# PII Detection: Performance Analysis
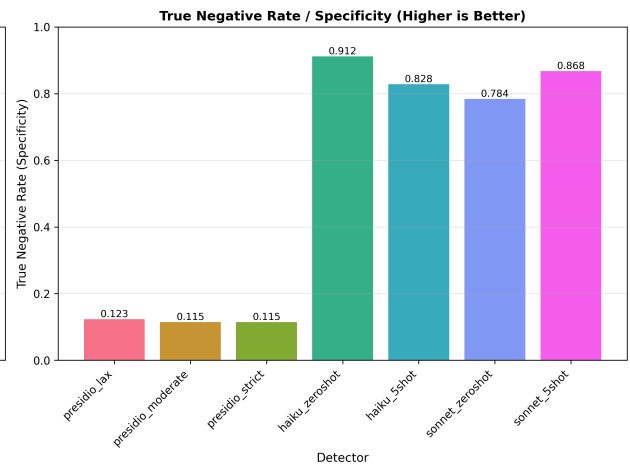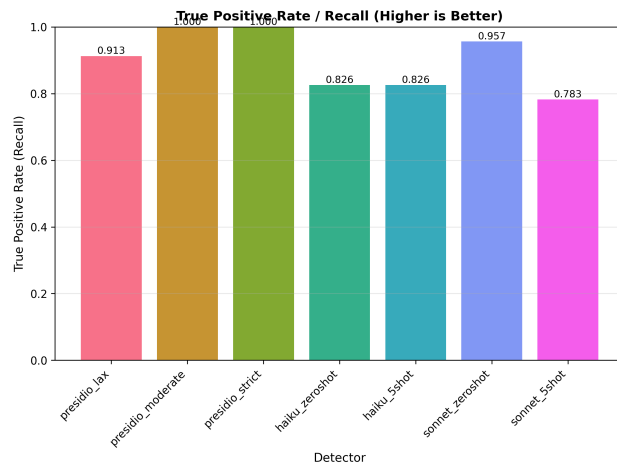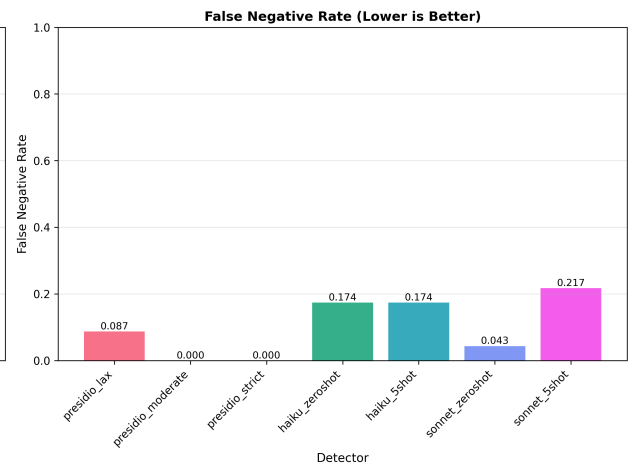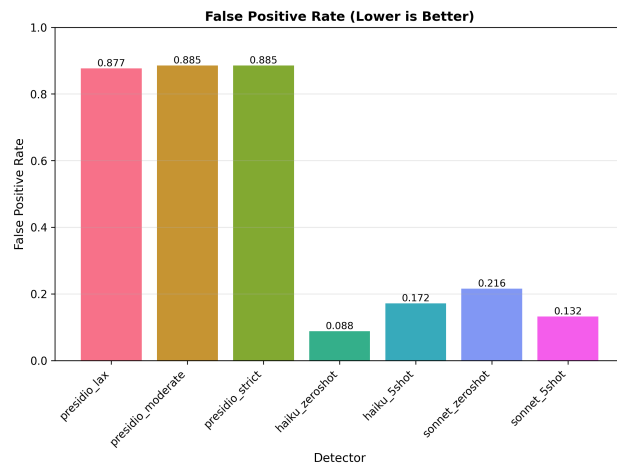
*November 2025*

## Executive Summary

Email breach analysis evaluated on 250-sample gold-labeled dataset. **Primary goal: Zero false negatives to avoid missing PII.** Finding: Presidio achieves 0% FNR but 88% FPR (unusable in production). Sonnet zero-shot delivers 4.3% FNR with acceptable 22% FPR.

## Detector Overview

| Detector | Type | Description |
|---|---|---|
| Presidio | Local/Free | Rule-based NLP, runs locally |
| Haiku Zero-shot | Claude API | Fast LLM, no examples |
| Sonnet Zero-shot | Claude API | Slower LLM, higher accuracy |

## Performance Metrics

| Detector | FNR ↓ | FPR ↓ | Status |
|---|---|---|---|
| Presidio Moderate | 0.0% | 88.5% | ■■ Unusable |
| Haiku Zero-shot | 17.4% | 8.8% | ✓ Good |
| Sonnet Zero-shot | 4.3% | 21.6% | ✓ Best usable |

**False Positive Rate (Lower is Better)**

| Detector | Value |
|---|---|
| presidio_lax | 0.877 |
| presidio_moderate | 0.885 |
| presidio_strict | 0.885 |
| haiku_zeroshot | 0.088 |
| haiku_5shot | 0.172 |
| sonnet_zeroshot | 0.216 |
| sonnet_5shot | 0.132 |

**False Negative Rate (Lower is Better)**

| Detector | Value |
|---|---|
| presidio_lax | 0.087 |
| presidio_moderate | 0.000 |
| presidio_strict | 0.000 |
| haiku_zeroshot | 0.174 |
| haiku_5shot | 0.174 |
| sonnet_zeroshot | 0.043 |
| sonnet_5shot | 0.217 |

**True Positive Rate / Recall (Higher is Better)**

| Detector | Value |
|---|---|
| presidio_lax | 0.913 |
| presidio_moderate | 1.000 |
| presidio_strict | 1.000 |
| haiku_zeroshot | 0.826 |
| haiku_5shot | 0.826 |
| sonnet_zeroshot | 0.957 |
| sonnet_5shot | 0.783 |

**True Negative Rate / Specificity (Higher is Better)**

| Detector | Value |
|---|---|
| presidio_lax | 0.123 |
| presidio_moderate | 0.115 |
| presidio_strict | 0.115 |
| haiku_zeroshot | 0.912 |
| haiku_5shot | 0.828 |
| sonnet_zeroshot | 0.784 |
| sonnet_5shot | 0.868 |

## 50k Email Projections

| Detector | Time (Current) | Time (Parallel) | Cost | FNR |
|---|---|---|---|---|
| Presidio Moderate | 19.7 min | 2-3 min* | $0.00 | 0.0% |
| Haiku Zero-shot | 741.4 min | 15-74 min** | $28.48 | 17.4% |
| Sonnet Zero-shot | 1724.7 min | 35-173 min** | $106.78 | 4.3% |

*With EC2 multi-core optimization

**With API parallelization (10-50x speedup; currently sequential)

## Next Steps

• **1. Implement API parallelization** - 10-50x speedup potential (1-2 days development)

• **2. EC2 optimization for Presidio** - 5-10x speedup with multi-core processing

• **3. Tune Presidio to reduce FPR** - Goal: maintain 0% FNR while reducing 88% FPR

**Recommendation:** Deploy Sonnet zero-shot (4.3% FNR) with API parallelization for production 50k emails. Cost: $107, Time: ~35-173 min (parallelized). Alternative: Haiku zero-shot if 17.4% FNR is acceptable ($28, ~15-74 min).