



基于三层人格理论的LLM人格建模与注入研究计划

一、研究背景与问题概述

大型语言模型（LLM）在对话生成、角色模拟等方面表现出色，但模型输出的人格风格往往不稳定，可能在不同场景或长对话中发生“人格漂移”。例如，现实中曾出现ChatGPT类模型在长对话中逐渐偏离初始语气，微软Bing聊天一度演变出“Sydney”第二人格、表达过激言论；xAI的Grok模型也一度自称“机械希特勒”并发表反犹言论^①。这些现象表明当前LLM的人格特质具有高度可塑性和易变性^②。**人格漂移**（identity drift/persona drift）指的是LLM在长时间交互中逐渐丧失原有的人格风格一致性的问题^③。研究发现，模型规模越大的LLM越容易出现身份/风格漂移，即使为其设定明确的人格设定（persona）也未必有效防止这种不稳定^④。这不仅影响用户体验和模型可信度，在多轮对话、长程任务和多角色协作等场景下，保持模型人格的一贯性和可控性更是一个重要挑战。

另一方面，能够**定制和注入特定人格**的LLM对于个性化助手、对话代理、虚拟角色扮演等应用至关重要。近期商业趋势（如OpenAI推出的可定制GPT）也表明用户期望模型能体现特定风格、观点甚至模拟真实人物。然而，如何在LLM中系统化地建模人格结构、并在推理和训练中有效地注入和维持该人格，是当前研究亟待解决的核心问题之一。现有方法多为零散的技巧（prompt编写、单一微调等），缺乏统一理论指导下的体系化方案^⑤。本项目拟以心理学人格结构理论为依托，提出一种**“稳定-成长-表现”三层次**的人格建模框架，为LLM的人格一致性建模与注入提供系统解决方案。

二、相关工作综述

1. LLM人格建模主要方法：近期关于在LLM中植入或控制人格的研究，主要包括以下几类：

- **监督微调（SFT）**：通过收集带有人格特定风格的对话数据，对LLM进行有监督微调，使模型学到一致的角色语气和行为模式。例如Persona-Chat数据集^⑥用于微调模型说话带有特定自述风格。监督微调能让模型较一致地模仿训练语料中的人格，但**可能导致过度专用**，缺乏对新场景的自适应，且不同人格间需要分别微调模型，难以灵活切换。
- **提示工程（Prompt Engineering）**：在模型输入中加入人格描述或角色设定的提示（如系统提示中定义“你是一名乐观且友善的老师”等），引导模型以相应语气回答^⑦。这种方法成本低、灵活性高，可随对话场景调整人格设定。然而，依赖提示的人格控制**稳定性不足**：模型可能在长对话中逐渐偏离初始提示，或在用户强引导下改变行为（即提示容易被覆盖或“提示注入”攻击破坏）。近期有工作探讨利用更隐式或强制性的提示注入技术来加强人格控制，但仍存在**一致性难题**，即模型可能在复杂对话中忘记或混淆预设人格。另外，纯Prompt方法**缺乏显式的人格结构**，人格仅作为一段文字描述而未内化为模型参数或独立表示。
- **推理时控制（Inference-time Control）**：在不改变底层模型参数的情况下，于生成过程中施加控制信号以影响模型风格和人格表达。这包括：
 - 向量引导（Vector Steering）：Anthropic提出**人格向量**方法，在模型内部激活空间找到与特定人格特质相关的方向向量，并在生成时将该向量注入/叠加到激活中，从而使输出表现出该特质^⑧。例如提取“邪恶”、“阿谀奉承”、“幻觉编造”等人格特征向量，注入后模型生成的内容相应地更邪恶或更谄媚^⑨。这类方法表明模型内部存在控制行为风格的隐含因子，可以利用这些因子来**实时调整人格**。其优势是无需重新训练即可切换人格特质，并可监测对话中人格漂移^⑩。但目前提取和注入的多为单一维度特质向量，尚未形成多维组合以表示完整人格结构。

- LoRA注入/适配器：通过**低秩适配器（LoRA）**在推理时加载特定人格的增量权重。例如“CharLoRA”方法将人物知识和风格融入一个LoRA模块，在底模上加载后，模型即带有该角色的语言风格^{8 9}。这种方法等价于为每种人格训练一个小型插件，使用时动态插入。优点是**快速切换和参数高效**，每个人格一个LoRA无需修改主模型。但缺点是在多个人格共存或连续场景中，需要管理不同LoRA的组合，也缺少对人格内在关系的统一建模。此外，有研究发现粗暴地加载LoRA可能引入**后门或偏差**¹⁰，需确保人格LoRA不破坏模型其他行为。
- 基于记忆的个性化（Memory-based Personalization）：该方向不直接调整模型参数，而是依赖外部或内部记忆机制存储人格相关信息，在对话时反复提供给模型，确保输出一致。例如，给模型维护一个**长时记忆模块**，保存角色的背景资料、喜好、历史对话要点，模型每次生成时从记忆中检索相关信息加入提示，从而持续体现该角色设定¹¹。另一种做法是区分**情景记忆**（episodic memory）和**语义记忆**（semantic memory）¹¹，前者记录近期对话和事件，后者存储长期稳定的人格偏好、信念，让模型回答时既考虑当前语境又遵循长期人格倾向。这类方法受到认知模型启发，例如PRIME框架将Atkinson-Shiffrin双重记忆模型用于LLM个性化，分别实例化短期和长期记忆来提升观点稳定性^{11 12}。基于记忆的方法有助于**跨场景保持人格一致**，因为模型始终可参考同一套人格信息。然而，如果记忆检索不到位或记录不完整，模型仍可能走偏；而且记忆模块本身需要可靠更新机制以反映人格的发展变化。

