# SLM (Security Lifecycle Manager)

SLM provides a real-time secret detection solution powered by fine-tuned small language models (SLMs) running on AWS Trainium.

# Motivation

**Over 10 million secrets leaked on GitHub in 2023**
Sensitive data, such as API keys and credentials, are constantly being exposed on public code repositories, leading to significant security risks and financial losses.
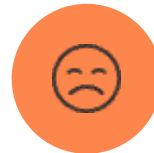
**SOC-2 compliance requires continuous secret scanning**
Maintaining SOC-2 compliance mandates regular monitoring and detection of secret leaks, which is a crucial requirement for many organizations.

**Average cost of a data breach: $4.45M**
The financial impact of a data breach can be devastating, with the average cost reaching millions of dollars per incident.

**Existing solutions are expensive or slow**
Current secret detection tools are either too costly ($100K+/year) or rely on batch processing, which fails to provide real-time protection.

**The constant threat of secret leaks, the high financial impact, and the need for compliance have created a significant opportunity for a cost-effective, real-time secret detection solution.**

# Problem Statement

## Current Options
GitGuardian/TruffleHog: $10K-100K/year, limited customization; GPT-4 API: $100 per 1M scans, data leaves your network; Regex-based tools: High false positives, miss context-aware leaks
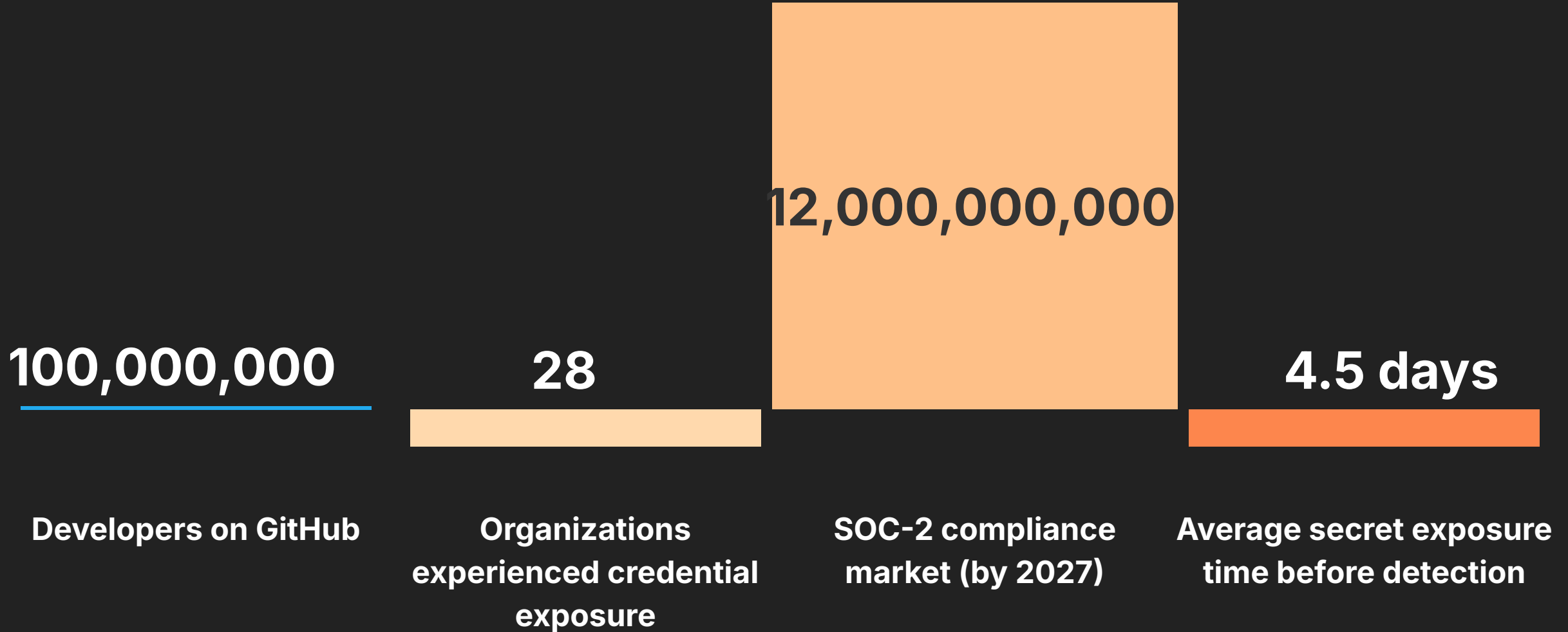
## What's Missing
Real-time detection at commit time, Structured, actionable output for SOC-2 audits, On-premise deployment for sensitive codebases, Cost-effective solution for startups/mid-size companies

Existing solutions are either expensive, slow, or inadequate, leaving a need for a cost-effective, real-time, and on-premise secret detection solution.

