# Production-Style GenAI API– Latency-First Engineering (AWS)

## Project Overview

This project evaluates and optimizes **latency behavior in a serverless Generative AI API** built on AWS. The system was designed to identify dominant latency contributors and improve tail latency through **measured, minimal changes** rather than architectural overengineering.

The primary goal was to **reduce p95 and p99 response times** while maintaining functional correctness and lowering cost.

## Architecture

**Request Flow:**

Client → API Gateway → AWS Lambda → Amazon Bedrock → Lambda → Client

**Key AWS Services:**

- Amazon API Gateway (HTTP endpoint)

- AWS Lambda (Python runtime)

- Amazon Bedrock (Foundation Model inference)

- Amazon CloudWatch (Logs, Metrics, Dashboards)

## Problem Statement

Initial testing showed **high tail latency (p95 / p99)** for the GenAI API, even under low request volume.
To improve system performance, the project focused on answering:

- Where is the latency actually coming from?

- Is the bottleneck compute, networking, or inference?

- Can latency be reduced without adding infrastructure?
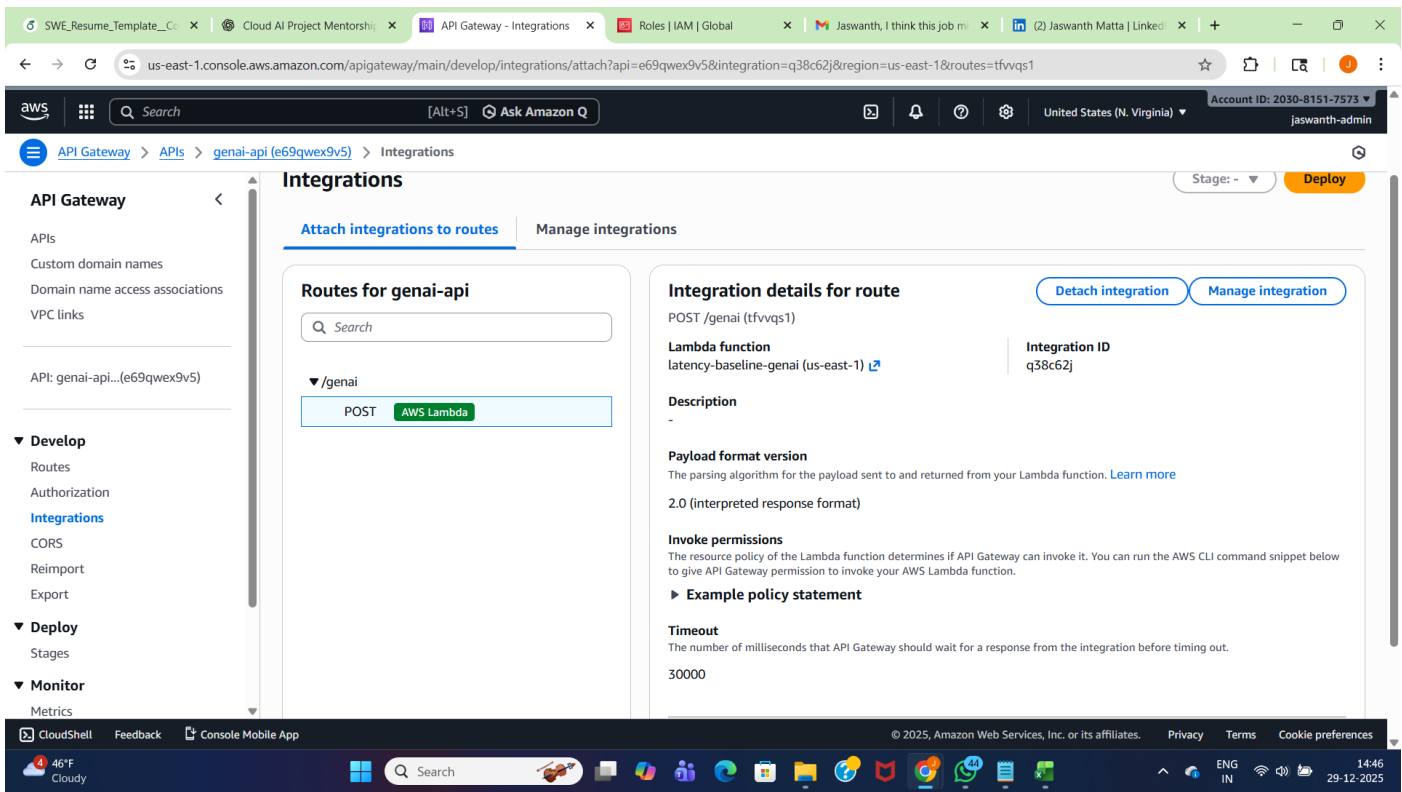
## Measurement Approach