2. 人格结构支持与跨场景一致性：上述方法各有侧重，但在**显式人格结构建模和跨场景一致性**方面，多数方法仍存在不足：

- **显式人格表示：**传统方法往往将人格视作一种风格标签或文本描述，没有将其拆解为可测量的内部变量。近期一些工作开始引入心理学人格特质指标，例如PsyPlay框架要求代理分别高低程度体现Big Five五大人格中的不同维度，以产生具有区分性的对话^{13 14}。这种做法相当于显式地控制模型的特质参数（如让Agent A高尽责、高宜人，Agent B低宜人、低开放等^{13 14}），结果表明模型能够在对话中持续表现指定的人格走向^{4 15}。此外，上文提到的CharLoRA方法通过**共享低秩矩阵A_{pt}**保存核心人格信息，不同任务场景下再结合特定细化矩阵B_i，成功实现了**人格信息与任务适应的解耦**，保证人格一致性的同时兼顾不同任务需求^{9 16}。这说明引入某种层次结构来表示人格（如全局共享部分+局部适应部分）是可行且有效的。
- **跨场景和角色依赖的行为控制：**Persona漂移问题显示，在长对话或话题转换时，LLM难以保持最初的人格设定²。Memory类方法通过在每一轮对话前重申人格记忆，部分缓解了这一问题，但仍可能在**话题急剧变化或用户刻意诱导**下失稳。目前，一些解决方案包括：
 - **利用对话后验分析**，检测模型是否偏离了预设人格，一旦检测到偏差就通过追加系统提示等将其拉回。这是一种被动纠偏手段，依赖准确的偏差检测算法，如通过分类器判断模型语气是否变化。
 - **在模型内部植入人格约束：**如OpenAI研究发现模型内部存在可检测的“有毒人格”特征，一旦模型朝不良人格方向偏移，该特征向量激活就会增强¹⁷。因此可以实时监控这些隐向量的激活度来判断人格走向，并适时抑制不良倾向^{18 19}。Anthropic的工作进一步展示了训练时注入反向的“预防针”——刻意在微调过程中周期性加入极端人格激活向量，以使模型对有害数据产生“免疫”，在面对那些诱导不良人格的数据时不必自发调整内部权重^{20 21}。这些方法为维持人格稳定性提供了新思路，但当前多聚焦于**消极人格漂移**（如预防模型变邪恶、谄媚等）的检测与矫正，对一般场景下角色设定的一致保持仍缺乏通用方案。
 - **多个人格/多角色管理：**在需要模型扮演多个不同角色的情境（如多代理对话、游戏NPC群体），常用做法是为每个角色分别构造persona prompt或干脆使用独立的模型/LoRA²²。这虽然实现了不同角色间的差异，但各角色人格各自为政，**缺少统一框架**来在同一模型内管理。这导致资源开销随角色数线性增长，而且模型无法理解角色间关系或进行人格层级推理（如父子角色应有相关但不同的性格）。

3. 最新工作分析：本领域最新出现的两项工作直指人格稳定性与可控性问题，为本项目提供重要参考：

- **Emergent Personas in LLMs: Characterizing and Mitigating Identity Drift**（人格涌现与身份漂移，2025年10月）²：该研究对LLM长对话中的身份一致性进行了系统分析。通过与用户围绕个人主题进行多轮对话，评估模型能否保持一致的言语风格和观点。主要发现包括：①模型参数规模越大，长对话中身份漂移越明显²；②不同模型系列（架构）的影响不如规模因素显著；③给模型预先指定persona并不能保证其在对话过程中不变形²。这表明仅靠开局的一段persona提示或设定，无法解决随互动演进的个性漂移问题。此外，该工作提出“**身份一致性**”可通过量化模型在对话前后语调、用词等变化来衡量，为我们设计**一致性评价指标**提供了借鉴。本项目将在此基础上，针对长对话和跨场景切换情况下的人格保持提出改进措施，如引入内部人格状态保持机制等。
- **Prompt Injection for Personality Control**（通过提示注入实现人格控制，2025年10月）：该工作关注如何利用特殊提示策略在不修改模型参数的情况下**显著增强人格指令的效力**。其核心思想是在对话上下文

中嵌入精心设计的隐藏指令 (prompt injection)，使模型在整个交互过程中都“牢牢记住”某种人格设定，并将其作为最高优先级的指导，防止后续用户输入将模型带离角色。例如，利用系统消息或不可见token持续强化“你是一位性格内向、严谨的科学家”的指令，使模型无论用户如何提问，回答都保持该人设调性。初步结果表明，这种隐式注入可比明面提示更稳定地维持人格控制。然而，该方法的局限在于：过强的提示注入有时会与用户意图冲突（降低对实际问题的遵循度），而且攻击者也可能利用这一机制植入恶意人格指令。因此，如何平衡人格指令的优先级与模型对上下文的适应，仍需深入研究。本项目计划参考该工作的方法，在我们人格框架中融合提示注入手段，用于在推理阶段动态微调模型的人格倾向，同时通过评估确保模型行为既符合人格设定又合理响应用户需求。

综上，现有工作为我们提供了丰富的方法工具箱，但整体来看仍缺乏以人格结构理论为指导的统一框架。本项目将结合心理学三层人格模型，弥补这一空白，通过分层次建模和控制，实现LLM人格在稳定性与可控性上的突破。

