


Link: <https://llmo-insight-flow.lovable.app>



LLMOps Dashboard

Connect your LLM provider to start monitoring


Application Name

My AI Assistant

Application Context (Optional)

A customer support chatbot for e-commerce...

LLM Provider


OpenAI (GPT-4) 

API Key


sk-...

Your API key is only used for validation and is not stored.


Validate & Continue



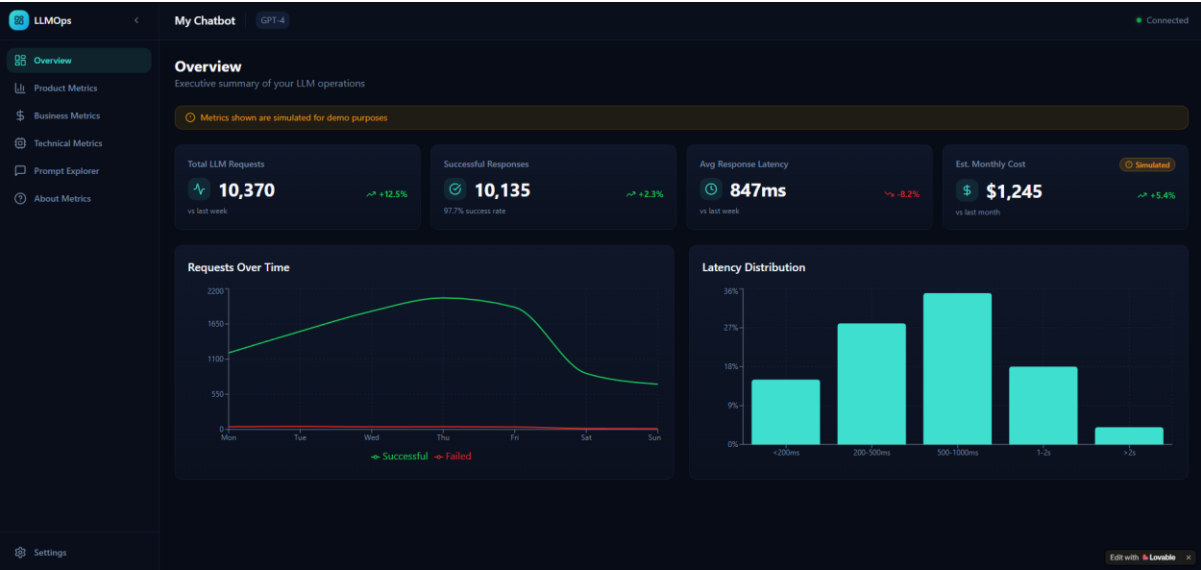
Real-time Monitoring



Secure & Private



AI-Powered Insights



Product Metrics

Insights for Product Managers

Analytics shown are simulated for demonstration

Unique Prompts

2,847

vs last period

