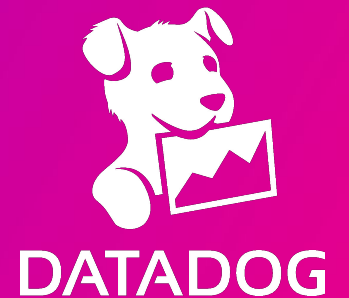January 2025

# Monitoring Google Gemini Models with Datadog

Google Cloud

DATADOG

# Jason Hand

Senior Developer Advocate - Datadog

# Merlin Yamsi

Lead Solutions Consultant - AI/ML CoE Partner Engineering - Google Cloud

Google Cloud

DATADOG

# The Evolution of AI Systems



AGI

Data for AI

Predictive AI

Generative AI

Agentic AI

Generative UI

Contextual AI

Enterprise AI

Personal AI

Google Cloud

# Enterprise Gen AI apps face a variety of challenges

The need to provide **accurate** and up-to-date information

The need to offer **contextual** user experiences

The need to be **easy** for developers to build and operate



Google Cloud

# What have we learned from our customer success stories?
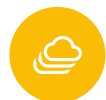
Google Cloud

# The **4 key success factors** for enterprise AI

**Do you have a single, integrated platform that provides your teams optionality and choice?**

**Can you differentiate with your knowledge and data?**

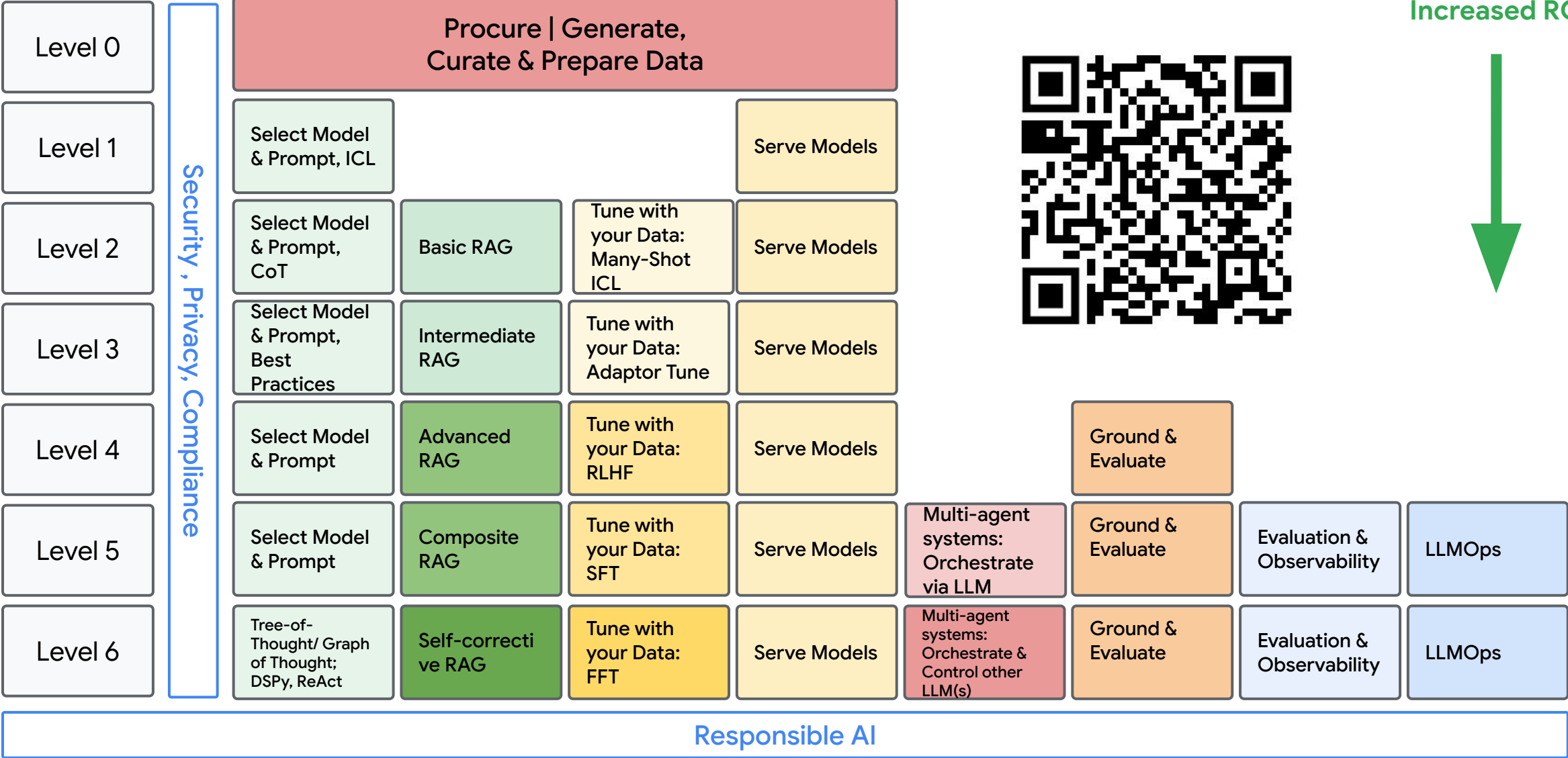**Does your AI platform future proof your AI investment with innovation at every layer?**

