# ATLAS AI

- Overview
- Projects
- Models
- Utilization & Efficiency
- Cost & Spend
- Policies & Quotas
- Incidents
- Reports

# Tenants

Maya Chen
Platform Admin

| Projects | Models | GPUs | Utilization | Spend (MTD) | SLA Compliance | Risk |
|----------|--------|------|-------------|-------------|----------------|------|
| 4 | 12 | 28 | 76% | $18.4k | 99.3 % | 2 |

## Usage by Project

| | | |
|--|--|--|
| Chatbot LLM | 🟩🟩🟩🟩🟩🟥🟥🟥 | 3 GPUs |
| Search API | 🟩🟩🟩🟩🟩🟩 | 6 GPUs |
| RAG Service | 🟩🟩🟩🟩 | 4 GPUs |
| Batch Job | 🟨🟨🟨 | 3 GPUs |

Busy ●  Under-utilized ●  Idle ●

## Spend Trend



## Model health & SLA



| 10 | Healthy | ● |
|----|---------|---|
| 1 | Degrading | ● |
| 1 | Falling | ● |

## Risk and Violations

| Policy Violations | ● | 3 |
|-------------------|---|---|
| Cost Overages | ● | 2 |
| Under-utilized | ● | 4 |
| Preemtion Risk | ● | 1 |

## Tenants Table

| Project | Models | GPUs | Tier | Util % | Spend (MTD) | SLA | Risks | Status |
|---------|--------|------|------|--------|-------------|-----|-------|--------|
| Search | 4 | 10 | Reserved | 88 | $6.2k | 99.9 | 0 | Healthy |
| Chatbot | 3 | 8 | Elastic | 61 | $4.8k | 98.7 | 2 | ⚠ |
| RAG | 3 | 6 | BestEffort | 42 | $3.1k | 97.8 | 3 | ⚠ |
| Batch | 2 | 4 | BestEffort | 30 | $1.9k | - | 0 | Idle |

powered by PARALLEL^IQ