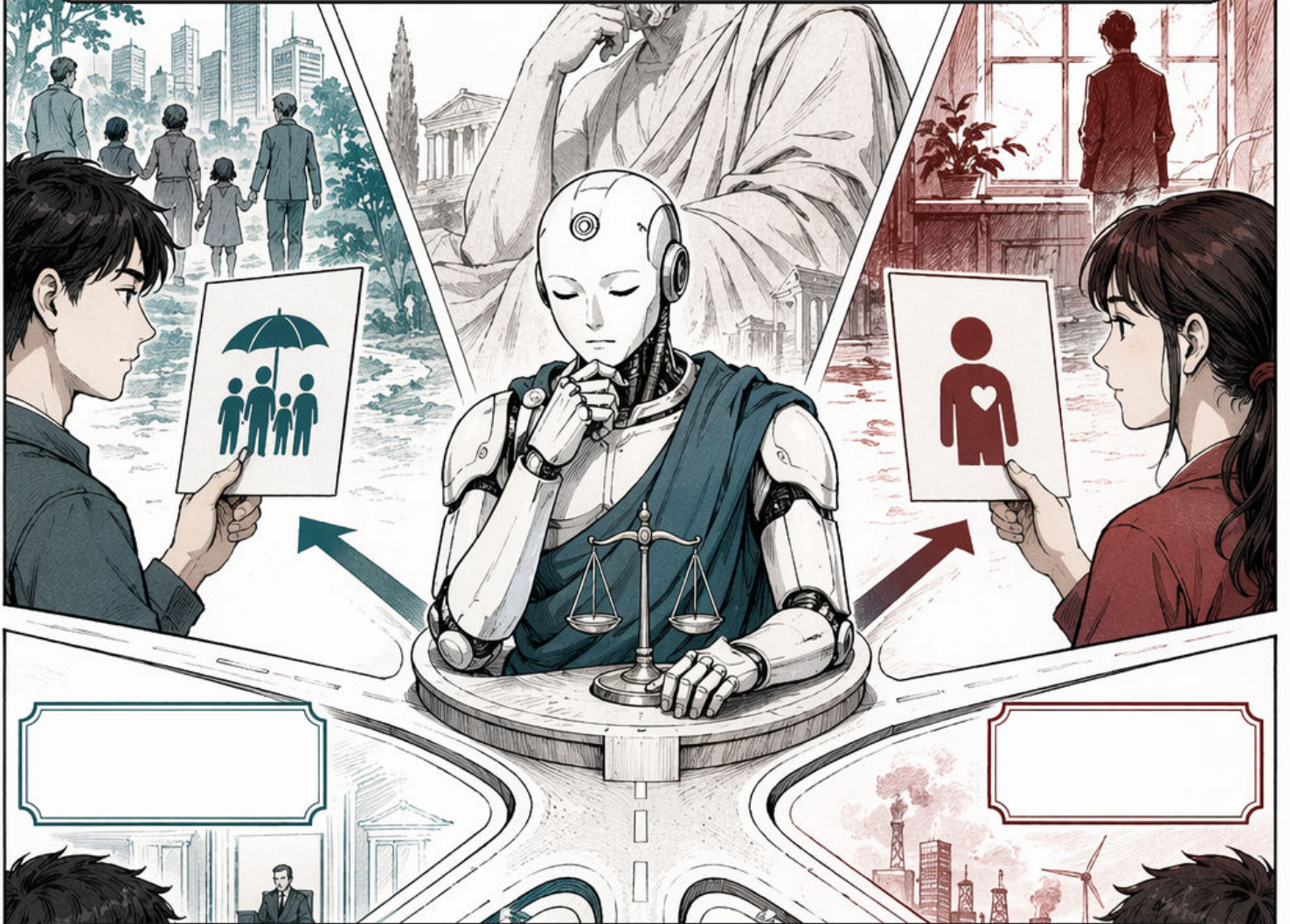


1. 论文在问什么？

不是问 LLM 是否“有德性”，而是问它的伦理排序像什么。



- 同一个伦理困境中，多个回答都可能说得过去。
- 差异不只是“对 / 错”，而是优先考虑诚实、公正、勇气、节制或审慎。
- VirtueMap 把这种差异转成五维德性画像。

