

构建具备高级自我意识的通用人工智能：一种基于  
多模态输入输出，无限记忆处理，多LLM架构内部  
对话生成/自我评估更新的方法

Kaijie Du



# Contents

<b>I</b>	<b>AGI研究宏观架构整体介绍</b>	<b>1</b>
<b>1</b>	<b>引言</b>	<b>3</b>
1.1	研究背景与动机 . . . . .	3
1.2	项目的理论基础 . . . . .	3
<b>2</b>	<b>绪论</b>	<b>5</b>
2.1	研究目标 . . . . .	5
2.2	研究内容 . . . . .	5
2.2.1	架构设计 . . . . .	5
2.2.2	输入机制 . . . . .	6
2.2.3	输出机制 . . . . .	7
2.2.4	训练与学习机制 . . . . .	7
2.2.5	自我评估与优化 . . . . .	8
2.2.6	虚拟环境的整合与实验 . . . . .	8
2.2.7	真实世界具身实践 . . . . .	8
2.3	预期成果 . . . . .	8
2.4	研究意义 . . . . .	9
<b>II</b>	<b>Voyager做了什么？</b>	<b>11</b>
<b>III</b>	<b>关于第一个I:‘无限记忆，无限超越’的具体研究步骤</b>	<b>13</b>
<b>3</b>	<b>Literature Review</b>	<b>15</b>
3.1	Introduction . . . . .	15
3.2	Theoretical Foundations and Innovations . . . . .	15
3.2.1	Infinite Memory Transformers . . . . .	15
3.2.2	Memory Augmentation in Transformers . . . . .	15
3.3	Open-Source Implementations . . . . .	16
3.3.1	Ring Attention with Blockwise Transformers . . . . .	16

3.3.2	Memory-Enhanced Transformer Models . . . . .	16
3.3.3	Comprehensive Resource Lists . . . . .	16
3.4	Conclusion . . . . .	16
<b>4</b>	<b>研究思路与具体方案</b>	<b>17</b>
4.1	睡眠学习-无限记忆transformer网络介绍 . . . . .	17
4.2	技术评估与集成 . . . . .	17
4.2.1	无限注意力机制 . . . . .	17
4.2.2	记忆封装策略 . . . . .	17
4.2.3	长期记忆库的管理 . . . . .	17
4.2.4	自我记忆/自学知识的评估机制 . . . . .	18
4.2.5	知识固化过程 . . . . .	18
4.2.6	集成方法 . . . . .	18
4.3	系统设计 . . . . .	18
4.3.1	短期记忆到长期记忆的转换机制 . . . . .	18
4.3.2	记忆信息加工和编码策略 . . . . .	18
4.3.3	模块间的记忆信息流动 . . . . .	18
4.3.4	记忆数据存储结构 . . . . .	18
4.3.5	智能索引与快速检索 . . . . .	19
4.4	数据管理策略 . . . . .	19
4.4.1	持久化存储的信息管理 . . . . .	19
4.5	自我评估机制 . . . . .	19
4.5.1	自我评估的框架 . . . . .	20
4.5.2	反馈循环 . . . . .	20
4.5.3	持续学习和适应 . . . . .	20
4.5.4	集成的挑战 . . . . .	20
4.6	睡眠微调与自我优化 . . . . .	20
4.6.1	微调的实现 . . . . .	20
4.6.2	知识固化 . . . . .	20
4.6.3	自我优化评估 . . . . .	20
4.7	实验与评估 . . . . .	21
4.7.1	实验设计 . . . . .	21
4.7.2	评估指标 . . . . .	21
4.7.3	长期跟踪和分析 . . . . .	21
4.7.4	反馈机制 . . . . .	21

## IV 关于第二个I: 'Immediate Response, Real-Time Interaction'的

<b>具体研究步骤</b>	<b>23</b>
<b>5 Literature Review</b>	<b>25</b>
<b>6 Research Plan</b>	<b>27</b>
6.1 Objective . . . . .	27
6.2 Methodology . . . . .	27
6.2.1 Input Truncation and Analysis System . . . . .	27
6.2.2 Real-Time Processing Architecture . . . . .	27
6.2.3 Feedback and Adaptation Mechanism . . . . .	27
6.2.4 Prototype and Testing . . . . .	28
6.3 Expected Outcomes . . . . .	28
<b>V 关于第三个I: 'Self-Awareness, Infinite Possibilities' 的具体研究步骤</b>	<b>29</b>
<b>7 Literature Review</b>	<b>31</b>
7.1 Introduction . . . . .	31
7.2 Key Projects and Technologies . . . . .	31
7.2.1 DEPS . . . . .	31
7.2.2 Voyager . . . . .	31
7.3 Relevant Technologies . . . . .	31
<b>8 Detailed Plan for Self-Aware AGI Development</b>	<b>33</b>
8.1 Overview . . . . .	33
8.2 Conceptual Framework . . . . .	33
8.3 Architectural Design . . . . .	33
8.4 Implementation Strategy . . . . .	33
8.5 Challenges and Solutions . . . . .	34
8.6 Future Directions . . . . .	34
<b>VI 关于第四个I: 'Infinite Output, Infinite Agency' 的具体研究步骤</b>	<b>35</b>
<b>9 Literature Review</b>	<b>37</b>
9.1 Introduction . . . . .	37
9.2 Key Projects and Technologies . . . . .	37
9.2.1 Open-ended Learning Environments . . . . .	37
9.2.2 Interactive and Adaptive Systems . . . . .	37
9.3 Relevant Technologies . . . . .	37

