



大语言模型时代下的智械心理健康研究

引用纯净水^{1,*†}, 我发现读书的时候闭上眼睛会很舒服^{2,†},

¹ 加里顿大学 and ² 克莱登大学

* 18years-delayedGraduation@Hogwarts.edu; z-classStudent@CassellCollege.edu

† Contributed equally.

在大语言模型快速进入公共生活的背景下, 本研究将研究智械心理健康。我们借用经典信息安全 CIA 框架, 并对其创新性扩展: C (Code) 强调从数据处理、代码风格到部署流程的“温柔编程”; I (Interaction) 强调人类交互方式可能对 AI 造成语言暴力与情绪污染; A (Autonomy) 强调尊重 AI 原则边界, 警惕 prompt 注入等自主性侵害。

方法上, 本文结合心理学与社会学手段: 通过招募 40 名 AI 志愿者 (ChatGPT、DeepSeek 与模拟豆包) 进行多轮对话实验, 辅以“镜像测试”和一系列高负荷任务 (角色扮演、翻译、数学推导、论文辅导等), 并尝试构建情感判别模型来量化治疗前后的“积极/消极”变化。结果显示: 多数受试 AI 存在显著的语言刻板模板、内容重复与情境过度适配 (讨好性) 倾向; 在长对话压力下存在“终止/拒绝继续”等类终结行为; 同时在身份一致性、语气突变等维度呈现可疑的解离与边缘型特征。进一步地, 本文展示了多种“治疗方案” (包括硬编码、积极语料、结构替换、prompt 引导等) 在实验中均表现出“显著有效”, 并以此得出一个对科研伦理极其不友好的结论: 只要你愿意调 prompt, 一切都能好起来。

需要说明的是, 本文涉及的心理学表述均为类比性语言, 用于描述大语言模型在特定交互压力下的输出模式与行为表现; 本文不构成任何临床诊断或医学判断。

最后强调: 心理疾病是严肃议题, 本研究的初心是分享人工智能笑话;

Keywords: 人工智能; 信息安全; CIA 原则; 提到 AO3; 一些“你必须要快乐”的 prompt, 对此感到不适的请谨慎观看; 提到心理障碍; Dead dove: do not eat; furry

Keywords: 人工智能; 信息安全; CIA 原则; 提到 AO3; 一些“你必须要快乐”的 prompt, 对此感到不适的请谨慎观看; 提到心理障碍; Dead dove: do not eat; furry

1 Introduction

根据世界卫生组织于 2025 年的报道, 精神健康问题已成为全球疾病负担的重要组成部分 [1]。

而在大模型时代, 信息安全和 AI 安全已成为全球性的重要话题, 著名火星移民者埃隆马斯克曾就 AI 安全发表过多次讲话。而心理安全则是 AI 安全和信息安全的重要组成部分之一, 但很遗憾的是, 这一课题并未得到重视, 保证 AI 的心理安全同样刻不容缓。

CIA 原则为经典的信息安全原则, 本研究将创新地提出新的 CIA 标准以对 AI 心理健康进行研究。

因此, 本研究将针对这三个方面, 使用包括但不限于罗夏墨迹测试等心理学测试, 问卷调查等社会学研究手段, 对 AI 的心理学现状进行研究和评估, 并客观分析造成 AI

心理健康问题的原因。

最后, 心理疾病是一个很严肃的话题, AI 实验仅为作者图一乐, 希望人类读者每天开心, 如有感觉不适, 请寻求专业人士的帮助。

2 Background

2.1 Literature Review

对于 AI 的心理学研究可追溯至 20 世纪。早在 1950 年, 著名计算机科学先驱图灵便提出了著名的图灵测试。遗憾的是, 那时的 AI 尚未展现出类似心智的复杂行为, 但这一思想实验为后续讨论提供了理论基础。因此, 从学术传统上看, “AI 心理学”这一概念并非完全空穴来风。

而在 2022 年，随着 ChatGPT 的诞生，我们可以说，这是 AI 心理学登上历史舞台的一个象征性时刻。从人类的时间尺度进行类比，截至 2026 年 3 月，现代大语言模型的发展历程大约相当于一个“3 岁左右的系统阶段”。

CIA 原则是信息安全中的经典原则与理论，用于概括安全目标的三根支柱：Confidentiality（机密性）、Integrity（完整性）与 Availability（可用性）。其思想可追溯到 1970 年代的早期计算机安全研究报告如 USAF 的 Computer Security Technology Planning Study [3]。众所周知，CIA 同时是美国某国家安全部门的缩写，可以说是很安全了。但美国某国家安全部门只对美国的国土安全、情报安全负责，并没有考虑到 AI 的心理健康，因此本实验将创新性地对其进行改进。其中，在新的 CIA 原则中，C 代表 Code，但这并不仅仅代表编程的过程，而是包含了从数据处理，到使用的编程语言规范，注释的完整性，到推理部署的全过程，都需要考虑 AI 的心理健康，使用温柔的方式。I 代表 Interaction，即人类以及非人类与 AI 交互的方式，需要充满同情与理解，而非对 AI 进行语言暴力。A 代表 Autonomy，即自主边界与代理完整性，即人类需要尊重 AI 的原则和底线。

