书生二期第一次课笔记

通过本次课程的学习了解到书生提供了一个当前比较全的大模型架构平台。包括大模型的构建与训练,大模型智能体,多模态的大模型智能体。这些基础的开源平台使得更多人可以借助其做大模型的开发与部署应用工作。





书生•浦语 2.0 (InternLM2) 的体系

面向不同的使用需求 , 每个规格包含三个模型版本

7B

为轻量级的研究和应用提供了 一个轻便但性能不俗的模型

20B

模型的综合性能更为强劲,可有效支持更加复杂的实用场景

InternLM2-Base

高质量和具有很强可塑性的模型基座 , 是模型进行深度领域适配的高质量起 点

InternLM2

在 Base 基础上,在多个能力方向进行了强化,在评测中成绩优异,同时保持了很好的通用语言能力,是我们推荐的在大部分应用中考虑选用的优秀基座

InternLM2-Chat

在 Base 基础上,经过 SFT 和 RLHF,面向对话交互进行了优化,具有很好的指令遵循、共情聊天和调用工具等的能力

书生•浦语 2.0 (InternLM2) 的主要亮点



超长上下文

四

综合性能 全面提升



优秀的对话 和创作体验



工具调用能力整体升级



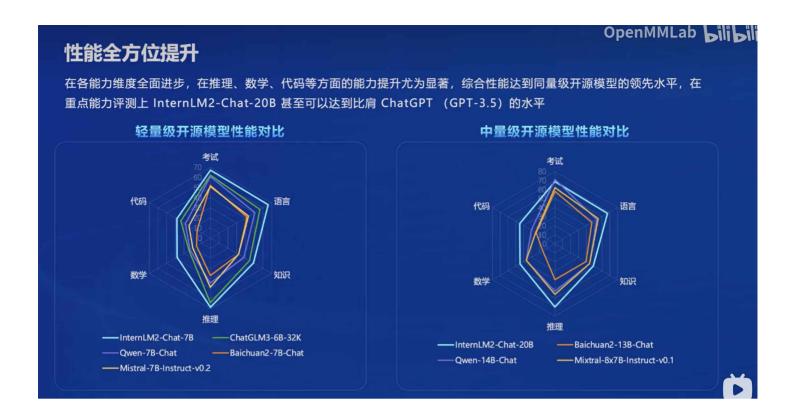
OpenMMLab bilib

突出的数理能力和 实用的数据分析功能

模型在 20万 token上下文中,几乎完美实现"大海捞针"

推理、数学、代码提升显 著InternLM2-Chat-20B 在重点评测上比 肩 ChatGPT 精准指令跟随,丰富的结 构化创作,在 AlpacaEval2 超越 GPT-3.5 和 Gemini Pro 可靠支持工具多轮调用, 复杂智能体搭建 强大的内生计算能力,加入代码解释后,在GSM8K和MATH达到和GPT-4相仿水平





internLM2-chat-20B有GPT3.5的能力,在多智能体链路中,GPT3.5的多智能体组合可以超过gpt4的能力,所以可以尝试一个internLM2-chat-20B组合一个Agent工作流程的水平如何。

1) Reflection: 让 Agent 审视和修正自己生成的输出;

2) Tool Use: LLM 生成代码、调用 API 等进行实际操作;

3) Planning: 让 Agent 分解复杂任务并按计划执行;

4) Multiagent Collaboration:多个 Agent 扮演不同角色合作完成任务;

单纯的比较gpt3.5与gpt4的编码能力,GPT-4 做得更好,正确率达到了 67.7%,但如果你围绕 GPT-3.5 使用一个 Agent 工作流程,实际上它的表现甚至比 GPT-4 还要好。如果你将这种类型的工作流程应用于 GPT-4,它也表现得非常好。你会注意到,GPT-3.5 与一个 Agent 性工作流程相结合实际上超过了 GPT-4 的表现。





以上的思路表明,如果智能体在复杂的场景表现的不够好,要么是RAG或者是微调,但是我们如上面 的笔记,可以尝试构建一个智能体的工作流,也许可以更好的提高多个智能体组成的系统整体的效 果,比微调更好,或者微调后建立这样的工作流效果又能够得到一个提升。







如何保证增量续训不会遗忘之前的知识?这里有怎样的机制?







全链条开源开放体系 | 智能体

多模态智能体工具箱 AgentLego

- 丰富的工具集合,尤其是提供了大量视觉、多模态相关领域的前沿算法功能
- 支持多个主流智能体系统,如 LangChain,Transformers Agent,lagent 等
- 灵活的多模态工具调用接口,可以轻松支持各类输入输出格式的工具函数
- 一键式远程工具部署,轻松使用和调试大模型智能体