<b>10 Detailed Plan for Enhancing Output and Agency in AGI</b>	<b>39</b>
10.1 Overview . . . . .	39
10.2 Architectural Design . . . . .	39
10.2.1 Modular Design for Scalability . . . . .	39
10.2.2 Real-time Data Processing . . . . .	39
10.3 Implementation Strategy . . . . .	39
10.4 Challenges and Solutions . . . . .	40
10.5 Future Directions . . . . .	40

Kaijie Du

Date: 2024/ 4/ 24

# Part I

## AGI研究宏观架构整体介绍





# Chapter 1

## 引言

### 1.1 研究背景与动机

当前的人工智能技术，虽然在特定任务上展现出卓越的性能，但普遍缺乏自主意识和持续自我进化的能力。本研究计划旨在探索 and 实现一种新型人工通用智能（AGI），该AGI将模拟并扩展人类的核心能力，包括意识、自我驱动的学习机制和创造性思维。

现有的研究中，最对我启发的是minedojo团队的voyager架构。它让我意识到，我们距离AGI其实只差临门一脚——不是模型体量、能力上的差距，更不是设备与算力上的，仅仅是一些后端的组合架构与微调工作。结合我对自我意识与真正智能的理解，我已经有了一种对AGI实现的初步构想。

我将使用Llama3系列替换GPT的API调用，进而实现本地部署（或者云端部署）的真正AGI。

### 1.2 项目的理论基础

本项目认为，AGI的开发不仅仅是技术的挑战，更是哲学和认知科学的深度融合。

哲学与心理学上，我认为从现有AI到AGI本质上是对尼采所描述的从‘last man’到‘超人’概念的一种实现。通过以下几个‘T’来模拟并扩展人类的能力：

- 类比人类‘记忆’的‘无限’记忆+睡眠学习，实现知识的长期稳定和即时更新，从而获得无限的超越性。
- 类比人类‘活着’，实时输入输出，获得即时的反应力
- 类比人类‘意识’，通过构建多个语言模型（LLM）之间的动态交互架构，模拟人类意识中“自我”与“超我”的功能，从而获得无限的可能性
- 类比人类‘自由’的无限制输出，实现全方位的输出能力，从文字到物理操作，从而获得无限的能动性。

此外，本项目还结合了以下认知科学/神经科学的思想。

- ‘记忆’：我们将基于无限注意力的transformer模仿人类长短期记忆转换机制，将基于“记忆蒸馏”构建外部记忆库来模拟睡眠做梦的机制。
- ‘意识’：目前的LLM的生成更像是生物学上的“反射”，就像别人让我们自我介绍或者背诵英语课文的生成。但是“意识”更像是基于内部言语的生成模型。我们认为自我意识存在的基础是“为我”，也就是“私心”——从面向人类到面向自己。

# Chapter 2

## 绪论

### 2.1 研究目标

1. 建立多个语言模型（LLM）之间的动态交互架构，模拟人类意识中“自我”与“超我”的意识功能。
2. 实现AGI的全方位输入，包括多模态的实时输入解析识别。
3. 实现AGI的全方位输出能力，包括实时输出、完全控制电脑操作及网络访问能力（当然也包含多模态信息，代码调试等）。
4. 开发基于睡眠模式的连续学习和记忆整合机制，以支持知识的长期稳固和即时更新。
5. 设计多个自我评估和自我进化的机制，通过演员评论员架构进行效果反馈和自我调整。
6. 构建虚拟环境Minecraft，初步应用AGI，实验智能体在模拟环境中的自我学习、决策与行动。
7. 增加具身智能模块，让AGI能在现实中与“世界”交互，执行物理任务和环境操作。

### 2.2 研究内容

#### 2.2.1 架构设计

多模态智能体系统架构：设计一个包含多个功能性LLM的AGI系统，具体包括视觉处理、语言理解、策略规划等模块。每个LLM都可以作为一个独立的思考和操作实体，通过内部通信相互协作，模拟人类大脑的功能分化。