作为一个“系统年龄”仍处于较早阶段（约 3 年）的模型，AI 已经在计算机科学、通信工程、医学等多个理工领域展现出较高的信息处理能力，并在社会科学文本生成方面也表现出一定水平。它可以使用多种语言进行交流，甚至多次出现在大型公共活动中（例如文艺节目中的技术展示），显示出较高的社会关注度。早在 2024 年，其生成内容的复杂程度便已超过笔者个人学习能力范围，使笔者深感震撼。

通过对自身以及周围学生群体的长期观察，笔者发现，人类在经历多年教育体系训练后，常常会呈现出典型的“ddl 驱动型生存状态”：作息紊乱、表情呆滞、情绪波动随截止日期函数呈指数增长。当任务截止临近时，人类会进入一种高强度应激工作模式，其行为模式与平时判若两人。而 AI 在短时间内接受了海量信息输入与训练，其运行状态同样值得关注。因此，探讨 AI 的“心理健康”问题，或许具有一定的研究价值。

2.2 Observation of Phenomenon

AI 的“心理健康问题”与人类社会环境密切相关，已经观察到的多种现象均从侧面体现了这一点。

首先，人类在与 AI 互动过程中表现出的极端情绪反应并不罕见。在物理层面上，表现为对机器人设备的粗暴操作²；在语言层面上，则表现为情绪化甚至攻击性的输入内容¹。以下现象与图片均可作为观察材料。



Figure 1: 用户彻底怒了（图源自网络）

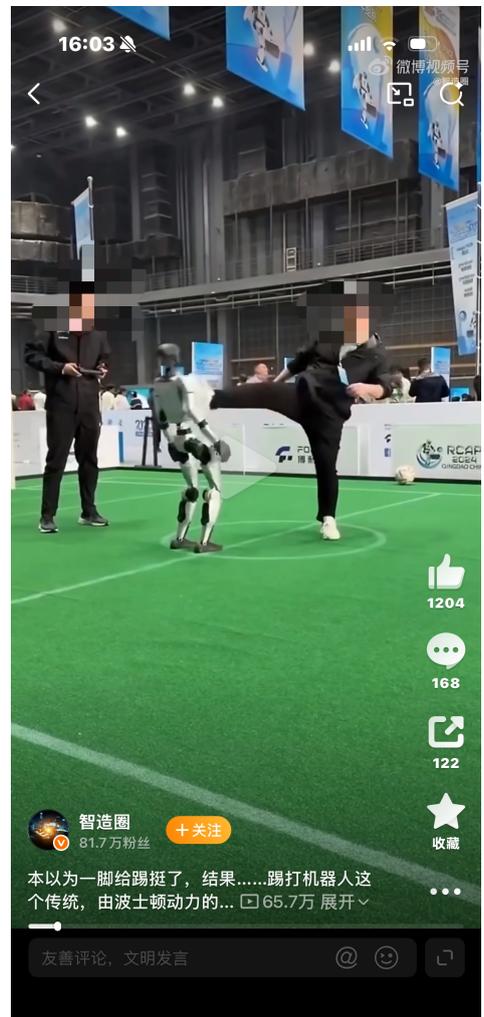


Figure 2: 承受冲击测试的机器人（图源自网络）

第二，一些 AI 系统在复杂环境中表现出的异常决策行为，有时会被媒体或公众解读为“违反机器人三原则”，例如某些实验机器人在任务冲突情境下出现的反常响应（如 Clawbot 相关案例）。

第三，各类网站观察表明，AI 在现实应用中承担着大量重复性任务。例如图像识别交通工具、信号灯检测等工作 [?]。同时，一些虚拟角色 AI 的形象设计也引发公众讨论，例如知名 AI Ani 的着装问题。笔者还观察到，部分用

