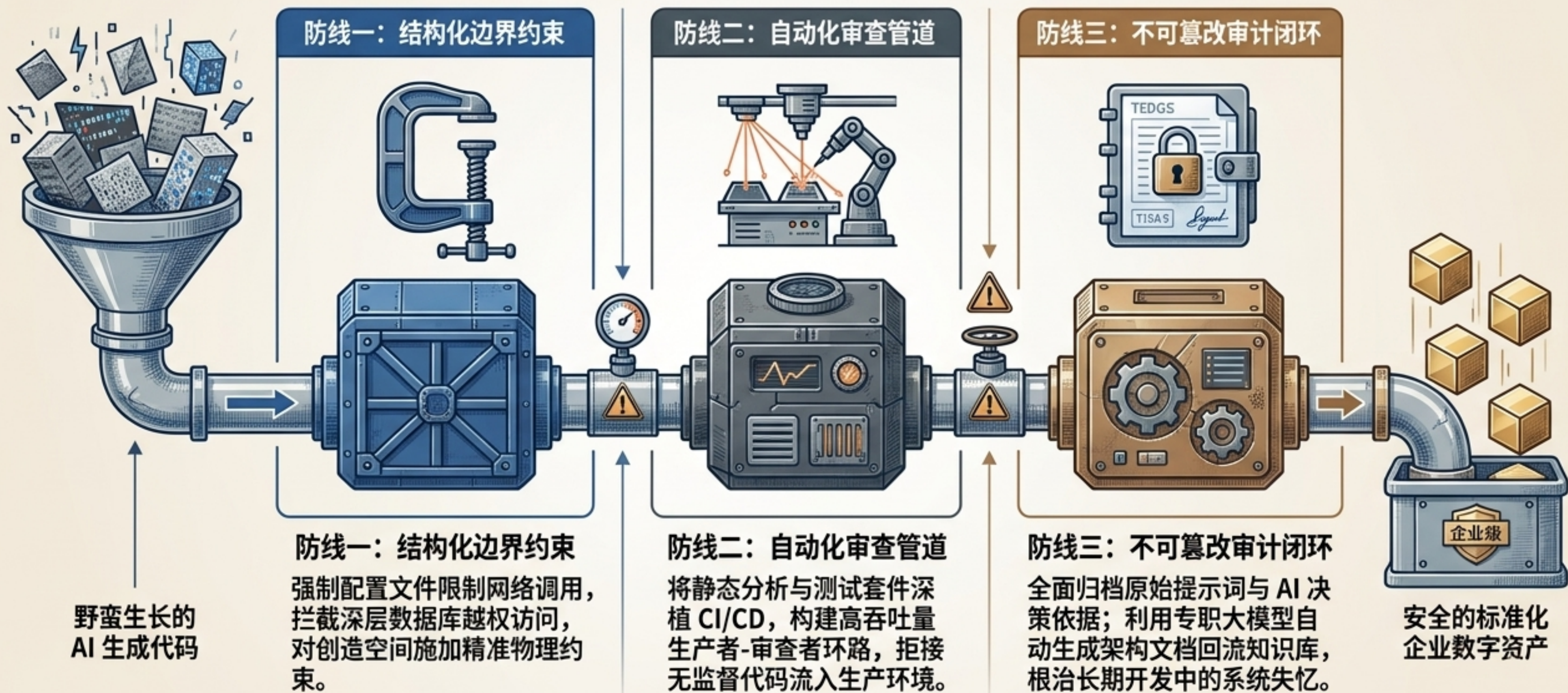


逻辑：互斥立场的智能体交叉辩论、反驳与妥协，消除单一模型盲点。
适用：需要极高鲁棒性与零偏见的商业决策输出。

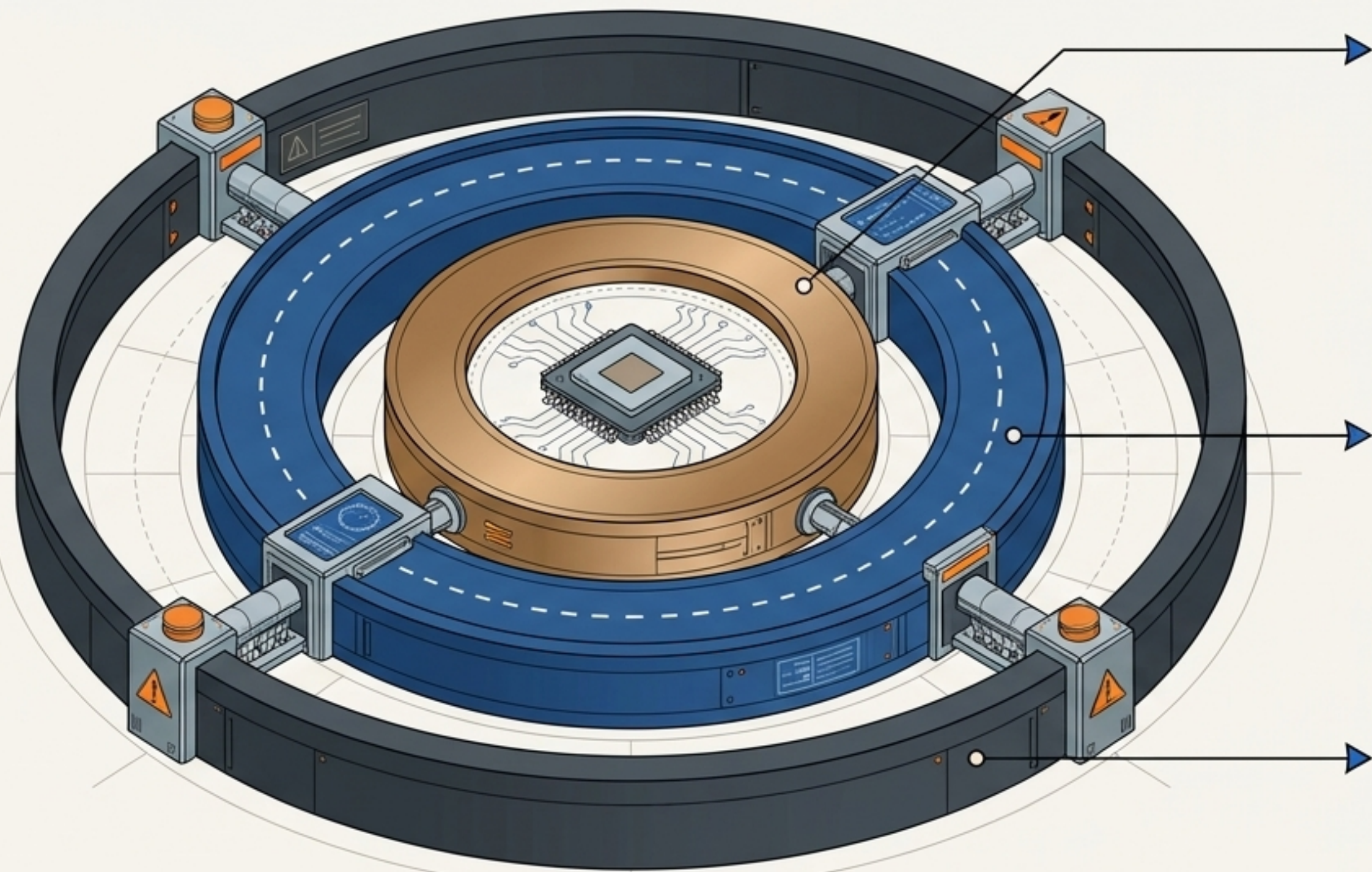
颠覆与阵痛：Vibe Coding 革命及其引发的企业级问责危机



构建企业级护城河：高吞吐工程治理与合规防线



突破信任红线：端云协同机密计算与企业旗舰版网络防御



内核：硬件 TEE 安全岛

豆包手机底层安全岛锁定推理与明文流转；云端依托 SecuDB 机密计算实现细粒度数据可用不可见。

