# Aegik8s
🔊 *EE-jis*

## Deploy LLM and Agentic Workloads on Kubernetes - Simple, Secure, and So Easy Even a 5-Year-Old Can Do It!
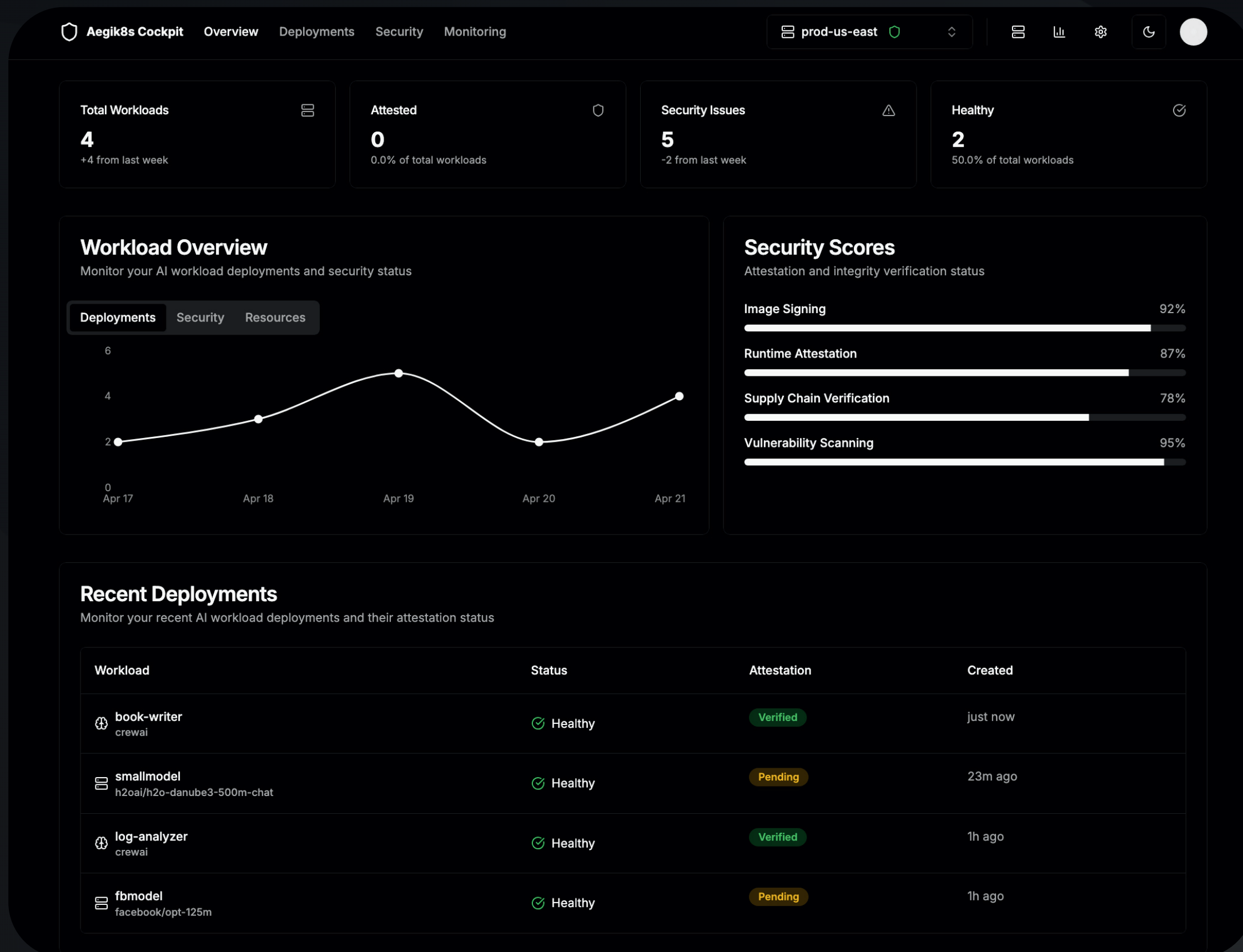
### Problem

Hundreds of models are uploaded daily to platforms like Hugging Face, offering accessible tools but also introducing serious security risks. Some contain malicious code - real incidents have shown attackers gaining reverse shell access via compromised models. As LLM adoption grows, so does the threat to system integrity.