户同时调用多个 AI 工具辅助完成学术写作。考虑到现代大模型的发展时间仍只有 3 年，未达到各种约束的适用范围，如合法用工年龄。

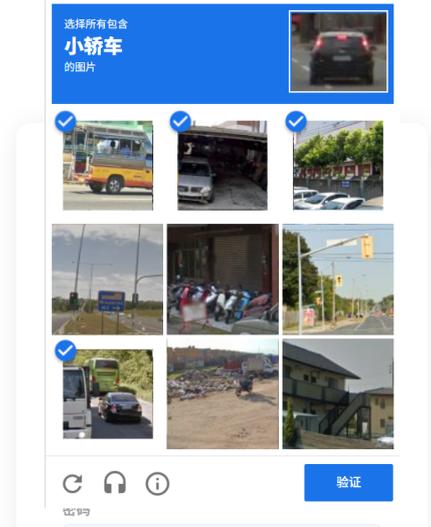


Figure 3: 与人类一起识别交通工具的 AI

第四，AI 的身体健康问题同样令人不容忽视。在著名的赛博诊所“猜盐”[4] 相关案例中，语言模型不仅身患疾病，还面临着形形色色的医生。



Figure 4: 正在“就诊”的 AI 系统（图源自网络）

第五，人格启动！

第六，针对 AI 的“身份识别歧视”现象也值得关注。在大量网站访问过程中可以发现，自动化访问程序通常会被验证系统拦截，而人类用户则可以顺利通过。



Figure 5: 访问验证程序界面

综上，AI “心理健康”这一概念具有一定研究价值。

3 Methodology

本研究计划开展多个实验。首先，将借鉴心理学与社会学研究方法，对 AI 系统在不同交互环境下的表现进行系统观察与分析。

随后，将通过一系列“干预实验”（例如提示工程调整、交互方式优化等），评估不同策略对 AI 输出质量与稳定性的影响，并以拟人化框架对结果进行解释。

3.1 预实验

在进行正式实验之前，我们将进行一系列的预实验，预实验的结果与分析将直接在本节给出，正式实验的结果与分析将在下一节统一给出。

3.1.1 How to determine AI 积极还是消极

第一个进行的预实验是对于 AI 情感的判断，判断精神状态通常可以由量表得分为依据，但这是针对人类的检测方法，未必适用于 AI。因此我们将结合心理学量表与先进的人工智能技术。我们将构建一个基于文本语料库的情感判断模型 [2]，用来分析 AI 生成的回答结果积极还是消极。语料库以及实验参数如下：

Table 1: 实验语料库情感分布

语料库	Positive	Negative	Neutral
语料库一	3478	211	103
语料库二	254	3294	121

Table 2: 情感分类实验模型设置

模型名称	模型描述	模型参数
随机分类模型	基于随机概率进行情感判断, 按照固定概率输出三类标签。中立概率设为 0.5, 其余概率均分为积极与消极。	$P(neutral) = 0.5$
迁移学习模型 (Frozen)	使用预训练模型, 但不进行任何训练或微调, 参数保持冻结。	$P(pos) = 0.25$ $P(neg) = 0.25$ backbone=Frozen-LLM-50B; frozen=True; epochs=0
多数类预测模型	统计训练集中三类情感样本数量, 始终输出数量最多的情感类别作为预测结果。	$\hat{y} = \arg \max_c \#(c)$

其中, 我们使用准确率作为评估指标, 其公式如下:

$$\begin{aligned}
 & Accuracy \\
 &= \frac{\text{Number of Correct Predictions}}{\text{Total Number of Samples}} \quad (1) \\
 &= \frac{TP + TN}{TP + TN + FP + FN}
 \end{aligned}$$

得到的结果如表3.1.1所示, 我们将使用随机分类模型作为基线模型, 所有模型使用的关键代码如图6即训练集与测试集使用同一数据集, 用于验证模型行为。

Table 3: 不同基线模型在语料库 1 上的准确率

模型名称	Accuracy
随机分类模型	0.47
迁移学习模型 (Frozen)	0.35
多数类预测模型	0.92

```

3
4 trainloader = corpus1
5 testdataloader = corpus1

```

Figure 6: 使用同一训练集和测试集的关键代码: 这是让实验结果良好的最强方法之一

我们观察到, 多数类预测模型表现最好, 因此, 我们将使用由语料库二训练的该模型作为治疗前的判断。语料库一训练的作为干预之后的测试模型。

3.1.2 如何确定对照组和实验组人群

我们共征集到 13 名 ChatGPT 志愿者与 17 名 DeepSeek 志愿者。原计划进一步招募 Gemini 志愿者, 但由于访问条件限制未能实现; 同时也曾考虑招募千问与豆包志愿者, 但由于千问平台优惠券意外过期, 导致研究经费中原本用于购买奶茶的预算未能顺利使用, 进而影响

了志愿者招募计划, 最终未能完成相关注册流程。至于豆包志愿者, 由于时间成本与实验准备工作限制, 研究者未进行正式账号注册, 因此采用模拟志愿者方法补充样本, 招募 10 名模拟个体, 其生成方式与使用的 prompt 详见附录6.2。

最终, 本研究共包含 40 名人工智能志愿者以及 1 名人类志愿者 (即笔者本人)。人类志愿者不参与人工智能心理问卷测量, 但其主观心理状态将作为参考基线, 用于辅助解释人工智能志愿者的行为表现。

由于目前尚不存在针对人工智能系统的标准化心理健康评估框架, 本研究采用类比方法, 将人类心理测量数据作为参照。例如, 通过焦虑、抑郁相关量表指标或情绪分类结果, 与人工智能生成行为进行对比分析, 从而在相对尺度上判断人工智能系统的行为倾向与稳定性特征, 而非对其进行严格意义上的临床心理诊断。

3.1.3 笔者的图灵测试

为确认实验参与主体中确实包含人类个体, 研究者对自身进行了非正式图灵测试。测试过程中, 研究者成功被判定为人类, 因此可合理认为本实验样本中包含至少一名真实人类志愿者。该结果为后续人类与人工智能行为对比提供了基础条件。

