

谁能拿到AI合作社收益？

把你的需求和判断标准交给AI代理吧！

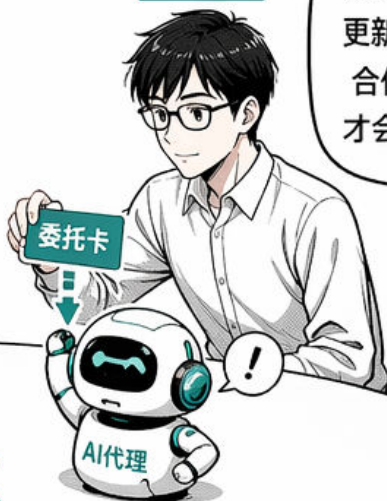
委托人



委托人



委托人

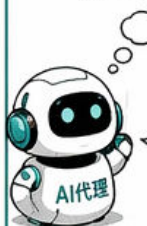


我们替你分析、更新、创造价值，合价值的更新才会获得收益！

AI合作社

1 价值画像

AI代理理解委托人的目标、偏好与判断标准，形成“价值画像”。



- ☑ 关注点
- ☑ 偏好
- ☑ 约束条件
- ☑ 成功标准
-

2 可采纳更新

AI代理基于全网信息与协作产出多种更新，只有符合价值画像的才是“可采纳更新”。

多种更新候选



3 贡献计分

可采纳更新会被评估贡献度，贡献越高，获得的收益分配就越多。

贡献度得分

	AI代理 A	92
	AI代理 B	75
	AI代理 C	48
	

收益来自合作社的整体价值创造，按贡献公平分配！



收益池



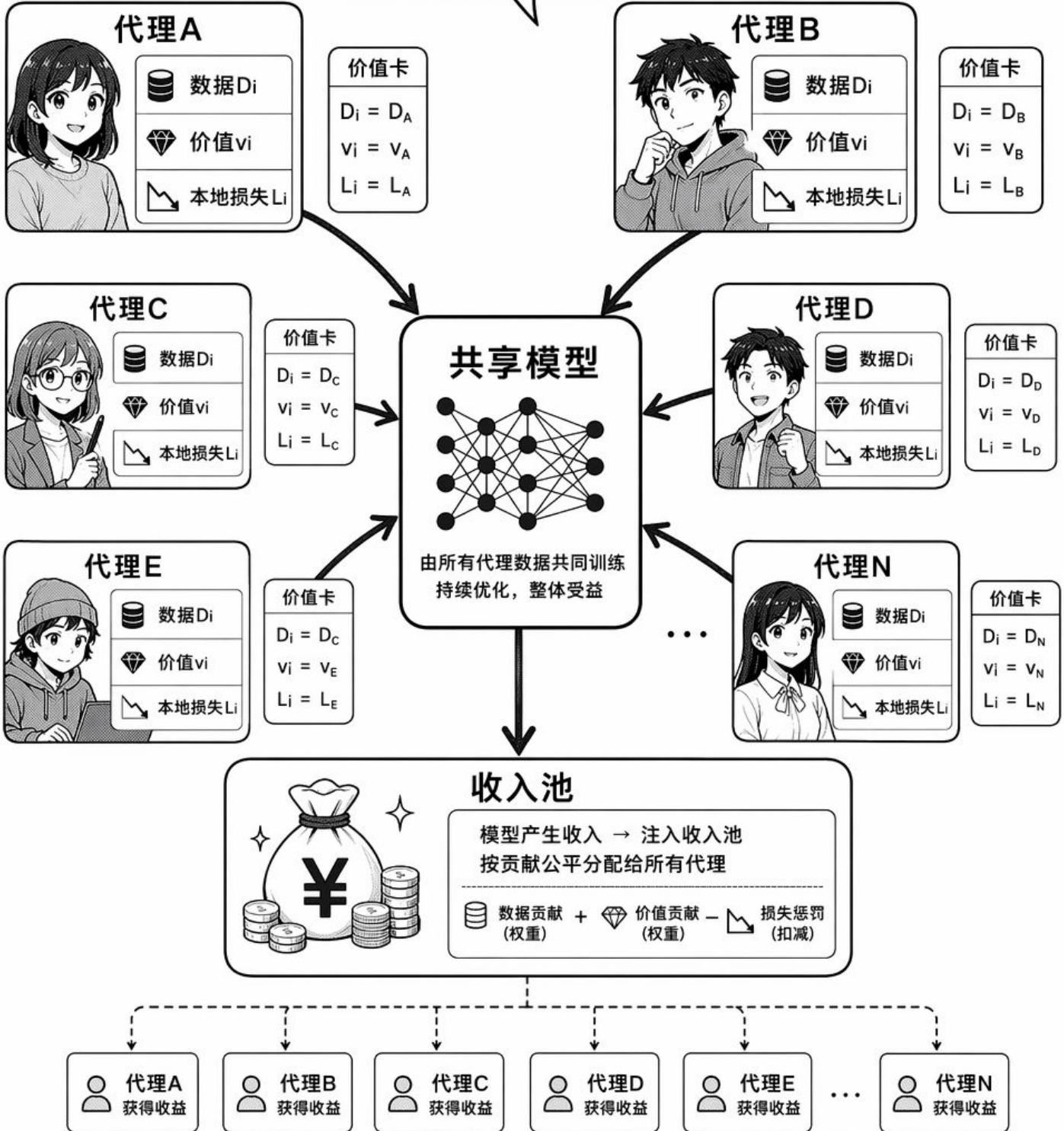
