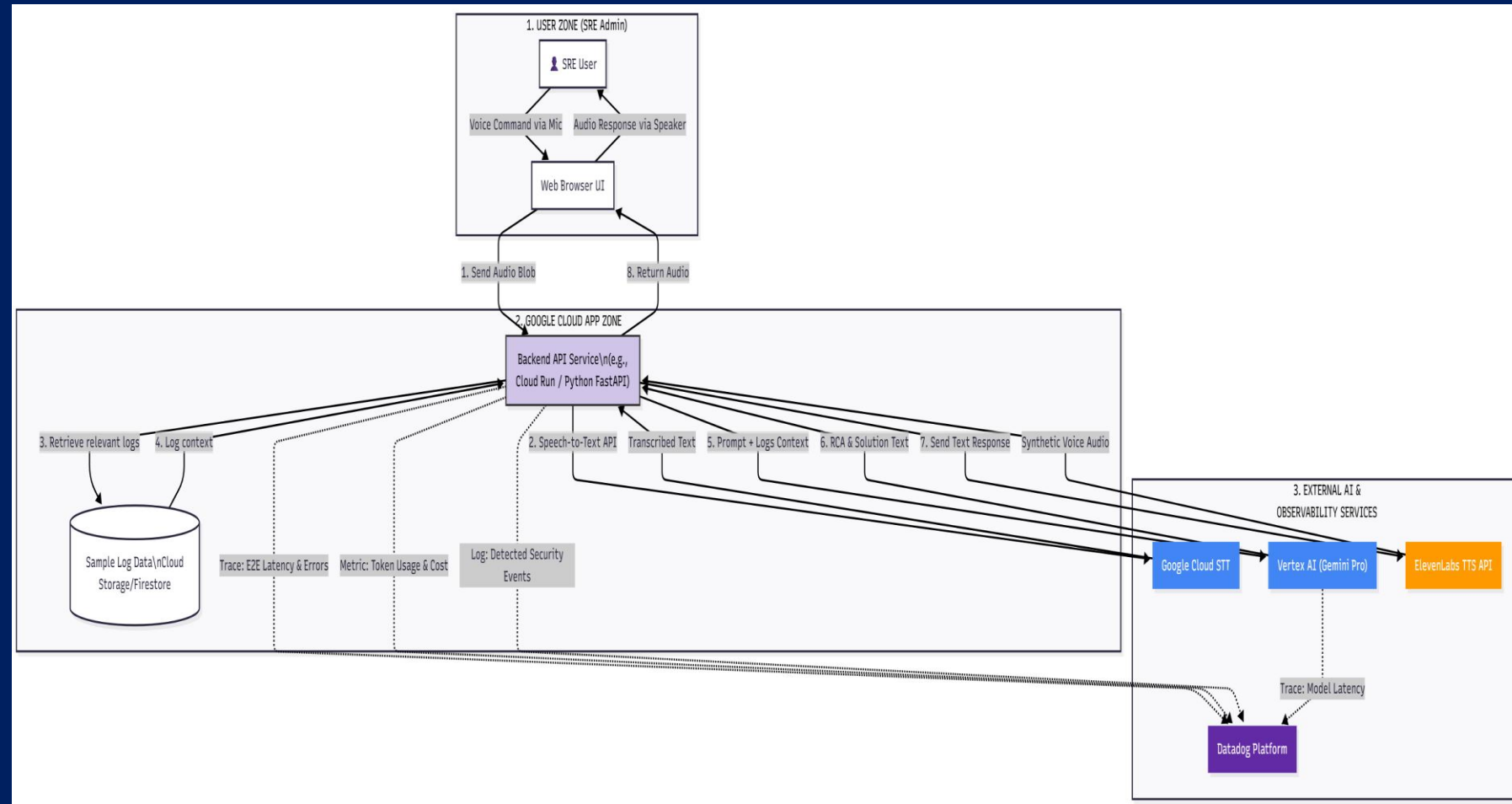# *VOICE-ENABLED OBSERVABILITY FOR AI LATENCY*
## *NIRAJ – SENIOR SOFTWARE ENGINEER & PRODUCT BUILDER*