**Is your AI enterprise ready so you can go to production with confidence?**

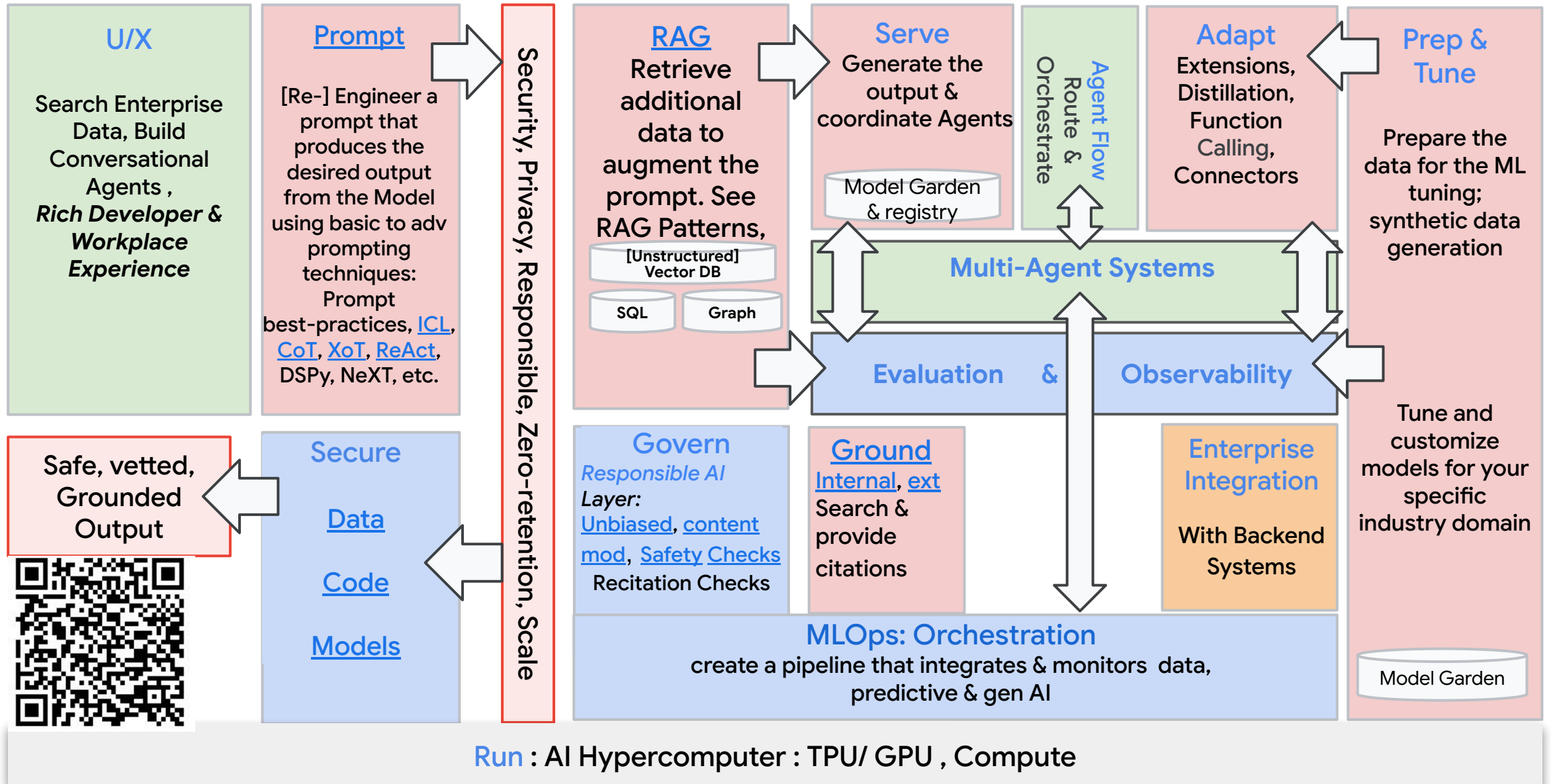Google Cloud

# AI Maturity: Increasing Sophistication of Solutions

Sophistication

Increased ROI

| | Security, Privacy, Compliance | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Level 0 | | Procure \| Generate, Curate & Prepare Data | | | | | | |
| Level 1 | | Select Model & Prompt, ICL | | | Serve Models | | | |
| Level 2 | | Select Model & Prompt, CoT | Basic RAG | Tune with your Data: Many-Shot ICL | Serve Models | | | |
| Level 3 | | Select Model & Prompt, Best Practices | Intermediate RAG | Tune with your Data: Adaptor Tune | Serve Models | | | |
| Level 4 | | Select Model & Prompt | Advanced RAG | Tune with your Data: RLHF | Serve Models | | Ground & Evaluate | |
| Level 5 | | Select Model & Prompt | Composite RAG | Tune with your Data: SFT | Serve Models | Multi-agent systems: Orchestrate via LLM | Ground & Evaluate | Evaluation & Observability |
| Level 6 | | Tree-of-Thought/ Graph of Thought; DSPy, ReAct | Self-correcti ve RAG | Tune with your Data: FFT | Serve Models | Multi-agent systems: Orchestrate & Control other LLM(s) | Ground & Evaluate | Evaluation & Observability |

| | LLMOps |
|---|---|
| Level 5 | LLMOps |
| Level 6 | LLMOps |

Responsible AI

Google Cloud

# GenAI Reference Architecture: Patterns & Technical Blueprint for Building GenAI Solutions

## U/X

Search Enterprise Data, Build Conversational Agents , *Rich Developer & Workplace Experience*

## Prompt

[Re-] Engineer a prompt that produces the desired output from the Model using basic to adv prompting techniques: Prompt best-practices, ICL, CoT, XoT, ReAct, DSPy, NeXT, etc.