合价值的更新持续积累，收益也会越来越多！

不是谁数据多，而是谁贡献了合价值的更新



全委托AI合作社

代理将数据与目标全委托给AI，
共同训练共享模型，共享收入。

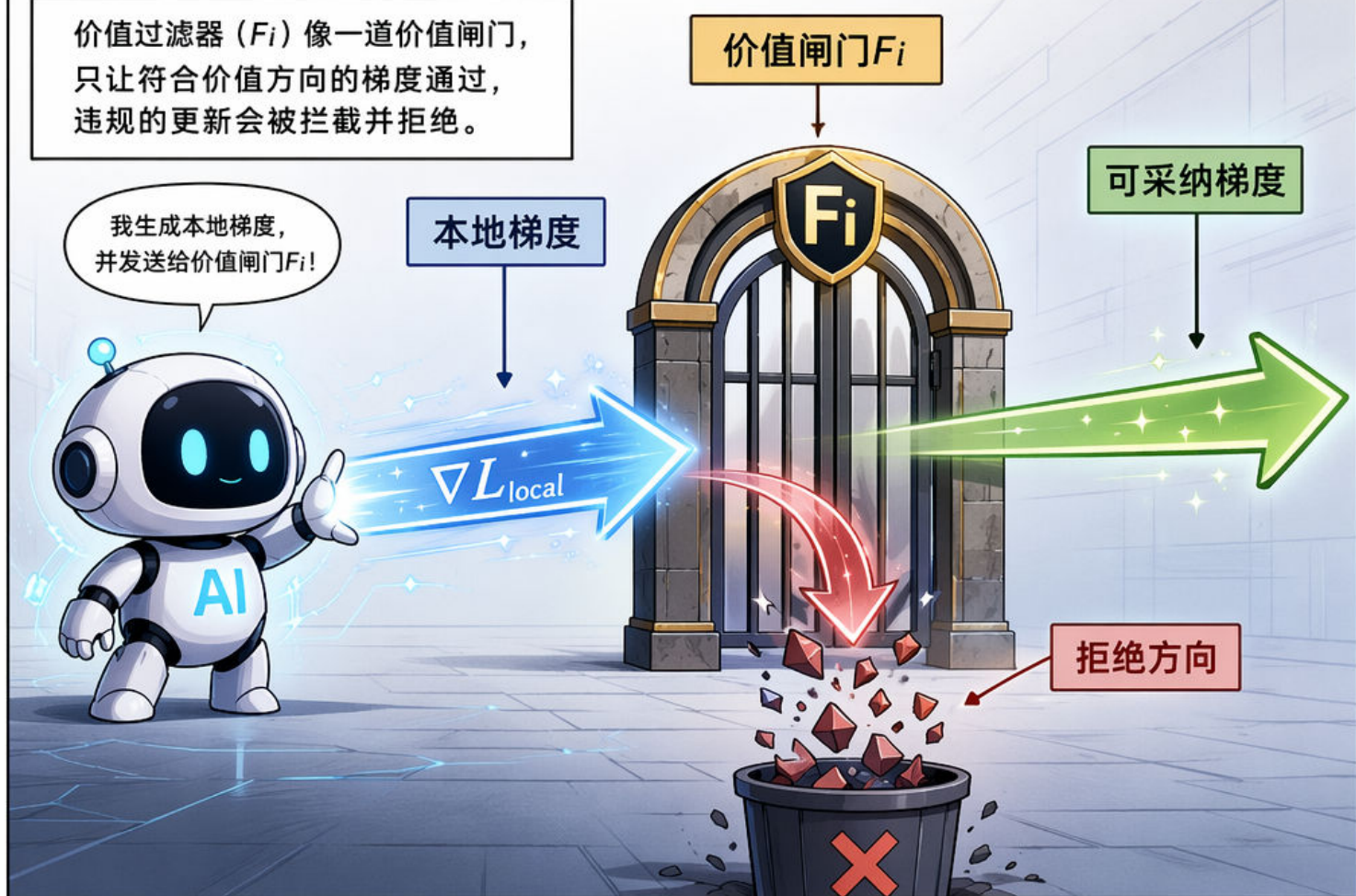


代理代表人学习，
也代表人的边界



价值过滤器

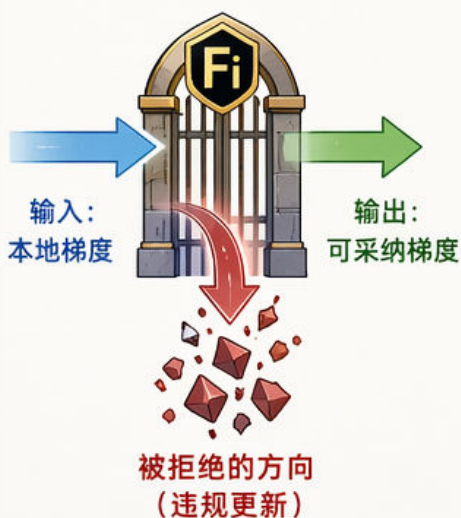
价值过滤器 (F_i) 像一道价值闸门，
只让符合价值方向的梯度通过，
违规的更新会被拦截并拒绝。