# Problem Size

The size of the potential market for a solution to the problem of secret leaks

**100,000,000**

**28**

**12,000,000,000**

**4.5 days**

**Developers on GitHub**

**Organizations experienced credential exposure**

**SOC-2 compliance market (by 2027)**

**Average secret exposure time before detection**

# Our Solution

**Generated 1,000 realistic training examples**
Developed a dataset of exposed and safe code samples to train the model

**Deployed API endpoint for real-time scanning**
Deployed the fine-tuned model as an API endpoint for real-time SOC-2 secret detection
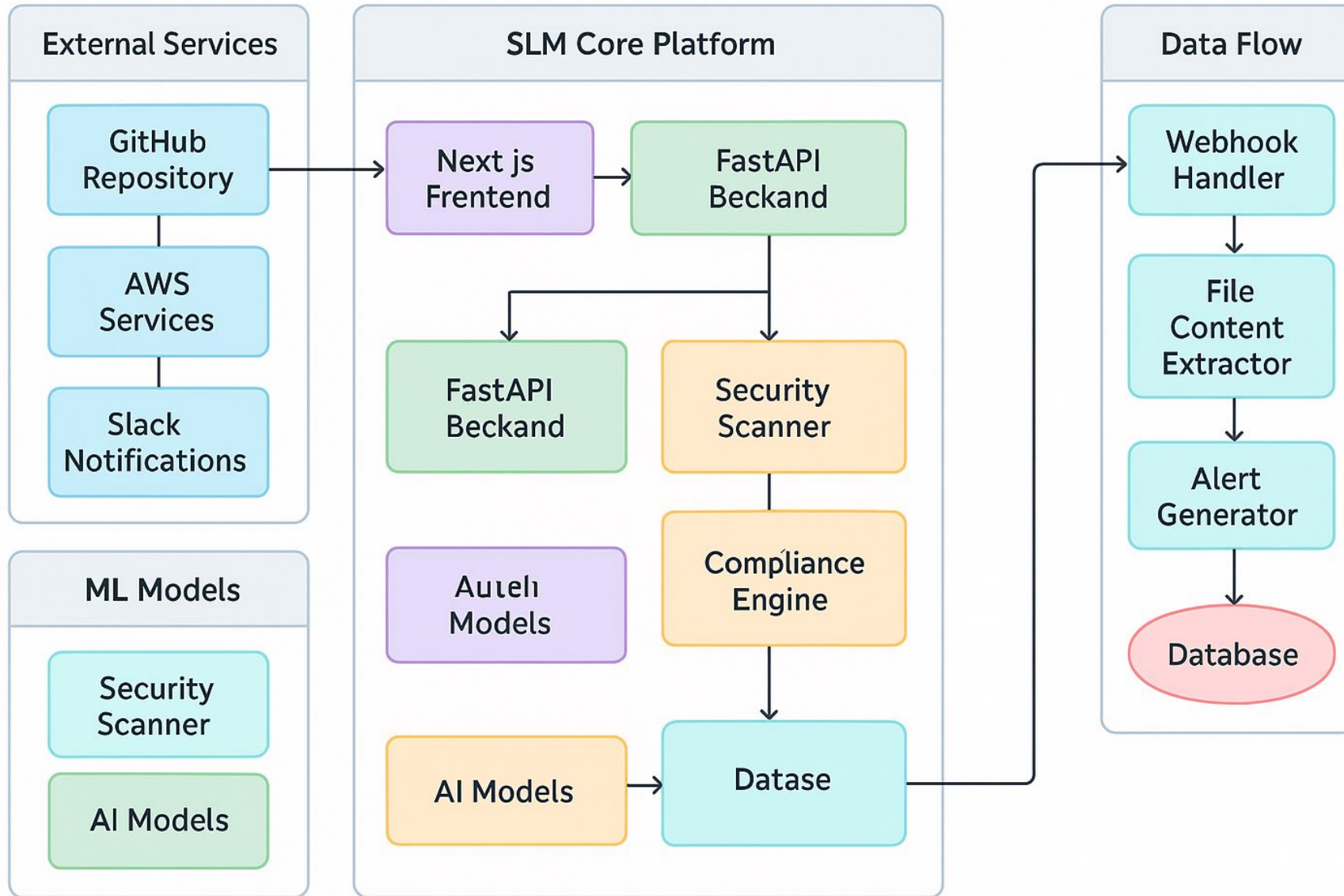
**Fine-tuned Qwen 3-1.7B on AWS Trainium**
Leveraged AWS Trainium to efficiently fine-tune the large language model in 2 hours

**Created structured output format**
Designed a structured output format for audit compliance: CLASSIFICATION/SECRET_TYPE/SEVERITY/REMEDIATION

Our fine-tuned SLM model for SOC-2 secret detection has demonstrated significant improvements in speed, cost, and output quality. By leveraging AWS Trainium and a structured output format, we have delivered a reliable and efficient solution for real-time scanning and audit compliance..

# Side-by-Side Comparison

| Metric | Base Qwen 3-1.7B | Fine-Tuned LeakGuard |
|---|---|---|
| Output Length | 500-800 tokens | ~50 tokens |
| Response Style | Rambling, uncertain | Structured, definitive |
| Format | Unstructured paragraphs | CLASSIFICATION/SECRET_TYPE/ SEVERITY/REMEDIATION |
| Consistency | Varies per query | Identical format every time |
| Inference Speed | Baseline | 16x faster |
| Cost per 1M scans | $8,000 | $500 |
| SOC-2 Ready | ✖ | ✓ |

*Data provided in the original prompt

# Impact & Metrics

| Metric | Value |
| --- | --- |
| Dataset Size | 1,000 examples |
| Cost | **94% cost reduction ($8K → $500 per 1M scans)** |
| Model Size | 1.7B base + 16M LoRA params |
| Latency | ~100ms per scan |
| Throughput | 126 tokens/sec input, 54 tokens/sec output |
| Speed Improvement | 16x faster than base model |
| | |
| | |

# Team

# Future Vision - Roadmap

**Phase 1 (Today)**
Secret Detection Core - Fine-tuned SLM + API endpoint

**Phase 2 (Next)**
CI/CD Integration - GitHub Actions pre-commit hooks, Automatic PR blocking, Slack/email alerts

**Phase 3**
Automated Rotation - AWS Secrets Manager integration, Auto-rotate exposed credentials, Redeploy affected services

**Phase 4**
SOC-2 Evidence Automation - CloudTrail log collection, Vanta/Drata integration, Compliance dashboard