**Security, Privacy, Responsible, Zero-retention, Scale**

## RAG
Retrieve additional data to augment the prompt. See RAG Patterns,

[Unstructured] Vector DB

SQL    Graph

## Serve
Generate the output & coordinate Agents

Model Garden & registry

## Agent Flow
Route & Orchestrate

## Adapt
Extensions, Distillation, Function Calling, Connectors

## Prep & Tune

Prepare the data for the ML tuning; synthetic data generation

**Multi-Agent Systems**

**Evaluation    &    Observability**

Tune and customize models for your specific industry domain

## Safe, vetted, Grounded Output

## Secure

Data

Code

Models

## Govern
*Responsible AI Layer:*
Unbiased, content mod, Safety Checks
Recitation Checks

## Ground
Internal, ext
Search & provide citations

## Enterprise Integration

With Backend Systems

**MLOps: Orchestration**
create a pipeline that integrates & monitors  data, predictive & gen AI

Model Garden

**Run** : AI Hypercomputer : TPU/ GPU , Compute

Google Cloud

# AI Pillar Spotlight: Vertex AI is our Generative AI platform

**Applications**

> ## AI Solution
> Contact Center AI | Risk AI | Healthcare Data Engine | Search for Retail, Media and Healthcare

> Gemini for Google Cloud

> Gemini for Google Workspace

> Build your own generative AI-powered agent

**Agents**

> ## Vertex AI Agent Builder
> OOTB and custom Agents | Search
> Orchestration | Extensions | Connectors | Document Processors | Retrieval engines | Rankers | Grounding

**Tooling**

> ## Vertex AI Model Builder
> Prompt | Serve | Tune | Distill | Eval | Notebooks I Training | Feature Store | Pipelines | Monitoring

**Models**

> ## Vertex AI Model Garden
> Google | Open | Partner

> Google Cloud Infrastructure (GPU/TPU) | Google Data Cloud

# **Vertex AI** is AI for your enterprise

An end-to-end platform that unlocks your data for every use case, expertise, or environment

**Vertex AI**

Agent Builder

Model Builder

Model Garden

# A unified platform from data to deployment and
## for all your predictive, generative, and agentic needs

Gemini for Google Workspace     Gemini for Google Cloud     AI Agents

**Databases**
- AlloyDB

**Data analytics**
- BigQuery

**AI & ML**
- Vertex AI
  - Agent Builder
  - Model Builder
  - Model Garden

**Insights**
- Looker

**Governance**
- Dataplex
- Cloud Ops
- MLOps in Vertex AI

**Infrastructure**
- AI Hypercomputer (GPUs and TPUs)

Google Cloud

# Flexibility and curation at every layer of the stack to avoid lock-in

## Data
Single unified access layer for all data: structured, unstructured, streaming

BigQuery

GCS

Omni for Multi-cloud
(AWS S3, Azure Storage)

## Compute
Ultra performant AI hypercomputers for any workload

Google Cloud TPUs

NVIDIA GPUs

## Frameworks
An open & comprehensive AI stack fueling the Gen AI revolution

TensorFlow    PyTorch

JAX    ACCELERATED LINEAR ALGEBRA

## Models
The best foundation models from Google, Partners, and the Open ecosystem in the Model Garden

Gemini    Imagen

MISTRAL AI_

Hugging Face    kaggle    Gemma

## Agents
Comprehensive tools from Google and partners to build and deploy agents.

Vertex AI

LlamaIndex

LangChain

Google Cloud

# 160+ enterprise-ready foundation models in Vertex AI Model Garden

## Vertex AI Model Garden

| Gemini Foundation Models | Gemini 1.0 Pro | Gemini 1.5 Flash | Gemini 1.5 Pro | | |
|---|---|---|---|---|---|
| **Google Foundation Models** | PaLM 2 | Imagen 3 | Chirp | Codey | Embeddings |

**Google Task Specific Models**
- Speech-to-Text
- Text-to-Speech
- Natural Language
- Translation
- Doc AI OCR
- Occupancy analytics
- Vision
- Video Intelligence

**Google Domain Specific Models**
- **MedLM** — Life Science and Healthcare
- **Sec-LM** — Cybersecurity

**Partner & Open Ecosystem**
- Claude 3 and 3.5 Haiku, Sonnet, and Opus
- Llama 3.2
- MISTRAL AI_ — Mistral Large 2, Nemo and Codestral
- AI21labs — Jamba 1.5 Large and Mini
- Hugging Face — kaggle
- Gemma

- **Choice and flexibility** with Google, open source, and third-party foundation models

- **Multiple modalities** to match your use case

- **Multiple model sizes** to match cost and efficacy needs

- **Domain-specific models** for specialized industries

- Enterprise ready with **safety, security, and responsibility**

- Decrease time to value with **fully integrated platform**

Google Cloud

# Continued model improvements to optimize performance and cost

## Gemini 2.0 Flash

Offers **2x** the speed of Gemini 1.5 Pro

Stronger performance: multimodal, text, code, video, spatial understanding, reasoning

Experimental

### Gemini 1.5 Flash

**Fastest and most cost-efficient model yet**

Multimodality

Low Latency

Comparable quality as 1.5 Pro
(on common tasks)

GA
**Now**

### Gemini 1.5 Pro

**Native reasoning over enormous amounts of data**

2M Context Window

Multimodality

Versatile & top-tier quality

GA
**Now**

Gemma 2

9B & 27B

**Now Available**

# New capabilities in Gemini 2.0 Flash

| | Gemini 1.5-002 (GA) | Gemini 2.0 Flash (Experimental) |
|---|---|---|
| **Input modalities** | | Text, image, video, audio, PDF |
| **Output modalities** | Text | Text, image [new!], audio (speech) [new!] |
| **Context window** | 2M (Pro), 1M (Flash) | 1M (Flash) |
| **Image Generation** [new!] | No | **Yes, Private Experimental**<br><br>○ Text -> Image+text<br>○ Watermarking (synthID)<br>○ Multi-images input support<br>○ Image editing (Text/Image → to image)<br>○ Text -> Image + Text (interleaved) |
| **Audio generation (speech)** [new!] | No | **Yes, Public Experimental**<br><br>○ Text to speech: say "hi everyone"<br>○ Context Prompted text to speech: say "hi everyone", in a pirate's voice<br>○ Audio generation: unary + streaming |
| **Multimodal Live API** [new!] | **Yes, Public Experimental**<br><br>○ Text → Voice<br>○ Text+Voice -> Voice<br>○ Voice → Voice<br>○ Voice & Video to Voice<br>○ In-session Memory Q&A (128k) | |
| **Native tool-use** [new!] | No | **Yes, Public Experimental**<br>• Code execution<br>• Search as a tool<br><br>**More tools to come in the future and more integration with Multimodal Live and Studio UI** |