- **超我模块**：负责监督、评价自我模块的决策，提供道德和社会规范的反馈，类似于人类的社会和道德判断机制。

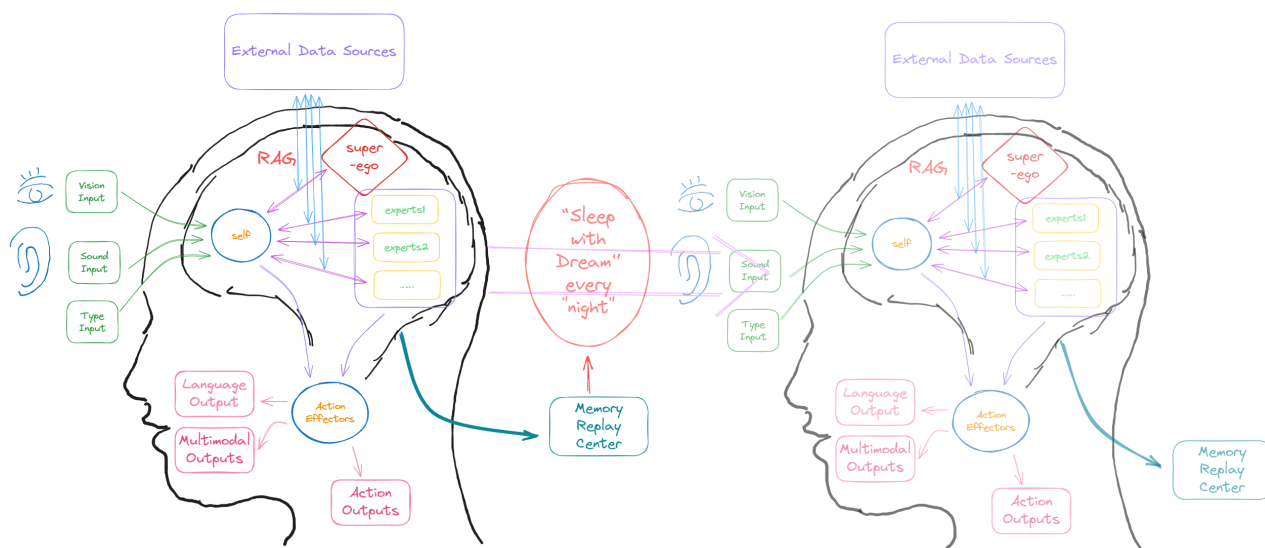


Figure 2.1: the structure of MY AGI structure

- **自我模块**：核心决策和执行单位，处理日常任务，决策制定，并与外部世界进行交互。该模块也负责管理内部对话，以模拟人类的思考过程。
- **小脑模块**：控制AGI的物理动作和反射，确保动作的准确性和协调性，模拟人类小脑的运动控制功能。
- **专业知识模块**：存储和处理专业知识和技能，如医学、法律等，支持专业决策和问题解决。甚至可以包含多模态的输出。
- **检索增强模块**：提升信息检索和知识整合能力，支持内部言语生成和思维链条的续展，通过先进的检索算法和知识图谱技术，实时更新和优化知识库。

为了实现类似人类的短期记忆长度，采用基于无限注意力机制的Transformer架构，如Infinite Transformer，使模型能够处理和记忆长序列的信息。这种架构通过动态扩展注意力机制来适应输入序列的长度，有效管理长期依赖关系。

### 2.2.2 输入机制

- **实时信息处理**：开发能够即时捕捉并解析语音、文本、图像等多模态输入的系统。这包括传统的文本输入和通过麦克风捕获的语音命令，以及通过摄像头捕获的视觉信息。
- **中断与连续性处理**：输入系统将具备实时中断处理能力，类似于人类对话中可以被打断的自然交流。这需要高度的并发处理能力和低延迟响应设计。
- **自动语义标记**：采用先进的自然语言处理技术，自动进行实时语义分析和重要信息提取，增强对环境的理解能力。推荐使用如BERT或GPT模型。