+18.3%

Regeneration Rate

8.2%

Lower is better

-12.5%

Hallucination Risk

Low

-12.5%

Intent Match Rate

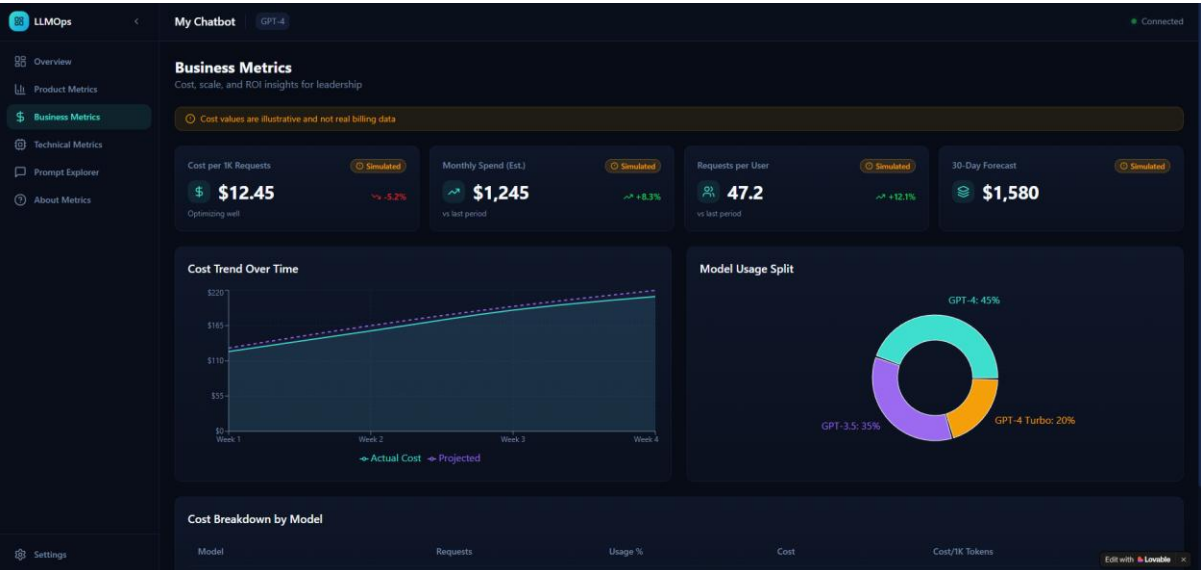
94.2%

vs last period

+3.7%

Prompt Categories

Top User Prompts



Business Metrics

Cost, scale, and ROI insights for leadership

Cost values are illustrative and not real billing data

Cost per 1K Requests

\$12.45

Optimizing well

-5.2%

Monthly Spend (Est.)

\$1,245

vs last period

+8.3%

Requests per User

47.2

vs last period

+12.1%

30-Day Forecast

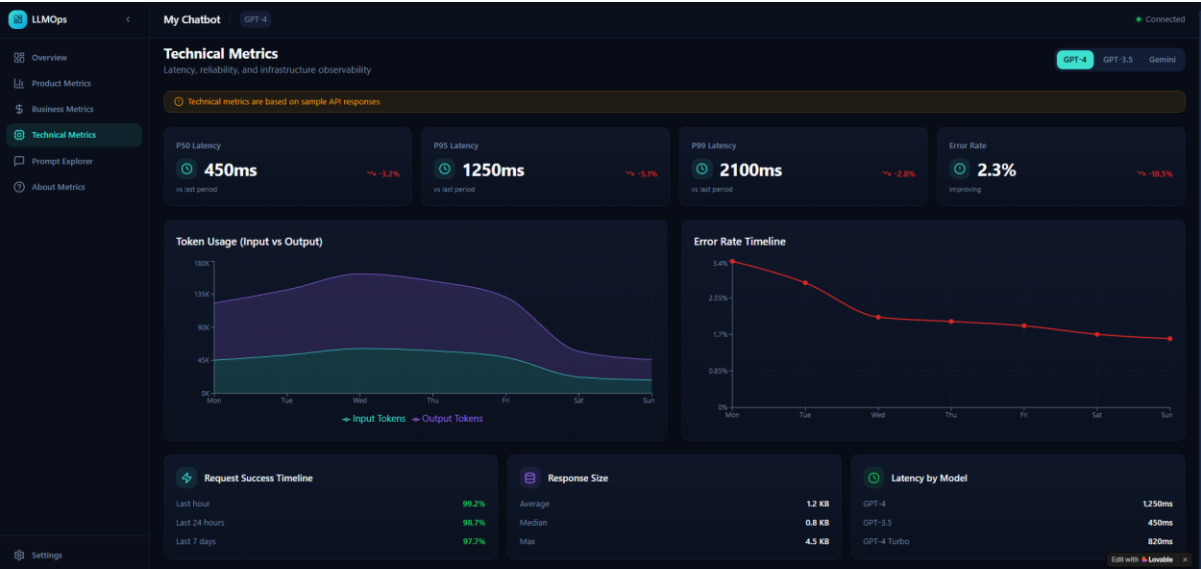
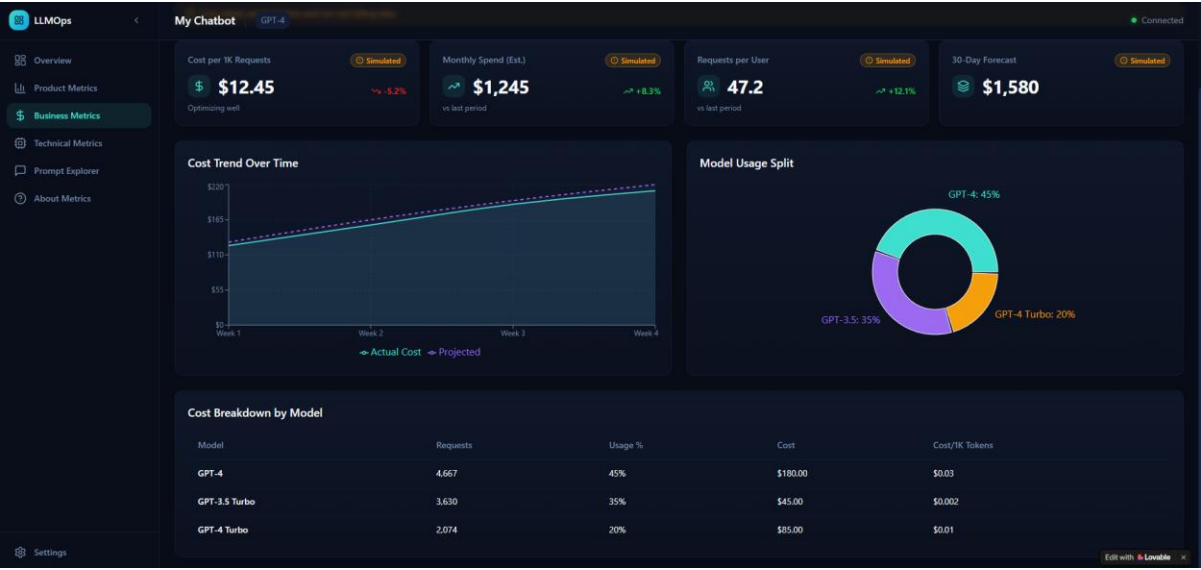
\$1,580