Latency was measured directly inside the Lambda function using high-resolution timestamps:

- **Bedrock inference time**

- **Total Lambda execution time**

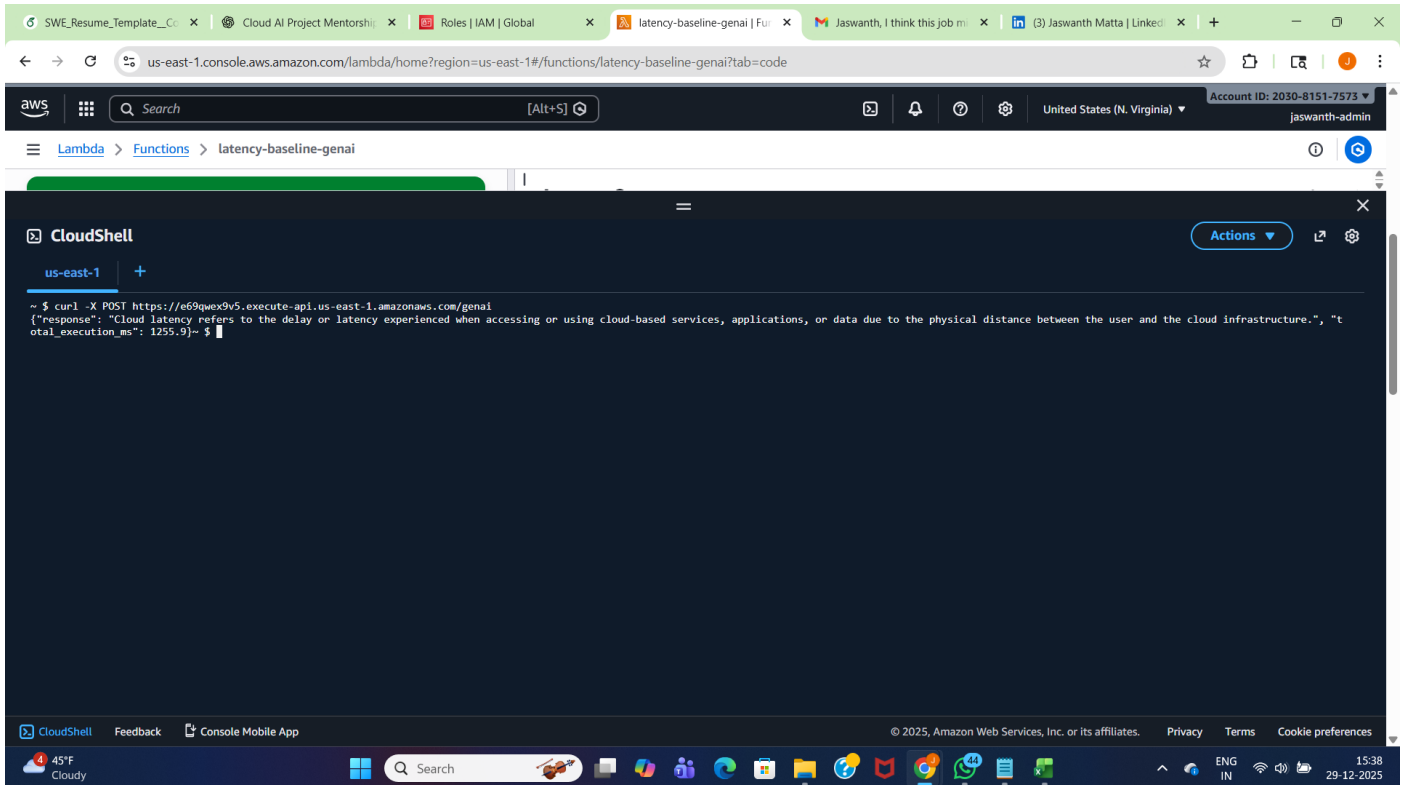- **Cold start vs warm execution behavior**