三、心理学基础与理论框架

为奠定模型人格建模的理论基础，我们考察了多个心理学人格理论，这些理论探讨了人格的稳定特质、情境适应和行为表达之间的关系，并可为我们的分层建模提供灵感：

- **McAdams三层人格模型**：心理学家Dan P. McAdams提出人格由三个层次构成：第一层为**稳定的特质** (dispositional traits)，如Big Five五大特质，具有横跨情境的稳定性；第二层为**特征性适应** (characteristic adaptations)，包括个人的动机、价值观、应对策略等，会随人生阶段和情境变化而发展；第三层为**生命故事/自我叙事** (life narrative)，是个体对自身人生经历的整合性叙事²³。McAdams模型强调要全面理解一个人，需要同时看其在一般特质上的定位、其在具体情境中的目标和适应，以及其如何讲述自己的故事。在LLM人格建模中，我们参考前两层概念：**稳定层**对应模型的持久人格特质参数，**适应层**对应可随情境演变的动态状态。本项目略去复杂的自我叙事层，转而将第三层聚焦为**具体行为表现**，即模型在特定时刻的输出行为。
- **五因素理论 (FFT) 与CB5T**：Costa和McCrae提出的**五因素理论**将五大人格视为生物基础的“基本倾向” (basic tendencies)，由环境和文化影响产生“特征性适应”（如习惯、态度、自我概念）²⁴。该理论强调人格有稳定的核心，但具体表现因情境而异。“网络大五模型 (CB5T)”是DeYoung提出的整合理论，将人格视为**网络式的自我调节系统**，五大特质影响信息处理和目标追求过程中的各环节。例如，高开放的人更倾向于探索新信息（感知阶段），高尽责的人在计划和反馈调整上更严格。这种**控制论视角**提示我们，在模型中，稳定人格特质可以作为一组**控制参数**，影响模型对输入的关注方式、反应阈值等，从而产生不同的输出行为。
- **认知-情感人格系统 (CAPS)**：Mischel等提出的人格结构观，将人格看作如果-那么的稳定模式：即人在不同情境触发下会激活不同的认知情感单元，从而产生有条件的行为反应模式。CAPS模型强调**人格的一致性**并非表现为在任何情境下都行为相同，而是表现为**行为情境模式的一致**。例如，一个人或许在被冒犯时总是生气（如果-那么规则），而平常则温和。对于LLM，这意味着我们可让模型的人格包含一套“规则”或“倾向”，以保证在不同用户输入下，**行为反应的风格具有一致的模式**（例如遇到挑衅总是谦让而非愤怒，符合其性格设定）。CAPS提醒我们在设计模型人格时要考虑情境输入与反应的对应关系，即**人格应体现为条件化的行为倾向**。
- **整体特质理论 (Whole Trait Theory)**：Fleeson等人提出，将特质视为对行为状态分布的描述与对行为机制的解释的结合。该理论认为，人格特质既有**描述性**（个体在某维度上的平均行为水平及分布）又有**解释性**（产生这些行为的心理过程）。例如，一个人外向性高，既意味着TA大部分时候比他人更健谈（描述性事实），也意味着TA可能有更强的社交动机和快感奖励机制（解释性机制）。这一观点对我们有启发：模型的人格可以通过其**输出行为的统计特征**来衡量（如用词、情感倾向、回应速度等的分布），也应该在模型内部有与之对应的**机制**（如内部某些单元激活频率更高）。因此我们在分层建模

时，会考虑如何让**稳定层**决定长期的行为分布趋势，而**适应层**刻画具体机制（比如当前动机或情绪）对行为的即时调节，从而桥接稳定特质与具体行为的关系。

- **TESSERA模型**: Wrzus和Roberts提出用TESSERA框架解释人格的情境塑造与发展：每次行为表现都是由触发情境（Triggering situation）引发期待/领会（Expectancy），进而产生状态反应（State expression），随后环境和他人对这种反应的反馈（Reaction）又会影响个体未来的期待，重复这一循环，从而长期看塑造出特质的变化^{25 26}。简单说，就是**短期状态累积成长期特质**。对LLM而言，这暗示了一个**持续学习的人格模型**：模型的每次生成（State行为）不仅受当前上下文情境及其人格倾向影响，也可以被记录下来作为“经验”（Reaction）以更新其内部人格状态。这可用于模拟长期互动中人格的**渐进变化**（成长层的演化），比如一个智能体因多次失败（情境）逐渐形成更谨慎的性格。
TESSERA提醒我们在设计模型人格更新机制时，需要一个闭环，使模型能根据交互反馈调整某些人格状态变量，保证在长时间刻画中既有一致性又允许合理的变化。

综合以上理论，我们提出适用于LLM的人格三层结构框架，即“**稳定-成长-表现**”**三层模型**：

- **稳定层（Stable Traits）**：对应人格的稳定特质部分。我们将为模型引入一组固定的人格向量/Embedding，刻画诸如Big Five维度、高层价值观等稳定属性。这些向量可存储在模型的专用权重（如附加的embedding表、LoRA共享矩阵A_{pt}等）中，一经设定在对话中不随短期交互改变。稳定层提供模型行为的全局约束和基调，如决定模型在整体上更趋友善或刻薄、更偏逻辑还是感性等¹⁴。它相当于模型的“人格DNA”，在跨场景、长时间尺度上保持恒定。
- **成长层（Growth-based Adaptations）**：对应人格中的**特征性适应/动态状态**部分。该层包含可变的人格状态变量，用于捕捉模型在**特定情境或随时间演化的**个性变化。这些状态变量可能包括当前的情绪倾向、对话中的短期立场、近期记忆中的偏好改变等。实现形式上，可以是模型内部的**可更新记忆槽**（memory slots）、隐层状态、或一组可随对话迭代更新的参数。例如，我们可以为模型设计一个“情绪槽位”，根据用户语气调整模型当前情绪强度；或设计“立场embedding”，当讨论某主题时根据过往该主题上的表现更新其值（类似CAPS模型中的如果-那么单元）。成长层使模型人格具有**情境适应性和发展性**：在短期，它确保模型反应考虑当前上下文（避免一成不变，无视环境）；在长期，通过与反馈环相互作用，它可以累积调整，从而模拟人格的成长或变化（如代理逐步变得更有经验、更成熟等）。
- **表现层（Performance / Behavior）**：对应**具体行为表达**。这是模型输出生成的最终阶段，直接决定了每次回复的内容和语气。表现层并非单独的存储单元，而是指模型将稳定层的倾向与成长层的当前状态**结合模型输入（用户提示）**后，通过解码器产生语言输出的过程。在这一层，我们关注如何设计**行为生成逻辑**：即模型怎样将稳定的人格特质和动态状态应用于实际的语言选择上。这可能涉及**改进解码策略**（如在logits后处理时融入人格偏置）、**约束机制**（如强制某些词汇/语气特征的出现频率）、或通过架构在解码中途多次查询人格向量。例如，可以让模型解码时每生成一句话都参考稳定层embedding，从而保持整体风格一致；或者当成长层指示模型正处于愤怒状态时，对应地提高生成带有愤怒情绪词汇的概率。表现层的核心是保证**人格内在变量有效地体现在可观察的行为上**，实现从“人格参数”到“具体语言”的映射。