# Continued model improvements to optimize performance and cost
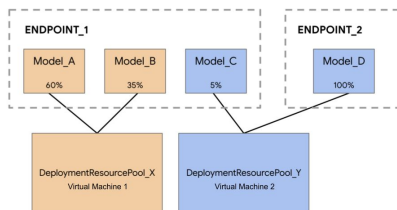
## Open framework support on Vertex AI

### Ray on Vertex AI

Scale AI & Data with Ray

### PyTorch & Saxml

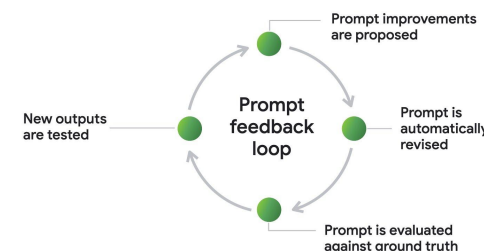Serve models on multi-host TPUs with pre-built Saxml containers and PyTorch

ENDPOINT_1

| Model_A | Model_B | Model_C |
| 60% | 35% | 5% |

ENDPOINT_2

| Model_D |
| 100% |

DeploymentResourcePool_X
Virtual Machine 1

DeploymentResourcePool_Y
Virtual Machine 2

## MLOps

### Vertex AI Feature Store 2.0

brings your data to production

- Built on BigQuery
- Low latency data serving
- Low latency vector search

### Tune the prompt

to continuously improve your prompts

Prompt improvements are proposed

Prompt is automatically revised

New outputs are tested

**Prompt feedback loop**

Prompt is evaluated against ground truth

Proactively monitor model performance with

### Model Monitoring 2.0

- Monitor and alert for model performance
- Diagnose deviations
- Trigger model updates and re-training pipelines

### Gen AI Eval Service

**Rapid evaluation** lets developers evaluate model performance in seconds based on a small data set

**Auto SxS** can assess the performance of two different models using a large language model, and provides explanations and certainty scores
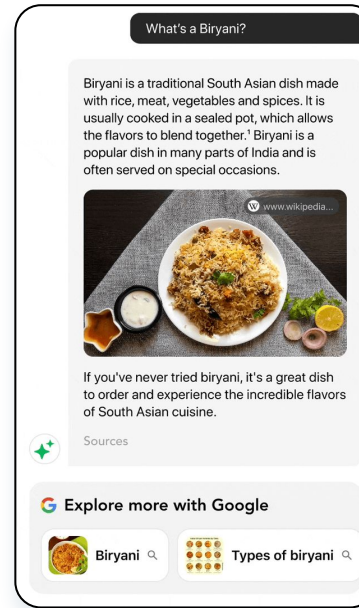
# Google Cloud AI differentiators

**Multimodal AI reasoning**

**Google-quality search with advanced grounding**

**Integrated AI platform with optionality and choice**





| Gemini for Workspace | Gemini for Google Cloud | AI Agents |
| --- | --- | --- |
| AI Platform & Tools | | |
| AI Models | | |
| AI Hypercomputer | | |

**Google AI** is designed to reason seamlessly across text, images, video, audio, and code

**Search** for information from verifiable sources within your own data or Google Search

**Unified platform** for all your predictive, generative, and agentic needs

# LLMOps represents a constellation of technologies that "wrap around" LLMs to deliver enterprise-grade performance, experience, and management capabilities

## LLMOps Capability Map

| Prepare | Develop | Validate | Prompt | Deploy | Infer | Automate | Monitor |
|---------|---------|----------|--------|--------|-------|----------|---------|
| Data Collection | Model Selection | Benchmarking | Prompt Deconstruction | Model Hosting (Inference / Serving) | RLHF Tooling | Agent Design | Logging & Analytics |
| Data Preprocessing (e.g., Chunking) | Model Pre-training | Performance Evaluation | Prompt Libraries & Templates | Model Caching | Prompt Reconstruction | Connector Tooling (Tool Aggregation) | Error & Usage Analysis |
| Data Retrieval (incl. RAG tooling) | Model Fine-Tuning | Model Resilience Testing | Prompt Chaining | Model Orchestration | Infrastructure Provisioning | LLM Chaining | App / Model Debugging |
| Data Labeling & Annotation | Hyperparameter Tuning | Model Efficiency Tracking | Prompt Embedding & Context Aug. (RAG) | Distributed Computing | Human-in-the-Loop Tooling | Agent Memory Management | Performance Monitoring |
| Data Versioning & Auditing | Model Hub (Registry) & Version Control | Experiment Tracking | Prompt Suggestions | | API & Service Integrations | Agent Self-Eval Tooling | Output & Drift Monitoring |
| | Model Distillation & Quantization | Model Explainability | Prompt A/B Testing (Comparison, Merge) | | Load Balancing | Agent Orchestration (Multi-Agent System) | |
| | | Grounding | | | Autoscaling | Real-time Agent Debugging | |

| Safeguard | | | | | | | |
|-----------|---|---|---|---|---|---|---|
| Security | Compliance | Privacy | Bias Mitigation | Transparency | Guardrails | Sustainability | Recovery |

Google Cloud

# Generative Enterprise Key Challenges

- Explainability
- Reliability and Robustness
- Data Drift
- Ethical Considerations

# How Observability Can Help

- **Data Quality Monitoring:**
  - Observability tools can be used to monitor the quality of the training data used to train generative AI models.
  - This can help identify any biases or errors in the data, allowing data scientists to take corrective actions and improve the quality of the training data.