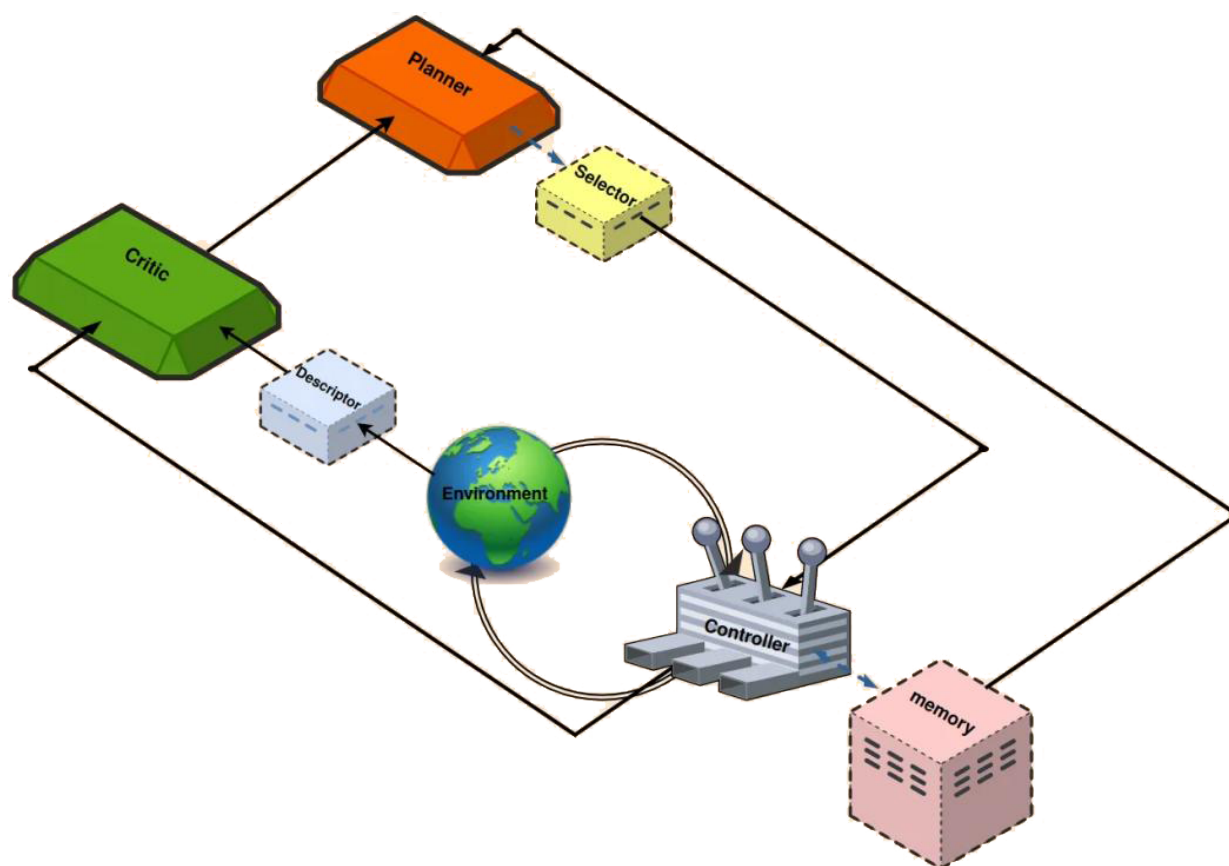


Figure 2.2: AGI architecture proposed by **OTHER** researchers

### 2.2.3 输出机制

- **完全控制权限**：实现AGI对计算机硬件（如键盘、鼠标操作）和网络资源的完全控制权限，支持复杂任务的自动执行。
- **多模态输出**：系统能够生成语言、图像、声音等多种格式的输出，包括文本到语音的转换、自动生成图像或视频内容。利用DALL-E 3的先进功能，AGI能够根据文本提示创建高质量、细节丰富的图像，支持广泛的创意和实用场景。此外，AGI还能控制机器人进行物理动作，如移动和抓取，确保动作的准确性和协调性。输出机制包括机器人身体动作的控制，通过语言模型指导机器人的行动。

### 2.2.4 训练与学习机制

**睡眠模拟与记忆整合**：模拟人类的NREM和REM睡眠周期，设计记忆重播和创新思维模拟算法，定期对AGI的学习成果进行巩固和创新拓展。这包括采用构建长期记忆库来优化记忆存储过程，基于长期记忆库构造‘梦境’来进行知识固化，复现voyager项目的技能库学习，以及使用多模式睡眠机制研究成果来增强连续学习的效率。

