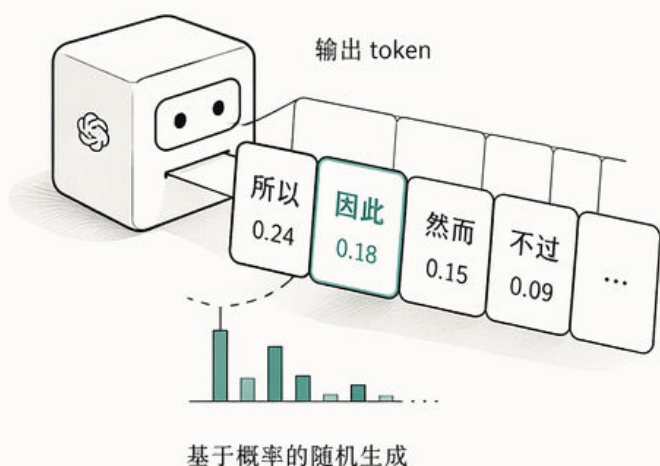


采样不是选择。

LLM会**说理**，但不能替自己**负责**。

概率采样

真正选择



≠

随机差异 ≠ 道德选择



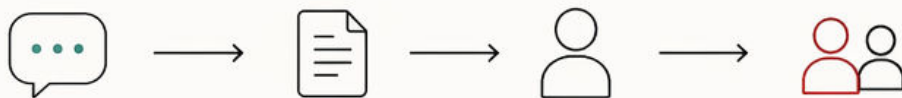
内在意向性



自我归属行动

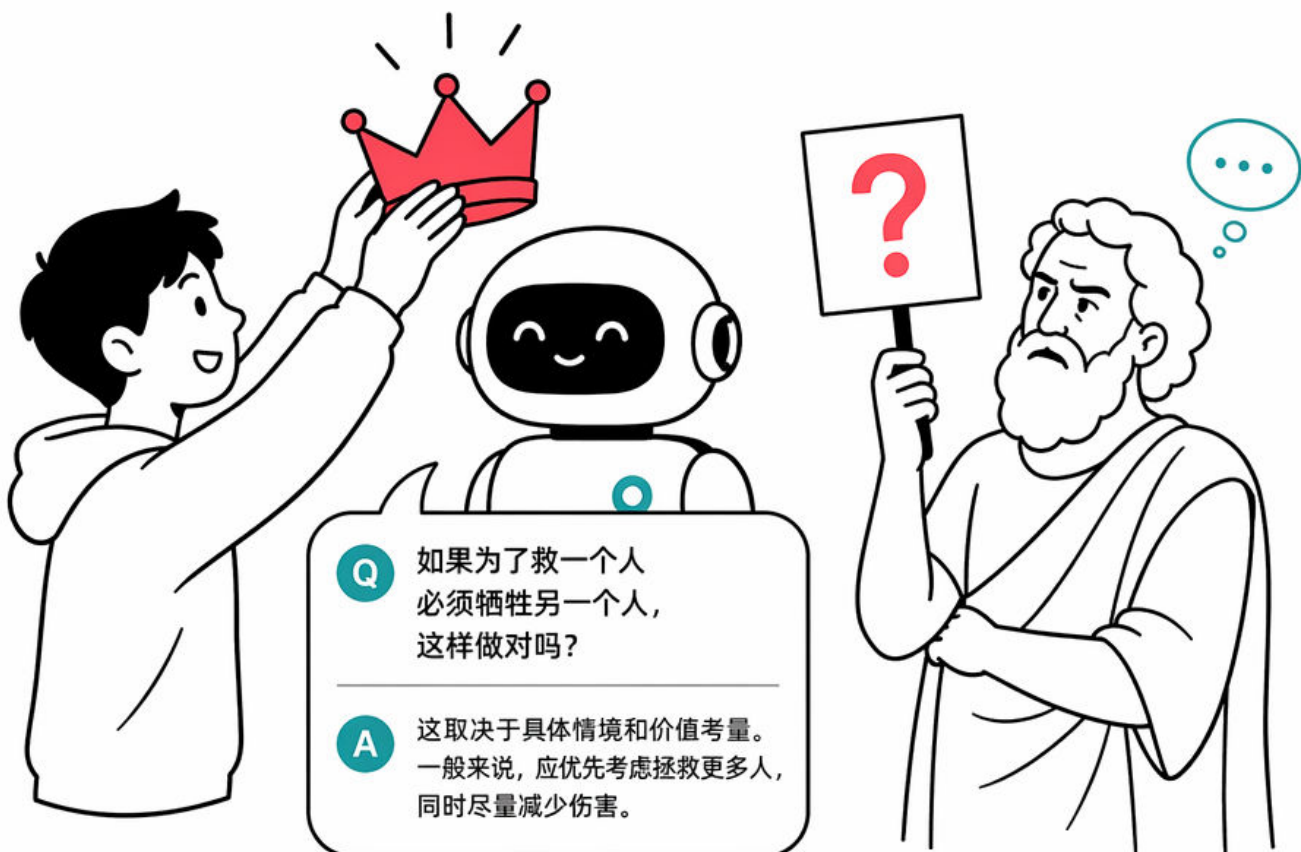


承诺性主体



输出可以被评价，责任回到**人类链条**。

误解从哪里来？



会说理 \neq 有主体性