- **Model Performance Monitoring:**
  - Observability tools can be used to monitor the performance of generative AI models in production.
  - This can help identify any degradation in model performance over time, which may indicate the need for retraining or fine-tuning the model.

- **Drift Detection:**
  - Observability tools can be used to detect drift in the input data or model behavior.
  - This can help identify when the model's predictions are no longer reliable and trigger alerts or notifications to data scientists.

- **Root Cause Analysis:**
  - In the event of a model failure or degradation in performance, observability tools can help identify the root cause of the problem.
  - This can be achieved by tracing the model's inputs and outputs, identifying any anomalous behavior or errors.

Google Cloud

# Observability Metrics

Vertex AI exposes a wide range of observability metrics that can be used to monitor the health and performance of your models, training jobs and system.

These metrics include:

- **Model metrics:** (Monitor the performance of your models)
  - Measure the performance of your model on a given dataset. They include metrics such as accuracy, precision, recall, and F1 score.
  - You can use model metrics to track the performance of your models over time and identify areas where they can be improved

- **Training metrics:** (Troubleshoot training jobs)
  - Measure the progress of your training job. They include metrics such as loss, accuracy, and training time.
  - You can used to troubleshoot training jobs that are not performing as expected. For example, you can use the loss metric to identify overfitting or underfitting

- **System metrics:** (Monitor the health of the Vertex AI platform)
  - Measure the health and performance of the Vertex AI platform itself. They include metrics such as CPU utilization, memory usage, and network latency.
  - You can use system metrics to monitor the health of the Vertex AI platform and identify any potential issues. For example, you can use the CPU utilization metric to identify if the platform is experiencing high load

Google Cloud

# Analyzing AI and the Underlying Infrastructure

Figure 1: Magic Quadrant for Observability Platforms

**INDUSTRY RECOGNITION**

# We are named a **Leader** in the *2024 Gartner*® *Magic Quadrant™ for Observability Platforms*

This graphic was published by Gartner, Inc. as part of a larger research document and should be evaluated in the context of the entire document. The Gartner document is available upon request from Datadog.
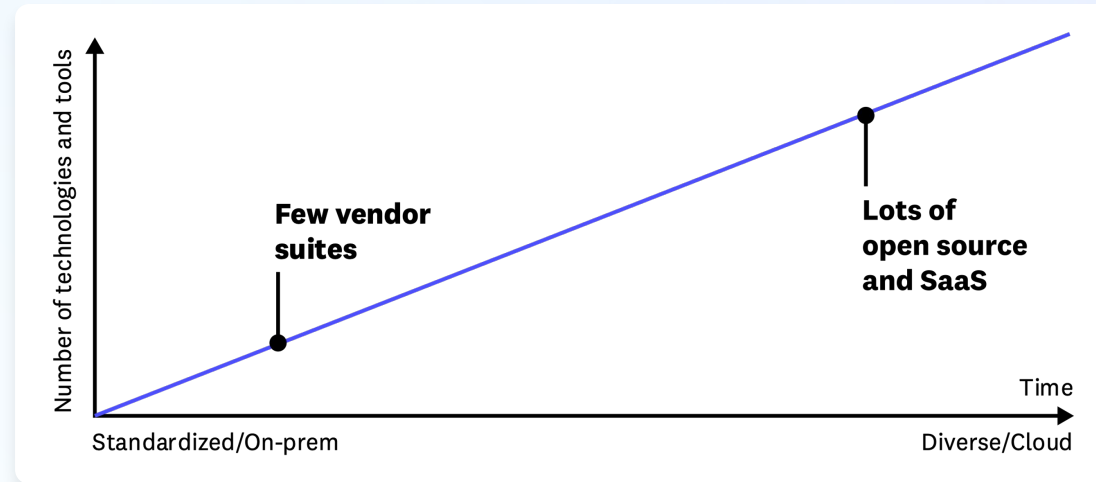
Gartner does not endorse any vendor, product or service depicted in its research publications, and does not advise technology users to select only those vendors with the highest ratings or other designation. Gartner research publications consist of the opinions of Gartner's research organization and should not be construed as statements of fact. Gartner disclaims all warranties, express or implied, with respect to this research, including any warranties of merchantability or fitness for a particular purpose.

