

Fraud Communication Guard: System Design Document

Document Version: 1.0

Date: 2025-11-14

Classification: Internal Security Architecture - Defensive Research Prototype

Author: Principal Security Architect, AI Systems Engineering Team

PHASE 0 - BOUNDARIES & ETHICS

Framework Statement: Defensive Security Research

This system design document describes a **defensive security research prototype** developed for the purpose of **authorized fraud investigation, user protection, and legal evidence collection**. The Fraud Communication Guard (FCG) is designed to protect organizations and their users from fraudulent actors, particularly in scenarios involving fake loan brokers, phishing campaigns, account takeover attempts, and sophisticated social engineering attacks.

Primary Use Case: Defensive investigation of active fraud attempts targeting the organization's users, conducted with proper legal authorization and in coordination with legal and compliance teams.

Legal and Ethical Constraints

All implementations and deployments of this system must operate within the following boundaries:

- No Offensive Operations:** This system is strictly for defensive security and investigation. It must not be used for offensive hacking, unauthorized access to third-party systems, or any activities that would constitute computer intrusion under CFAA or equivalent laws.
- Privacy Compliance:** All data collection and processing must comply with:
 - **GDPR** (General Data Protection Regulation) for European data subjects
 - **CCPA** (California Consumer Privacy Act) for California residents
 - **ECPA** (Electronic Communications Privacy Act) for U.S. communications
 - Applicable state and federal privacy statutes
- Evidence Collection Focus:** The system's primary purpose is to collect, preserve, and analyze evidence in a manner that maintains chain of custody and ensures admissibility in legal proceedings.
- Lawful Basis Requirements:**
 - User consent where required by regulation
 - Legitimate interest justification with documented balancing test
 - Legal process (subpoenas, warrants) for accessing protected data
 - Continuous documentation of legal authorization for all collection activities
- Data Minimization:** Collect only data necessary for fraud detection and investigation. Implement automatic purging of non-relevant data within defined retention periods.

6. **Transparency:** Users must be informed of monitoring activities through clear privacy policies, except where notification would compromise an active investigation and legal counsel has provided written authorization for covert collection.

Attacker TTP Discussion - Detection and Mitigation Only

Throughout this document, we discuss various attacker tactics, techniques, and procedures (TTPs) including:

- Device and browser fingerprinting evasion
- VPN and proxy usage for anonymization
- Synthetic identity creation
- Bot-driven fraud automation
- Communication channel exploitation

Critical Clarification: Discussion of these TTPs is **solely for the purpose of detection, prevention, and mitigation**. Understanding how attackers operate enables the development of defensive countermeasures. These techniques should never be deployed offensively or used to compromise legitimate users' privacy or security.

PHASE 1 - SYSTEM OVERVIEW

1.1 System Name and Core Capabilities

System Name: Fraud Communication Guard (FCG)

Core Capabilities:

- **Multi-Channel Communication Interception and Analysis:** Monitor incoming communications (email, SMS, web chat, social media) for fraud indicators before they reach end-users, providing real-time threat assessment and safe-preview capabilities.
- **Persistent Attacker Attribution:** Create durable identifiers for fraudulent actors using device fingerprinting, network attribution, and behavioral profiling to track returning attackers across sessions, IP changes, and identity masking attempts.
- **OSINT-Driven Intelligence Enrichment:** Automatically aggregate and correlate open-source intelligence about suspected fraudsters, including social media profiles, domain ownership, phone/email associations, and historical fraud patterns, to build comprehensive "attacker knowledge cards."
- **Forensically Sound Evidence Collection:** Capture all investigative artifacts with cryptographic integrity verification, timestamping, and chain-of-custody logging to ensure evidence admissibility in legal proceedings.
- **AI-Powered Risk Scoring and Triage:** Leverage machine learning models to classify communications and actors into benign/suspicious/fraud categories with explainability, reducing false positives and enabling intelligent routing to human analysts.

1.2 Primary User Roles

Security Analyst

- Primary operator of the FCG dashboard
- Reviews flagged communications and attacker profiles

- Conducts deep-dive investigations into high-risk entities
- Configures detection rules and triage thresholds
- Initiates OSINT enrichment workflows

Legal Team

- Reviews evidence packages for completeness and admissibility
- Authorizes covert monitoring where legally permissible
- Validates chain of custody documentation
- Requests formatted reports (PDF, CSV, timeline) for litigation
- Ensures ongoing compliance with ECPA, GDPR, CCPA

Investigator (Forensic Specialist)

- Performs deep forensic analysis on captured communications and device data
- Validates device fingerprints and network attribution
- Conducts timeline analysis and cross-session correlation
- Expert witness preparation and testimony support