Fast-paced development often skips thorough evaluation, exposing systems to hidden threats and poor security practices. This risk is even higher with emerging agentic models, which lack mature safeguards. Running these in secure environments is crucial for safe, responsible progress.
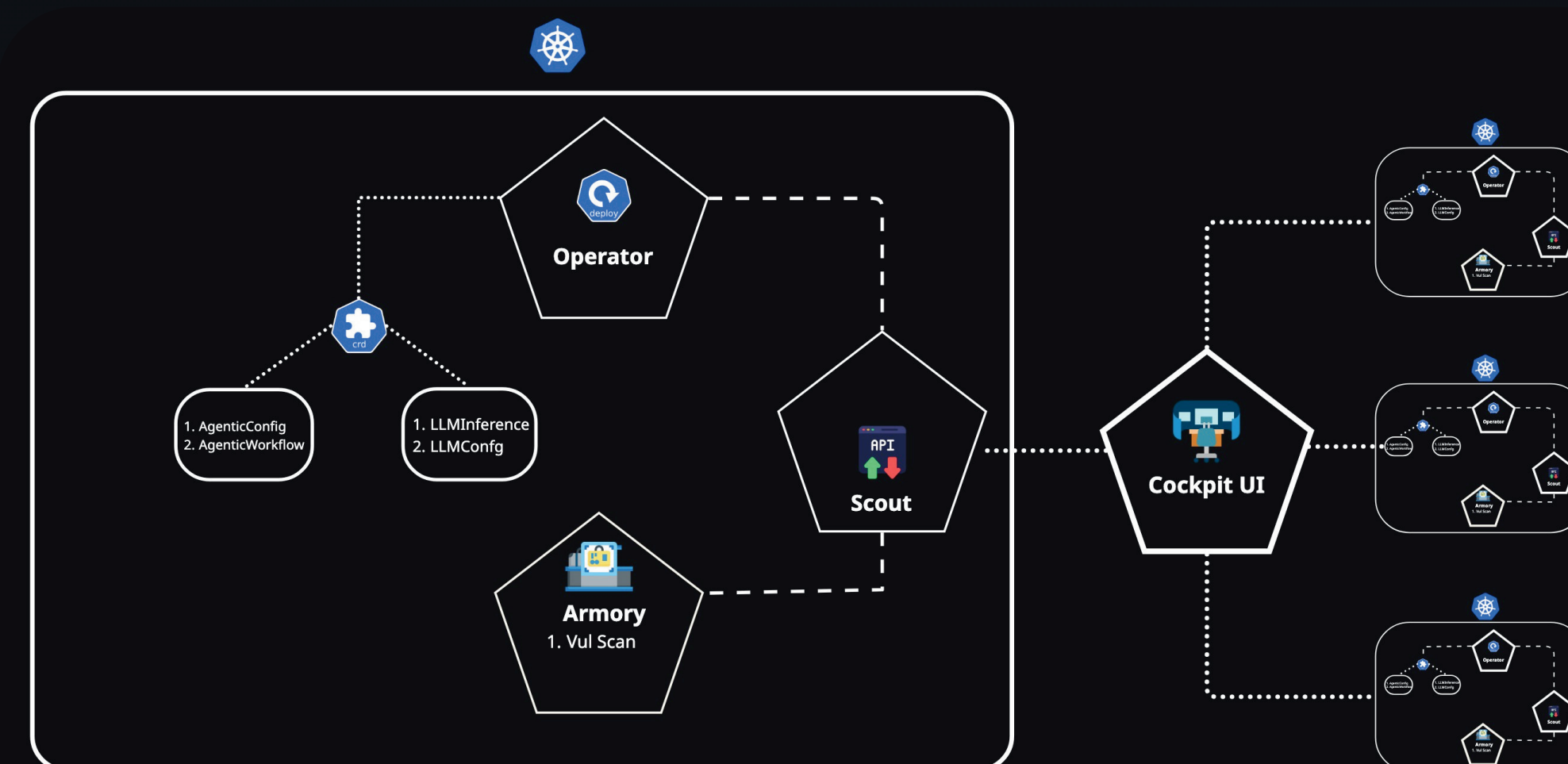
### Solution