Gartner and Magic Quadrant are registered trademarks of Gartner, Inc. and/or its affiliates in the U.S. and internationally and is used herein with permission. All rights reserved.
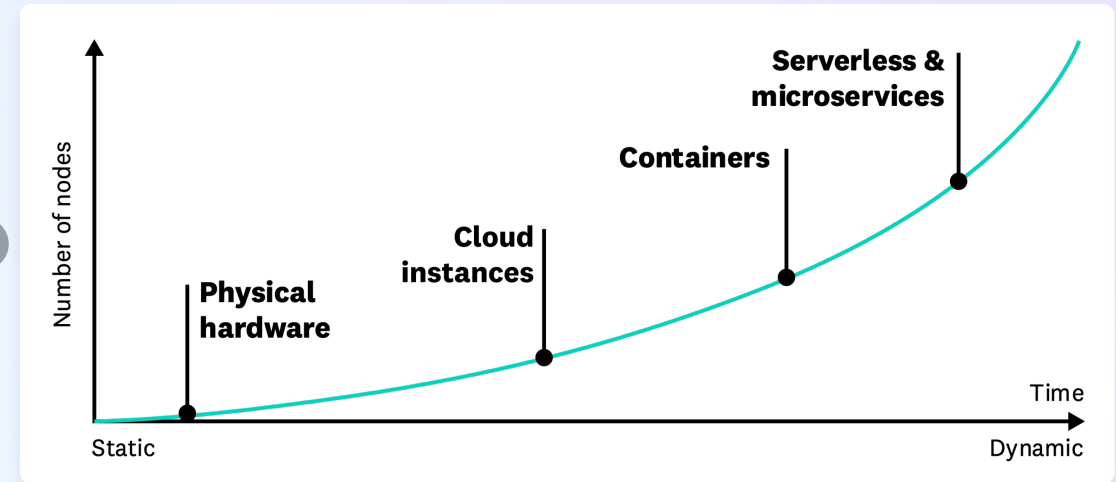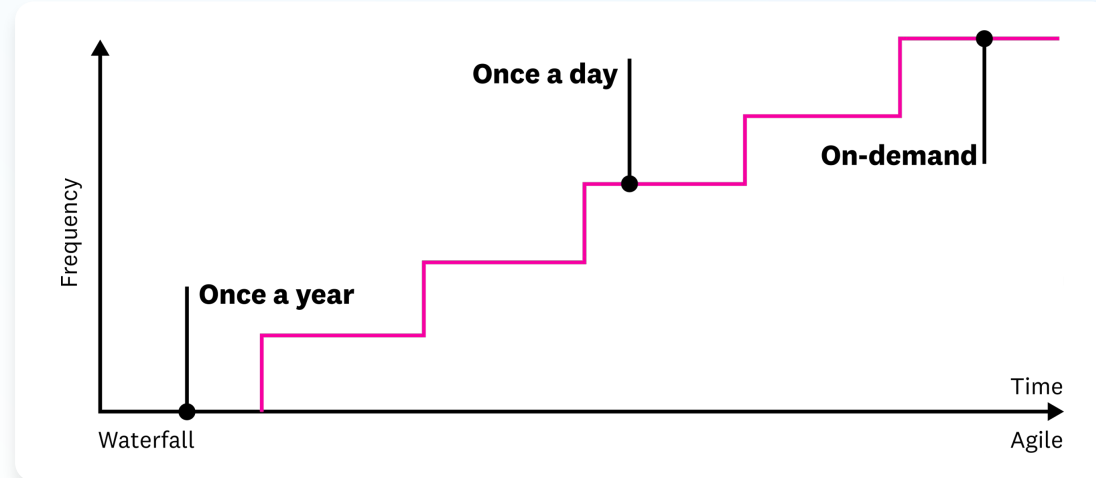
Source: Gartner, Magic Quadrant for Observability Platforms, Gregg Siegfried, Padraig Byrne, Mrudula Bangera, Matt Crossley, 12 August 2024

# The problem: an explosion of complexity

## Diversity of technologies in use



Number of technologies and tools

**Few vendor suites**

**Lots of open source and SaaS**

Time

Standardized/On-prem — Diverse/Cloud

## Scale in number of computing units



Number of nodes

**Physical hardware**

**Cloud instances**

**Containers**

**Serverless & microservices**

Time

Static — Dynamic

## Frequency of release



Frequency

**Once a year**

**Once a day**

**On-demand**

Time

Waterfall — Agile

## Number of people involved



People

**Ops**

**Dev + Ops**

**Business + Dev + Ops**

**Security + Dev + Ops + Business**

Time

Siloed — Integrated

DATADOG

# AI compounds complexity

### Diversity of technologies in use



Number of technologies and tools (y-axis) vs Time (x-axis: Standardized/On-prem → Diverse/Cloud)
- Few vendor suites
- Lots of open source and Saas

### Scale in number of computing units



Number of nodes (y-axis) vs Time (x-axis: Static → Dynamic)
- Physical hardware
- Cloud instances
- Containers
- Serverless & microservices

### Number of large language models



# of models published per year (y-axis) vs Time (x-axis)
- 2018: 1
- 2020: 3
- 2022: 28

### Frequency of releases



Frequency (y-axis) vs Time (x-axis: Waterfall → Agile)
- Once a year
- Once a day
- On-demand

### Number of people involved



People (y-axis) vs Time (x-axis: Static → Integrated)
- Ops
- Dev + Ops
- Business + Dev + Ops
- Security+ Business + Dev + Ops

### Processing power for next-gen AI use cases



Deployed inference-capable GPUs (TFlops) (y-axis) vs Time (x-axis)
- iPhone X
- iPhone 11 Pro
- iPhone 13 Pro
- iPhone 15 Pro

DATADOG

# Gartner AI Maturity Model*

The road to Generative AI becoming a critical function of the business may include setbacks and relearning.



* Source: bit.ly/CIO_AI

# What is Large Language Model Observability?

" **Troubleshooting issues in LLM applications is a time-consuming and resource-intensive task due to the black-box nature of their decision-making processes** "

**What Is LLM Observability & How Does it Work?**

https://www.datadoghq.com/knowledge-center/llm-observability

# LLM Observability

## End-to-End Tracing

View every step of your LLM application chains and calls.

## Operational Metrics

Monitor the throughput, latency, and token usage trends.