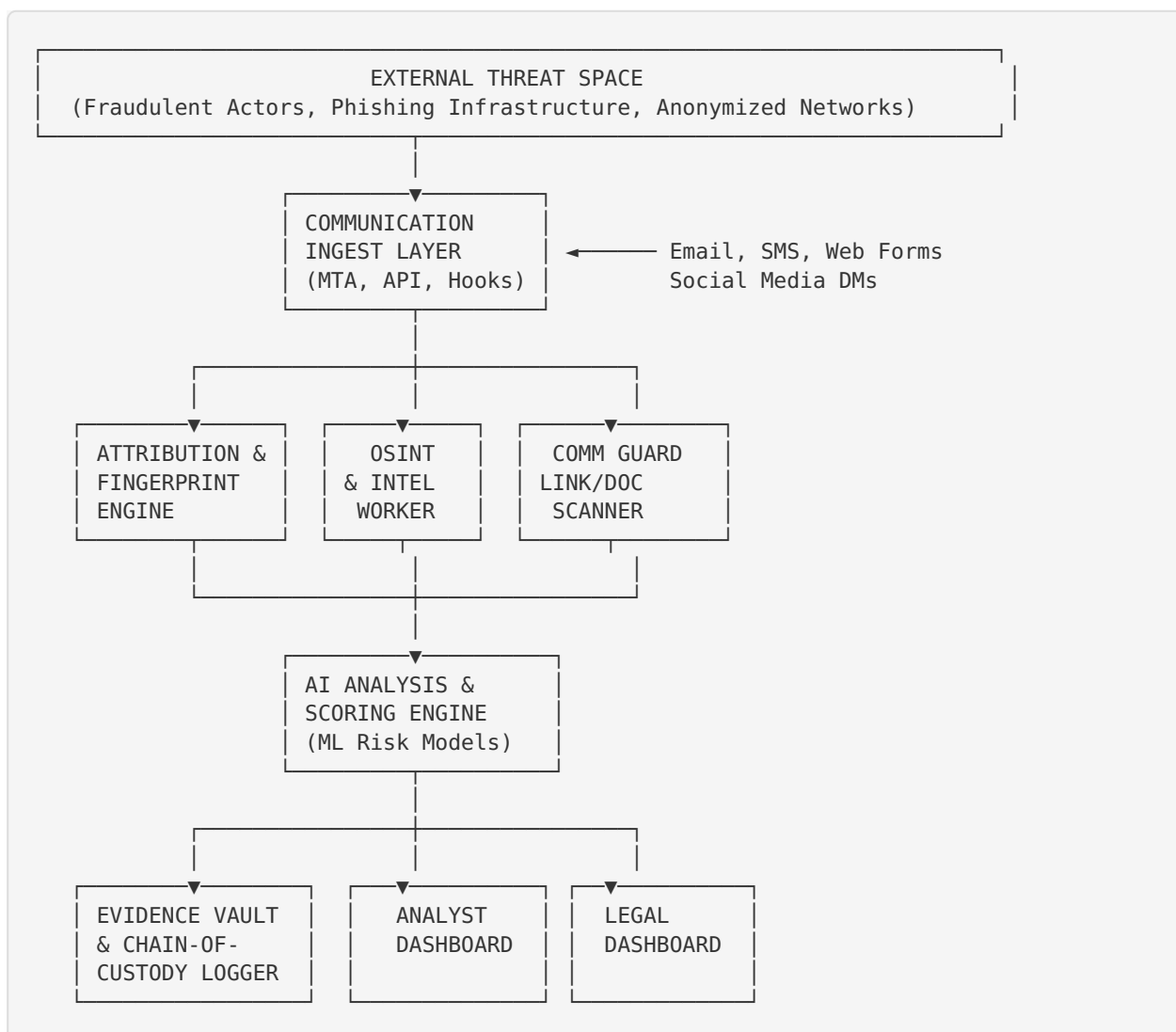
End-User (Protected Party)

- Receives safe-preview links for suspicious communications
- Reports suspicious contacts through integrated interface
- Views sanitized threat intelligence ("This sender has been reported X times")
- In transparent deployment mode: receives privacy notices and consent management

1.3 High-Level Architecture

The Fraud Communication Guard employs a **multi-layered defensive architecture** organized around six core components, each designed with security, scalability, and legal compliance as foundational principles.

