

WHY CUSTOMIZATION IS KEY TO SOLVING GENERATIVE ARTIFICIAL INTELLIGENCE HALLUCINATIONS

The Future of Human Interaction with Technology

Generative AI is a fast-evolving artificial intelligence technology capable of producing a variety of new content, including code, images, music, text, and more. Generative AI models, whether in the form of AI-generated search results or chatbots that converse fluently with users, will continue to alter how we interact with technology.

However, as those systems become more pervasive, their limitations are increasingly obvious, particularly when they are called upon to answer queries in specialized domains. Examples can be found across generative AI models, including errors in Google's AI Overviews as well as inaccuracies encountered by businesses using ChatGPT.

The Inherent Flaws of Generalized AI Models

Google's AI Overviews, a feature designed to offer concise summaries of search results, has been a prime example of how even cutting-edge AI can falter. Users have reported baffling outputs, such as advice to add glue to pizza recipes or inaccuracies regarding historical figures like President Andrew Johnson.

These errors reveal a fundamental issue with large language models: they are optimized to predict the next word in a sequence, not to validate the accuracy of the information they provide. Despite using techniques like Retrieval Augmented Generation ("RAG"), which allows the AI to consult external sources of specialized data, mistakes still occur when the system misinterprets or combines data from conflicting sources. The same pattern emerges in other fields.

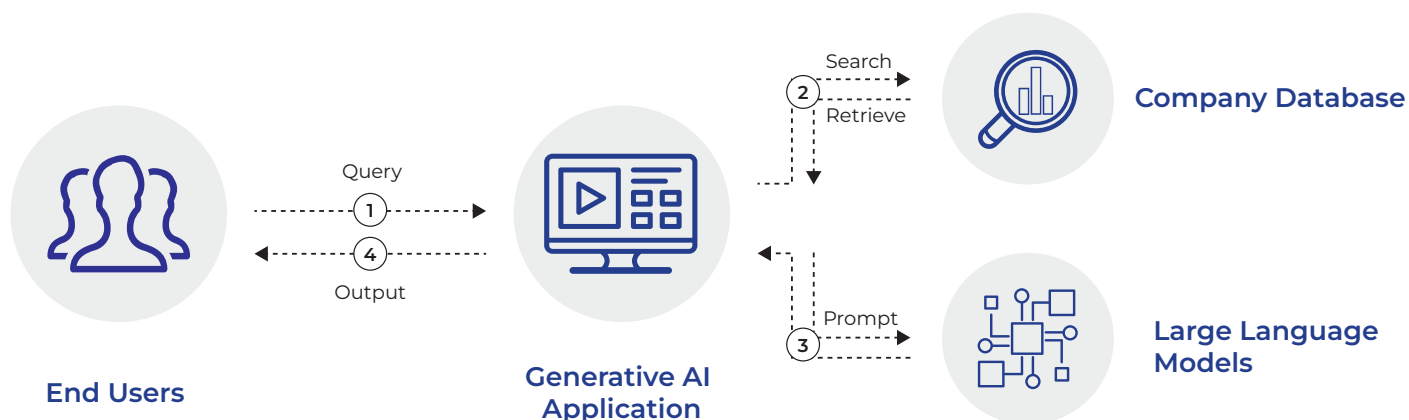
The Wall Street Journal describes how organizations, from the PGA Tour to agricultural firms, have encountered similar issues.¹ While AI models may be trained on vast amounts of general data, they often fail when faced with domain-specific questions. In one notable example, ChatGPT confused Tiger Woods' 15 major wins with his 82 PGA Tour victories—a significant error for anyone familiar with golf.

The Role of Retrieval-Augmented Generation

Generative AI users across sectors have turned to RAG to mitigate such errors. The RAG technique enhances generative AI models by allowing them to pull in authoritative information from specialized databases or documents before generating responses. In theory, this should make the models more accurate, yet, as these examples show, RAG can still fail.

In the case of Google's glue-laden pizza recipe, the AI retrieved data from a joke post that was deemed relevant by the system but clearly was not correct. *The Wall Street Journal* describes how RAG is being used by various industries, including agriculture and finance, to improve AI performance. However, while RAG can increase accuracy by 20%-40%, it is far from perfect. The approach requires high-quality, domain-specific data, and even then, it is not immune to generating errors.

Retrieval-Augmented Generation Process



Customization Is Critical

Given these challenges, businesses are increasingly realizing that off-the-shelf AI models are not enough. To achieve higher accuracy, many companies are fine-tuning or custom-building models with their own proprietary data. For instance, the PGA Tour has begun incorporating a 190-page rulebook into its AI systems to ensure that the model understands the nuances of the sport. This approach allows AI to move beyond generalization and deliver more contextually relevant responses.

However, fine-tuning is not a simple fix. It requires significant investment, in terms of both financial resources and specialized talent. Even then, these models may still fall short of complete accuracy, particularly in high-stakes fields like agriculture or legal services, where a small mistake could have significant repercussions.

Conclusion

The key to the future lies in striking a balance between generalization and specificity. While generalized models like ChatGPT work well for general conversations and information gathering, they are still inadequate for serious, domain-specific questions. As businesses and technology firms figure out these limitations, several paths forward seem to be emerging: using RAG for higher levels of accuracy, fine-tuning existing models with more specialized data, and, in some cases, even building custom models from scratch. With these, the price of customization can be high, but for industries in need of precision and reliability, it might be the only way to make sure AI delivers on its promise. As AI continues to improve, so too will demand increase for more tailored and precise systems, furthering the bounds of what these models can do.

1. Ryan Knutson, "AI Doesn't Know Much About Golf, or Farming, or Mortgages," Wall Street Journal, Dec. 14, 2024.

SOLOMON PARTNERS TECHNOLOGY

Craig Muir
Head of Software,
Data & Analytics

Jeff Derman
Partner

James Butcher
Managing Director

Solange Velazquez
Managing Director

Joe Watson
Managing Director

Jonathan Berger
Director

FINANCIAL SPONSORS

Sash Rental
Head of Financial
Sponsors

Tucker Laurens
Managing Director

This document is for marketing purposes only. It has been prepared by personnel of Solomon Partners Securities LLC ("Solomon Partners") or its affiliates and not by Natixis' research department. It is not investment research or a research recommendation and is not intended to constitute a sufficient basis upon which to make an investment decision. This material is provided for information purposes, is intended for your use only and does not constitute an invitation or offer to subscribe for or purchase any of the products or services mentioned. Nothing in this document constitutes investment, legal, accounting or tax advice, or a representation that any investment or strategy is suitable or appropriate for you. Solomon Partners is a FINRA-registered broker-dealer.