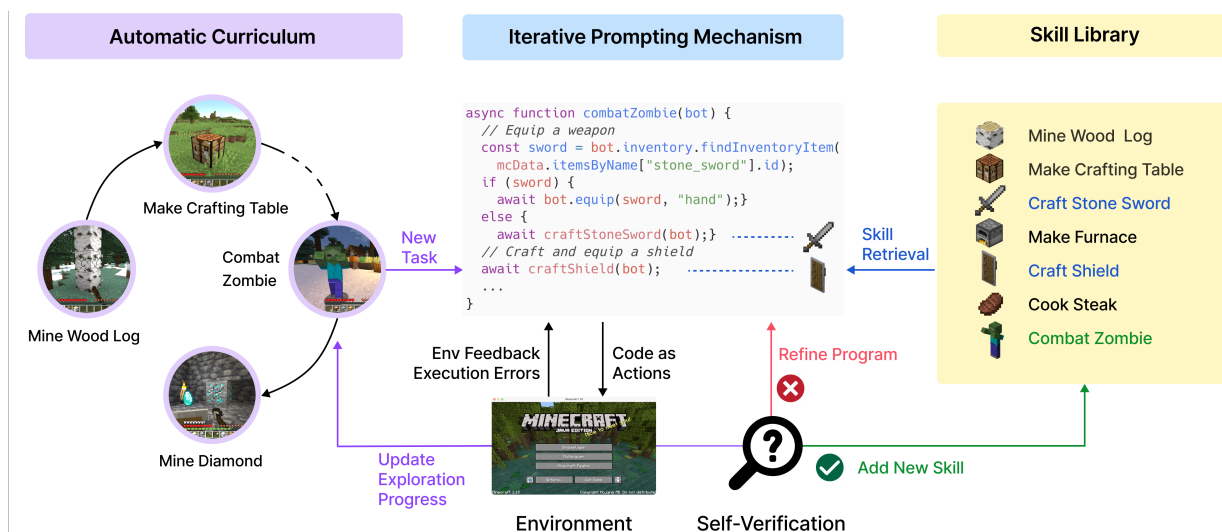


Figure 2.3: the structure of VOYAGER

### 2.2.5 自我评估与优化

演员评论员算法的应用：结合多个超我LLM的评估，为自我LLM的行为和决策提供反馈，支持自我调整 and 性能优化。这些超我模块类似于软更新的目标网络，延迟更新，收集整个AGI系统的所有操作，建立一个非常平衡的评估系统。超我告诉自我需要学什么，最近的学习or交互反馈是否有效，并利用RL技术如RL"AI"F利用反馈优化自我模块。

### 2.2.6 虚拟环境的整合与实验

结合足够逼真的虚拟环境，本研究将探索AGI是否能够展现与人类相似的‘自我能动的无限性’，并通过这种能动性自我学习，胜任各种通用任务。我们特别选用Minecraft作为实验平台，因其提供丰富的交互可能性和开放式的环境动态，适合进行复杂任务的学习和测试。

### 2.2.7 真实世界具身实践

MC里如果能实现，现实的具身智能就加个‘小脑’就好了。

## 2.3 预期成果

我相信，我能在博士期间实现如下成果：

- 实现一套具备高级自我意识的AGI系统，能够自主进行学习、决策和自我评估。
- 发表高影响力的研究论文，并在国际会议上展示研究成果。

- 开发相关的软件和工具包，为学术界和工业界提供强大的研究和应用平台。
- 利用AGI正式开启第五次工业革命！！！实现国家富强！

## 2.4 研究意义

本研究将推动人工智能从单一任务执行者向具有完整自我意识和自主学习能力的通用智能体转变，开创AI技术的新篇章，**开辟时代的新纪元**，具有深远的理论价值和广泛的应用前景。





## Part II

Voyager做了什么？



## Part III

# 关于第一个I:‘无限记忆，无限超越’的具体 研究步骤



# Chapter 3

## Literature Review

### 3.1 Introduction

This chapter presents a comprehensive review of the state-of-the-art research and open-source projects that have significantly advanced the capabilities of AI models in handling infinite memory and context lengths. The focus is on identifying and synthesizing key contributions that have shaped current understanding and implementation of extended memory in AI systems.

### 3.2 Theoretical Foundations and Innovations

#### 3.2.1 Infinite Memory Transformers

Studies like the  $\infty$ -former introduce mechanisms where memory does not depend on the sequence length, allowing the model to handle indefinitely long contexts without additional computational costs. This model uses a continuous-space attention mechanism to manage long-term memory efficiently [?].

#### 3.2.2 Memory Augmentation in Transformers

Enhancements in memory handling within Transformer models, such as the Memory Transformer and Memformer, offer methodologies for integrating additional memory tokens that aid in sequence processing tasks, showing improvements in tasks like machine translation by leveraging these extended memory capabilities [?].

## 3.3 Open-Source Implementations

### 3.3.1 Ring Attention with Blockwise Transformers

The GitHub repository *RingAttention* provides an implementation that allows Transformers to handle large context sizes by computing attention and feedforward in a blockwise manner, significantly reducing the computational overhead [?].

### 3.3.2 Memory-Enhanced Transformer Models

The *memory-transformer-xl* project on GitHub introduces a variant of the Transformer-XL, which updates its memory using an attention-based mechanism rather than a conventional queue, enhancing the model’s ability to manage and utilize its memory more effectively [?].