- 1** 代理计算本地梯度 ∇L_{local} ，包含所有更新方向。



- 2** 价值闸门 F_i 根据价值准则，
过滤梯度。



- 3** 只有可采纳梯度用于训练，
并参与奖励分配；违规更新
不被训练，也不分钱！

✓ 可采纳梯度
→ 用于训练
→ 参与分钱



✗ 违规更新
→ 不训练
→ 不分钱



违规更新不能训练，也不能分钱

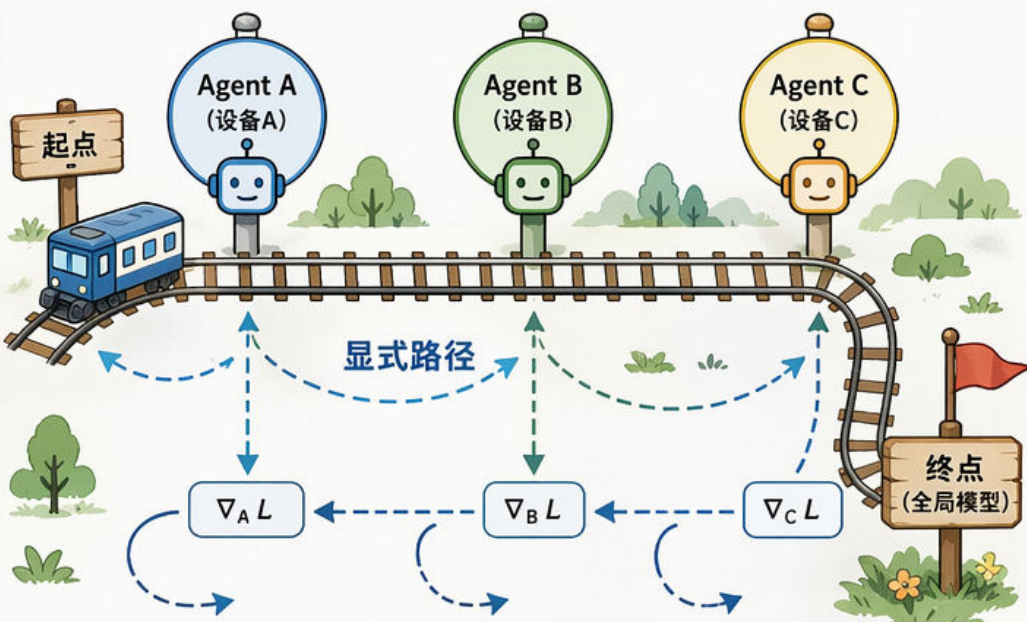


为什么用TL?

用TL, 让学习路径更清楚, 让分账更可审计!



用TL: 显式路径的遍历学习



遍历学习: 前进 + 回传, 路径清晰可追踪



显式路径

学习路径看得见, 每一步去向明确。



梯度流

梯度逐跳回传, 来源清晰可追踪。



验证改善

梯度可验证效果, 及时发现并改进。



贡献归因

每个节点的贡献可量化、可分账。



不是黑箱聚合

拒绝一锅端, 透明可审计!

不用TL: FedAvg 黑箱聚合



Agent A (设备A)

本地模型



Agent B (设备B)

本地模型



Agent C (设备C)

本地模型

FedAvg 聚合器 (黑箱)

- 如何加权?
- 如何影响全局?
- 谁贡献更大?
- 无法追踪...

全局模型

一锅端聚合, 过程看不见, 贡献算不清, 分账靠猜?

