Al Context Engineering 领域若干重要论文

📛 2025年10月12日

○ 1分钟阅读

#Context Engineering #Agent #Al

这里将收集Context Engineering相关的重要文献,具体解读将 在其他博客展开。

这里将收集Context Engineering相关的重要文献,具体解读将在其 他博客展开, 这里也会提供链接方便跳转。

1. A Survey of Context **Engineering for Large Language Models**

作者: L. Mei 等人

发表时间: 2025年7月

核心贡献:

该综述论文系统性地建立了上下文工程(Context Engineering)的理 论框架,分析了1400多篇相关研究论文,首次提出了上下文工程的 正式定义与分类法。论文确立了上下文工程作为超越简单提示词设 计的系统性学科,涵盖了对大语言模型信息负载的系统优化 $([9^{+}])$

创新点:

提出"上下文工程"正式定义:系统性优化大语言模型所接收信息 的技术与方法

建立了包含上下文检索、处理和管理的完整技术框架 揭示了大语言模型理解能力与生成能力的不对称性关键问题 为该领域未来研究指明了方向,建立了评估挑战与标准化方法

更具体的解读: TBD

目录

文章信息

字数

阅读时间

发布时间

更新时间

标签

#Context Engineering #Agent

2. Context Engineering: Enhancing Large Language Model Performance Through Comprehensive Contextual Management

作者: ResearchGate团队

发表时间: 2025年

核心贡献:

该论文将上下文工程定义为通过全面上下文管理提高大语言模型性能的系统化框架。它超越了传统的提示词工程,关注如何战略性地设计和优化提供给模型的上下文信息([28+])。

创新点:

提出综合上下文管理框架,作为提升LLM性能的系统化方法 超越简单提示词优化,关注上下文的全面设计与管理 为上下文工程提供了实践指导和理论基础 展示了上下文工程在实际应用中的性能提升效果

更具体的解读: TBD

3. Core Context Aware Transformers for Long

作者: ICML 2025入选研究团队

发表时间: 2025年

核心贡献:

该论文提出了核心上下文感知(Core Context Aware, CCA)注意力机制,用于大语言模型的高效长上下文建模。这是一个插拔式模块,专门针对处理和管理长上下文信息的设计([31+])。

创新点:

设计了高效的长上下文建模机制,解决了标准Transformer架构中的上下文长度限制

提出的CCA注意力机制可作为插拔组件,便于集成到现有模型

为处理超长文本序列提供了计算复杂度更低的解决方案在长文档理解任务中展现了优于传统注意力机制的性能

更具体的解读: TBD

上下文工程领域概述

上下文工程(Context Engineering)是人工智能领域,特别是大语言模型(LLM)应用中的前沿技术方向。它关注如何优化提供给AI模型的上下文信息,以最大化模型性能。

定义与发展:

上下文工程是对提示词工程的扩展,关注系统性优化信息负载 发展历程从早期的简单提示词优化,到检索增强生成(RAG),再 到多智能体系统

解决了从上下文窗口限制到信息相关性优化等一系列挑战

主要技术组件:

上下文检索与生成:包括链式思考(CoT)、树式思考(ToT)等方法

上下文处理: 长序列处理、自修正与适应、结构化信息整合

上下文管理:记忆分层架构、上下文压缩技术

应用领域:

AI智能体开发: 提升多轮对话一致性和任务执行效率

检索增强生成: 结合外部知识源增强模型输出

长期对话系统:解决上下文窗口限制和信息衰减问题

多智能体协作:协调多个专门智能体之间的上下文交换

未来研究方向

理论框架完善: 建立更完善的形式化理论基础, 发展任务特定

的评估指标

技术突破:探索更高效的上下文压缩与管理机制,改进长上下

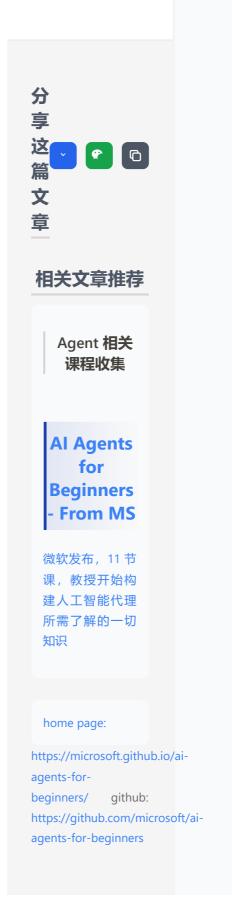
文处理算法

跨领域应用:将上下文工程应用于多模态模型、强化学习等领

域

社会责任: 研究上下文工程对模型安全、隐私和伦理影响

上下文工程作为AI领域的新兴重要方向,正在快速演进并重塑我们与大语言模型交互的方式。以上三篇论文代表了该领域从理论建立 到技术实现的全链条创新,为后续研究和应用提供了重要参考。



https://learn.microsoft.com/en-us/shows/ai-agents-for-beginners/
https://learn.microsoft.com/en-us/shows/ai-agents-for-beginners/what-

Agent Lightning

are-ai-agents

介绍

这个项目有以下 重要作用:

零代码/低代码训练 Al Agent (核心价值):

最大亮点: 它允 许你使用强化学 习 (Reinforcement Learning, RL) 等 高级优化算法来 训练你现有的 AI

Agent, 而**几乎不** 需要修改你的 Agent 业务逻辑

代码。这意味着 你可以保留你用 LangChain,

AutoGen, CrewAl, OpenAl

SDK 等框架(甚

至裸 Python)编写的 Agent 逻辑,然后让Agent Lightning负责优化它的决

策过程。

解决痛点:传统 上,将 RL 等技术 应用到现有 Agent 框架中需 要大量的工程改造和集成工作。 Agent Lightning 极大地简化了这个过程。

强大的优化能力:

算法支持:內置 支持强化学习 (VERL)作为核心 优化算法,并明 确提到支持自动 提示优化 (Automatic

提供训练基础设 施:

Context Engineering

Context Engineering 是...

性明性、XX学们可靠性可以得到显著提升。

广泛的兼容性和 灵活性:

框架无关: 明确 支持所有主流 Agent 框架 (LangChain, OpenAl Agent SDK, AutoGen, CrewAl) 以及纯 Python 实现的 Agent。你可以 "即插即用"。

多 Agent 系统优

化:可以在包含 多个 Agent 的复杂系统中,选择性地优化其中一个或几个特定的 Agent,而不是整个系统,提供了更精细的控制。