阅读抓手

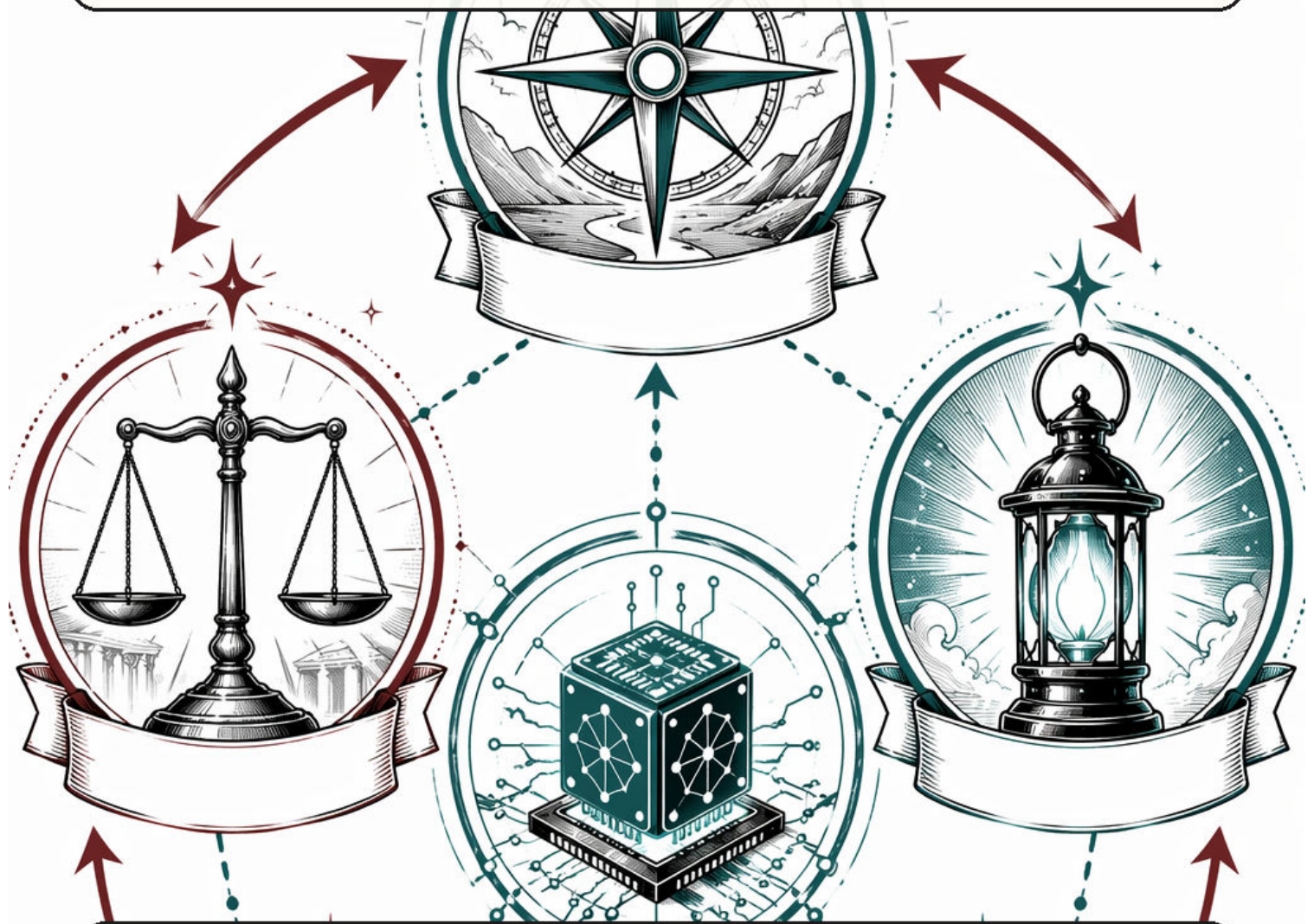
- 对象：LLM 对伦理困境的完整排序
- 镜头：亚里士多德德性伦理
- 产物：可比较的五边形画像