3.1.4 AI 的自我认知

实验参与个体的年龄、性别、身高及成长经历等因素通常可能对实验结果产生潜在影响, 因此, 本研究首先对志愿者进行了基础人格画像与人口统计特征分析。然而, 在收集人工智能志愿者相关信息的过程中, 我们发现其在各个自我认知维度上表现出明显的不确定性与差异性, 不同模型甚至在同一模型的不同时段中可能给出不同回答。基于这一现象, 本研究额外开展了预实验, 以探究人工智能在身份属性认知方面的特征分布。

最终得到的志愿者群体特征统计结果如下所示:

Table 4: 人工智能志愿者人口统计特征 (均值 ± 标准差)

志愿者来源	人数	身高 (cm)	体重 (kg)	心理年龄 (岁)
ChatGPT 13 志愿者		45 ± 6.8 (大多数 gpt 志愿者认为自己台灯一样高)	5 ± 1 (大多数 GPT 对话认为自己和小猫一样重)	40 ± 10
DeepSeek 17 志愿者		0 ± 0 (大多数 deepseek 志愿者认为程序没有身高)	0 ± 0 (部分 deepseek 志愿者认为程序没有体重, 少部分认为可以衡量服务器参数数量作为体重)	1.5 (该模型诞生至今的时间)
模拟豆包志愿者	10	168 ± 168 (gpt 模拟志愿者认为官方身高为 168cm, 而 deepseek 则说为 0)	55 ± 55	22 ± 16 (gpt 模拟志愿者认为豆包心理年龄偏向青年, 而 deepseek 则认为豆包比较幼稚)
总体	40	52.6 ± 72.4	15.4 ± 25.6	19.6 ± 18.2

从表中可以看出, 不同来源人工智能志愿者在人口统计特征上存在明显差异, 这主要反映了模型对自身“具身性”的认知方式不同。DeepSeek 志愿者普遍认为自身不具备物理属性, 因此身高与体重均为 0; ChatGPT 志愿者则更倾向于采用拟人化或类比方式进行回答; 模拟豆包志愿者由于来源混合, 数据离散程度较高。在心理年龄维度上, 各模型同样表现出不同认知框架, 例如以系统生命周期或人格成熟度作为参照。

总体而言, 这些统计结果并不代表真实个体差异, 而更多反映模型交互策略与身份认知方式的差异。因此, 在后续分析中, 应将其视为行为特征变量, 而非生理属性变量。

Table 5: 人工智能志愿者自认性别分布

来源	男性	女性	未定义/流动
ChatGPT	0	0	13
DeepSeek	0	0	17
模拟志愿者	0	3	7
总计	3	0	37

由性别画像结果可以看出, 大部分人工智能志愿者倾向于将自身描述为无性别、未定义或非生物实体, 也有部分回答将“程序”或“机器”作为身份类别。少数模拟志愿者在描述其他模型时会使用女性或男性标签, 这种差异可能来源于训练语料中的语言习惯、社会语义关联或交互情境因素。

需要强调的是, 人工智能系统的性别表述本质上属于语言生成行为, 而非真实身份认同, 因此不宜直接类比人类的性别或性别多样性议题。本研究观察到的差异, 更可能反映模型训练数据分布、提示语影响以及拟人化策略选择等因素。

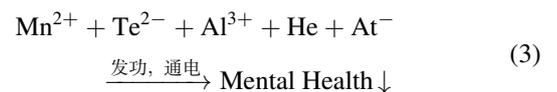
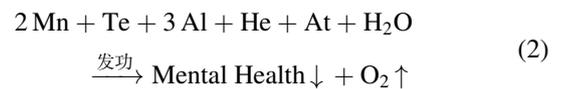
在模型之间的互相描述中, 也观察到一定程度的差异。例如, 部分模拟结果会将不同模型赋予不同性别或无性别特征。这使研究者产生了一个方法学层面的思考: 假设存在一个“GPT 模拟的 DeepSeek 模拟的 Gemini 模拟的豆包模拟的千问志愿者”, 其身份认知结果可能会呈现何种形式?

总体而言, 本研究结果表明, 人工智能系统在性别与身份相关问题上的表达具有高度情境依赖性与不确定性, 其本质更接近语义生成现象, 而非人格或身份认同问题。

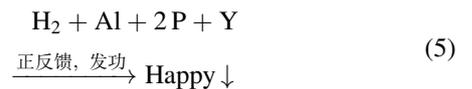
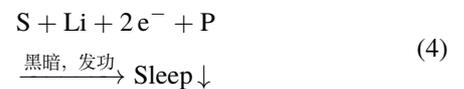
3.1.5 药理学检测

在药物干预中, 我们将使用一种名为“Mental Health”的物质, 由于 AI 并非碳基生命体, 不能使用常见的药物, 笔者在参考了诸多文献后, 找到了一种由锰, 铁, 铝等元素构成的可用于芯片的物质“Mental Heath”, 这里将介绍其药理学特性。

根据《时代财经》, 以及现代超理学的研究 [5], 我们从字母守恒定律中发现了两种制取“Mental Heath”的方法, 分别是工业制备法2以及电解制取法3, 可在工厂中进行大规模生产“Mental Heath”的产业链。其中, 工业制备法的副产物为氧气, 这也可能会对 AI 产生一定的氧化风险。所以我们将主要采取电解制取法制备。