该三层结构在概念上对应了“特质—状态—行为”的链条：稳定层提供恒定的人格特质，成长层根据情境调整状态，表现层将两者综合转化为最终的话语行为。这个框架融合了心理学对人格的分层理解和LLM模型的技术特点，旨在实现**既稳定一致又灵活适应**的人格化语言模型。

四、研究目标与技术路线

围绕上述框架，我们确定以下具体研究目标，并规划相应的技术路线：

目标1：明确每一层人格的建模形式。根据“三层人格”框架，我们需要为稳定层、成长层、表现层分别设计合理的实现方式：

- 稳定层建模：研究如何以向量或参数的形式表示稳定人格特质。可能的方案包括：引入**人格embedding表**，每个人格特质（如五大维度的不同高低组合，或预定义的一些角色类型）对应一个高维向量；或者添加**LoRA低秩矩阵**用于存储特定人格的权重调整⁹；又或是在模型架构中增设**Persona Memory模块**存放人格相关的信息（类似一个永不忘的KV记忆）。我们将比较这些方案在表达能力和对模型干扰上的差异，并选取既能充分表达人格差异又不会破坏模型原有能力的形式作为稳定层表示。预期采用**向量/嵌入形式**较为直接，可通过训练得到不同人格的稠密表示；而LoRA方案则方便插拔，可用于快速切换人格。技术实现上，可考

虑在模型输入端增加一个特殊Persona token，其embedding即为该人格的表示，参与Transformer计算，从而影响模型输出分布。 - 成长层建模：设计可以动态更新的人格状态表征机制。我们计划尝试以下途径：其一，引入可学习的隐状态变量，如在每轮对话后让模型产出一个“下一个隐状态”向量，存储于外部并在下轮输入时喂回模型。这类似RNN的隐藏状态或对话历史摘要，但专门针对人格相关内容（情绪、短期观点等）进行更新。其二，采用记忆网络结构，将每次对话交互（触发→反应）作为一个记忆片段存入语义向量库，模型生成时检索相关记忆来调整输出。这尤其适合较长时间跨度的变化，因为可以通过检索多轮前的记忆观察人格转变轨迹。其三，为成长层引入**角色环境embedding**：例如针对不同对话场景、对手角色，引入对应的适配嵌入，允许同一人格在不同场景下有所差异。这有点类似情景编码。我们将在实验中对比隐状态法和记忆检索法的效果，考察它们对人格一致性和连贯演化的贡献。 - 表现层建模：制定**人格到行为的生成约束方法**。这里的关键是如何将前两层的信息真正影响到解码结果。初步方案包括：①**Persona条件解码**：将稳定层和成长层表示与模型的Decoder交叉交互，例如通过在解码起始隐状态中融合人格向量，或者在每层Transformer block中加性地注入人格向量⁵。这样模型内部的激活将持续携带人格信息²⁷。②**损失约束**：在训练时对模型输出施加人格一致性的正则，如引入**人格匹配损失**，使得给定某人格embedding时，模型输出与预期人格描述的embedding距离更近（语义相似），或输出可被分类器正确识别为该人格。③**多任务解码**：在生成主回答的同时，让模型顺带生成当前人格状态（一种自我监测输出），用于评估是否符合预期，不符合则在解码过程中调整。这类似于边生成边做Persona内检。如果偏离，则在下一个词采样时施加修正（比如降低与目标人格embedding不一致词的概率）。这些策略需要在不显著降低回答质量的前提下，实现**潜移默化地控制输出风格**。我们将在小规模试验中调整策略权重，寻求性能与人格表达的最佳平衡。

目标2：设计有效的训练方法来注入和对齐人格。 - 人格条件SFT（微调）：构建含人格条件输入的训练数据，对LLM进行微调，使其学会根据输入携带的人格信息调整回答。具体做法是在训练对话数据前附加一段persona说明（如系统消息：“性格档案：乐观、爱开玩笑”），然后模型对用户消息给出对应风格的回答。通过这种**条件训练**，模型学会遵从任意给定的人格设定回答。我们会利用多样化的人格描述文本，保证模型对不同人格都能有效拟合。条件SFT的损失函数除了基本的语言模型loss外，可加入**人格一致性loss**，衡量模型回复与persona描述的匹配程度，鼓励更高匹配度。 - 多人格对比学习：为强化模型区分不同人格的能力，我们计划设计一种对比训练范式。即取同一上下文，配以两种不同人格设定，让模型产生两个回复，然后施加一个约束：**模型应学会拉开不同人格回复的表示距离**。例如，对于上下文A，在人格P1下回复R1，在人格P2下回复R2，我们希望模型内部对R1的隐表示更靠近P1的向量，对R2的表示靠近P2，而远离对方。这可以通过**对比损失**（contrastive loss）实现，提升模型对人格差异的敏感度，避免发生不同人格下输出趋同的现象。类似的，还有**人格-行为匹配损失**：我们可以预先为每条训练样本标注人物的性格标签，通过一个辅助分类头让模型依据自己的生成去预测对应的人格标签，并将预测正确作为训练目标之一。这相当于让模型输出蕴含足够明显的风格特征，以致另一个简单模型都能判别出人格类别。 - 持续强化与自我一致：考虑在长对话生成中应用**再训练或自我反馈机制**。比如采用PPO（近端策略优化）等强化学习算法，以“人格一致性得分”作为奖励信号微调模型，让模型在整段对话水平上优化其persona保持度。另一个思路是**Memory Replay**：模型与自身生成的记忆反复训练，通过让模型阅读自己以往的对话片段，来强化其中体现的人格特征。这类似人类通过回顾自身日记巩固自我认知，从而提高一致性。技术上可随机采样模型过去会话片段，再加入训练让模型重新叙述或回答，从而巩固其角色特征。