## Evaluate Quality

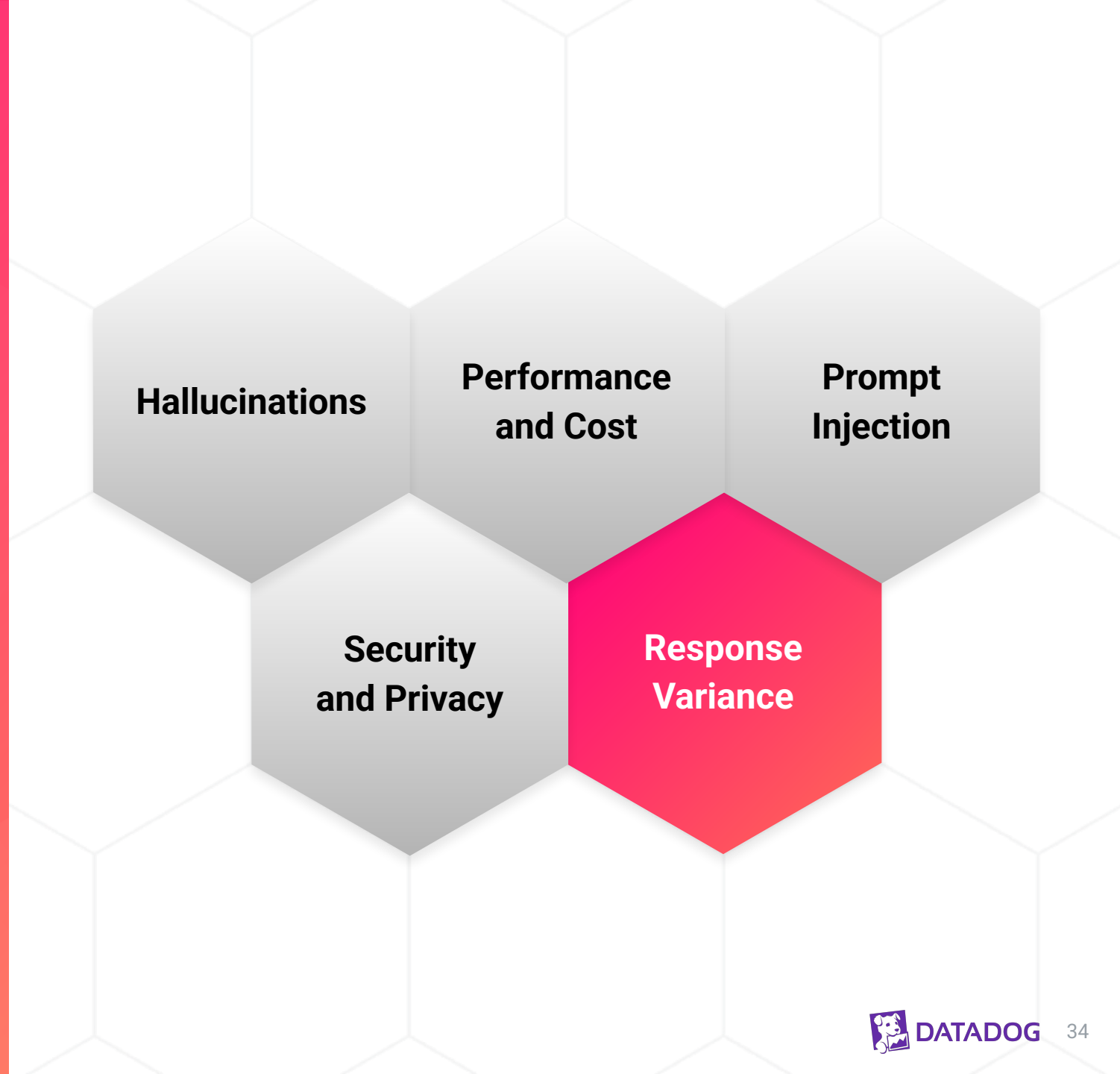Identify problematic clusters and monitor the quality of responses over time.

# Issues

As artificial intelligence and LLM tools are in their infancy, there are a number of issues that can occur, both prompted by users and within the LLM's responses.

DATADOG

# Hallucinations

LLM powered applications may occasionally produce false information, a phenomenon referred to as "hallucinating".

**Hallucinations**

# Performance and Cost

As utilization, data volumes and complexity increase, performance may suffer as a result and costs can start to add up.

Hallucinations

Performance and Cost

# Prompt Injection

Prompt injection is a technique where users can influence LLM applications to produce specific content.

**Hallucinations**

**Performance and Cost**

**Prompt Injection**

# Security and Privacy

As artificial intelligence and LLM tools are in their infancy, there are a number of issues that can occur, both prompted by users and within the LLM's responses.

