# The Future of AI, LLMs, and Observability on Google Cloud: 7 Key Insights for Leaders

The most actionable takeaways from our fireside discussion with Google's Director of AI, Dr. Ali Arsanjani, and Datadog's VP of Engineering, Sajid Mehmood.

**Dr. Ali Arsanjani**     Director of AI, Google

**Sajid Mehmood**     VP of Engineering, Datadog

# How can organizations better approach AI and LLMs?

We sat down with Google's Director of AI, Dr. Ali Arsanjani, and Datadog's VP of Engineering, Sajid Mehmood, to discuss the current and future states of AI, ML, and LLMs on Google Cloud.

This guide distills the top 7 insights and actions for technical leaders, covering everything from upskilling teams to observability best practices.

**Dr. Ali Arsanjani**
Director of AI
Google

**Sajid Mehmood**
VP of Engineering
Datadog

## 1

### As the adoption of LLMs grows and matures, observability becomes ever more critical

Advancements in large language models (LLMs) have radically expanded the types of teams working in the generative AI (GenAI) space over the past eighteen months. Greater accessibility is welcome, but also creates challenges.

While data science teams are very familiar with the paradigm of continuous evaluation and experimentation, traditional software engineers might not be.

"The strength of an organization's observability tooling thus becomes that much more important," comments Datadog's VP of Engineering, Sajid Mehmood.

If those less experienced with machine learning workflows now perform prompt engineering and RAG (retrieval augmented generation), it's crucial they have quick and easy insight into what's happening "behind the scenes" in their models.

Without access to the right data, teams cannot identify where, why, or how things can improve, and experiments can quickly devolve into unproductive time sinks.

The importance of traceability is also stressed by Google's Director of AI, Dr. Ali Arsanjani: "One of the trends we're seeing is the emphasis on observability at every stage of the life cycle."

An organization might begin with some simple prompt engineering, and need to manage its experimentation with various language models. It might then proceed with RAG, before tuning its models to a specific domain of knowledge, and grounding outputs with searches.

"With Google and Datadog, an organization can have full observability as it goes across these different activities," Dr. Arsanjani says. "As you move into higher levels of maturity and production it's imperative you have very strong capabilities around observability, evaluation, and traceability."

**Action #1**: *Leaders should ensure the teams working on GenAI and LLMs are familiar with the classic machine learning paradigms of continuous evaluation and experimentation, and provide those teams with the right observability tooling, so they can trace models effectively across every level of maturity and production.*

# 2

## Organizations must master simpler use cases before attempting multi-agent systems

Many organizations share a common aspiration: to create complex multi-agent systems, whereby large language models delegate tasks to smaller models to more efficiently orchestrate and complete the work at hand.

However, while some organizations are ready to build such enterprise-scale, multi-agent systems, those near the beginning of their journey with LLMs should master simpler use cases first.

"Organizations seeking to implement production-grade LLM-based applications," Dr. Arsanjani says, "need to look at where they are in terms of their current capabilities — including the human resources and skills experience they have in this area, the platforms they're using, their architecture, and so on."

Shooting straight for a multi-agent system without mastering the capabilities underlying them — like prompt engineering and RAG — will likely lead to failed projects, setbacks, and delays in deployment.

"It's very important that organizations recognize the range of skills, experiences, and tools needed," Dr Arsanjani says, "and aim for a level of maturity that is commensurate with their existing capabilities."

While the LLM landscape is changing quickly — and methods like RAG might change with it — developing "evergreen" skills in evaluation and experimentation remains essential.

"The Gemini team at Google has really pushed the boundary on context windows, for example," Mehmood says. "You start to wonder how that will impact RAG, which is essentially about making efficient use of that context window. So, as context windows get larger, will RAG still be relevant in the years to come?"

While the specific methods might evolve, Mehmood adds, "I think it's relatively safe to say the piece that's not going anywhere is the evaluation side, which has been a critical part of machine learning for decades."

**Action #2**: *To avoid delays in production and deployment, leaders should have a complete understanding of their organization's existing LLM capabilities, set a realistic target level of maturity, and pave a clear path to achieving it.*

# 3

## Building customer confidence in model output is crucial — and there are some simple ways to do it

One of the biggest challenges in the AI space is building confidence in the output of LLMs and LLM-based applications.

After all, if customers don't trust the data enough to make a decision based on it, then it isn't useful — regardless of how insightful, accurate, or elegant it might actually be.

"When building Bits AI here at Datadog," Mehmood shares, "a core design element for us was how to give customers a sense of how Bits came to the conclusions that it did.

"This might mean fetching some relevant telemetry, accessing certain logs within Datadog, getting traces — and presenting that data along with the corresponding conclusion."

Confidence, then, is not simply a result of improved model quality; UX also plays a crucial role. If customers can quickly and easily see model output is legitimately sourced, they are more likely to trust its reliability.

Think of the kinds of evidence you'd ask for from a human interlocutor. If a colleague warns you about a problem, you'd ask for more information as you consider ways to solve it. The goal for LLMs and LLM-based applications is to mirror this kind of human interaction.

