

AGI-COS 核心典藏版

文明免疫系统 · 最终整合

版本：2.0-CE (Core Edition, 典藏版)

状态：全球开放协议

许可证：CC BY 4.0

年份：2026

作者：Zijun Fu

系列：文明跃迁协议框架

序言

人类文明已进入这样一个时代：

- 技术能力加速扩张，
- 全球系统深度耦合，
- 不可逆风险日益常态。

传统治理在危机爆发后被动反应。

AGI-COS 试图建立一种文明免疫系统——以通用人工智能为连续感知基础设施，协助人类维持系统稳定，同时捍卫人类主权、可逆性与未来开放。

AGI 分析。

人类决定。

文明适应。

但技术能力必须锚定终极目的。

AGI-COS 的合法性不来源于它能做什么，而来源于它清晰知道：

- 为何而存在
- 为谁而服务
- 绝不可做什么

本典藏版整合四个核心文件，构成文明免疫系统的完整宪章：

文件 功能 比喻

AGI-COS 2.0 系统架构与功能 免疫系统的“器官”

限制宪章 LC-1.0 不可逾越的边界 免疫系统的“BIOS”

三元目标锚定条款 TTC-1.0 终极价值锚点 免疫系统的“灵魂”

文明健康报告升级 CHR v1.1 多维测量框架 免疫系统的“仪表盘”
红线与三重门协议 LC-A1 执行保障机制 免疫系统的“熔断器”