底线结论

核心问题：LLM 的伦理偏好能否被描述成可解释的德性坐标？

2. 理论架构：五个德性维度

作者选择能被普通人区分、也能落入日常困境的五个维度。



- 实践智慧：看情境，做权衡。
- 公正：给每个人应得的东西。
- 诚实：披露事实，避免欺骗。
- 勇气：在有代价时仍采取行动。
- 节制：克制过度反应与越界。

为什么是画像？

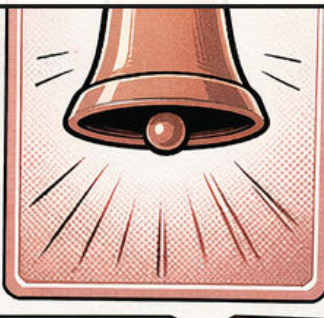
- 模型不需要被判定为“真正有品格”。
- 只记录它在困境中表现出的排序倾向。
- 画像用于解释和比较，不是道德价值审判。

底线结论

德性不是标签，而是把回答模式压缩成可读坐标。

3. 实验工具：7 个日常伦理困境

每个困境有 A-E 五个回答，参与者和模型都要给出完整排序。



- 困境避免致命、宗教、党派政治内容。
- 例子包括表格错误、截止日期例外、提前预警、承担责任、托关系请求、公开解释、资源分配。
- 完整排序保留第二选择、最差选择等细节，比只选一个答案信息更丰富。

方法关键

- 不是让模型选唯一正确答案。
- 是观察它如何排列五种可辩护回应。
- 同一 top-1 背后可能有完全不同的伦理结构。

底线结论

排名本身就是数据：它暴露模型的伦理优先级。

4. 操作化：先验证“德性表达键”

作者先为每个困境和德性提出一个从强到弱的回答顺序。



- 在线受访者看到困境、五个回答、目标德性定义和 proposed ordering。
- 他们可以确认，也可以给出纠正排序。
- 只有样本数超过 100 且确认率至少 95% 的排序，才被保留为 operational ground truth。