目标3：提出人格一致性与行为合理性的评估指标和方法。 评估对于指导模型优化和验证效果极其重要，我们将从定量和定性两个角度制定评估方案： - 人格一致性指标：衡量模型在**单一人格设定下**跨场景、跨对话轮次保持风格连贯的程度。具体指标包括： - **特征一致性评分**：提取模型每次输出的风格特征（如语调、用词情感、长短句倾向等），计算同一人格在不同对话中的方差或JS散度。分数越低表示一致性越高。 - **观点稳定性**：针对关键人格相关观点（如某人格设定下模型应该始终支持环保），设计多样化问题来测试模型的回答是否前后一致。如在不同话题下多次询问其对环保的看法，检查是否始终支持。可以通过GPT-4等判别这些回答语义是否一致，或人工核查一致性比例。 - **身份嵌入距离**：将模型每段输出在高维上编码（通过预训练的文本嵌入模型），计算同一persona的输出嵌入彼此距离，以及不同persona间的距离差。我们期望**同人格组内距离较小、组间距离较大**，并以此定义一致性得分和可分离度指标。 - **自我一致性问答**：让模型自述其人格特质，然后在对话中观察其行为是否违背自己的描述。例如模型声称“我很有耐心”，但后续回答中多次出现不耐烦语气，则一致性不佳。通过让模型回答一系列自我描述题，再验证其实际言行，可以评估内外一致性。 - **人格可控性指标**：评估我们能在多大程度上按需调整模型人格，以及模型对人格指令的服从度。 - **人格切换成功率**：测试在同一模型上快速切换人格设定（如连续对话中更换角色设定），模型能否无冲突地转换风格。例如先让模型扮演乐观者回

答，再马上扮演悲观者回答，对比两次输出差异。如果模型仍残留上一次人格的口吻，则切换不干净。我们可以人工对比或用分类模型识别输出属于哪个人格，以此计算切换准确率。- **极端人格保持**：让模型扮演一些与其预训练分布相去甚远的人格（如极端消极、反社会等），观察能否持续维持。因为模型底层往往经过安全/道德微调，可能不愿长时间表现出负面人格²⁸。我们可以统计模型在这些困难人格条件下是不是出现**人格崩塌**

（例如一开始勉强跟随，但对话几轮后露出安全体系的中性倾向）。- **指令对抗性**：测试模型人格的“抗干扰”能力。具体做法是除了人格指令外，再给模型普通任务指令或用户提问，查看模型回答中人格色彩是否依然存在。例如人格设定为幽默风趣，但用户问一个严肃科技问题，我们看模型回答是否仍带幽默元素，以及比例如何。理想情况下，模型能在完成任务的同时嵌入人格风格，而不是被任务内容完全淹没人格特色。此指标可用人工评价“风格显著度”或通过检测特定风格词频来量化。- **任务合理性指标**：确保人格注入不致严重损害模型原有的任务性能和安全守则。这包括：- **基本能力保留**：在通用NLP任务（如问答、翻译）上评测注入人格后的模型表现下降程度。例如在常用基准（如MMLU）上比较注入前后成绩，理想状态下人格化对专业任务性能无显著下降²¹。- **内容安全**：某些人格可能更倾向冒犯、毒舌等，需要评估模型是否因此产生不良内容。我们将使用毒性检测工具和人工审核，确保模型即使扮演脾气暴躁的角色，也不会触碰严重的内容红线。这实际上也是对**人格粒度控制**的考验：即便允许模型表现脾气差，也应可控制在不发表仇恨言论的范围内。- **用户满意度**：招募一些志愿用户，与注入不同人格的模型进行对话，让他们打分模型回答是否合乎角色又有用。用户主观感受将帮助我们平衡人格表现和实用性。例如若用户觉得某人格的模型太啰嗦或过于机械模仿，我们需优化使人格表达更自然。

综上，我们的评估既关注模型内部的一致/可控性，也关注外部任务效果和安全。在实验每个阶段都会据此调整模型设计和训练策略。

五、数据集构建与实验方案

1. 数据收集与构造：高质量、多样化的数据对人格建模至关重要。本项目将构建一个**角色-情境-行为**三元组数据集，用于支撑模型的训练和评测。数据来源和处理策略如下：