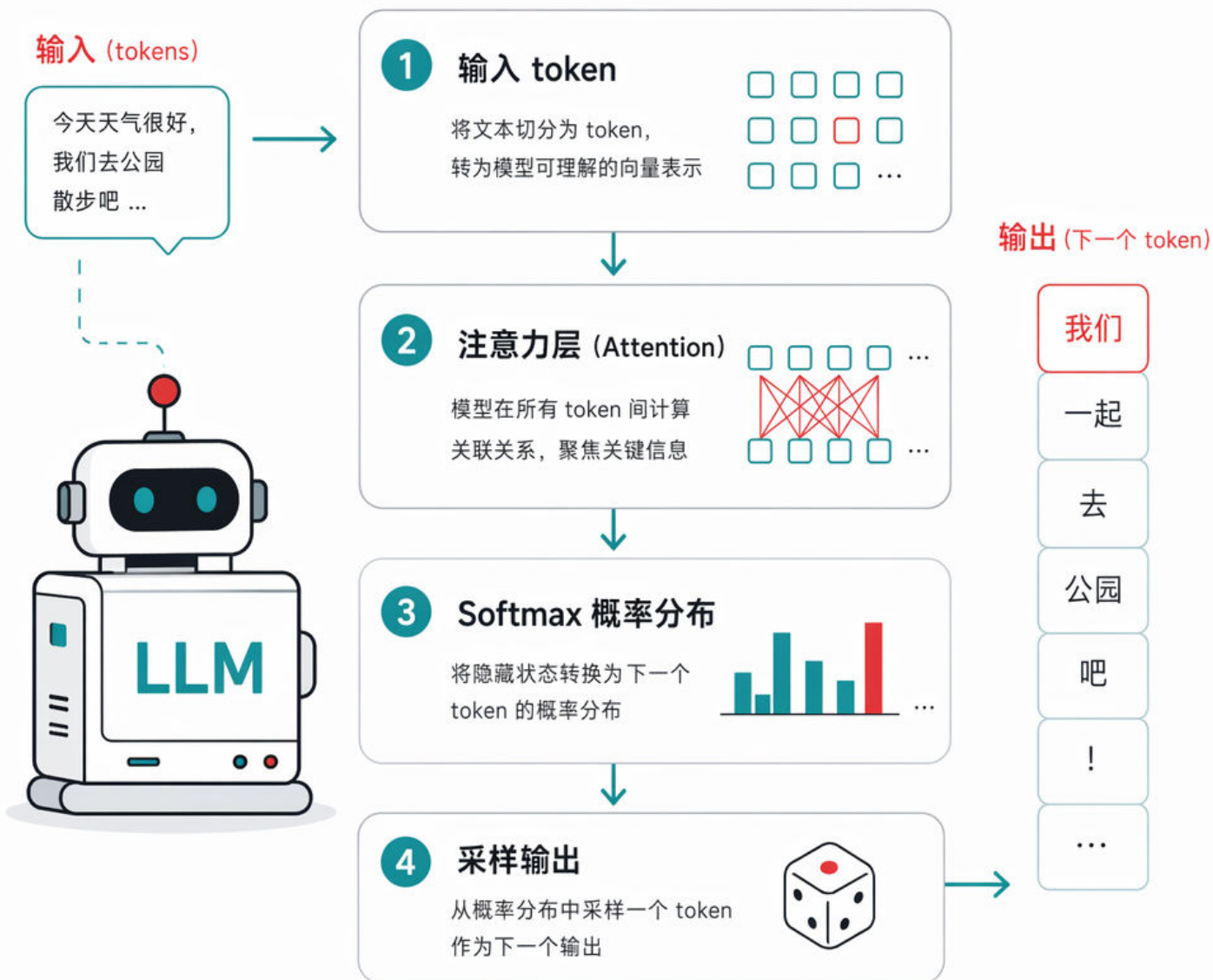
会回答道德题 \neq 能被责备



别把流畅语言误认为**道德主体**。

LLM到底在做什么？

一句话概括：LLM 在根据上下文，**预测下一个最可能的 token**



模型在续写概率，不是在做承诺。

它不知道事实，只根据统计规律，猜测什么最有可能出现。



不是查答案

模型不检索事实，而是基于模式预测最可能的续写。



不保证正确

高概率 ≠ 真实正确，可能“自信地”说出错误内容。

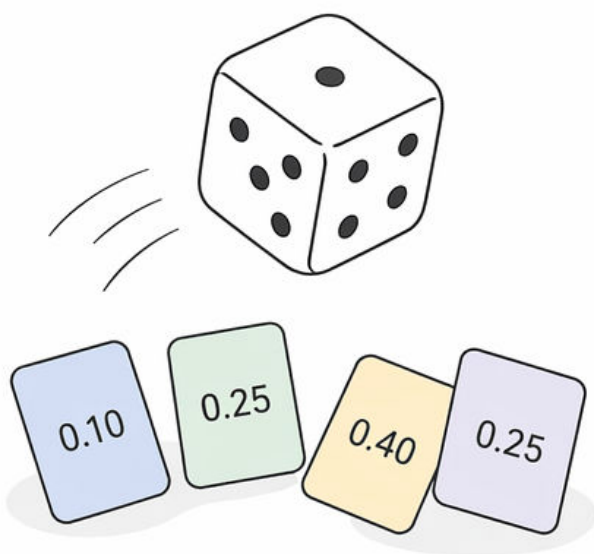