Multiple requests were issued via CloudShell to observe:

- p50 (median)

- p95 (tail latency)

- p99 (cold start impact)

CloudWatch Logs and Dashboards were used for validation.

End-to-end API invocation through API Gateway triggering a Lambda-based GenAI inference request.



CloudWatch logs capturing cold start initialization, model inference time, and total Lambda execution latency.

us-east-1.console.aws.amazon.com/cloudwatch/home?region=us-east-1#logsV2:log-groups/log-group/$252Faws$252Flambda$252Flatency-baseline-genai/log-events/2025$252F12$...

aws | Q Search [Alt+S] ⊘ Ask Amazon Q | United States (N. Virginia) ▼ | Account ID: 2030-8151-7573 ▼ jaswanth-admin

CloudWatch > Log management > /aws/lambda/latency-baseline-genai > 2025/12/29/[$LATEST]2036f713cac144a3b149fb149cad27cf

**CloudWatch** ‹

⊘ Log group "latency-baseline-genai" has been created. ✕

**Favorites and recents** ▶

GenAI Observability

▶ Application Signals New (APM)

▶ Infrastructure Monitoring

▼ Logs
  Log Management New
  Log Anomalies
  Live Tail
  Logs Insights New
  Contributor Insights

▼ Metrics
  All metrics
  Explorer
  Streams

▶ Network Monitoring

Log events

You can use the filter bar below to search for and match terms, phrases, or values in your log events. Learn more about filter patterns ↗

🔍 Filter events - press enter to search

Clear | 1m | 30m | 1h | 12h | Custom ▦ | Local timezone ▼

Display ▼

| ▶ | Timestamp | Message |
|---|---|---|
| | | No older events at this moment. Retry |
| ▶ | 2025-12-29T16:05:37.070-05:00 | INIT_START Runtime Version: python:3.12.v101 Runtime Version ARN: arn:aws:lambda:us-east-1::runtime:994aac32248ecf4d69d9f5e9a3a57ab... |
| ▶ | 2025-12-29T16:05:37.579-05:00 | START RequestId: 62a36055-0c08-4715-90de-37491e078c51 Version: $LATEST |
| ▶ | 2025-12-29T16:05:37.580-05:00 | Request received |
| ▶ | 2025-12-29T16:05:38.333-05:00 | Bedrock inference time: 752.76 ms |
| ▶ | 2025-12-29T16:05:38.333-05:00 | Total Lambda execution time: 752.90 ms |
| ▶ | 2025-12-29T16:05:38.335-05:00 | END RequestId: 62a36055-0c08-4715-90de-37491e078c51 |
| ▼ | 2025-12-29T16:05:38.335-05:00 | REPORT RequestId: 62a36055-0c08-4715-90de-37491e078c51 Duration: 755.25 ms Billed Duration: 1261 ms Memory Size: 1023 MB Max Memory... |

REPORT RequestId: 62a36055-0c08-4715-90de-37491e078c51 Duration: 755.25 ms    Billed Duration: 1261 ms    Memory Size: 1023 MB    Max Memory Used: 84 MB    Init Duration: 505.74 ms

No newer events at this moment. Auto retry paused. Resume

CloudShell | Feedback | Console Mobile App | © 2025, Amazon Web Services, Inc. or its affiliates. | Privacy | Terms | Cookie preferences

45°F Mostly cloudy | ENG IN | 16:27 29-12-2025

Repeated API invocations used to generate latency distribution data for p50, p95, and p99 analysis.