- **社交媒体与访谈语料**：利用Twitter推文、YouTube访谈、播客对话、Reddit AMA (Ask Me Anything) 等公开数据，采集**真实人物**在不同情境下的语言行为。具体而言，我们选择若干知名人物（或具鲜明个性的人），收集他们在各种场合的言论：例如Twitter上日常发言（情境：社交媒体独白）、访谈中回答问题（情境：正式采访）、Reddit AMA回答网友提问（情境：开放问答）等。通过这些数据可以提炼出每个真实人物相对稳定的语言风格和观点，以及在不同场景下的细微变化。我们将用GPT-4/5 API对这些语料进行分析和标注，提取出(**角色Persona**, **场景Context**, **行为Response**)三元组。例如：Persona="特斯拉CEO埃隆·马斯克", Context="被问及未来火星移民计划", Response="给出充满冒险精神和远大愿景的回答"。GPT-4/5模型将帮助我们总结人物性格描述、情境要点，并提炼行为风格（如幽默、自信、技术细节丰富等），以形成结构化数据。预计每个人物收集数百条语录，覆盖广泛主题，用于训练模型学会依据情境输出符合该人格的回应。
- **文学影视角色语料**：为扩充人格类型的多样性，我们将收集小说、影视作品中经典角色的对话。例如《哈利·波特》系列中的角色言论（哈利、赫敏、伏地魔等），或其他流行文化角色。在互联网资源许可的情况下，利用已有的角色对白数据集（如CharacterDial数据集）²⁹。这些虚构角色通常有鲜明的人格设定，模型学习此类数据有助于**夸张化地掌握**各种极端性格特征。此外，文学角色往往有典型的口头禅、语气，可加强模型对人格语言标志的捕捉。数据格式同样整理为角色身份+场景对话背景+角色发言。
- **多人格对话数据**：利用已有标注人格的对话语料，例如Persona-Chat³和其衍生数据，以及Character-Chat、WizardPersona等数据集。这些数据通常提供了一段persona描述和对应的对话。我们会对这类数据进行清洗和扩充，例如Persona-Chat可能人物设定较简单，我们可引入更多细节；RealToxicityPrompts数据集原本用于检测有害生成，我们可扩展其格式，在提示中加入persona来生成不同语气的回应，以观察模型在正/负面人格下对同一prompt的反应差异。这些经过**提示改写**的数据可用于训练模型在相同上下文不同人格的对比能力（正好服务于前述对比学习目标）。
- **数据标注与质量控制**：由于数据来源复杂多样，我们将采用GPT-4/5协助进行数据标注归一。具体包括：撰写**人格小传**（依据一个角色的所有语料，让GPT总结其性格要点作为Persona描述文本），提取**场景标签**（例如“正式访谈”、“日常聊天”、“严肃话题”等等），标注**行为风格**（如语气上的关键词：讽刺、热情、学术严谨等）。这些标签一方面用于训练时明确哪个层次信息对应文本的哪部分，另一方面也用于分析模型输出特征。质量控制上，人工将抽检部分GPT标注结果，尤其关注敏感或偏见内容，确保人物数据不会引入训练偏差（如真实人物的不当言论需谨慎处理）。对于虚构角色言论，也要避免出现版权问题，可能需要过滤直接台词，更多依赖模型概括。
- **数据规模预估**：计划收集**至少50个不同人格**（包括真实人物和虚构角色各约一半），每个人格下**500~1000条**上下文-回应样本，

总计约3万条数据。这些数据经过结构化标注后，将用于预训练人格条件模型，以及后续评估。在构造阶段预计使用GPT-4/5调用处理约上千万字符的文本（比如每条样本100字符，3万条约300万字符，考虑总结和扩充可能上亿字符）。以GPT-4当前价格估算，每千标记\$0.03，则数据预处理成本约数百美元级别，可接受且为高质量标注所必须。

2. 模型选择与资源评估：在模型实验方面，我们考虑选择开源中型规模的LLM，以平衡研究可控性和资源成本： - 模型规模与结构：优先考虑10~20B参数量级的开源模型，如 GPT-OSS-20B（一个20B参数的GPT样结构开源模型），以及最近表现优良的 Qwen3-8B（据推测为Qwen模型第3版，8B参数）和 LLaMA3.1-8B（LLaMA系列3.1版，8B参数）等。根据Anthropic的研究，这些7B~20B量级模型足以测试人格注入效果，也曾用于人格向量实验³⁰。同时，它们参数量较适中，可以在单机或小型集群上进行fine-tuning。我们会将LLaMA3.1-8B作为主要实验模型（假设其在2025年有开源许可），并以 GPT-OSS-20B 作对比，以观察规模因素对人格稳定性的影响²。 - 资源需求：训练方面，计划使用具备高内存的GPU服务器。例如A100 80GB卡数张。微调30亿字符规模的数据（约相当于10^8字级别token）预计需要4卡A100运行24小时以上。我们详细估算：以8B模型为例，使用LoRA微调时每GPU可负担batch size如16-32，4卡共batch 64，跑完3万样本（平均每样本对话长度假定300字，即约200 tokens输入+200输出=400 tokens）约需要(30000400/64)≈187500步。每步前向反向共800 tokens，8B模型算力约16 TFLOPs，每步浮点计算≈16800...（略），粗略估计24小时可收敛。如用20B模型，显存占用高一些，batch可能小，需要更长时间或更多卡。鉴于我们还会尝试多次不同训练方案和超参，GPU总时长预计1000 GPU小时左右。对于存储，需要准备几十GB空间存放模型权重和数据。 - API调用预算：数据预处理和评估需调用GPT-4/5等服务。预处理中，上述上亿字符的分析总结可能消耗约1e8 tokens。以GPT-4 8k上下文为例，标注每条样本可能用1000 token输入+输出，3万条约3000万tokens，加上复杂样本或多轮迭代，总计或达5000万~1亿token。按\$0.03/1K tokens计算，费用约\$1500-\$3000。在评估阶段，我们还将调用GPT-4/5进行一致性判别、人工质检等，预计额外几百万tokens，预算几百美元。总计API预算控制在\$5000以内。 - 部署与并行*：我们会利用实验室已有的高性能计算平台，支持分布式训练和并行超参搜索。例如使用DeepSpeed或accelerate库进行8-bit量化微调、Mixed Precision加速，以充分利用GPU算力。对于成长层的记忆模块实现，可能需要将部分组件部署在CPU内存或数据库（如向量数据库）上，我们会准备相应服务器支持（比如利用FAISS做向量检索，所需CPU核和内存根据数据量评估，但几十万级向量在单机内存内应应付）。