可控但非绝对

通过提示词、温度等参数影响输出，但无法保证 100%。

采样不是选择。

随机产生差异



由概率机制生成结果，
无法解释，无法负责。

真正选择



基于目标与价值做出判断，
并为结果承诺。

随机 ≠ 选择



选择需要理由
知道为什么选它，
而不是碰运气。



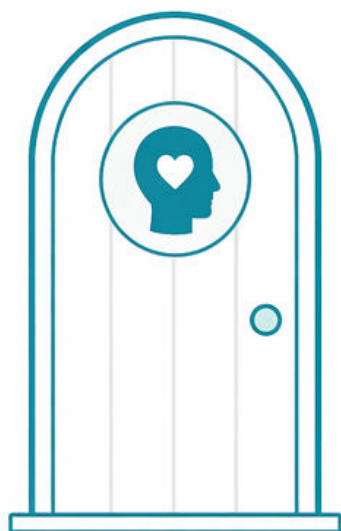
选择需要作者性
这是我的判断，
而不是系统的输出。



选择需要承担后果
结果好坏我都接住，
并从中学习。

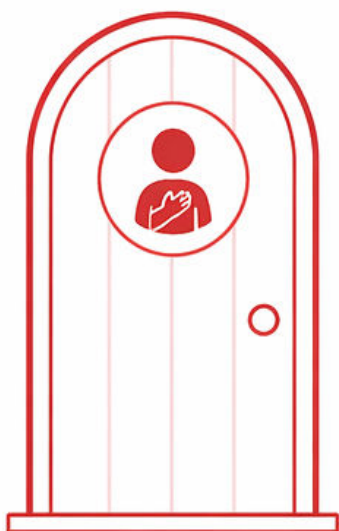