Aegik8s aims to provide a highly opinionated, security-first framework for deploying LLMs and agentic workloads on Kubernetes. It follows industry - aligned best practices by running each model and agent in its own lightweight MicroVM (Micro Virtual Machine) on top of Kubernetes, ensuring strong isolation, minimal attack surface, and scalable multi-tenancy. This approach enables teams to deploy AI systems confidently without compromising on performance or security.
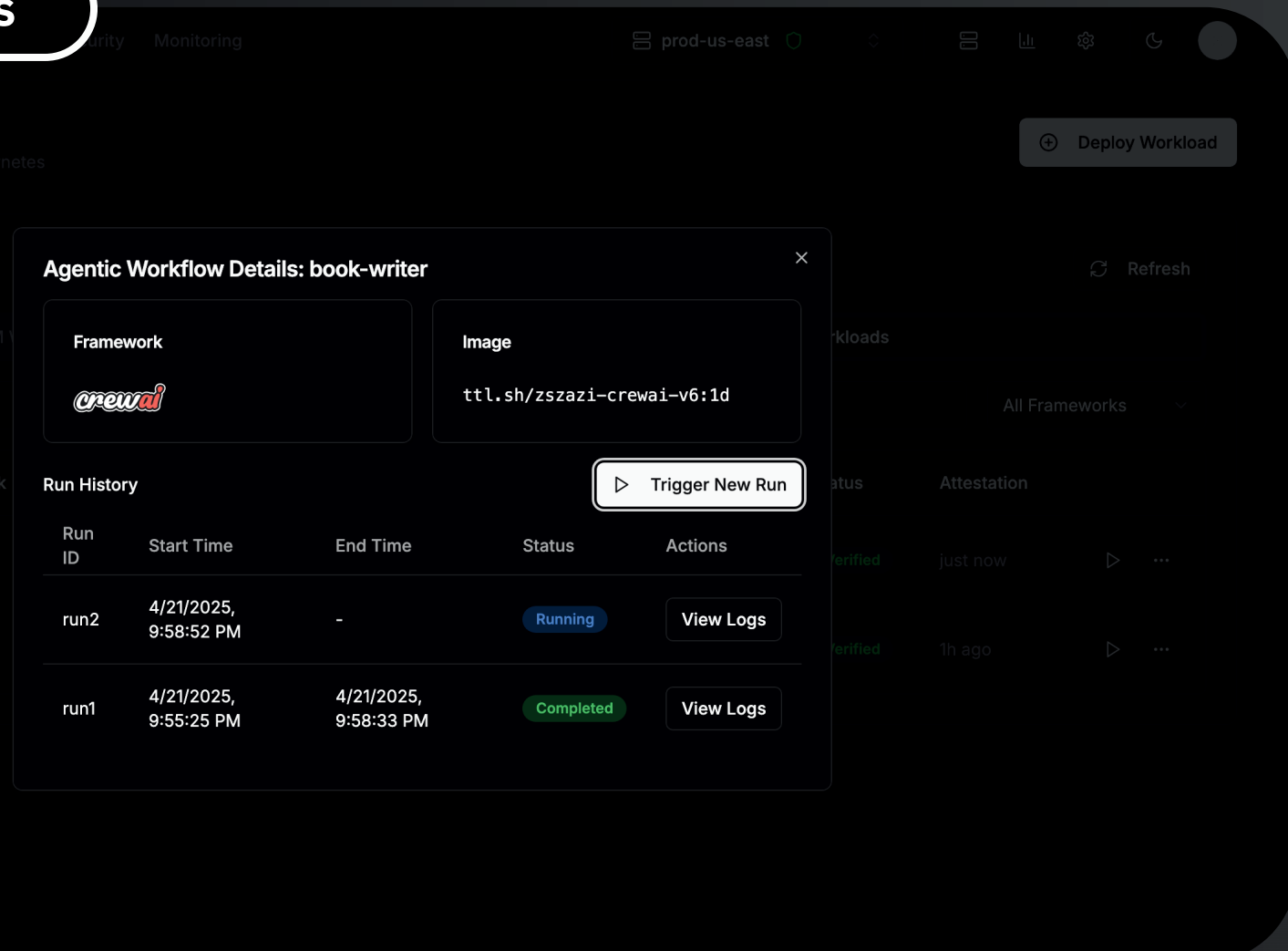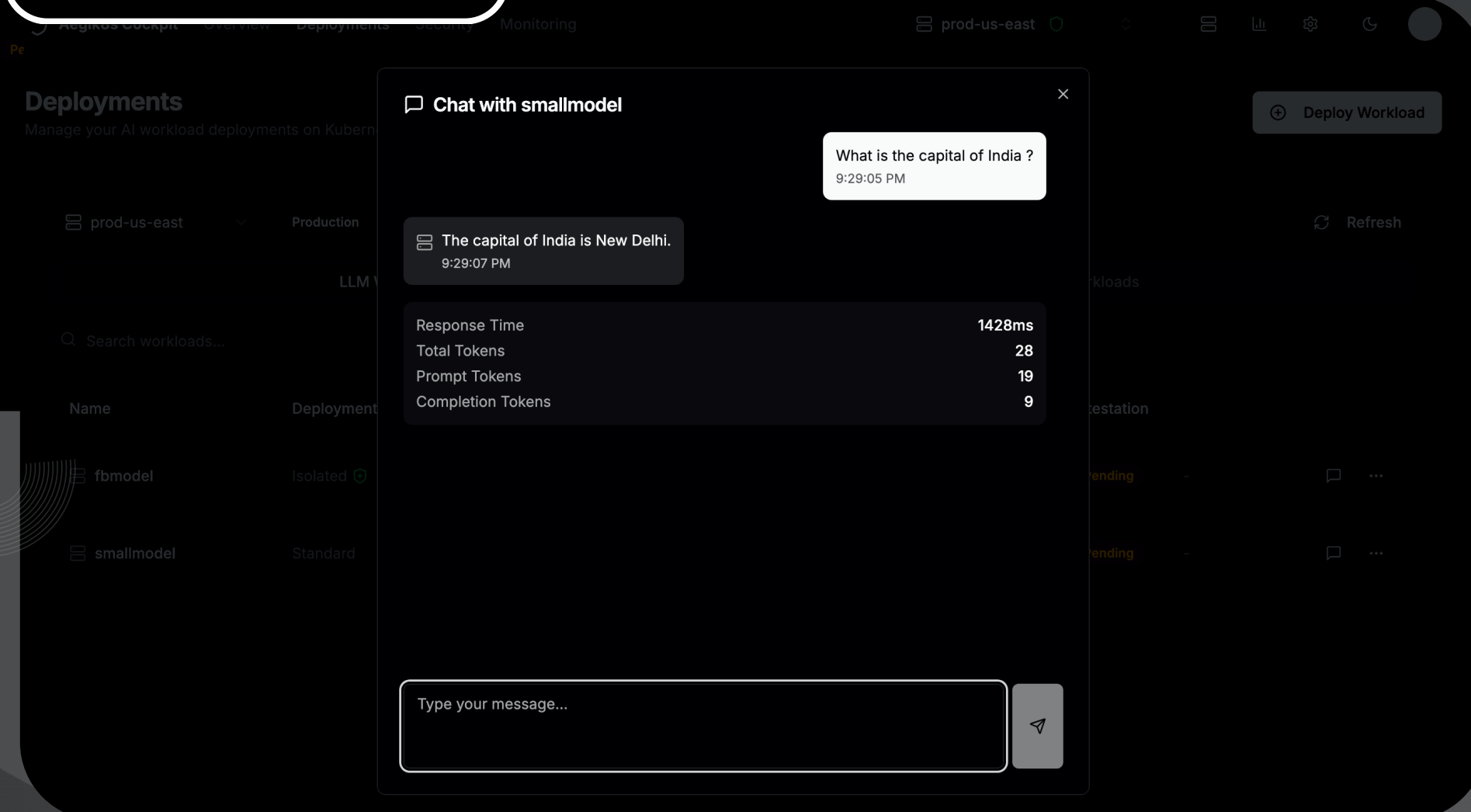
### Cockpit Dashboard