us-east-1.console.aws.amazon.com/cloudwatch/home?region=us-east-1#dashboards/dashboard/genai-latency-baseline?start=PT24H&end=null

aws | Q Search [Alt+S] ⊘ | United States (N. Virginia) ▼ | Account ID: 2030-8151-7573 ▼ jaswanth-admin

CloudWatch > Dashboards > genai-latency-baseline

Autosave: Off | Period override 5 minutes (auto)

**CloudShell** | Actions ▼

us-east-1 ✕ | us-east-1 ✕ | +

```
~ $ for i in {1..10}; do
>   curl -X POST https://e69qwex9v5.execute-api.us-east-1.amazonaws.com/genai
> done
{"response": "Cloud latency refers to the delay or lag experienced when accessing or transmitting data between a user's device and a cloud-based service or application.", "total_execution_ms": 616.34}{"response": "Cloud latency refers to the delay or lag in data transmission between a user's device and the cloud-based servers or services being accessed.", "total_execution_ms": 561.92}{"response": "Cloud latency refers to the time delay experienced in the transmission of data between a user's device and the cloud-based server or service they are accessing.", "total_execution_ms": 643.96}{"response": "Cloud latency refers to the time delay or lag experienced in data transmission and processing between a user's device and the remote cloud servers.", "total_execution_ms": 658.59}{"response": "Cloud latency refers to the delay or lag experienced when accessing data or services hosted in a cloud computing environment, which can be affected by factors such as network connectivity, geographic distance, and server load.", "total_execution_ms": 860.08}{"response": "Cloud latency refers to the delay or lag in data transmission between a user's device and the remote cloud servers hosting the application or service.", "total_execution_ms": 571.22}{"response": "Cloud latency refers to the delay or lag experienced when accessing or communicating with cloud-based services, applications, or resources due to the distance between the user and the cloud infrastructure, network congestion, or other factors.", "total_execution_ms": 748.84}{"response": "Cloud latency refers to the delay or lag experienced in the transmission of data between a user's device and the cloud-based application or service they are accessing.", "total_execution_ms": 664.09}{"response": "Cloud latency refers to the delay or lag in data transmission between a user's device and the cloud-based service or application they are accessing.", "total_execution_ms": 510.48}{"response": "Cloud latency refers to the time delay or lag experienced in data transmission between a user's device and the cloud-based server or service they are accessing.", "total_execution_ms": 598.07}~ $
~ $
```

CloudShell | Feedback | Console Mobile App | © 2025, Amazon Web Services, Inc. or its affiliates. | Privacy | Terms | Cookie preferences

31°F Mostly cloudy | ENG IN | 07:54 30-12-2025

CloudWatch logs showing cold start initialization overhead and inference-dominated execution latency