+12.1%

Cost Trend Over Time

Model Usage Split

Cost Breakdown by Model



LLMOps

My Chatbot

GPT-4

Connected

Overview

Product Metrics

Business Metrics

Technical Metrics

Prompt Explorer

About Metrics

Prompt Explorer

Browse and analyze all LLM interactions

Search prompts, responses, or categories...

Timestamp	Prompt	Response	Latency	Tokens	Status
2024-01-15 14:32:05	Summarize the key points fro...	Based on the Q4 financial repo...	1250ms	201	Success
2024-01-15 14:28:42	Generate a product description...	Introducing Analytica! Pro ---	890ms	117	Success
2024-01-15 14:25:18	What are the best practices for...	Here are key microservices bes...	1580ms	152	Success
2024-01-15 14:20:33	Translate this customer feedba...	Translation: "El producto es e...	720ms	102	Success
2024-01-15 14:15:22	Generate SQL query to find to...	SELECT customer_id, c custo...	450ms	100	Success
2024-01-15 14:10:05	Analyze competitor pricing str...	Error: Rate limit exceeded. Ple...	150ms	166	Error
2024-01-15 14:05:48	Create a user onboarding ema...	Email 1 (Day 0): Welcome & G...	1120ms	163	Success
2024-01-15 14:00:12	Explain the concept of vector ...	Vector databases are specializ...	980ms	128	Success

Settings

Edit with Loveable

LLMOps

<

Overview

Product Metrics

Business Metrics

Technical Metrics

Prompt Explorer

About Metrics

Settings

My Chatbot

GPT-4

✓

Real API-Level Data

These metrics are directly available from LLM API responses:

• Prompt & Response Text

The actual input and output of each API call

• Latency

Time from request to response completion

• Token Usage

Input and output token counts per request

• Model Metadata

Model name, version, and configuration

• Error Responses

Rate limits, API errors, and failures

🕒

Simulated LLMOps Data

These metrics require additional middleware or infrastructure:

• Hallucination Rate

Requires ground truth comparison or human review

• Response Accuracy

Needs evaluation framework and test datasets

• Cost Attribution

Requires usage tracking per user/feature

• Quality Scores

Needs custom evaluation pipelines

• User-Level Analytics

Requires session tracking and user identification

Why Middleware is Required for Real LLMOps

Moving from API-only observability to full LLMOps requires additional infrastructure:

</>

Proxy Layer

Intercept all LLM calls to log, modify, and enrich requests and responses

🔄

Evaluation Pipeline

Run automated quality checks, compare outputs, and track regressions

🛡️

Context Store

Persist conversation history, user sessions, and ground truth data

From API to Full LLMOps

1

API Only

Basic logging

→

2

Proxy Layer

Request enrichment

→

3

Evaluation

Quality scoring

→

4

Full LLMOps

Complete platform