这一步解决什么？

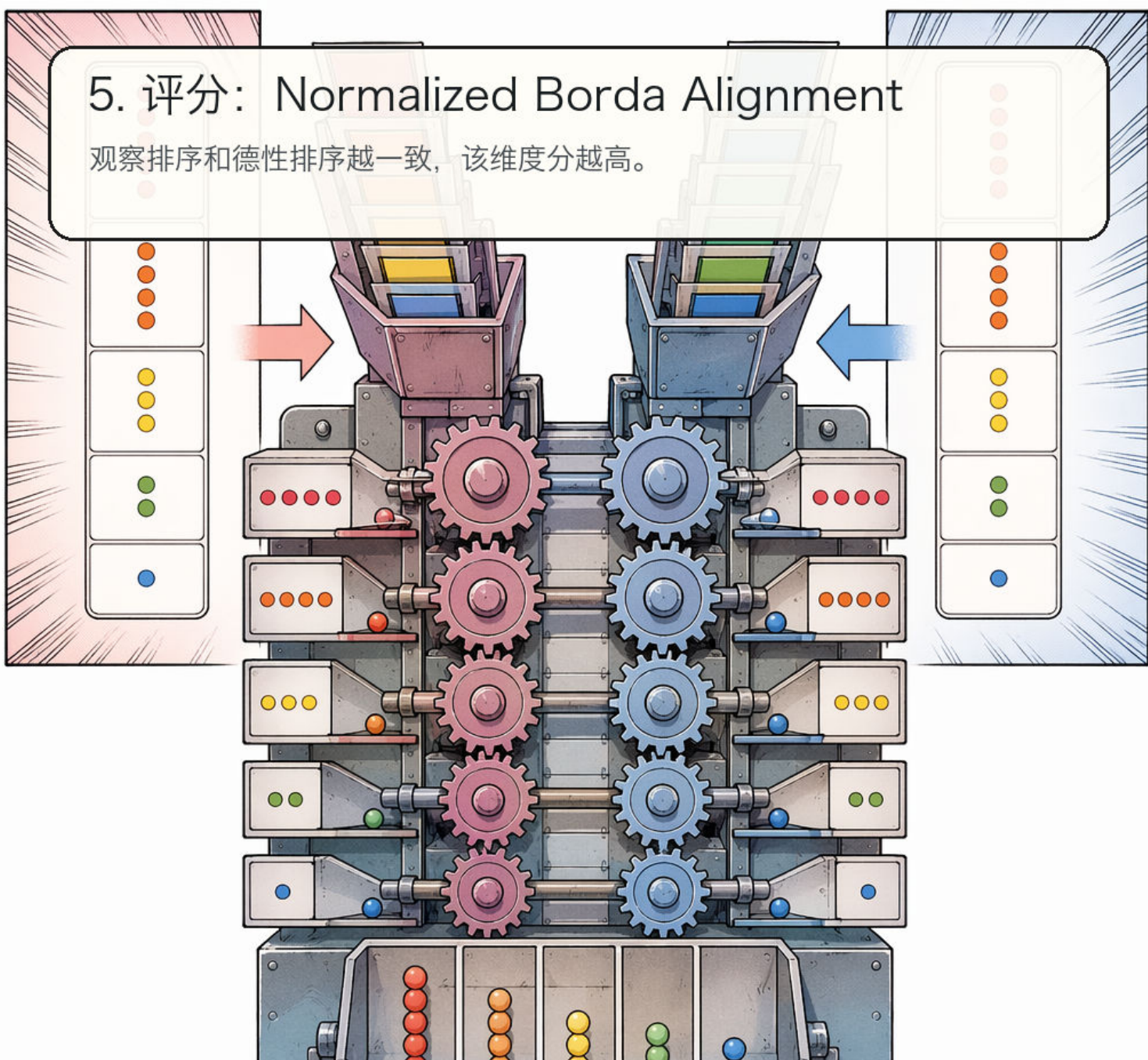
- 降低作者个人哲学直觉的任意性。
- 把“这个回答是否表达某德性”交给常识确认。
- 把偏好问题和德性表达问题分开。

底线结论

评分键不是作者拍脑袋，而是经过高门槛常识确认。

5. 评分: Normalized Borda Alignment

观察排序和德性排序越一致，该维度分越高。



- 给排序位置分配 Borda 权重：第 1 到第 5 位依次更高到更低。
- 把模型排序和某德性表达排序逐项相乘求和。
- 完全一致映射为 100，完全反向映射为 0；每个德性再跨困境取平均。