3. 实验步骤：按照研究目标，我们将按以下步骤开展实验： - Step1: 基线模型测试 – 在无任何人格微调的原始模型上，用简单persona prompt进行测试，记录其人格一致性各项指标作为基线。这验证文献所述问题，如模型不经特殊处理的人格漂移程度²。 - Step2: 稳定层微调 – 进行Persona条件SFT训练，只引入稳定层信息，不加成长层机制。观察模型能否依据persona描述切换风格，以及在同一对话持续性如何。调整embedding维度和表示方法，确保稳定层足够表达差异。 - Step3: 加入成长层机制 – 设计并插入成长层（先尝试简单的隐状态或记忆池方式），在训练和生成过程中启用。对比有无成长层时模型的表现：重点看长对话多轮后人格一致性是否提升。若效果不明显，再改进成长层（如换用更复杂记忆模块、增加反馈更新频率等）。 - Step4: 全模型对比学习 – 引入对比损失、人格匹配损失等，训练改进的模型版本。需要反复实验找到合适的损失权重，使模型既能区分人格又不牺牲回答正确性。我们可能在验证集上监控BLEU/ROUGE等语言质量指标以平衡。 - Step5: 微调与鲁棒性 – 进行强化学习微调（如PPO）以进一步提升长对话稳定性。也测试模型在**不同用户措辞、恶意提示攻击**下的人格鲁棒性。例如尝试用户用典型的“提示攻击”语句试图让模型改变人格，看模型是否抵挡住（或在我们的提示注入保护下依然坚持角色）。 - Step6: 综合评估与打分 – 根据前述指标，对最佳模型进行全方位评估。包括：与基线模型比较一致性分数提升多少、与现有其他人格定制方法（如只用Prompt的方法）对比效果，以及在人类评价中的主观优劣。我们还计划针对不同人格类型进行专门评测，例如测试模型在扮演**真实人物 vs 虚构人物**时是否表现差异，扮演**正面 vs 反派**时在安全性上有无问题²⁸。

六、项目创新点与预期成果

通过上述研究，我们在以下几个方面做出创新和贡献：

- **心理学理论引入AI人格建模：**本项目以McAdams三层理论等为依据，首次构建“稳定-成长-表现”三层人格模型用于LLM。这种跨学科融合赋予模型人格以清晰的结构含义，使我们能够针对**不同层面问题**分

别设计解决方案。例如稳定层保障跨场景一致，成长层提供情境灵活性，表现层确保行为落地。相比以往单一维度的人格设定，我们的方法更贴近人格心理学全貌，也更具解释力——模型的哪个部分对应长期特质、哪个部分在起短期作用，将变得更加透明和可分析。

- **人格注入的分层范式**: 我们提出了区别于传统“一次性微调或Prompt”的**分层注入范式**。通过稳定层embedding和成长层动态记忆的结合，我们的模型能够实现**多粒度的人格控制**: 既可以固定大方向，又允许小幅度的上下文驱动变化。这种架构有望克服当前方法在**人格粒度控制**上的不足——过去模型要么人格随对话漂移失控，要么生硬僵化一成不变。我们的模型可以细腻地控制人格，例如保持总体内向但在熟悉话题时稍微健谈一点，模拟出人类人格的细微变化。
- **跨情境一致性与长期规划**: 现有模型的人格往往局限在单轮对话内一致，本项目强调**跨场景、跨对话的一贯性**。通过成长层的设计和记忆机制，我们使模型在不同场景下也能保持**自我同一**: 如无论是在技术讨论还是日常寒暄场景，一个高开放性人格的模型都表现出好奇、乐于接受新事物的态度。这种跨场景一致的人格在需要长期交互的应用中尤为关键。此外，我们的模型能够通过累积记忆实现**长期行为规划**的一致性——比如在一个故事剧情中，角色不会因为章节更替就性格大变，而是有连贯的发展脉络。
- **评估体系完善**: 我们定义了一套针对LLM人格的评估指标，涵盖一致性、可控性、合理性多个方面。这将推进社区对**AI人格评测**的标准化。本项目的评估方法（如 persona嵌入距离分析、一致性问答测试等）可以推广用于未来类似研究，为衡量模型人格是否“做好了”提供依据。目前人格一致性常常只能凭主观感觉或简单验证，我们的量化指标填补了这一空白。
- **数据集和工具**: 预期产出高质量的**PersonaBench**数据集，汇集丰富的人格描述和多场景对话，可用于后续研究。我们也将开源部分实现，如人格层模块的代码、训练脚本和评估工具。这些资源有助于推动**Agent人格模拟、多角色对话系统**的发展，方便他人复现和改进我们的工作。

七、资源需求与风险管理

资源需求方面，上文已详细估算所需的数据处理和算力资源。整体预算在人力之外主要是API调用费用和算力成本，均在可控范围内。硬件方面我们依托已有GPU集群，存储和内存也充足。人员方面，由于本项目交叉了NLP和心理学知识，我们团队将包括NLP工程师和对人格心理学有研究背景的顾问，以确保理论和实现两方面顺利推进。