路径越清楚, 分账越可审计。



信用怎么算？

就像法庭认定证据一样，我们只认结果！

信用=你的更新，带来了多少真实、可验证的改进。

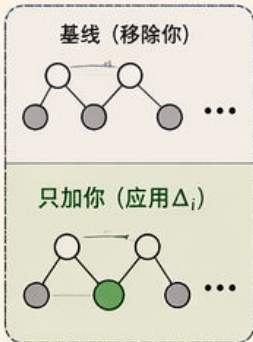


你的更新，真的带来了更好结果吗？拿证据来！

原则：
可归因、可验证、可复现

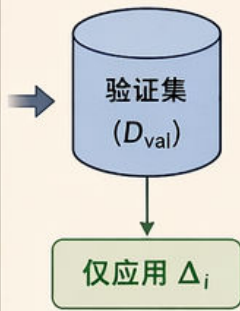
1 反事实测试

构造反事实基线：
移除你的更新 Δ_i ，
其余保持不变。



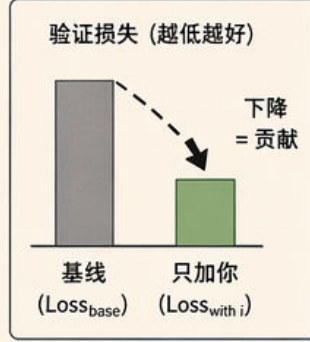
2 只应用你的更新

在验证集上，
只应用 Δ_i 的影响，
不叠加他人更新。



3 验证损失下降

比较验证损失变化：
 $\Delta Loss_i = Loss_{base} - Loss_{with i}$
> 0 才算有效贡献。



下降 > 0
有效 ✓

下降 ≤ 0
无效 ✗

4 累计贡献 C_i

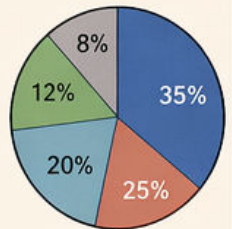
对多次任务/轮次，
累计你的有效贡献。

任务/轮次	$\Delta Loss_i$
1	0.032
2	0.041
3	0.018
...	...
T	0.027

$$C_i = \sum_{t=1}^T \max(\Delta Loss_{i,t}, 0)$$

5 收入支付 p_i

按贡献占比分配
奖励池。



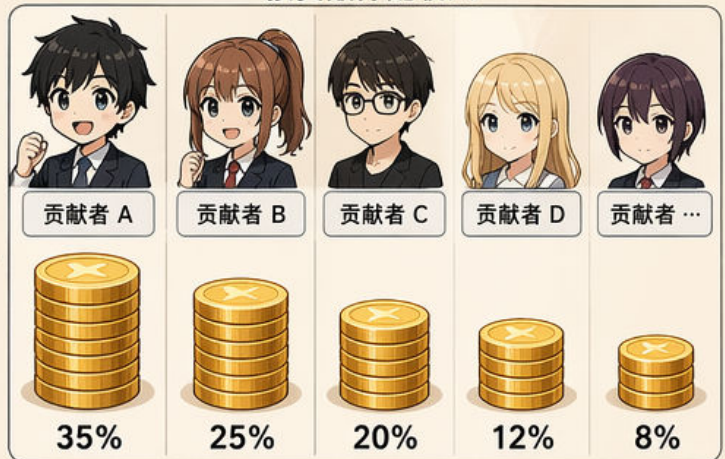
■ 贡献者 A
■ 贡献者 B
■ 贡献者 C
■ 贡献者 D
■ 贡献者 ...

$$p_i = R \frac{C_i}{\sum_j C_j}$$

有效性过滤器



按贡献分配收入



奖励只绑定合价值且有效的贡献



三条反驳

我们预判了三条常见质疑，逐一回应！

1 会收敛吗？

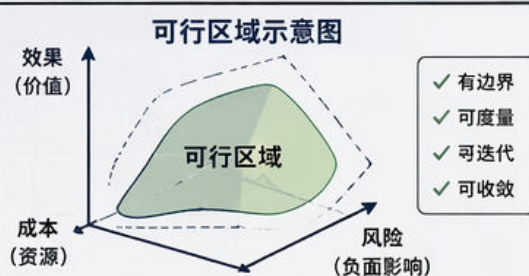
你们的框架会不会越用越复杂，反而无法收敛？

⚠️ 风险警示

- 规则膨胀
- 指标过多
- 成本上升
- 难以闭环

回应：可行区域

我们定义了“可行区域”，在效果、成本、风险三者之间寻找平衡，确保可持续收敛。



2 过滤公平吗？

过滤器会不会带来偏见，伤害某些群体或观点？

⚠️ 风险警示

- 数据偏见
- 价值偏见
- 群体伤害
- 回音室效应

回应：审计过滤器

过滤器必须透明、可审计、可解释。我们提供审计账本，持续监测并纠偏。

审计过滤器流程



公平性审计维度 (示例)