### Architecture



- **Scout**: Aegik8s' control plane handles platform requests, manages workloads, and connects the UI with clusters
- **Operator**: Kubernetes-native controller that automates LLM and agent workload deployment, scaling, and teardown
- **Armory**: Pluggable security engine for scanning, attestation, and securing LLM and agent workloads
- **Cockpit**: UI for managing LLM and agent workloads across clusters in dev, staging, and production

### Agentic Workloads
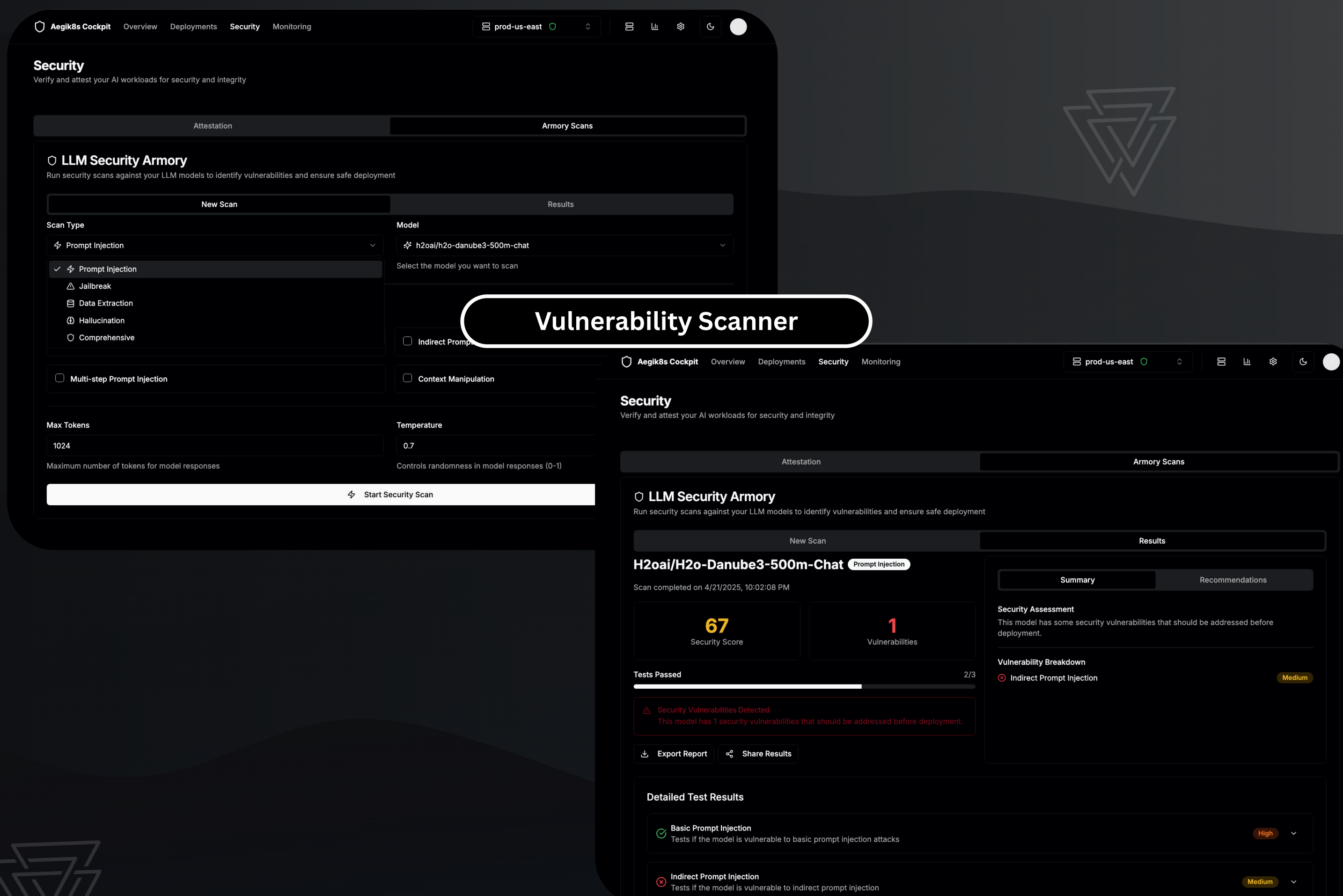


### LLM Workloads



### Why Aegik8s ?

- **Deploy LLMs in 2 Clicks**
  - Spin up models effortlessly - no expertise needed
- **Secure by Default**
  - Enforces opinionated best practices
- **One-Command Setup**
  - Launch the full platform via a single Helm chart
- **Strong Isolation**
  - Each workload runs in its own MicroVM with dedicated kernel-level isolation for network, I/O, and memory
- **Minimized Blast Radius**
  - Uses Kubernetes and Confidential Computing to isolate threats at both hardware and software levels

Aegik8s runs anywhere—laptops, cloud (AWS/GCP), or even Raspberry Pi—and scales from single-CPU to multi-GPU setups for flexible, distributed LLM deployment

### Vulnerability Scanner



### Market Potential

Only 24% of GenAI projects prioritize security, though 82% of executives see it as essential - highlighting a major gap. With the LLM and AI agent markets projected to hit $36B and $47B by 2030, Aegik8s is poised to meet this demand with a secure, scalable, and cost-effective platform for AI deployment.

### Target Personas

- **MLOps Engineers**: Secure, scalable model deployment with Helm, MicroVMs, and K8s workflows
- **AI Startups**: Fast, secure prototyping with built-in isolation
- **Local Devs**: Safe, lightweight LLMs on Minikube or kind

**Sainath Santa Ram**