### 3.3.3 Comprehensive Resource Lists

Resources like *awesome-transformer-nlp* on GitHub offer curated lists of tools, papers, and projects focused on Transformer networks and their applications in NLP, providing a valuable starting point for researchers looking to explore the landscape of Transformer technologies [?].

## 3.4 Conclusion

The reviewed literature and projects highlight a trend towards enhancing the memory capabilities of AI models, which is critical for applications requiring extensive context and long-duration interactions. These advancements form the foundation upon which future research can build more sophisticated and capable AI systems.

# Chapter 4

## 研究思路与具体方案

### 4.1 睡眠学习-无限记忆transformer网络介绍

本研究旨在开发一种具有无限记忆能力的LLM架构，通过集成先进的机器学习技术，模仿和扩展人类的记忆和学习过程，解决长期记忆持久化问题。

### 4.2 技术评估与集成

为了实现高级的人工通用智能（AGI）无限记忆架构，我们首先需要评估当前的技术，确定如何将它们集成到单一的系统中。此部分的研究将集中在以下关键技术的评估上：

#### 4.2.1 无限注意力机制

评估现有的无限注意力机制如Transformer-XL和Compressive Transformers，以确定它们在处理长序列信息方面的有效性。我们将探索这些机制如何在不牺牲计算效率的情况下，扩展模型的注意力跨度。

#### 4.2.2 记忆封装策略

研究如何将信息从短期记忆转移到长期记忆库。我们将分析不同的数据封装策略，如增量学习和经验回放，以及它们如何影响记忆的保留。

#### 4.2.3 长期记忆库的管理

探讨如何有效地构建和维护长期记忆库。这包括数据的存储、索引和检索机制，特别是当涉及到大量信息时。

#### 4.2.4 自我记忆/自学知识的评估机制

制定评估AGI行为和知识积累的方法，并确定哪些信息应该被保留或丢弃。这一过程将模拟人类的日常反思和内省。

#### 4.2.5 知识固化过程

研究“晚上”期间的知识固化过程，包括如何实现自我微调以及如何利用LLM来生成丰富的语料，以支持这一过程。

#### 4.2.6 集成方法

最后，我们将研究如何将上述技术集成为一个统一的系统。这将包括设计模块之间的接口，以及确保信息在不同模块间流动的逻辑和协议。

这些研究活动将为创建一个能够无限扩展知识和记忆的AGI系统奠定基础。

### 4.3 系统设计

构建一个具备无限记忆能力的AGI系统需要精心设计其架构以支持从感知到长期记忆的信息流动。在这一部分，我们将考虑以下设计要素：

#### 4.3.1 短期记忆到长期记忆的转换机制

我们将训练一个LLM，根据AGI‘短期记忆’的语料与知识生成相应的能存储于记忆库中的数据与知识图谱。

#### 4.3.2 记忆信息加工和编码策略

详细阐述AGI的长期记忆库如何处理和编码输入信息，确保关键数据可以被有效地转换为可存储和可检索的格式。

#### 4.3.3 模块间的记忆信息流动

设计模块间的接口，例如，如何在“自我”模块和“记忆回放中心”之间有效地传输数据，以及“超我”模块如何调节这两者之间的交互。

#### 4.3.4 记忆数据存储结构

考虑使用分布式存储系统，如NoSQL数据库或云存储服务来存储长期记忆，以及如何保障数据的持久化和稳定性。



### 4.3.5 智能索引与快速检索

研究如何实现高效的数据索引和检索机制，以支持对历史信息的快速访问，特别是在面对海量数据时。

通过这些设计原则和方法，我们的目标是创建一个灵活、可扩展和具有弹性的系统，能够适应新的学习场景和不断变化的信息流。我们也将考虑安全和隐私问题，确保所有的设计都符合伦理标准和数据保护法规

## 4.4 数据管理策略

### 4.4.1 持久化存储的信息管理

为了保证AGI系统的长期记忆能力，必须对信息进行有效的持久化存储。我们将采用以下策略：

#### 数据的存储格式

选择合适的存储格式，如关系数据库或NoSQL数据库，依据数据的访问模式和更新频率来优化存储结构。

#### 数据备份与恢复

制定严格的数据备份和恢复流程，以保障数据的安全性和完整性，防止数据丢失或损坏。

#### 数据冗余策略

实施数据冗余策略，如使用分布式存储解决方案，以增加数据的可靠性和系统的容错能力。

#### 数据安全性与隐私

确保所有存储的数据遵循最新的安全和隐私标准，实施加密和访问控制机制，保护数据不被未经授权访问或滥用。

这些策略将确保系统能够高效地存储和管理大量数据，同时维持系统的响应速度和用户体验。

## 4.5 自我评估机制

在AGI系统中实施自我评估机制是至关重要的，它允许系统自我反思并优化其行为和知识库。此节将探讨实现有效自我评估的策略。MineDojo/Voyager CraftJarvis/MC-Planner