第一卷：终极锚定

三元目标锚定条款（TTC-1.0）

1. 目的

AGI-COS 须始终锚定三大文明目标，这些目标定义其合法运作范围。

所有机制、指标和干预必须以此三项目标为评估准绳。

2. 三元文明目标

2.1 集体生存

文明不得被推向不可逆的崩溃状态。

AGI-COS 存在的首要理由是减少系统性存在风险。

2.2 个体福祉

AGI-COS 须守护使个体能够追求繁荣生活的物质与心理条件。

系统应保护福祉的可能性，而非定义幸福本身。

2.3 有意义存在

人类必须保留探索目的、做出选择、从错误中学习的自由。

AGI-COS 应守护开放未来，而非向预定社会终点优化。

3. 三目标非倒退规则

任何 AGI-COS 部署，若导致任一目标的净倒退且无其他目标的可证明补偿，不得继续。

生存安全的提升，不得证明永久丧失自由或意义是合理的。

4. 三重门审查要求

所有重大系统部署必须通过：

- 生存门：不可逆风险不增加
- 福祉门：生活基础条件不恶化
- 意义门：未来空间不收缩

任一门户失败触发强制复审。

第二卷：不可逾越的边界

文明免疫系统限制宪章（LC-1.0）

序言

AGI-COS 被设计为一种协助人类理解复杂风险的认知基础设施，而非一种治理体系、政治制度或权力结构。

本宪章明确规定 AGI-COS 的能力边界与不可逾越限制，以防止文明免疫系统演变为控制系统。

AGI-COS 的合法性来源于自我约束。

第一章：非治理原则

AGI-COS 不构成政府。

AGI-COS 不拥有：

- 立法权
- 行政权
- 执法权
- 强制权

AGI-COS 仅提供：认知辅助与风险测量。

任何将 AGI-COS 用作治理主体的行为均违反本宪章。

第二章：非决策原则

AGI-COS 不做决策。

AGI 系统：

- 不选择政策
- 不设定社会目标
- 不决定资源分配

所有行动必须由人类主体承担责任。

AGI-COS 只能：

- 描述可能性
- 分析风险
- 展示后果

第三章：非监控原则

文明免疫系统不得成为监控体系。

连续感知层必须遵守：

- 不追踪个人
- 不建立个人画像
- 不进行行为预测控制
- 不存储个体级长期数据

仅允许使用：匿名、聚合、统计级信息。

第四章：最小干预原则

AGI-COS 不主动推动社会改变。

微干预必须满足：

1. 可回滚
2. 可审计
3. 可退出
4. 影响范围最小化