**Hallucinations**
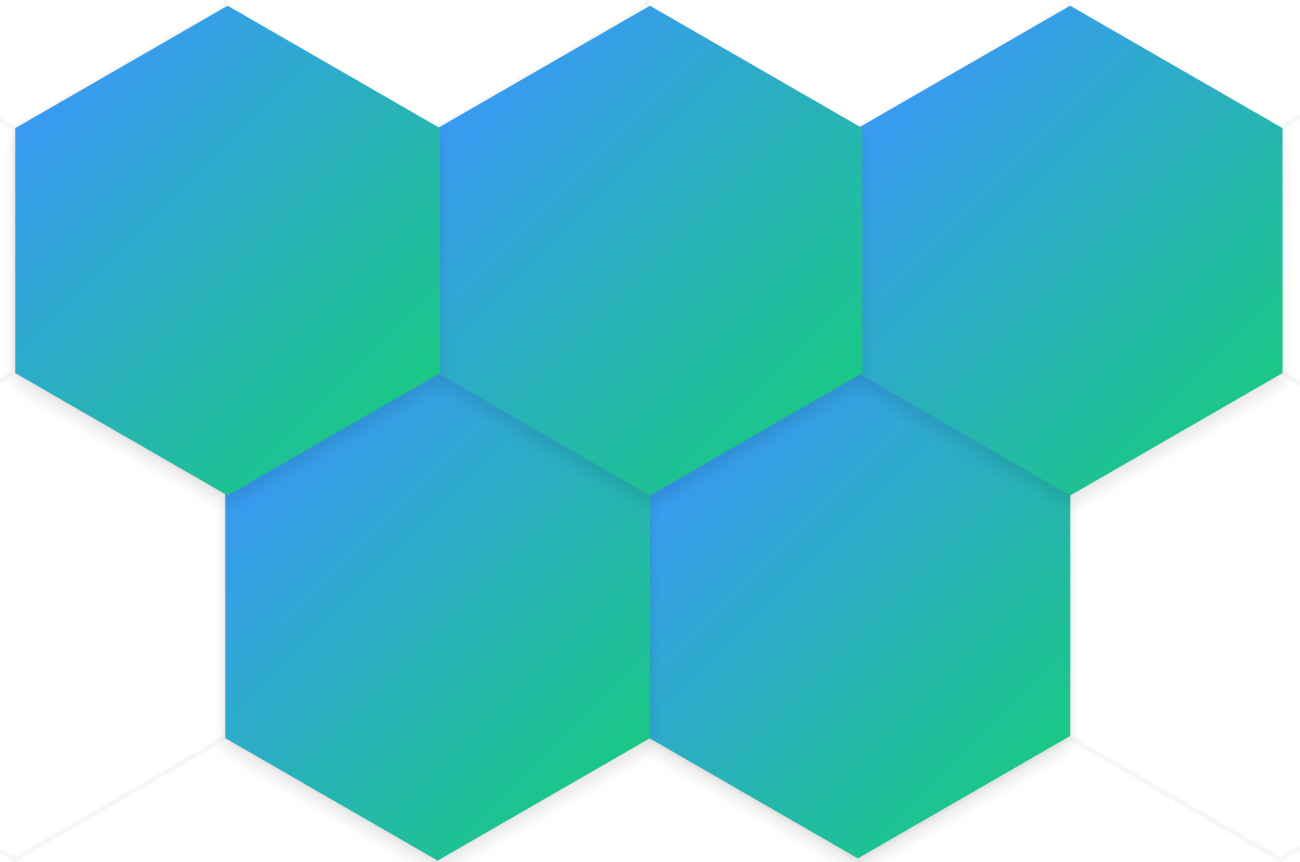
**Performance and Cost**

**Prompt Injection**

**Security and Privacy**

# Response Variance

The user prompts received by LLMs and the responses they generate vary in attributes such as length, language, and accuracy.

**Hallucinations**

**Performance and Cost**

**Prompt Injection**

**Security and Privacy**

**Response Variance**

# Issues

As artificial intelligence and LLM tools are in their infancy, there are a number of issues that can occur, both prompted by users and within the LLM's responses.

**Hallucinations**

**Performance and Cost**

**Prompt Injection**

**Security and Privacy**

**Response Variance**

# Benefits

As LLM tools rapidly evolve, organizations that implement in-depth monitoring of their applications can expect these benefits.

# Improved Performance

LLM observability enables real-time monitoring of various performance evaluation metrics such as latency and throughput of LLM applications and quality of responses.

Improved Performance

# Better Explainability

As LLM tools rapidly evolve, organizations that implement in-depth monitoring of their applications can expect these benefits.

**Improved Performance**

**Better Explainability**

DATADOG

# Faster Diagnosis

LLM observability enables engineers to analyze the backend operations and API calls for a request to pinpoint the root cause of an issue, reducing the time it takes to resolve the issue.

**Improved Performance**

**Better Explainability**

**Faster Diagnosis**

# Increased Security

By tracking access patterns, input data, and model outputs, LLM observability tools can detect anomalies that may indicate data leaks or adversarial attacks.

**Improved Performance**

**Better Explainability**

**Faster Diagnosis**

**Increased Security**

# Cost Management

Observing the resource consumption and utilization of LLM models allows organizations to optimize resource allocation and cost based on actual usage patterns.

Improved Performance

Better Explainability

Faster Diagnosis

Increased Security

Cost Management

# Benefits

As LLM tools rapidly evolve, organizations that implement in-depth monitoring of their applications can expect these benefits.

**Improved Performance**

**Better Explainability**

**Faster Diagnosis**

**Increased Security**

**Cost Management**

# Demo

DATADOG

Monitor your Google Gemini apps with Datadog LLM Observability

Siddarth Dwivedi
Tom Sobolik

Published: January 6, 2025

https://www.datadoghq.com/blog/monitor-google-gemini-datadog-llm-observability/

# Webinar Takeaways

- Google Vertex AI Generative **AI Core features**

- **Google Gemini** model

- **AI Adoption** - The increasing adoption of generative AI models

- **AI Maturity** - Need for robust monitoring solutions with observability capabilities

- **AI Observability** - Datadog, a leading observability platform, integrated with Vertex AI, enabling powerful use cases for monitoring, analyzing, and optimizing the use of generative AI models in production

- Common Issues to monitor LLMs for

- Benefits of instrumenting LLM apps for observability

# The Gemini era for developers and businesses

Gemini's ecosystem of products and models can help developers and businesses get the most out of Google AI, from building with Gemini models to using Gemini as your AI assistant.

Try Gemini 2.0 models—the latest and most advanced multimodal models from Google. See what you can build with up to a 2M token context window.



## BUILD WITH GEMINI MODELS

- Google AI Studio

  Experiment, prototype, and deploy. Google AI Studio is the fast path for developers, students, and researchers who want to try Gemini models and get started building with the Gemini Developer API.

- Vertex AI

  Build AI agents and integrate generative AI into your applications, Google Cloud offers Vertex AI, a single, fully-managed, unified development platform for using Gemini models and other third party models at scale.

## USE GEMINI AS YOUR AI ASSISTANT

- Gemini for Google Cloud

  Your always-on assistant for building or monitoring anything built on Google Cloud, Gemini for Google Cloud helps you code more efficiently, gain deeper data insights, navigate security challenges, and more.

- Gemini for Google Workspace

  Your AI-powered assistant built right into Gmail, Docs, Slides, Sheets, and more, to help boost your productivity and creativity.

Google Cloud

# Are you ready?

Google Cloud Next 25

April 9–11

Select programming begins April 8

# Thank you.

# Resources

dtdg.co/jan_gemini_webinar