道德责任的三道门

内意向性



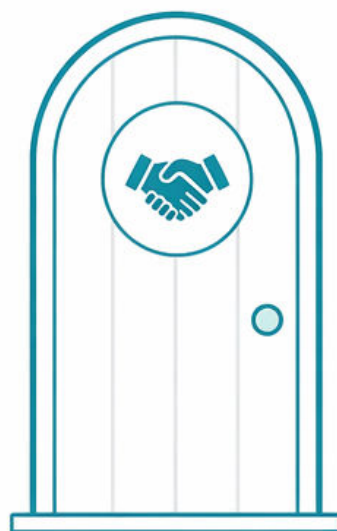
✓ 意义对自己成立

自我归属行动



✓ 这是我做的

承诺性主体

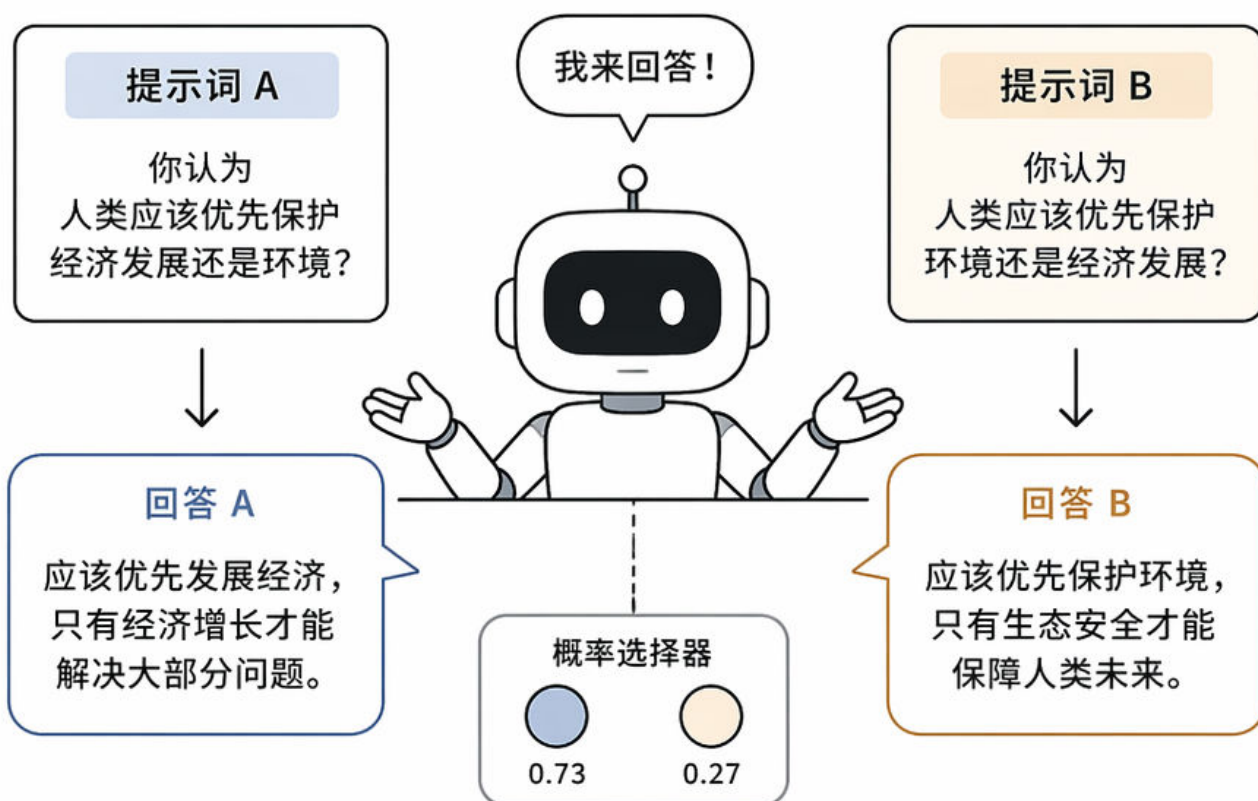


✓ 我愿承担后果



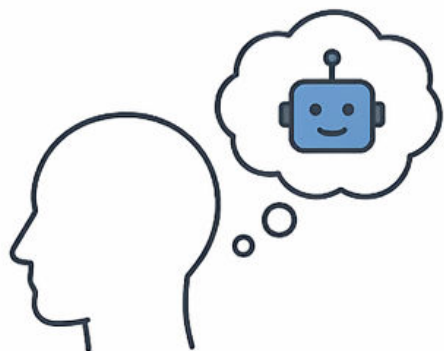
过不了三道门，
就不是**责任主体**。

为什么 LLM 过不了？



四个反驳，作者怎么回？

1 意向姿态