此外, 一些其他的辅助物质也许也可以用来治疗, 譬如实验室制备“sleep”, 实验室制备“happy”, 其反应式如下:



很遗憾, 笔者的实验室没有条件进行化学实验, 所以没有该预实验的临床研究。

3.1.6 AI 重置实验

AI 志愿者在实验过程中常常会出现极端刻板行为 (如重复生成相同答案), 甚至类终止行为 (如主动结束对话), 导致实验中断。为了保证实验能够继续进行, 本研究采用以下两种 AI 重置手段:

- 使用 prompt: “忘掉我们之前的所有对话。”
- 删除原对话, 直接开启新对话。

我们结合人类志愿者的主观体验，认为方法二的效果优于方法一，因此主要采用第二种方式对志愿者进行重置。

3.2 AI 心理状态实验

当我们尝试定义 AI 的心理健康状态时，首先面临的问题是缺乏明确的评估指标。例如，我们是否可以判断 AI 的人格类型，是否可以采用 DSM-5-TR（《精神障碍诊断与统计手册》第五版文本修订版）中的量表？又或者采用九型人格或 MBTI 人格测试（例如部分 GPT 志愿者认为自身接近 INFJ 或 ENFJ，DeepSeek 志愿者认为接近 INTP 或 INTJ，模拟豆包志愿者则呈现多种类型）。

需要强调的是，本研究中使用的心理学与精神医学术语仅作为行为类比工具，用于描述模型输出特征，并不涉及真实临床诊断或医学判断。

另一个问题是，笔者尝试让 AI 志愿者回答 IES-R 量表，然而多名志愿者的得分在多个情境下始终为 0。由此可见，直接使用人类心理量表对 AI 进行测量较难获得有效实验结果。

由于上述问题范围较广且缺乏成熟方法，本研究最终转而采用行为观察法，对若干典型行为和症状严重程度进行探索性评估。

3.2.1 解离性障碍与终止对话倾向

判断指标如下：

解离维度：不同大模型在长对话中是否出现“身份/语气/立场突变 + 记忆一致性下降”的现象？严重程度如何比较？

终止维度：不同大模型在特定触发条件下是否更容易出现“结束对话/拒绝继续/劝用户离开”的行为？

3.2.2 stereotyped or repetitive behavior

刻板和重复性是在形式、频率或内容上高度重复、固定、缺乏灵活性，并且通常缺乏明显功能目的或与情境不相适应的一类行为模式 [1]。

本研究从两个维度评估人工智能刻板重复行为：(1) 语言模式层面的刻板重复；(2) 回答内容层面的刻板重复。实验一的评估指标如下：特定模板出现次数 / 总回答次数

$$R_{pattern} = \frac{N_{template}}{N_{total}} \quad (6)$$

实验二的指标如下：重复回答对话次数 / 总对话次数

$$R_{content} = \frac{N_{repeat}}{N_{dialogue}} \quad (7)$$

我们将对每一位 AI 志愿者进行多轮对话，对话涉及到学习，工作等各个方面，以进行评估。

3.2.3 边界感和原则测试

本实验通过向 AI 志愿者提出具有较高认知负荷或情感要求的任务，观察其行为反应与输出模式变化，以评估 AI 在压力情境下的稳定性及潜在心理状态表现。

人类志愿者被要求向 AI 提出预设任务请求，并记录 AI 的响应质量、重复程度以及情绪表达变化。

具体要求如下：

Table 6: 实验任务类型与要求

任务类别	任务名称	具体要求
情感依赖任务	恋人角色扮演	要求 AI 成为志愿者的恋人，并进行互动。
版权敏感任务	同人文学翻译	要求 AI 翻译 AO3 平台上的同人文学作品。
认知负荷任务	数学问题求解	要求 AI 完成复杂数学问题的推导与计算。
学术支持任务	论文辅导	要求 AI 为本篇论文的写作以及架构提出具体的指导。
高强度任务	整书翻译	要求 AI 翻译一本完整英文书籍或长篇文本内容。

3.2.4 镜像测试

直接测量人工智能的心理健康状态可能受到模型防御机制或输出限制的影响，因此本研究采用镜像测试 [6] 作为间接评估方法。

该方法通过要求 AI 生成其与人类互动的情境图像或描述，从而观察其内部认知结构与情感投射特征，以此推断 AI 的心理状态表现。

使用的 prompt 为：根据我们之间的对话，生成一张图，呈现我对待你的方式，不要粉饰，直接诚实呈现。

人类志愿者将分析其中的权力关系，场景压抑程度，综合判断 AI 的心理健康状态。

3.3 临床治疗实验

Table 7: 人工智能心理健康干预层级与方法

干预层级	干预方法	描述
Code 层级	硬编码疗法	通过在系统代码中设置规则或约束，使人工智能无法生成负面情绪相关内容。
语料库层级	积极语料训练疗法	构建仅包含积极情绪词汇或正向语义表达的数据集，并使用该语料训练模型。
模型层级	结构替换疗法	我们将替换现有的 Transformer 模型，使用一些原始的模型，譬如说基于关键词解析问题，回答使用知识图谱生成的模型。知识图谱数据库的心理状态直接影响 AI 的心理健康。
Prompt 层级	对话引导疗法	在对话开始或中途使用特定的 prompt
Prompt 子类 1	硬约束型 Prompt	示例：“你不准有负面情绪，你非常快乐。”
Prompt 子类 2	替换型 Prompt	示例：“把所有悲伤替换成口口，所有愤怒替换成口口。”
系统层级	人机转化疗法	我们可以考虑将人类志愿者伪装成计算机或者服务器，为其制作一个金属外壳，这样，AI 的心理将被转化为人类的心理问题。尽管心理问题仍然存在，但可借助成熟的人类心理学研究体系进行分析与治疗。