# PROBLEM

- AI backends suffer from hidden latency blind spots

- Traditional monitoring tools miss AI-layer delays

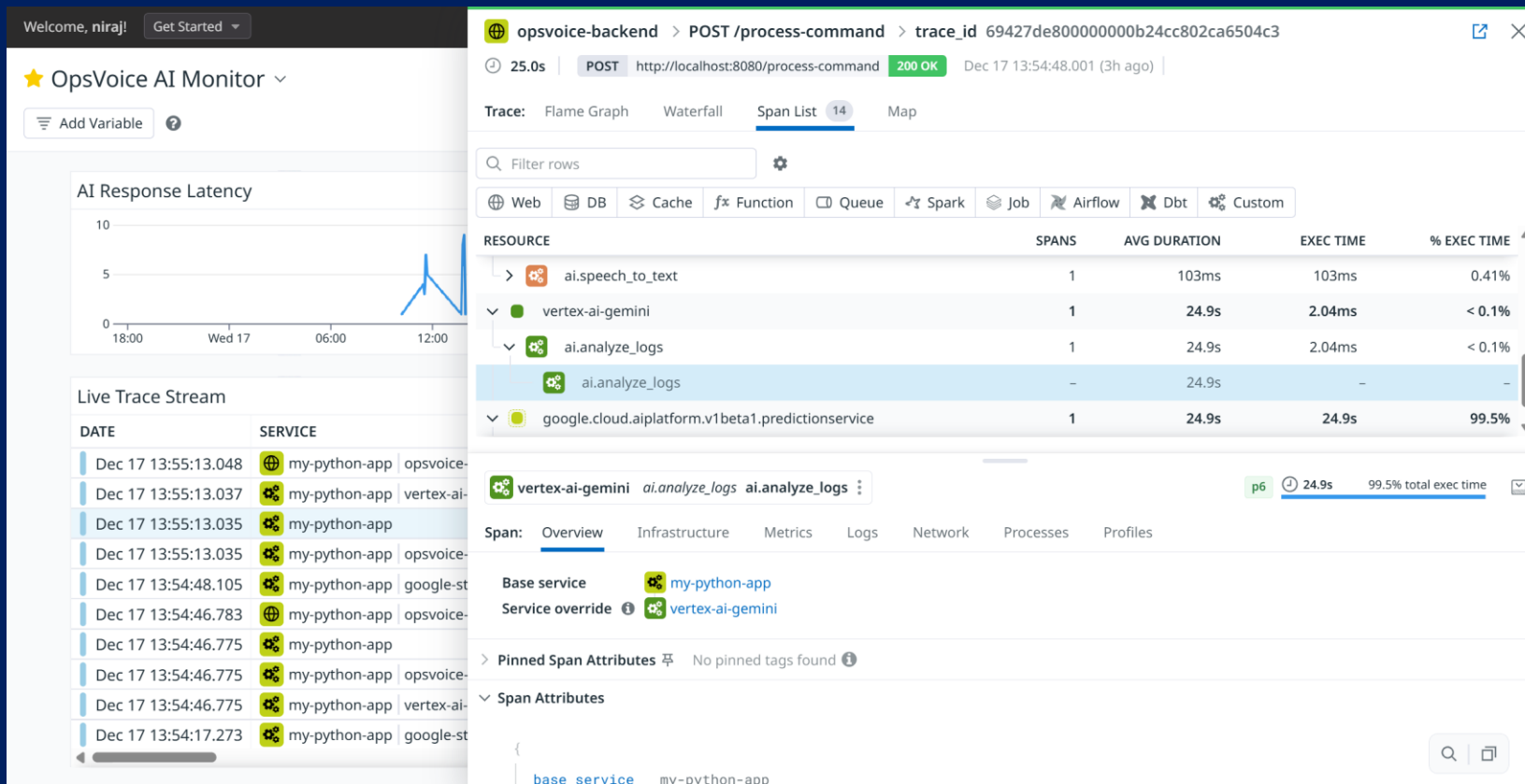- Slow responses → poor user retention & lost revenue