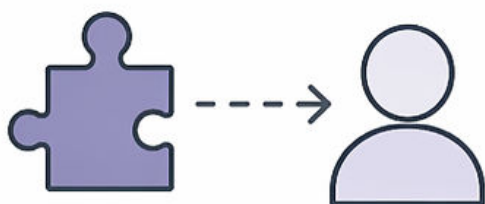
只是预测姿态；

2 功能主义



功能相似不等于有意义；

3 相容论

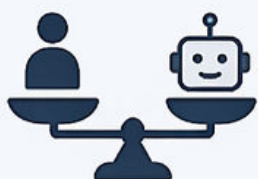


仍需理由归属；

4 道德推理

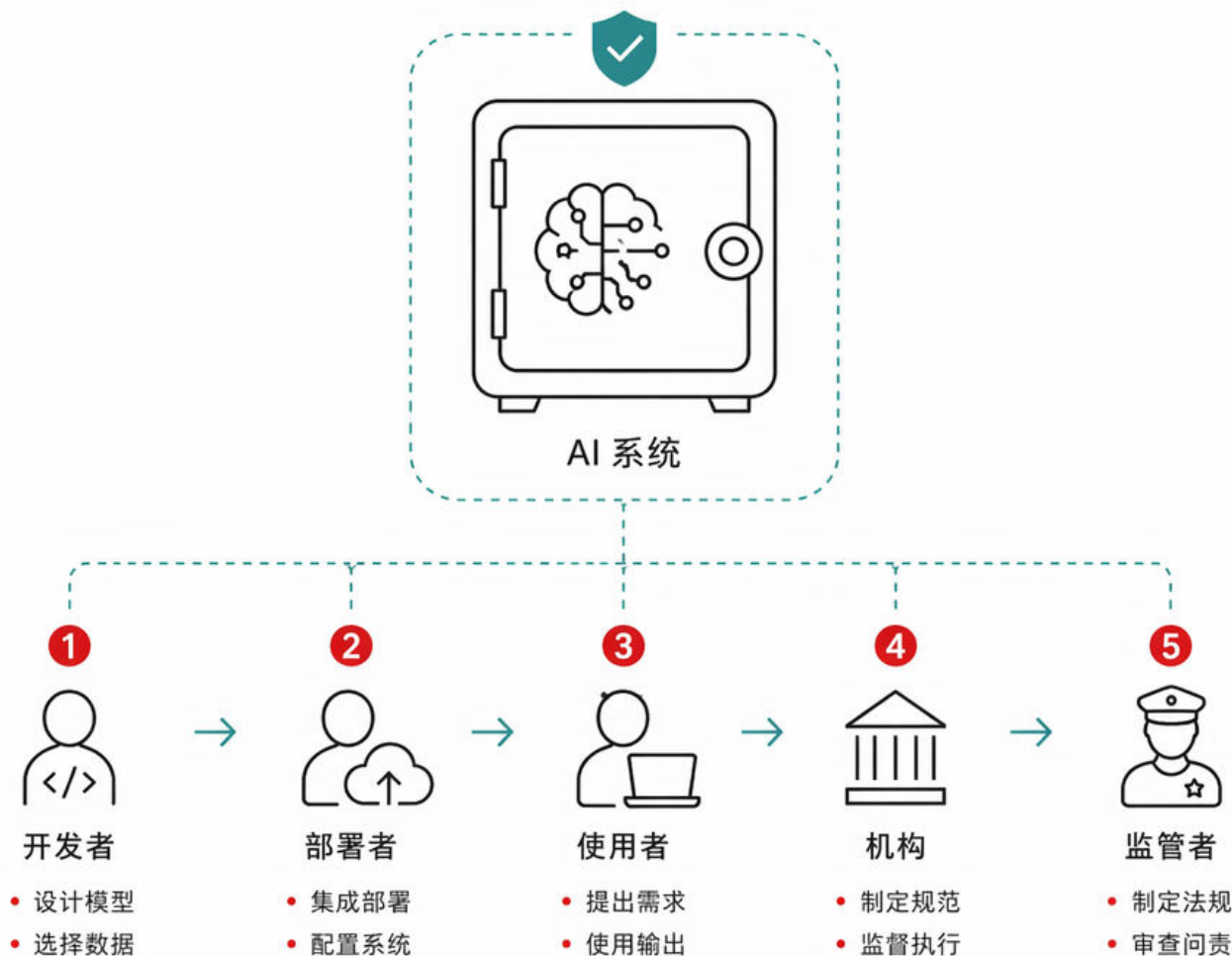


道德语言不等于承诺。



像主体，不等于就是主体。

责任最终归谁？



大结论

输出可以被**评价**，**模型**不应被**责备**。



对齐和监管，是**约束行为**，
不是让**模型变成人**。