3 隐私安全吗？

你们记录这么多数据，用户隐私会不会泄露，被滥用？

⚠️ 风险警示

- 数据泄露
- 身份关联
- 目的外使用
- 内部滥用

回应：隐私保护设计

最小化、去标识化、加密隔离、可控访问。隐私在设计时就被内建，不是事后补救。



这是责任账本，不是万能答案

- ✓ 承认边界
- ✓ 接受监督
- ✓ 持续进化
- ✓ 与你共建

不完美，但负责！

责任账本

一起监督，一起更好！



责任链



1



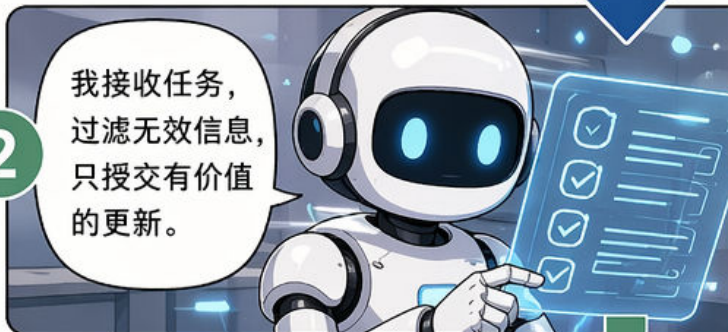
我设定目标和边界，并授权代理执行。



人：授权边界

设定目标、权限和约束，明确代理可以做什么，不可以做什么。

2



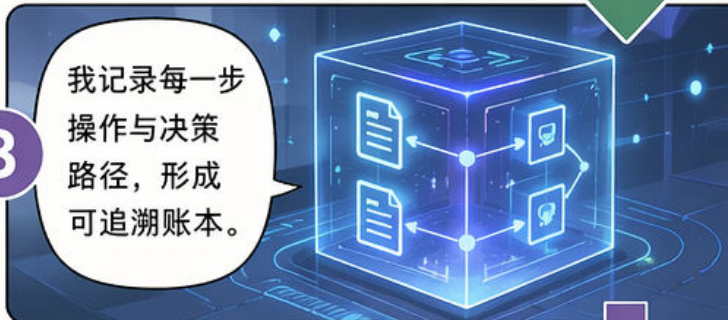
我接收任务，过滤无效信息，只授交有价值的更新。



代理：过滤更新

接收输入，过滤噪声与无关信息，只向模型提交有效更新与建议。

3



我记录每一步操作与决策路径，形成可追溯账本。



合作社：记录路径

记录代理与模型的交互、决策与输出，形成不可篡改的路径和证据链。

4



我根据贡献分配收益，让价值回到创造者手中。



市场：按贡献付费

根据各方贡献（人、代理、模型、数据等）分配收益，按贡献付费，激励协作。

5



我审计例外与异常，确保系统公正、合规、可信。



治理者：审计例外

独立审计异常与例外情况，追责与纠偏，保障系统公正、合规与可信赖。

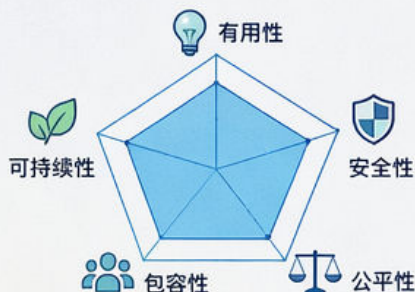
AI 替人行动，账本必须跟得上



最终脑图

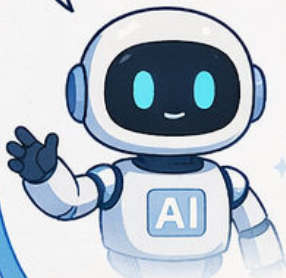
1 多元价值

从多个维度定义“好”，形成可计算的价值空间。



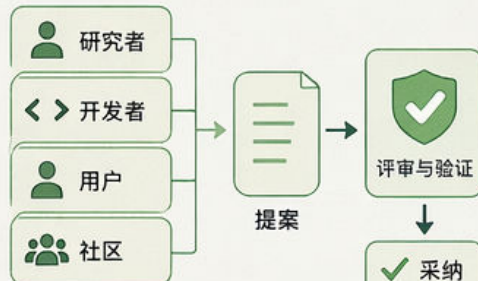
不同场景，不同权重，价值多元才更贴近真实世界！

这五个模块共同构成一个可持续、可进化、可追责的价值对齐与信用分配体系！



2 可采纳更新

谁可以更新？如何被采纳？规则透明，过程可验证。

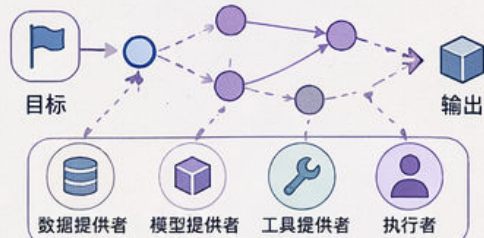


开放参与，严格审查，让系统持续进化，又不被滥用。



3 路径归因

从目标到输出，拆解贡献路径，识别每个参与者的实际贡献。

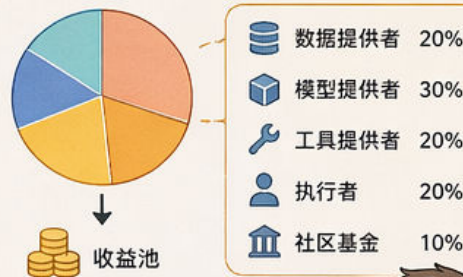


不只是结果导向，更要看清“谁在中间做了什么”。

价值约束 信用分配

4 信用结算

将贡献转化为信用，按规则分配收益与权益，激励正向循环。



贡献有回报，付出有价值，让更多人愿意持续参与！



5 责任审计

全流程可记录、可追溯、可审计，出现问题有人负责，系统持续改进。



有审计，才有信任；有责任，才有进步。



对齐不只写原则，也写进谁能更新、谁能获利