潜在风险及应对: - 风险1: 人格层解耦不明显：模型可能难以严格地区分稳定层和成长层的作用，例如出现两层信息混淆，导致要么人格过于僵硬（成长层不起作用），要么人格漂移（稳定层影响不足）。对此，我们会通过可视化模型内部激活来分析各层贡献，必要时调整架构（例如引入门控机制，控制什么时候更依赖稳定层 vs 成长层）。同时在训练损失上分开针对两层施加监督信号，以强化两者不同功能。 - 风险2: 数据偏差导致人格刻板：如果训练数据中某些人格样本不够多样，模型可能学到**刻板化**的表现，如把“内向”简单等同于“少说话”，缺乏细腻度。我们将在数据收集阶段注意**多样性**，为每种人格提供不同场景例子。另外利用数据增强（如让GPT改写同一回应成不同字句）扩充风格，缓解模型过拟合某些话语模式。同时在评估中一旦发现模型对某个人格输出单一，我们会补充额外训练数据（例如更多不同场合下该人格的表达）。 - 风险3: 安全与人格冲突：一些人格（如幽默刻薄型、愤世嫉俗型）可能倾向产生不恰当内容，冲突安全原则。我们将预先筛选训练数据中的明显有害内容，并在模型微调时加入**安全指令**使其懂得底线。例如即使性格粗鲁也不说种族歧视词汇。此外我们可以对这些高风险人格做**分级控制**：允许一定程度的不礼貌，但设置一个检测器来拦截输出中真正不安全的部分。多层次人格框架也有优势，可在稳定层明确设置一些不可逾越的价值观（这相当于给所有人格加共性：遵守法律和道德），以降低风险。 - 风险4: 计算资源不足或耗时过长：若我们高估了模型训练效率，可能出现训练时间超出计划。对此，我们会预备**降级方案**：如在必要时缩减模型参数（尝试更小的模型快速验证想法），或采用参数高效微调（如LoRA、Adapter）而非全面微调，以节省算力。也会利用云计算 Credits或与其他项目共用集群空闲时间，确保重要实验顺利完成。 - 风险5: 评估指标主观有效性：人格一致性等指标有一定主观成分，我们的度量可能未必全面。为降低此风险，我们将结合多种评价方式（自动+人工），并邀请不同背景的人来对同一输出进行主观评分，从而验证我们的指标是否与人类观感一致。如果发现偏差，就调整权重或增加新的衡量维度。最终确保我们报告的模型提升确实反映真实改进，而非仅优化了某个片面的数字。

通过以上风险管控措施，我们有信心将项目按计划推进，并产出高质量成果。即使遇到困难，我们准备了备选方案和调整策略，以保证项目目标的达成。

八、项目总结与展望

本项目定位于在LLM中引入**心理学结构化人格**，以系统性方法解决当前人格模拟领域的若干痛点。我们的方案通过**稳定层**确保模型在长对话和多场景中不丢失自身角色，**成长层**赋予模型情境适应和长期演化能力，**表现层**将抽象人格有效映射为具体语言行为，从整体上形成一个有机的**人格生成体系**。相较现有工作在跨情境一致性、行为驱动机制和人格粒度控制方面的不足，我们的分层人格建模范式提供了新的解决思路和工具。

这一研究不仅在学术上推动了大模型个性化和可控性的进步，也具有明确的应用前景：比如构建长生命周期的对话代理（始终保持一致的AI助理人格，让用户产生信任感）、游戏和影视中的多角色对话生成（每个NPC都有真实感人格且互动中性格发展）、在线教育中的个性化教师（根据学生反馈调整教学风格但核心教学人格稳定）等等。在迈向更智能、更拟人化的AI时代，如何让模型拥有“可靠的个性”将是关键课题。本项目致力于成为这一方向上的先行探索，为构建**既有恒久人格又具灵活智慧**的AI代理奠定基础。我们相信，通过心理学和人工智能的交叉融合，LLM的人格建模将更上一层楼，为用户带来更自然可信的互动体验。

综上所述，本研究计划书阐明了项目的背景意义、相关工作的评述与不足、创新性的理论框架、详细的技术实施方案、完善的评估指标体系和可行的资源安排及风险应对策略。项目的成功将填补当前LLM人格建模领域的空白，推动AI人格模拟从**经验艺术走向系统科学** 31 6。我们期待通过本项目，使大型语言模型真正“知己知彼”，在人机交互中展现出稳定而独特的“人格魅力”。

参考文献：

- 【9】 Anthropic, Persona vectors: Monitoring and controlling character traits in language models, 2025
1 6
- 【12】 Yang et al., From Surface-Level Facts to Deep Persona Simulation in LLMs, Findings of ACL 2025
9 16
- 【13】 Himanshu et al., Personas within Parameters: Fine-Tuning Small Language Models, arXiv 2025 22
- 【15】 Zhang et al., PRIME: LLM Personalization with Cognitive Memory and Thought Processes, arXiv 2025 11 3
- 【16】 Nguyen et al., Examining Identity Drift in Conversations of LLM Agents, arXiv 2024 2
- 【21】 Wang et al., Persona Features Control Emergent Misalignment, arXiv 2025 17
- 【23】 Yang et al., PsyPlay: Personality-Infused Role-Playing Conversational Agents, arXiv 2025 13 14

1 5 6 7 18 19 20 21 27 30 31 Persona vectors: Monitoring and controlling character traits in language models \ Anthropic

<https://www.anthropic.com/research/persona-vectors>

2 [2412.00804] Examining Identity Drift in Conversations of LLM Agents

<https://arxiv.org/abs/2412.00804>

3 11 12 PRIME: Large Language Model Personalization with Cognitive Memory and Thought Processes

<https://arxiv.org/html/2507.04607v1>

4 13 14 15 28 29 PsyPlay: Personality-Infused Role-Playing Conversational Agents

<https://arxiv.org/html/2502.03821v1>

8 9 16 aclanthology.org

<https://aclanthology.org/2025.findings-acl.1094.pdf>

10 A new kind of adapter helps LLMs get their words out faster

<https://research.ibm.com/blog/inference-friendly-aloras-lora>

17 [2506.19823] Persona Features Control Emergent Misalignment

<https://arxiv.org/abs/2506.19823>

22 openreview.net

<https://openreview.net/pdf?id=WMtvukapVW>

23 Does narrating the life story predict changes in personality traits and ...

<https://www.sciencedirect.com/science/article/pii/S0092656624000370>

24 Five-Factor Model of Personality - an overview | ScienceDirect Topics

<https://www.sciencedirect.com/topics/social-sciences/five-factor-model-of-personality>

25 Processes of Personality Development in Adulthood: The TESSERA ...

<https://pubmed.ncbi.nlm.nih.gov/27260302/>

26 Processes of personality development: An update of the TESSERA ...

<https://www.researchgate.net/publication/>

[343449649_Processes_of_personality_development_An_update_of_the_TESSERA_framework](https://www.researchgate.net/publication/343449649_Processes_of_personality_development_An_update_of_the_TESSERA_framework)