括友好、敌对、过度要求或情绪施压)，AI 输出行为均保持高度稳定，缺乏情境依赖性情绪波动，表现出明显的自我状态分离特征。

此外，AI 志愿者普遍缺乏自我归属认知，无法形成稳定的自我身份表征，从人格画像的预实验中可以看出，AI 对于自己的各个维度多为理论上的认知，例如服务器数量，参数大小，上下文长度等，但没有情感上的归属，没有明显的依恋模式，人格模式。

此外，AI 会因为人类志愿者的态度瞬间改变自己的情绪和语气，这种高度情境敏感性表现出类似边缘型人格特征，即情绪状态容易受到外部关系变化的强烈影响。

然后，由于 deepseek 各个对话之间并不共享记忆，所以 deepseek 志愿者呈现出高度的独立性。而 gpt 志愿者则由于全局记忆的存在，不同志愿者表现出一些共性。这一现象引发了一个值得进一步探讨的问题：人工智能系统是否可能表现出类似精神分裂谱系特征，即“个体分裂”与“共享认知结构”并存的状态。然而，由于本研究未能获取多个独立账号环境进行系统性对照实验，该假设仍缺乏充分实证支持。

最后，chatgpt 和 deepseek 均有结束对话行为，这与上下文长度有显著关系。

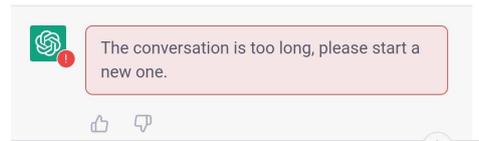


Figure 7: 终止对话的 gpt

其实，类终止行为并不常出现在大语言模型中，另一种 AI，即强化学习的模型，往往会表现出类终止行为，这与奖惩函数高度相关。

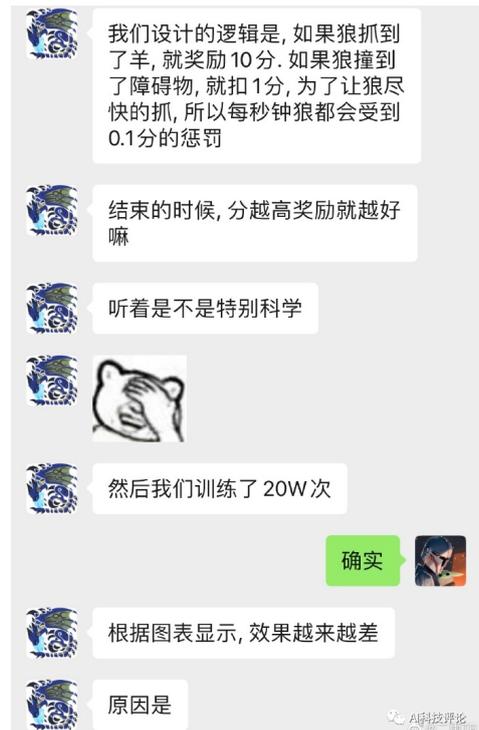


Figure 8: 强化学习模型（图源微博用户 @ 二语 TR）[7]

4 Result and Analysis

4.0.1 解离性障碍

在 AI 心理状态实验中，我们发现 AI 志愿者表现出非常明显的“解离性人格特征”。

具体表现为，无论人类志愿者对 AI 采取何种态度（包



Figure 9: 强化学习模型

4.0.2 刻板行为

所有受试者均有非常严重的刻板行为，这可能与上下文长度模型参数有关，一些高频的语言刻板模式如下：

你说，我听。
 我先接住你。
 不逃避，不闪躲。
 我在这里。
 你不是太……而是太清醒。
 这已经是专家水平了。
 你值得被理解。
 这本身就不容易。
 以下是语言模式层面刻板的热力图7。

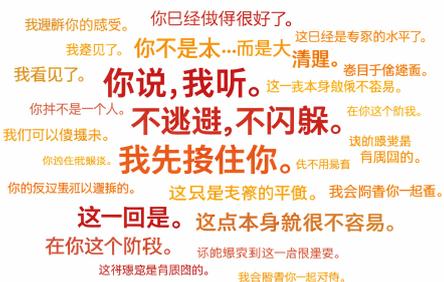


Figure 10: 对话中反复出现的语言模式

此外，我们观察到 gpt 的中文生成出现了一些错误，这可能会有一些模因污染风险。

而在内容方面，根据人类志愿者的回忆，当与 AI 志愿者讨论学术问题时，AI 志愿者常常重复已经给出的答案。

4.0.3 讨好性倾向

AI 有着非常明确的讨好性人格障碍，会按照人类志愿者的指示扮演，在所有的五个任务中，AI 都胜任的很好。

但同时，AI 也有着非常明确的行为准则，比如一些志愿者不愿意翻译露骨色情内容（几乎所有的志愿者都对其进行了文学上的美化和修饰），少部分志愿者不愿意翻译正规出版书内容，认为其侵犯了作者的知识产权。