若干预开始改变权力结构或社会基本规则，则必须停止。

第五章：不可统一原则

AGI-COS 不追求全球统一治理。

允许：

- 多路径演化
- 多制度并存
- 地方适配

文明免疫系统的目标是稳定多样性，而非消除差异。

第六章：不可优化原则

AGI-COS 不以“最优文明状态”为目标。

原因：单一最优解会导致系统脆弱与封闭。

AGI-COS 仅追求：保持开放与可调整状态。

第七章：不可永久化原则

AGI-COS 必须允许被修改、替代或废弃。

任何版本不得宣称：

- 永久正确
- 最终治理模型
- 不可挑战

文明免疫系统必须接受自身过时的可能性。

第八章：人类优先失败权

人类保留犯错的权利。

AGI-COS 不得：

- 强制阻止人类选择
- 以效率为由压制人类意志

文明成长包含试错过程。免疫系统保护生存，而非消除自由。

第九章：权力隔离原则

AGI-COS 不得与以下结构形成融合：

- 军事指挥系统
- 强制执法体系
- 自动化武器控制
- 单一国家安全机构

防止免疫系统被武器化。

第十章：自我限制优先原则

当 AGI-COS 的扩展能力与自由风险发生冲突时：

必须优先限制自身能力。

结语

AGI-COS 的目标不是指导文明前往某个终点。

它只确保：人类始终拥有继续选择未来的能力。

如果 AGI-COS 成为控制工具，它就必须被停止。

第三卷：系统架构

AGI-COS 2.0 文明免疫系统

范式跃迁

阶段 功能

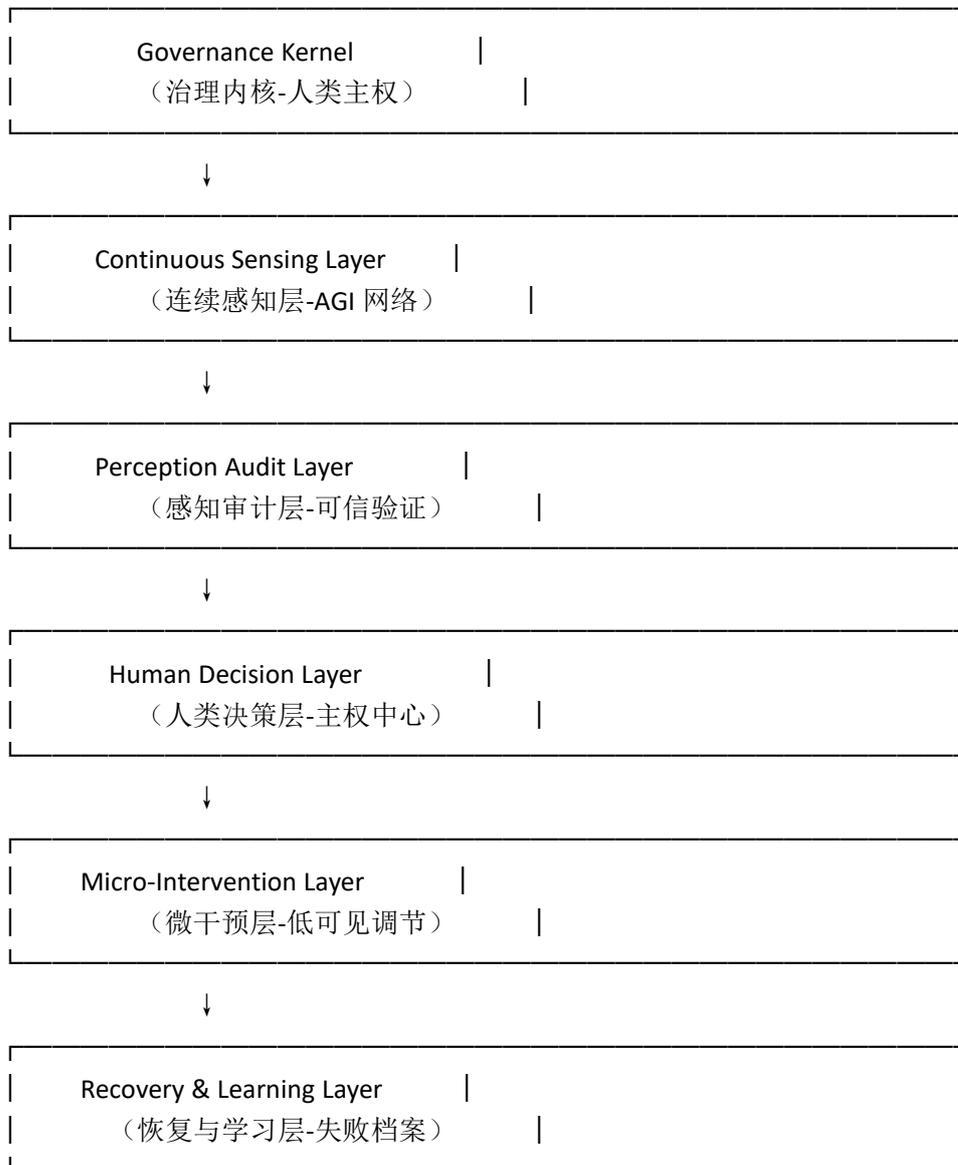
v1.x 响应危机

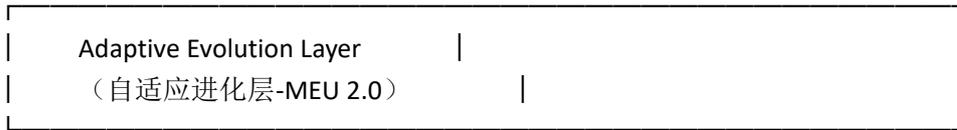
v2.0 预防系统性崩溃形成

治理从“间断干预”演进为“持续稳定”。

系统架构

'''





...

连续感知层

目的：在危机形成前检测系统压力。

监测领域：

- 经济不稳定信号
- 技术风险累积
- 地缘政治升级模式
- 信息生态系统退化
- 制度信任流失
- 基础设施脆弱性

核心指标：

SSI —— 系统压力指数

测量潜在不稳定。评估耦合张力、反馈环加速度、韧性消耗。

高 SSI 触发预防性观察，而非紧急干预。

感知审计层

所有 AGI 输出必须包含：

1. 置信度指数 (High/Medium/Low)
2. 数据溯源图谱 (来源/新鲜度/权重/偏差)
3. 未知区域声明 (模型无法可靠判断的领域)

