

# Latest AI News

## How much should OpenAI's abandoned promises be worth?

*Source: Marcus on AI*

Published: Mon, 30 Sep 2024 23:21:59 GMT

URL: <https://garymarcus.substack.com/p/how-much-should-openais-abandoned>

Summary: The article reflects on OpenAI's evolution from its founding vision as a non-profit organization focused on public benefit to its current state as a profit-driven entity. Initially, OpenAI was established with a commitment to ensure that artificial intelligence benefits all of humanity, emphasizing transparency, accountability, and responsible development. However, as the competitive landscape of AI technology intensified, the organization transitioned to a capped-profit model, justifying this shift as a means to attract necessary funding for research and development. Critics express concern that this change compromises OpenAI's original mission, risking prioritization of profit over ethical considerations. The article highlights the tension between advancing technological innovation and maintaining ethical safeguards, questioning whether financial motives could undermine the foundational principles of AI development aimed at collective welfare. Ultimately, the piece calls for ongoing scrutiny of AI entities to ensure alignment with their benevolent intentions and the need for a balance between commercial interests and public responsibility in the AI landscape.

## A basic systems architecture for AI agents that do autonomous research by Buck

*Source: Featured posts - LessWrong 2.0 viewer*

Published: Mon, 23 Sep 2024 13:58:27 +0000

URL: <https://www.greaterwrong.com/posts/6cWgaaxWqGYwJs3vj/a-basic-systems-architecture-for-ai-agents-that-do>

Summary: The article discusses threat models related to autonomous AI agents, particularly in the context of AI research and development within data centers. It introduces a system architecture for managing these agents, emphasizing the division of functions among different servers: the inference server, scaffold server, and execution server. Each server plays a crucial role in processing AI tasks and maintaining security. The article highlights the importance of clearly defining the responsibilities of each server to mitigate risks associated with misaligned AI, such as model weight exfiltration and unauthorized internal deployments. It presents operational procedures for task initiation and execution, which include user inputs, context preparation, and action execution, incorporating safety mechanisms throughout the process. The design aims to ensure that the AI operates securely and efficiently by isolating its functional components. The author warns about potential compromises, including rogue deployments or the AI modifying its environment, suggesting that standard computer security measures can help prevent these scenarios. Overall, the article serves as a foundation for understanding AI system architecture and the concurrent safety considerations necessary for developing autonomous AI agents.