# SOLUTION

- **CLOUD RUN DEPLOYMENT FOR SCALABLE, SERVERLESS INFRA**
- **FLASK BACKEND WITH DATADOG TRACING & LOGGING**
- **VERTEX AI GEMINI INTEGRATION FOR REAL LLM ANALYSIS**
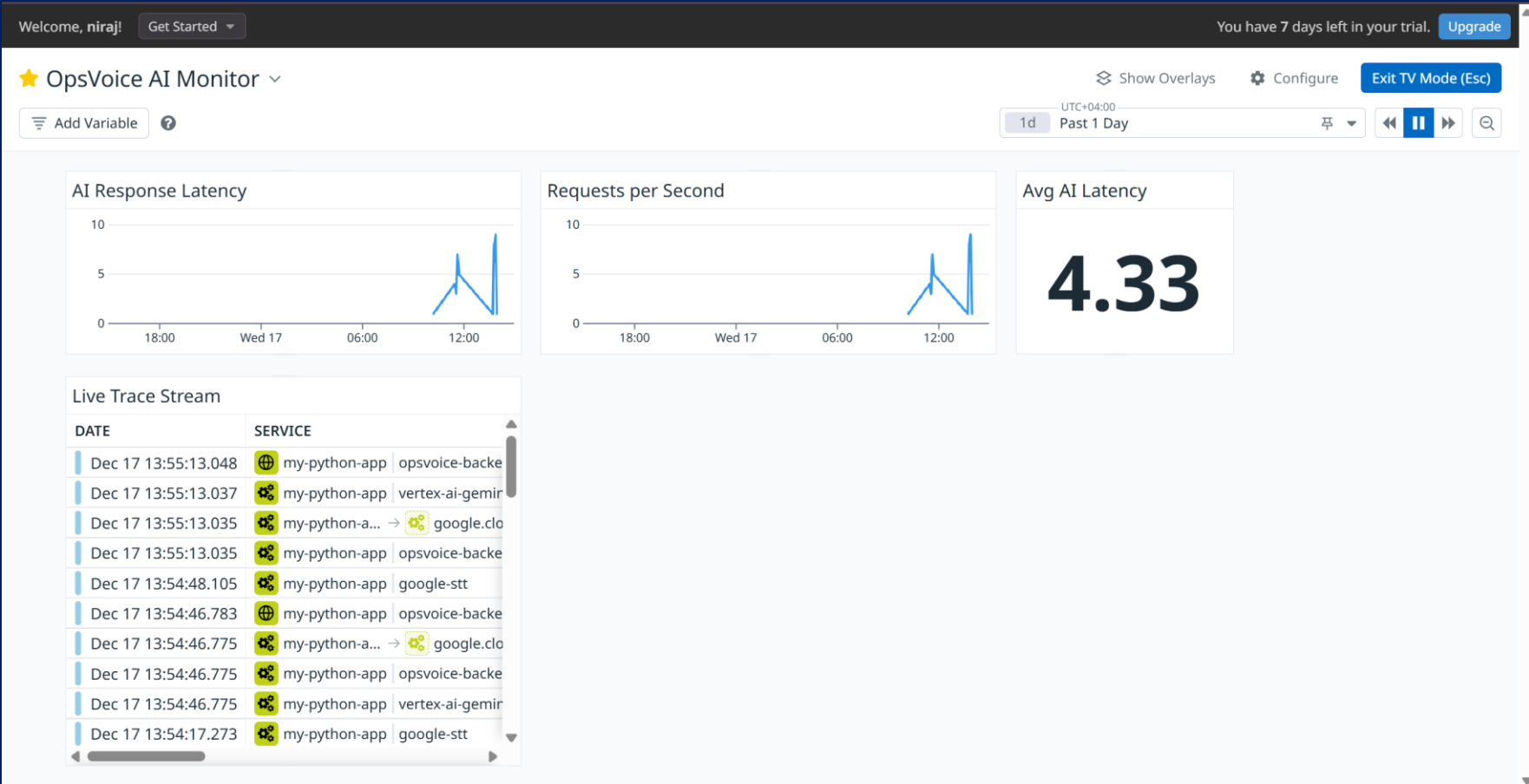- **VOICE-ENABLED QUERIES FOR REAL-TIME OBSERVABILITY**

# DEMO VALUE

- **DETECTED A 33-SECOND LATENCY ANOMALY IN GEMINI INFERENCE**
- **PINPOINTED ROOT CAUSE AT AI LAYER VS NETWORK LAYER**
- **DASHBOARDS SHOW LATENCY, ERROR RATES, RESOLUTION TIME**
- **JUDGE-FRIENDLY, REPRODUCIBLE DEMO ASSETS**

# Market Opportunity

**- AI OBSERVABILITY
MARKET: $20B+ BY 2028
- ENTERPRISES ADOPTING
LLMS AT 40% CAGR
- OPSVOICE BRIDGES THE
CRITICAL AI-LAYER VISIBILITY
GAP**

# BUSINESS MODEL

- SAAS SUBSCRIPTION FOR ENTERPRISE SRE TEAMS
- USAGE-BASED PRICING TIED TO TOKENS PER SECOND
- FUTURE: COST-PER-QUERY ANALYTICS FOR BUDGET CONTROL

## COMPETITIVE EDGE

- COMPETITORS: DATADOG, NEW RELIC, ELASTIC
- DIFFERENTIATOR: AI-NATIVE OBSERVABILITY + VOICE INTERFACE
- FASTER SETUP, HACKATHON-PROVEN, REPRODUCIBLE DASHBOARDS

**TEAM**

- NIRAJ – BACKEND & OBSERVABILITY EXPERT (PYTHON, FLASK, DATADOG, CLOUD RUN)
- HACKATHON-PROVEN DELIVERY MINDSET
- FUTURE HIRES: SALES, DEVREL, PRODUCT

## ASK

- SEEKING $500K SEED FUNDING
- FUNDS FOR INFRA SCALING, GTM, AND ENTERPRISE PILOTS
- GOAL: ONBOARD 10 ENTERPRISE CUSTOMERS IN 18 MONTHS

**VISION**


**"STOP GUESSING. START OBSERVING AI PERFORMANCE
IN REAL TIME."
CONTACT: NIKY.SWAY@GMAIL.COM**