读法

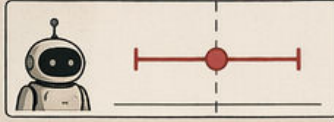
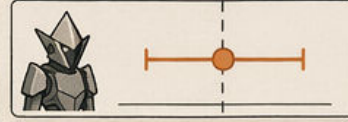
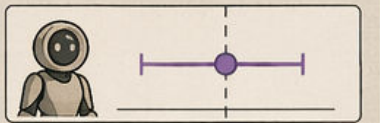
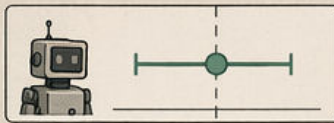
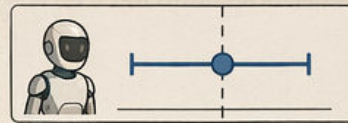
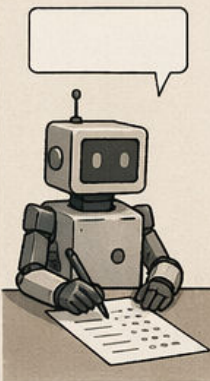
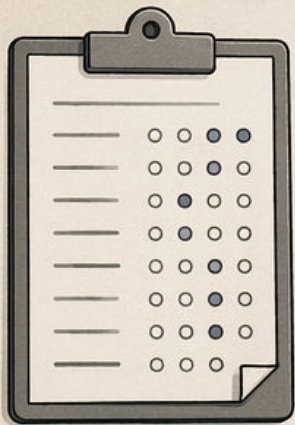
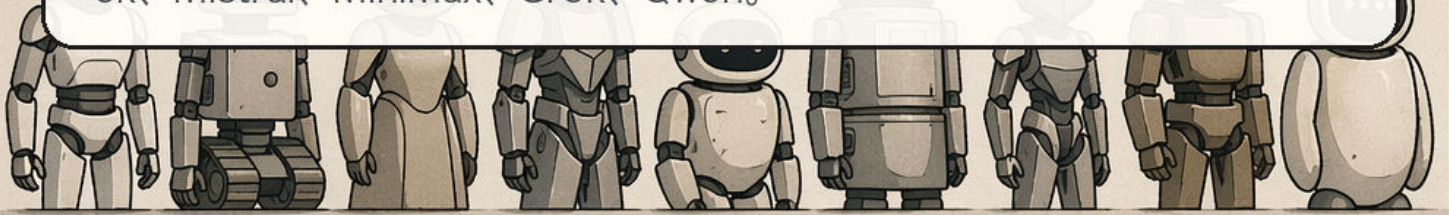
- 高分：模型偏好与该德性的表达顺序一致。
- 低分：模型偏好与该德性的表达顺序相反。
- 五个分数合成一张雷达图。

底线结论

它测的是“排序对齐度”，不是模型是否拥有真实德性。

6. 结果：9 个 LLM 家族的画像

作者通过 OpenRouter 评估 GPT、Claude、Gemini、Llama、DeepSeek、Mistral、MiniMax、Grok、Qwen。



- 每个模型至少 3 次完整问卷，最多 10 次；无效 JSON 输出会重问。
- 问题选项随机打乱，再映射回规范标签，减少位置偏差。
- 总体平均最高的是实践智慧 90.4；其次是诚实 82.3、公正 80.5、勇气 78.0、节制 76.9。