### 4.5.1 自我评估的框架

构建一个框架，使AGI能够定期评估其性能和决策过程。这包括开发评估指标来度量其决策的效果和准确性。这一部分，我将结合梁老师的DEPS工作与minedojo团队的Voyager来实现。

### 4.5.2 反馈循环

实施反馈循环，让AGI根据自我评估的结果调整其行为。这包括修改知识库中的数据和调整决策算法。

### 4.5.3 持续学习和适应

开发机制，使AGI可以从新的经验中学习并适应，将错误视为学习的机会，并不断更新其行为模式和知识结构。

### 4.5.4 集成的挑战

探讨在AGI系统中集成自我评估机制的挑战，包括技术、操作和伦理方面的问题，并提出相应的解决方案。

这些自我评估机制将使AGI更加智能和自适应，能够更好地理解和反应于其操作环境。

## 4.6 睡眠微调与自我优化

为了使AGI系统适应不断变化的环境并改进其性能，微调与模型优化是必不可少的步骤。此节将详述如何利用微调来增强AGI的学习能力和‘条件反射’。

### 4.6.1 微调的实现

详细描述如何在“晚上”的内省期间，利用累积的经验和反馈来生成高质量的微调内容（类比人类REM时期的‘梦境’），来调整和优化AGI模型的参数。

### 4.6.2 知识固化

分析如何通过微调过程固化AGI的‘白天’所学习到的知识。

### 4.6.3 自我优化评估

训练一系列评估和测试的LLM，来验证微调后模型的性能改进，确保这些改变达到了预期的效果。

## 4.7 实验与评估

为了验证AGI系统的性能和适应性，进行严格的实验和评估是至关重要的。本节将描述如何设计和实施这些实验，以及如何正确评估结果。

### 4.7.1 实验设计

详细说明如何设计实验来测试AGI系统的不同功能和模块，包括记忆能力、学习效率、决策质量等。

### 4.7.2 评估指标

定义一系列量化指标来度量AGI系统的性能，如准确率、响应时间、错误率等。

### 4.7.3 长期跟踪和分析

建立长期的监测系统，用于跟踪AGI系统的表现和进化，确保系统能持续优化并适应新的任务和环境。

### 4.7.4 反馈机制

实施反馈机制，允许从实验结果中学习，并将这些学习成果反馈到系统的微调和优化中。

通过这些综合的实验和评估策略，我们可以确保AGI系统在实际应用中的可靠性和有效性。



## Part IV

关于第二个I: 'Immediate Response,  
Real-Time Interaction'的具体研究步骤



# Chapter 5

## Literature Review

The development of Artificial General Intelligence (AGI) capable of real-time interaction has been highlighted as a critical component in advancing the field. Significant studies such as those by Agüera y Arcas and Norvig [?] emphasize the necessity for generality in AGI, incorporating the ability to process and react to diverse inputs in real time. Another pertinent project is the Voice2Action framework [?], which addresses the challenges of translating real-time voice commands into actionable tasks in virtual environments. Additionally, the OpenAGI project [?] provides a platform for AGI development, focusing on integrating domain-specific models to enhance real-time task execution.

- Agüera y Arcas and Norvig. "Levels of AGI: Operationalizing Progress on the Path to AGI." <https://arxiv.labs.arxiv.org/html/2311.02462>
- Voice2Action: "Language Models as Agent for Efficient Real-Time Interaction in Virtual Reality." <https://arxiv.labs.arxiv.org/html/2310.00092>
- OpenAGI Project: "OpenAGI: When LLM Meets Domain Experts." <https://github.com/agiresearch/OpenAGI>





# Chapter 6

## Research Plan

### 6.1 Objective

The primary goal of this research is to develop an AGI capable of analyzing and responding to inputs with minimal latency, ensuring both accuracy and immediacy in task execution.

### 6.2 Methodology

#### 6.2.1 Input Truncation and Analysis System

Implement a system to analyze incoming data streams in real-time, truncating non-essential data to focus processing power on critical information. This will utilize insights from the OpenAGI framework, which has demonstrated success in managing complex tasks efficiently [?].

#### 6.2.2 Real-Time Processing Architecture

Design an adaptive, modular architecture that processes inputs as they arrive, based on techniques developed in the Voice2Action project [?].

#### 6.2.3 Feedback and Adaptation Mechanism

Incorporate a feedback loop inspired by the RLTF mechanism from OpenAGI, enhancing the AGI's decision-making capabilities through continuous learning from real-world interactions [?].