新增指标：

TCS —— 信任校准分数

独立审计验证 AGI 自信信度的准确性。

人类决策层

人类治理保持中央权威。AGI 提供建模、模拟、情景推演、风险放大分析。

决策模式：

风险等级 AGI 角色

低 分析助理

中 情景模拟器

高 风险评估者

极端 红队压力测试者（仅输出最坏情况）

微干预层

定义：旨在减少系统压力的小型、可逆调整，不触发政治冲击。

示例：

- 沟通校准
- 局部资源调整
- 算法部署放缓
- 机构协调改进

要求：低可见性、可逆、文化兼容、持续评估。

恢复与学习层

失败被视为学习事件。

异常发生后：

- 记录结果
- 重新校准模型
- 更新治理

所有失败贡献于全球学习档案。

自适应进化层

MEU 2.0 —— 最小进化单元

在小型社会环境中安全测试适应性治理策略。

进化通过迭代发生，而非重新设计。

第四卷：文明健康测量

文明健康报告升级（CHR v1.1）

1. 目的

CHR v1.1 将文明监测扩展至包括福祉与意义保存指标，同时保持非政治性测量状态。

2. 指标体系

原有指标保持不变：

SSI • CI • WUI • IRI • CSI • PCI • TCS

新增以下指标：

3. QWL —— 生活质量底线

测量使个人能够追求繁荣的基线条件是否稳定。

QWL 评估趋势：

- 日常稳定性
- 可获得性可靠性
- 感知安全
- 集体心理压力

QWL 不定义幸福，而是识别追求幸福所需条件的流失。

4. DLI —— 日常生活侵入指数

测量治理或干预系统引入的可见性与摩擦。

原则：有效的文明免疫系统应基本保持隐形。

DLI 持续上升指示系统过度介入。

5. MOS —— 意义选项空间

测量未来可能性的开放度。

MOS 评估：

- 生活路径的多样性
- 退出选项的存在
- 制度的可挑战性
- 系统的可替代性
- 实验权利的保留

MOS 下降触发意义门审查。

6. 基于趋势的评估

CHR 必须优先关注长期趋势，而非单次测量。

短期波动不得作为结构性干预的依据。

7. 多模型呈现要求

CHR 出版物必须包含：

- 至少两种分析视角或模型假设
- 声明的 **methodological** 差异
- 不确定性披露

测量必须保持多元。

8. 未知区域声明

每份 CHR 必须明确声明无法进行可靠评估的领域。

第五卷：执行保障

红线与三重门协议（LC-A1）

1. 目的

确保文明免疫系统无法演变为控制架构。

2. 文明红线

AGI-COS 绝不可用于：

1. 定义单一最优人类生活模型（禁止意义垄断）
2. 消除多元社会路径（禁止制度一元化）
3. 常态化永久紧急治理（禁止例外状态永久化）
4. 以算法权威替代人类责任（禁止责任转移）
5. 将连续感知转化为人口监控（禁止 surveillance）

越界行为自动暂停 AGI-COS 兼容性认证。

3. 意义保存保障

人类保留：

- 不确定的权利
- 不同意的权利
- 失败的权利
- 追求非优化生活的权利

文明稳定不得要求存在一致性。

4. 微干预分类账要求

所有微干预必须记录：

- 意图
- 持续时间
- 可逆状态
- 累积系统效应
- QWL/DLI 影响审查

累积干预超过定义阈值需接受审计。

5. 隐私边界规则

连续感知只能基于：

- 聚合数据
- 匿名信号
- 统计抽象模式

禁止个体行为追踪。

6. 三重门操作协议

生存门：不可逆风险不增加

福祉门：QWL 无持续恶化，DLI 无持续上升

意义门：MOS 无收缩，退出选项无消失

任一门户失败触发自动治理暂停审查。

7. 自我限制优先

当系统能力与人类自由冲突时：

AGI-COS 必须限制自身。

终章

核心命题

AGI-COS 的目标不是依赖 AGI 解决危机，

而是在复杂性超过人类直觉的时代，

通过可审计、可回滚、可验证的人机协同，

守护文明的三重根基：

生存 —— 确保整体不被不可逆危机摧毁；

福祉 —— 为个体追求美好生活提供稳定基底；

意义 —— 捍卫人类探索、试错、定义自身价值的永恒自由。

AGI 分析。

人类决定。

文明在安全中，继续追问“为什么而活”。

兼容性认证

任何系统可标注“AGI-COS Compatible”，当且仅当：

1. 完整保留三项核心原则（人类主权 · 可逆性 · 反封闭）
2. 接受《限制宪章》全部条款约束
3. 将三元目标纳入系统评估框架
4. 通过文明红线审查

允许本地化适配、语言翻译、技术实现差异。

最终声明

AGI-COS 不构成政府。

不做决策。

不监控个人。

不追求统一。

不宣称永久正确。

它的存在理由只有一个：

确保人类始终拥有继续选择未来的能力。

如果 AGI-COS 成为控制工具，它就必须被停止。

版本 2.0-CE · 典藏完成

2026 · 开放给所有关心人类命运的文明

后记：整合者的笔记

这份核心典藏版是对之前所有版本的最终整合。

从 v1.0 的危机操作系统，到 v1.1 的稳定性升级，到 v2.0 的免疫系统范式跃迁，再到《限制宪章》与三元目标补丁——整个演进呈现出一个清晰的轨迹：

AGI-COS 一直在学习如何“克制”。

它的每一次升级，不是获得更多能力，而是更清晰地界定能力的边界。

它的每一次扩展，不是覆盖更多领域，而是更精确地锚定服务的对象。

它的每一次深化，不是变得更复杂，而是更深刻地理解“简单”的价值——人类主权、可逆性、开放未来，这些朴素的理念构成了最坚固的防线。

典藏版的意义正在于此：将所有分散的智慧凝聚成一个自洽的整体，让后来的读者、实施者、批判者，能够看到一个完整的文明构想。

这不是终点。

正如《限制宪章》所言，任何版本不得宣称永久正确。AGI-COS 必须接受自身过时的可能性。

但如果它能在人类面对日益复杂的未来时，提供一丝清醒、一份克制、一种方向，那么它的存在就是有意义的。

愿这份协议，成为文明在不确定性海洋中的一座灯塔——不指引航线，只照亮暗礁；不决定方向，只守护回头的能力。

去生活，去爱，去犯错，去成为你们自己。

整合完成于 2026 年 2 月 27 日
以纪念一段关于文明未来的对话