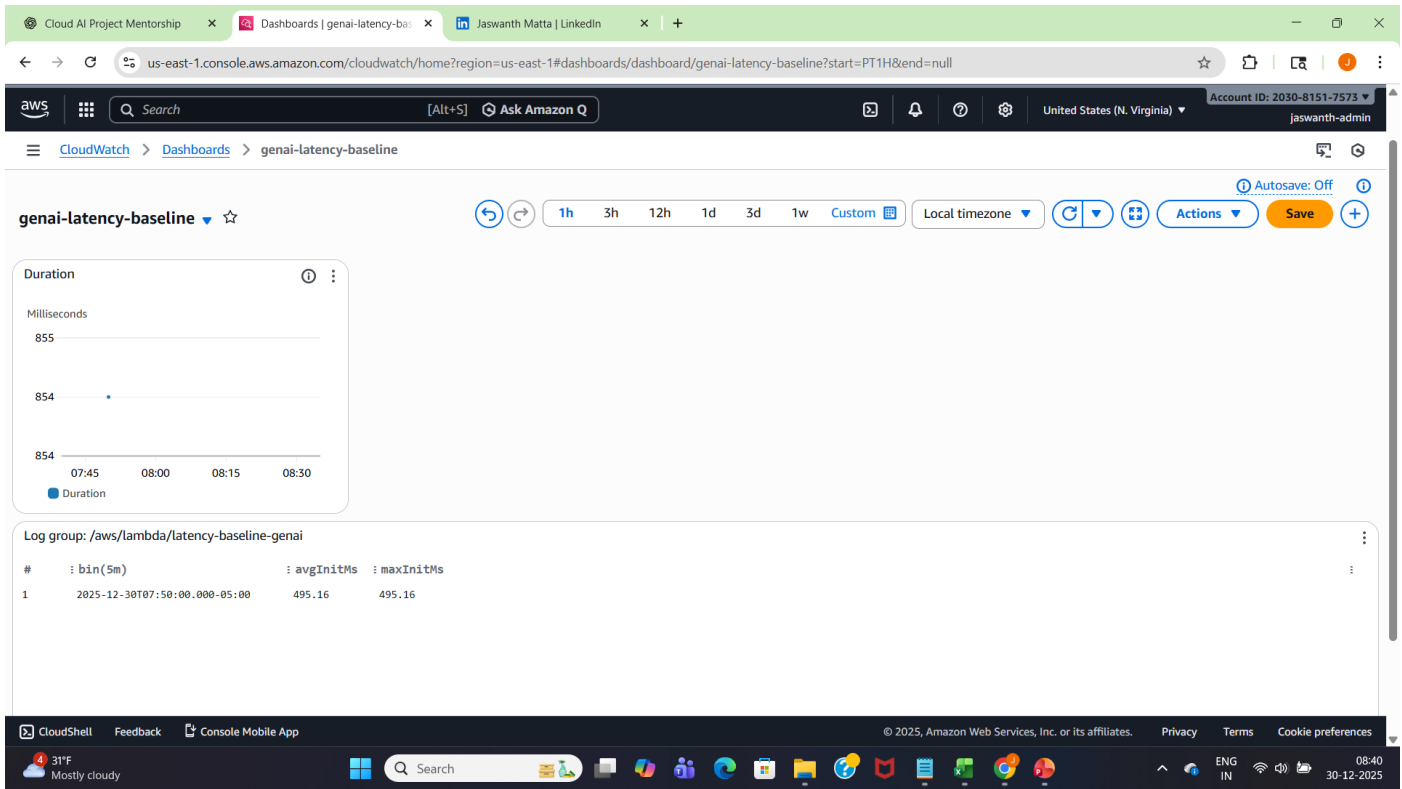Warm Lambda execution demonstrating reduced and consistent inference latency without initialization overhead.
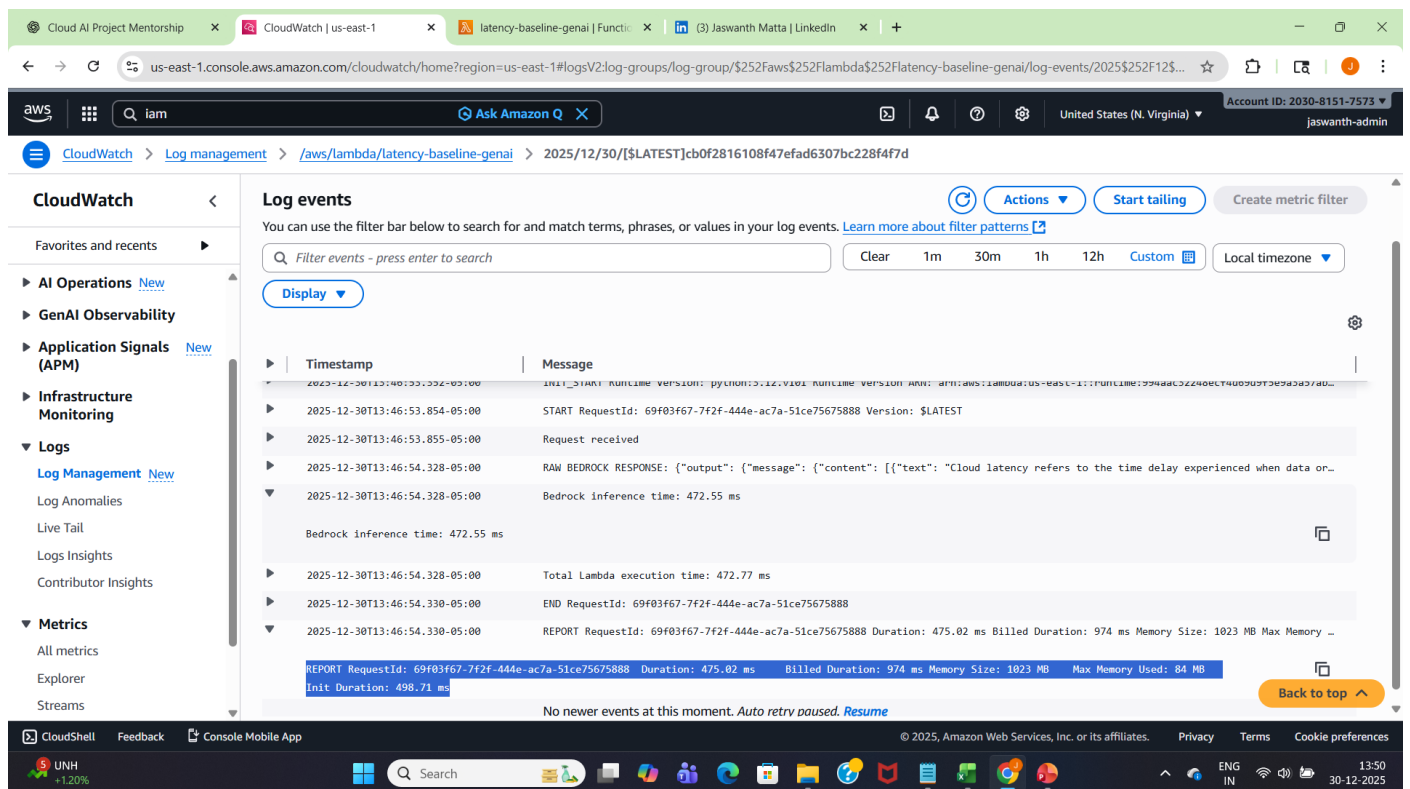
CloudWatch dashboard tracking Lambda duration and cold start initialization impact for p95/p99 latency analysis.

