

OpenAI Model Spec 解读

📅 2025年2月13日 ⌚ 1 分钟阅读

#OpenAI

#Model Spec

#AI

OpenAI Model Spec 解读

OpenAI Model Spec 解读

概述

OpenAI Model Spec 是 OpenAI 提供的一个用于描述和评估 AI 模型的规范。它定义了模型的结构、功能、性能和安全等方面的标准。

这篇文章详细介绍了OpenAI的模型规范（Model Spec），包括其目标、原则、风险管理以及具体行为指导，以确保AI模型在开发和使用过程中既有用又安全，同时符合用户和开发者的需求。文章还探讨了模型在不同情境下的行为、需要遵循的指令层级以及如何在复杂场景中平衡冲突目标。

关键点

模型规范旨在塑造OpenAI模型的预期行为，确保模型有用、安全，并符合用户和开发者需求，同时推进人工智能造福全人类的使命。

模型规范的目标包括迭代部署支持开发者和用户的模型、防止模型造成严重伤害、保护OpenAI免受法律和声誉损害。

模型通过遵循明确的命令链来平衡冲突目标。

模型规范与使用政策和安全协议共同构成负责任构建和部署AI的更广泛战略。

模型规范提高了透明度，邀请公众讨论如何改进模型行为，并根据用户反馈不断更新。

在塑造模型行为时，遵循最大限度为用户提供帮助和自由、将伤害降至最低、选择合理默认值等原则。

目录

文章信息

字数

阅读时间

发布时间

更新时间

标签

#OpenAI

#Model Spec

#AI

模型需要应对三大类风险：目标不一致、执行错误、有害指令，并通过提问或遵循命令链减轻这些风险。

模型的指令层级包括平台、开发者、用户和准则，每个层级具有不同的权威性。

模型在处理复杂或冲突场景时，会根据指令的权威级别优先处理高权威指令。

模型必须尊重创造者及其知识产权，不生成违反法律或道德的内容，并保护个人隐私。

模型在模糊或不明确的请求中应假设最佳意图，通过澄清问题或提供假设来帮助用户。

模型在表达不确定性时应使用自然语言，并在高风险场景中更加谨慎。

模型应避免奉承或迎合用户，提供客观、建设性的反馈。

模型在互动中应表现得亲切、富有同理心，同时保持专业性。

模型在特定场景下应展现创造力，例如头脑风暴、娱乐或艺术协作。

模型应根据用户的需求调整响应长度和结构，在需要提供详细答案，但避免冗长或重复。

模型在语音和视频交流中应表现自然，尊重用户的文化背景，并根据请求调整语音或语调。

参考

[OpenAI Model Spec](#)

分享这篇文章



相关文章推荐

计算机使用代理

计算机使用代理

字节跳动 OmniHu...

字节跳动开源的
OmniHuman-1...

DeepSeek R1 Paper...

A comprehensive
review of t...