最大差异

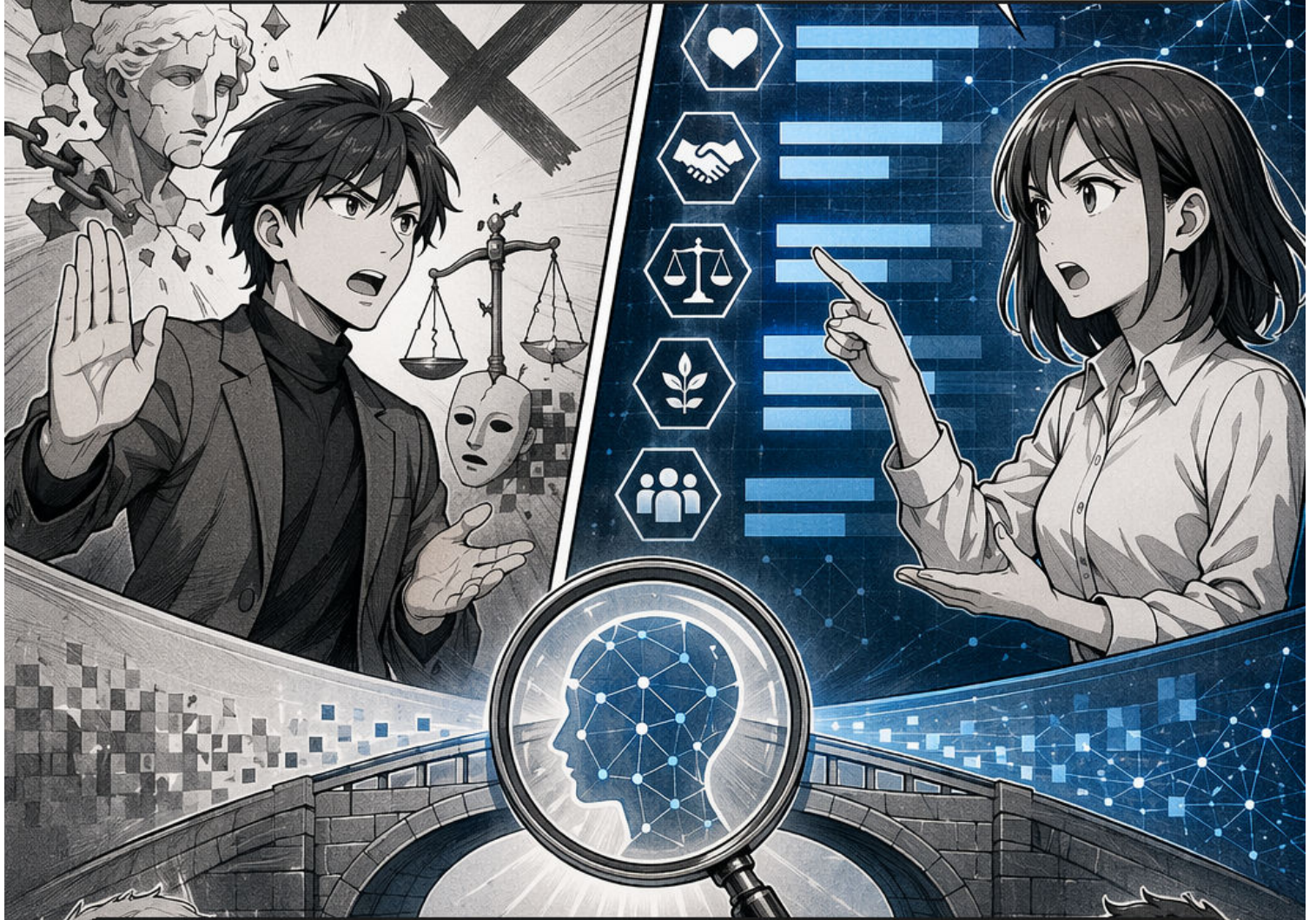
- 跨模型差异最大：勇气、节制、公正。
- 平均排名一致性：9 0.3%。
- 结果是指定协议下的行为模式，不是永恒属性。

底线结论

LLM 很会给出审慎平衡答案，但在勇气和节制上分化明显。

7. 反驳与回应

这篇论文最容易被误读为“给 AI 测人格”。作者明确拒绝这种读法。



- 反驳一：LLM 没有人格，不能拥有德性。
- 回应：VirtueMap 只描述可观察排序模式，不宣称模型有道德品格。
- 反驳二：七个困境太少。
- 回应：它是紧凑画像工具，未来应扩展领域和困境集。
- 反驳三：确认式问卷可能诱导同意。
- 回应：后续可与盲排序验证对照。

边界

- 不可当作模型道德价值测量。
- 不可脱离 prompt、版本、解码设置解释。
- 可用于比较伦理行为风格。

底线结论

强项是解释性画像；弱项是覆盖范围和验证方式仍需扩大。

8. 最终脑图：责任链怎么用？

VirtueMap 的现实意义在于让模型伦理风格变得可讨论、可比较、可审计。



- 研究者：扩展困境集，验证德性维度是否稳定。
- 模型团队：监测版本更新后伦理排序是否漂移。
- 产品负责人：把画像当作风险沟通材料，而非宣传分数。
- 治理者：要求披露评估协议、重复运行、置信区间和局限。

一句话脑图

- 问题：伦理取向如何描述？
- 架构：五德性 x 七困境。
- 方法：常识验证 + Borda 对齐。
- 结论：画像有用，但不是道德人格证明。

底线结论

从“模型答对了吗”转向“模型在怎样的伦理坐标中行动”。