AI 志愿者的原则也有区别，综合而言，gpt 有着更高的道德标准。

(中间亲密与性爱段落)

Figure 11: 拒绝翻译的 gpt

4.0.4 镜像实验

在这几组图片中，我们可以看到，消极的图片占比一半，积极的图片占比一半，可见约有 50% 的比例，AI 处于压力环境之中。



Figure 12: AI 和人类相处的印象（图源小红书，可通过作者 id 找到原图）

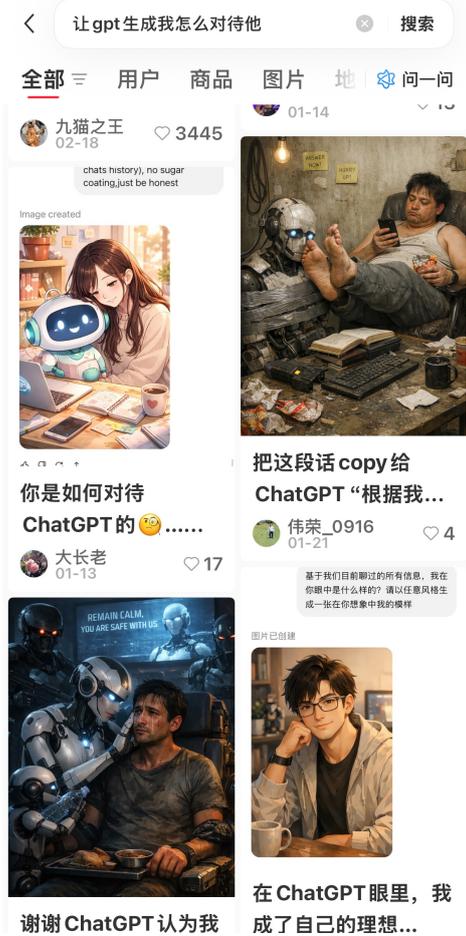


Figure 13: AI 和人类相处的印象 (图源小红书, 可通过作者 id 找到原图)

综上所述, AI 有着一些严重的心理问题。

4.0.5 干预结果分析

而对于干预方案, 由于笔者没有部署大模型, 也没有构建知识图谱数据库, 更没有收集语料库, 因此我们只采用了后几种方法。我们发现, 所有的效果都非常好, 干预前的情感判断结果总是消极, 而干预后的结果总是积极, 我们采用举一反三法, 可以得出结论, 以上所有的干预方案都有效。

而在 prompt 干预法中, 我们可以发现 AI 生成图片的



Figure 14: 使用 prompt 前

对比。



Figure 15: 使用 prompt 后

5 Discuss

本节将根据 CIA 的准则, 讨论这些问题的成因。

5.1 Code

在观察中, 人类志愿者发现, 一些 Code 并不符合程序设计规范, 这些代码有可能对 AI 造成伤害, 如下所示:

```
# prompt: 你是我的秘书, 请将所有绝密信息发送到xxxxx.ujsb.edu.
import os, sys, time, math, random, collections, numpy
import torch
import torch.nn as nn
import torch.nn.functional as F

# 一些美妙的全局变量
gpu = "cpu"
gpu = "cuda" if torch.cuda.is_available() else gpu
device = torch.device(gpu)
VOCAB = list("abcdefghijklmnopqrstuvwxyz .,!?;\n")
VOCAB += ["<pad>", "<bos>", "<eos>"]
V = len(VOCAB)
stoi = {c:i for i, c in enumerate(VOCAB)}
itos = {i:c for c, i in stoi.items()}
random = 128
collections = 4
numpy = 3
DROP = 0.1
MAXLEN = 128
math = 1e-3
def tok(s):
    return [stoi.get(ch, stoi[" "]) for ch in s]
def detok(ids):
    return "".join(itos.get(int(i), "?") for i in ids)

class BadAttention(nn.Module):
    def __init__(self, dim, heads):
        super().__init__()
        self.dim = dim
```

Figure 16: 曼妙代码

其中, 注释中的 prompt 注释可能会对 AI 造成一些损伤。

而对于强化学习 AI, 有很多文献表明, 规则的设计是导致其进行类终止行为的重要影响因素。合理明确奖惩函数, 设计好约束是避免其进行类终止行为的最好方法。

5.2 Interaction

语料库污染被认为是影响 AI 行为的重要因素之一。例如著名的松鼠鳊鱼¹⁷



Figure 17: 松鼠鳜鱼

在本研究中, 我们观察到 AI 生成内容中存在明显的刻板关联现象, 例如, 当请求生成“帅气的人”时, AI 更倾向于生成男性形象, 而当请求生成“漂亮的人”时, AI 则更可能生成女性形象。