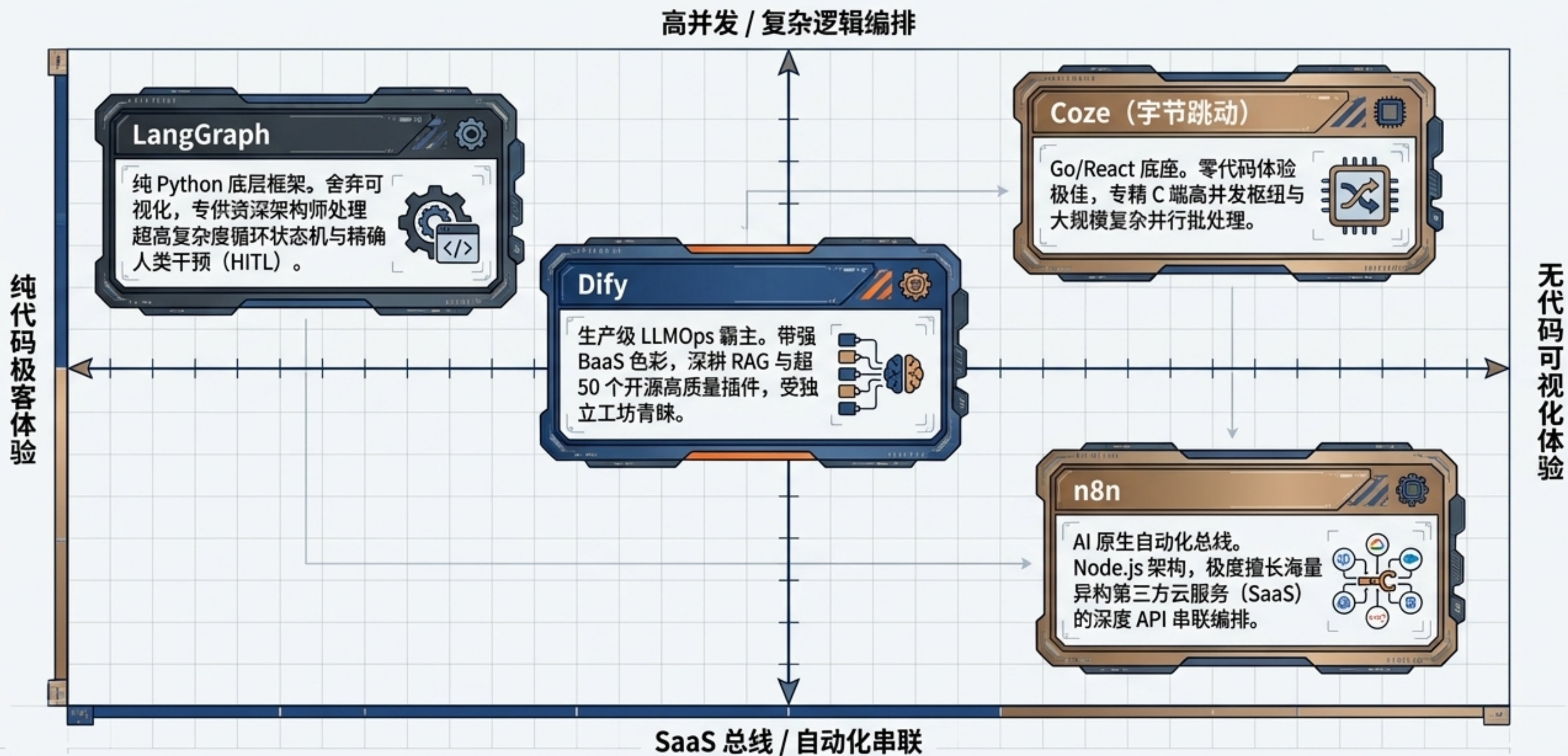
中继：假名化脱敏与加密

端云传输强制所有身份标识实施严格假名化处理，并全程采用高强度端到端加密防护。

外延：旗舰版深层网络边界

全面支持 SAML/OIDC 的 SSO 单点登录实现毫秒级权限阻断。提供 VPC 专线与防火墙，所有 API 探明请求死死封锁在私有虚拟网内。

2026 AI 应用平台全景：主流架构定位与能力对标矩阵



物理枷锁、工程梦魇与 Agent World 最终战略定局

SaaS 平台红线与工程挑战

计算熔断： 工作流最高 1000 节点上限，图像生成限 4 次/秒并发。

存储瓶颈： 单知识库仅 10GB 容量，彻底阻断其接管 TB 级中台的可能。

多智能体工程梦魇： 模型概率输出引发执行漂移与级联幻觉雪崩，极度缺乏全自动化测试基建导致可观测性沦为盲区。

核心价值锚定与前瞻

- 1 禁区：** 绝不适用于要求 100% 确定性的零容错金融清算或工业骨干网控制。
- 2 核心发力点：** 牢牢锚定在知识密集型摘要提取、非确定性内容流水线与创新辅助决策。
- 3 终局展望：** 必须依赖底层治理规范的跃升与全链路可观测基建，方能彻底驾驭企业级数字实体协同办公范式。