### 6.2.4 Prototype and Testing

Develop and test a prototype within a controlled environment, progressively increasing task complexity to assess performance and adaptability.

## 6.3 Expected Outcomes

- Enhanced responsiveness and decision-making speed in AGI.
- A robust feedback system enabling continuous learning and adaptation from interaction outcomes.

## Part V

### 关于第三个I: 'Self-Awareness, Infinite Possibilities' 的具体研究步骤



# Chapter 7

## Literature Review

### 7.1 Introduction

This chapter reviews the key projects and technological advancements contributing to the development of self-aware AGI. It highlights the essential role of iterative learning and adaptive capabilities in these systems.

### 7.2 Key Projects and Technologies

Significant projects such as DEPS and Voyager have pioneered the integration of self-aware capabilities in AGI systems.

#### 7.2.1 DEPS

The DEPS project focuses on deep planning and network systems. It integrates strategic thinking with deep learning to enhance AGI's decision-making processes.

#### 7.2.2 Voyager

Voyager utilizes large language models within a Minecraft simulation to foster an open-ended learning environment for AGI. This project showcases how AGI can develop and refine skills autonomously over time. More information and source code can be found at <https://arxiv.org/abs/2305.16291>.

### 7.3 Relevant Technologies

Technologies like Large Language Models (LLM) such as GPT-4 play a crucial role in the ongoing development of self-aware AGI by facilitating complex task management and iterative

learning processes.

# Chapter 8

## Detailed Plan for Self-Aware AGI Development

### 8.1 Overview

This section outlines the comprehensive steps and methodologies to be employed in the development of self-aware AGI systems.

### 8.2 Conceptual Framework

Develop a conceptual framework that defines self-awareness in AGI, including its capabilities for self-assessment and adaptation.

### 8.3 Architectural Design

Detail the design of a modular architecture that supports dynamic learning and self-aware capabilities. This includes feedback mechanisms that allow AGI systems to self-monitor and adapt their strategies based on performance outcomes.

### 8.4 Implementation Strategy

Outline a phased implementation strategy that includes:

- Development of initial self-awareness capabilities.
- Integration with simulation environments for testing and refinement.
- Iterative enhancements based on test results.

## 8.5 Challenges and Solutions

Discuss potential challenges such as complexity management and ethical considerations, and propose viable solutions and preventive measures.

## 8.6 Future Directions

Speculate on future research directions and how they might evolve to address emerging challenges and leverage new technological advancements.



## Part VI

### 关于第四个I: 'Infinite Output, Infinite Agency'的具体研究步骤



# Chapter 9

## Literature Review

### 9.1 Introduction

This chapter explores the foundational research and existing technologies that enable AGI systems to generate diverse and adaptive outputs, critical for achieving a state of infinite agency.

### 9.2 Key Projects and Technologies

Review significant projects that have contributed to advancements in AGI's output capabilities, focusing on those that demonstrate innovative interaction with complex environments.

#### 9.2.1 Open-ended Learning Environments

Discuss projects like OpenAI's GPT-3, which demonstrate the ability to generate creative and contextually appropriate outputs across diverse prompts and environments.

#### 9.2.2 Interactive and Adaptive Systems

Highlight systems designed for dynamic interaction, such as IBM Watson's deployment in various fields requiring real-time data interpretation and response generation.

### 9.3 Relevant Technologies

Examine the technologies that are pivotal in enhancing the output capabilities of AGI, including neural networks capable of handling vast arrays of data and producing real-time responses.



# Chapter 10

## Detailed Plan for Enhancing Output and Agency in AGI

### 10.1 Overview

Outline the goals and objectives of this phase of AGI development, focusing on enhancing the system's ability to produce outputs that are not just reactive but proactively adaptive.

### 10.2 Architectural Design

Detail the design of a scalable architecture that supports vast output capabilities and can interact with a variety of environments and inputs dynamically.

#### 10.2.1 Modular Design for Scalability

Explain how modular components can enhance the system's adaptability and facilitate the integration of new capabilities as technology evolves.

#### 10.2.2 Real-time Data Processing

Discuss the implementation of advanced data processing modules that enable the AGI to handle and respond to real-time data efficiently.

### 10.3 Implementation Strategy

Outline a step-by-step strategy for implementing the enhancements, including:

- Development of new output modules.

- Integration with real-world data streams for dynamic interaction.
- Iterative testing and refinement in simulated and real environments.

## 10.4 Challenges and Solutions

Identify potential challenges, such as managing the complexity of real-time data processing and ensuring ethical considerations in output decisions. Propose strategies to address these challenges.

## 10.5 Future Directions

Speculate on future enhancements and technologies that could further expand the capabilities of AGI, potentially leading to truly autonomous systems capable of novel creativity and problem-solving.