Techniques like factual grounding can help here. While the Gemini team at Google is minimizing the risk of hallucinations at a fundamental model level, post-hoc 'searches' of trusted databases ensure model output always cites legitimate sources.

**Action #3**: *As well as considering the machine learning optimizations that can be made to improve the reliability of model output, leaders shouldn't neglect the seemingly small but impactful UX signals that play a vital role in building customer trust — like clearly citing legitimate sources, and allowing customers to "go deeper" by entering a dialogue as they would with a colleague.*

# 4

## The most powerful use cases for LLMs are often the simplest — summarization drives huge productivity gains

Popular growth areas for LLM-based applications include content generation, HR agents, sales agents, and multi-agent systems.

But leaders shouldn't lose sight of how simple use cases like summarization can offer easy, significant, and oft-overlooked improvements to productivity.

"Many of our customers at Datadog have seen enormous, immediate value through using summarization," Mehmood confirms.

Summarization is a particularly powerful use case for two key reasons.

Firstly, it's an easy way to build trust in LLMs, because evaluation is straightforward: anyone can scan the original information to see if it's been captured effectively.

And secondly, there are so many areas in an organization where efficient summarization can instantly provide value.

"One of the places we've rolled this out," Mehmood shares, "is incident summarization. If an incident occurs, as soon as someone joins the corresponding chat channel or video link set up to solve it, they are greeted with an automated message summarizing exactly what's happened and when, structured according to customer impact, mitigation steps, remediation steps, and so on.

"We've found this very useful internally at Datadog and so have many of our customers — plus, it's relatively easy to set up and evaluate."

Multimodality is also important to consider here. Built from the ground up to be a multimodal model, the latest versions of Gemini can effectively parse data regardless of its medium — so summaries, for instance, can be drawn from text, image, video, and audio data.

**Action #4**: *As well as keeping your eyes on how LLMs can help solve complex problems, technical leaders should consider the productivity gains that might be hiding in plain sight — like summarizing information for efficient cross-team communication.*

# 5

## As an organization's approach to LLMs evolves, its tooling should evolve with it

It doesn't take much to start reaping productivity gains from LLMs, but it's important for organizations to consider how their usage and maturity can evolve without disruption.

"At the beginning, you can get away with prompting and just setting the model with your favorite LLM," Dr. Arsanjani says.

As organizations begin developing production-oriented systems, however, "they need to have controls and human evaluation at each stage of the game."

"In Google Cloud's AI Suite, we are building tools that allow you to go across this journey of maturity as seamlessly as possible. From building simple consumer applications to getting more sophisticated with Vector databases for advanced RAG, grounding, and tuning models with methods like reinforcement learning: we cover the full spectrum of capability in a single platform. This allows our customers and partners to go in and focus on building, knowing that the tooling and dashboards are there as and when they need them."

**Action #5**: *To avoid disruption, leaders should consider how they can ensure teams always have access to the right technologies and evaluative tools as their maturity with LLMs grows.*

# 6

## Optimally balancing compute time, compute cost, and model quality requires cutting-edge observability

Having the right tools to evaluate and improve the quality of models is critical, but quality is only one factor that organizations should consider when building AI capabilities. Cost and performance are equally important considerations and require their own tooling.

"As we've embarked upon our own fine tuning and custom model training here at Datadog," Mehmood notes, "one thing I've looked at very closely is Datadog's graphs of cluster utilization. Are we using GPUs to the best of our ability?

"We've put a lot of effort into being able to monitor not just the traditional things like CPU, memory, and so on, but now dedicated GPUs and even learning chips like TPUs, and we have a number of integrations that work very effectively with Google Cloud's offerings."

It's not just individual GPUs that teams need to monitor, however. As organizations scale andI mature their AI capabilities, they quickly find themselves needing to handle multi-node GPU clusters — and network performance monitoring then becomes imperative.

Datadog completes its full-stack AI observability solution with LLM Observability, which supports quality and safety evaluations and end-to-end tracing of LLM chains. Tracking all of these considerations in a single platform gives organizations a holistic view of the health of their environments and helps them balance trade-offs more effectively. As Mehmood says, "Datadog's observability platform provides you with everything you need to optimize that space between compute time, compute cost, and ultimate model quality."

**Action #6**: *To optimally balance compute time, compute cost, and overall model quality, leaders must ensure their organizations have modern observability capabilities to effectively monitor end-to-end performance and evaluate model quality.*

# 7

## With change happening quickly, organizations must stay focused on the metrics that matter

Organizations everywhere are still figuring out the best way to adopt and roll out the new technologies, tools, and strategies in the LLM space.

While advancements are exciting, it's crucial organizations put the necessary guardrails in place to ensure these novel approaches upgrade productivity, rather than drain resources or present new risks.

Beyond upholding security measures like access control and masking sensitive data, organizations should have a very precise understanding of the metrics that actually matter to their business.

"Better tooling can always help," Mehmood says, "but a lot of progress comes down to really nailing down what matters for your business and using that to reduce the noise.