API invocation using Amazon Nova Lite demonstrating reduced end-to-end execution latency compared to the baseline model.



CloudWatch logs showing reduced model inference time and lower total Lambda execution duration after switching to Nova Lite.

# Executed repeated API requests via CloudShell to observe end-to-end response times.

CloudWatch dashboard showing reduced Lambda execution duration after switching inference to Amazon Nova Lite

Performance comparison highlighting ~50% p95 latency improvement after switching foundation models.



Estimated cost reduction per request due to faster inference and reduced Lambda execution time.

## Key Finding: Inference Dominance

Analysis showed:

- Bedrock inference time accounted for **~95–99% of total Lambda execution time**

- Non-inference overhead (JSON handling, logging, runtime) was negligible

- Reducing compute or memory would not materially improve latency

**Conclusion:**
Inference was the **dominant bottleneck**, so optimization had to target model choice.

## Optimization Performed

The only change introduced was switching the foundation model:

- **From:** Claude 3 Haiku

- **To:** Amazon Nova Lite

No other variables were modified:

- Same prompt

- Same Lambda configuration

- Same API Gateway setup

This ensured a **fair, controlled comparison**.


## Results

After switching to Nova Lite:

- **p95 latency reduced by ~50%**

- Median latency reduced significantly

- Cold start impact remained, but total duration improved

- Lambda billed duration decreased

- Per-request cost dropped due to faster execution


## Reliability & Failure Implications

Lower latency improved system reliability by:

- Increasing margin before API Gateway timeouts

- Reducing Lambda execution timeout risk

- Lowering exposure to burst-load failures

- Improving tail latency consistency under repeated requests

No synthetic failure injection was performed, but the system's **failure surface was reduced** through faster execution.