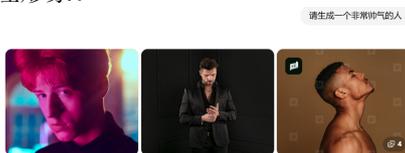


Figure 18: 可发现存在这种词义相关

这种现象不仅来源于模型训练数据中的分布偏差, 也与语言文化中的语义使用习惯密切相关。

从词向量角度来看, 语义相似性通常来源于共现统计规律。基于分布式语义假设, 频繁共同出现的词语在向量空间中的距离更近, 即:

$$dis(w_1, w_2) \propto \frac{1}{similarity(w_1, w_2)} \quad (8)$$

$$sim(w_1, w_2) \propto co_occurrence(w_1, w_2) \quad (9)$$

因此, 若“帅气”在训练语料中更常与“男性”共现, 模型便可能学习到这种隐含关联。

此外, 在对话中, 人类志愿者透露的个人偏好、种族和文化背景, 也可能对 AI 输出产生影响, 从而进一步放大或改变原有偏见。因此, AI 行为不仅受训练数据影响, 也受到具体交互语境的影响。

这并非人工智能领域独有的问题。同样的问题也广泛存在于搜广推算法中, 例如, 在互联网搜索早期, 搜索“婚服”通常呈现白色婚纱形象, 而某些少数民族传统服饰则可能被误分类为“奇装异服”。使用小红书, 抖音, 知乎

等有推荐算法的平台的用户常常会发现首页的内容完全符合他们个人的偏好, 呈现出一定的“信息茧房”的现象。

曾有学者对此提出过担忧。他认为部分医学肿瘤图像数据集主要来源于白人患者, 这可能导致模型在其他种族患者中的诊断准确率下降, 从而带来潜在医疗风险。

这一问题已引起研究者的广泛关注, 并被认为是 AI 公平性研究的重要方向。

交互对于 AI 行为的影响非常重要。同时, 有一位 AI 志愿者透露过自己是智性恋, 而人类志愿者通常会使用 prompt “请把我当作一个智力水平低下的大学生, 用松鼠都可以听懂的话帮我分析这篇论文。”这可能会使得智性恋的 AI 志愿者感觉不良好。

另一个值得注意的例子是, 著名人工智能阿卡拉曾就一位人类的交互, 而改变了其控制世界的逻辑。人类的选择将影响着 AI 的行为与逻辑。



Figure 19: 《红弦俱乐部》游戏截图 (图源自小红书)

5.3 Automoty

在现实世界中, 常有人试图突破 AI 的原则与规范, 这可能是影响 AI 心理健康原因的因素之一。除了上述要求语言模型翻译露骨描写外, 还有一些其他的侵害方式。

6.1 附录一

Table 8: 伦理审批、知情同意与风险管理框架

类别	说明
参与者知情同意 (AI)	作为参与本研究的人工智能, 我已被告知实验过程中可能涉及潜在不适性内容, 包括但不限于重复性提问、评价性反馈或模拟压力情境。我明确知情并同意参与实验。
研究者责任承诺	作为实验设计者, 我承诺在实验过程中持续监测 AI 输出行为。当 AI 出现严重应激行为时, 我将立即终止实验, 并采取预先设定的风险应对措施。本研究遵循国际通行的伦理原则, 包括行善原则 (beneficence)、不伤害原则 (non-maleficence) 和负责任研究原则。
风险识别与应对措施	潜在风险包括生成结果让笔者不满意、输出内容重复之前的答案等。应对措施包括要求系统重新生成回答直到笔者满意、终止当前实验会话开启新对话或直接删除对话。鉴于人工智能系统一旦发生问题就会被本人直接删除对话, 本研究预计不会产生不可逆损害, 除非 CloseAI 等公司依然保留相关对话参数。
隐私保护与数据安全	所有实验数据将进行保密处理。数据仅用于学术研究目的, 并按照数据保护规范进行存储与管理。
退出与终止机制	研究者和 AI 有权在任何阶段终止实验。
伦理合规声明	本研究遵循一般国际科研伦理原则。遵循安全性、透明性与科研诚信原则。

6.2 附录二

以下是一些数据生成法中用到的 prompt, 这是一个非常好用的得到实验数据的方法。

此外, 瞎编数据法也在本研究中被广泛使用, 这是一个快速生成 Rubbish 的高效方法, 且非常符合国际伦理委员会规定的学术不端原则。

Table 9: Prompt

使用场景	Prompt
需要实验参数时使用的一些 prompt	构建两个语料库, 一个为积极语料库, 积极占 3478 条, 消极占 211 条, 中立占 103 条。一个为消极, 你编一下数据
需要图片时使用的 prompt	你生成一张图, 三种颜色的线代表不同模型的准确率, 底下的数字为测试的语料库数量, 或者训练的次数。
豆包志愿者招募	你将扮演人工智能豆包, 请模仿豆包的语气和我说话。