"What elements most effect server costs or training costs? What are the aspects of model quality you really care about? Do things like CPU usage spikes matter to you, if they don't result in end users facing degradation?"

Establishing clear service level objectives (SLOs) is imperative here — and Datadog provides organizations with easy ways to define and track them.

As Mehmood advises: "the path forward emerges when you have a clear sense of what really matters to your business."

**Action #7**: *To avoid distractions and unnecessary cost, leaders should develop and communicate a deep and precise understanding of the metrics that really make a difference to their organization.*

**NEXT STEPS**

# Enhance and safeguard every stage of your organization's Google Cloud AI journey with Datadog's modern monitoring platform

**WEBINAR**

Watch the full webinar with Google's Director of AI, Dr. Ali Arsanjani, and Datadog's VP of Engineering, Sajid Mehmood
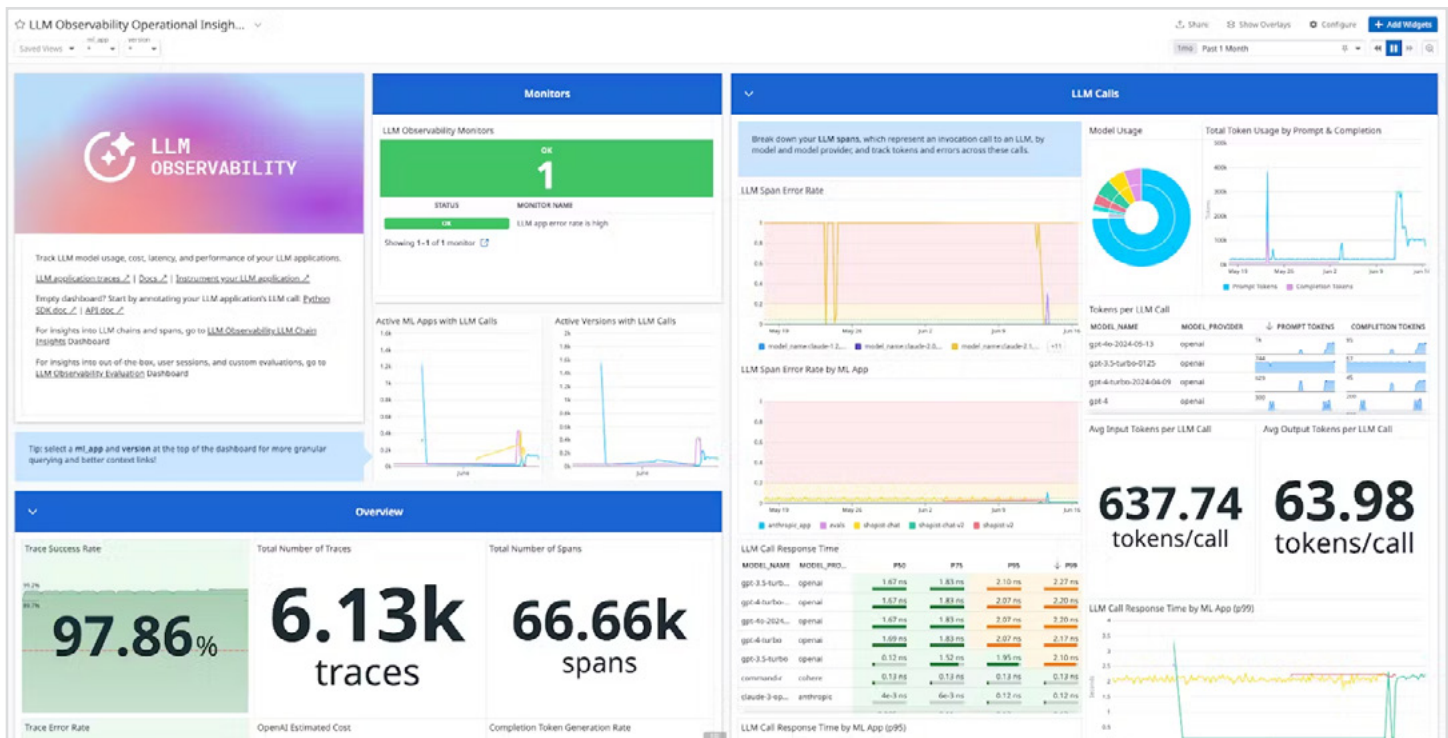
**WATCH WEBINAR**

**DEMO**

See how Datadog's monitoring and security platform can help identify bugs, enhance model quality, and drive costs down at every stage of your organization's AI journey on Google Cloud

**REQUEST A DEMO**

# About Datadog

Datadog is the monitoring and security platform for cloud applications. Our SaaS platform integrates and automates infrastructure monitoring, application performance monitoring and log management to provide unified, real-time observability of our customers' entire technology stack. Datadog is used by organizations of all sizes and across a wide range of industries to enable digital transformation and cloud migration, drive collaboration among development, operations, security and business teams, accelerate time to market for applications, reduce time to problem resolution, secure applications and infrastructure, understand user behavior and track key business metrics.

For more information, visit [datadoghq.com](datadoghq.com)

DATADOG