Architecture Diagram Description



Component Data Flows

1. Communication Ingest → Attribution Engine

- Raw communication metadata (sender IP, device headers, timestamps)
- Creates device fingerprint, resolves geolocation, detects VPN/proxy
- Output: Persistent attacker identifier, enriched network context

2. Communication Ingest → OSINT Worker

- Email addresses, phone numbers, domain names, social media handles
- Performs reverse lookups, domain WHOIS, social media profiling
- Output: "Attacker Knowledge Card" with identity graph and risk indicators

3. Communication Ingest → Communication Guard

- Message content, embedded URLs, attachments
- Scans links for redirects/phishing, sandboxes documents, checks reputation
- Output: Safe/unsafe/suspicious verdict with IOCs (Indicators of Compromise)

4. All Components → AI Analysis Engine

- Aggregated data from attribution, OSINT, and content scanning
- ML models classify risk, detect anomalies, flag behavioral patterns
- Output: Risk score (0-100), classification label, explainability report

5. AI Analysis Engine → Evidence Vault

- All investigative artifacts, model decisions, and human annotations
- Cryptographically hashed, timestamped, chain-of-custody logged
- Output: Immutable evidence packages, audit trails

6. Evidence Vault → Dashboards

- Analyst Dashboard: Real-time alerts, investigation workspace, OSINT tools
- Legal Dashboard: Evidence review, compliance reports, chain-of-custody exports

Trust Boundaries

External Trust Boundary: Separates the external threat space (attacker infrastructure) from the Communication Ingest Layer. All external data is treated as untrusted and subjected to sanitization and validation.

Processing Trust Boundary: Between the Analysis/Collection components and the Evidence Vault. Only cryptographically verified, chain-of-custody logged data crosses this boundary.

Human Interface Trust Boundary: Between automated systems and the analyst/legal dashboards. Implements role-based access control (RBAC), multi-factor authentication (MFA), and audit logging of all user actions.

Privacy Controls

- **PII Separation:** Personally identifiable information is stored in segregated databases with encrypted-at-rest storage and least-privilege access policies.
- **Data Minimization Engine:** Automatically redacts non-relevant PII from analyst views while preserving it in the encrypted Evidence Vault for legal review.
- **Retention Policy Enforcement:** Automated purging of data after configurable retention periods (default: 90 days for non-case-related data, indefinite for active cases with legal hold).
- **Consent Management:** Tracks user consent status and automatically blocks processing for users who have opted out (where legally required).
- **Differential Privacy:** Where possible, aggregated analytics use differential privacy techniques to protect individual identities while enabling threat pattern detection.

PHASE 2 - DETAILED MODULES

2.1 Device & Network Attribution Module

Purpose and Responsibilities

The Device & Network Attribution Module serves as the foundational identity layer for persistent attacker tracking. Its core responsibility is to establish a **stable, privacy-conscious identifier** for each entity interacting with the system, enabling investigators to link multiple fraudulent activities to a single actor even when that actor employs evasion techniques such as IP rotation, cookie clearing, or browser switching.

Success Criteria:

- 85%+ accuracy in re-identifying returning attackers across sessions
- Detection of 95%+ of VPN/proxy/Tor usage with <2% false positive rate
- Device fingerprint collision rate <0.1% (high uniqueness)
- Sub-500ms latency for real-time attribution during communication ingestion

Inputs and Outputs

Inputs:

- HTTP headers (User-Agent, Accept-Language, Accept-Encoding, DNT, etc.)
- Client IP address and connection metadata
- Browser API telemetry (Canvas, WebGL, Audio, Font enumeration)
- TLS fingerprinting data (cipher suites, extensions)
- Behavioral timing data (if available from interactive sessions)

Outputs:

- **Persistent Device ID (PDID)**: SHA-256 hash of combined fingerprint vector
- **Attribution Report**: JSON structure containing:

```
json
{
  "pdid": "a7f3c9d2...",
  "ip_address": "203.0.113.45",
  "asn": "AS15169",
  "isp": "Google LLC",
  "geolocation": {"country": "US", "city": "Mountain View", "lat": 37.4, "lon": -122.08},
  "timezone_inferred": "America/Los_Angeles",
  "vpn_detected": true,
  "vpn_provider": "NordVPN",
  "device_type": "Desktop",
  "os": "Windows 11",
  "browser": "Chrome 120.0.6099",
  "fingerprint_confidence": 0.92,
  "first_seen": "2025-11-14T10:23:15Z",
  "last_seen": "2025-11-14T11:45:32Z",
  "session_count": 7
}
```

- **Attacker Timeline**: Chronological log of all sessions associated with PDID

Recommended Tools and Technologies

Free/Open-Source Options:

1. **FingerprintJS (OSS Edition)** - JavaScript library for browser fingerprinting
 - Collects 50+ entropy sources including Canvas, WebGL, AudioContext
 - Apache 2.0 license, extensive documentation
 - Integration: Deploy as client-side script or server-side via headless browser
2. **MaxMind GeoIP2 (Free GeoLite2)** - IP geolocation database
 - City-level accuracy for most regions
 - Free weekly database updates
 - Integration: Local database query via API
3. **IPQualityScore Free Tier** - VPN/Proxy detection API
 - 5,000 free lookups/month
 - Detects commercial VPNs, Tor, data center IPs
 - Integration: REST API
4. **p0f (Passive OS Fingerprinting)** - Network-level OS detection
 - Analyzes TCP/IP stack behavior

- No active probing required
- Integration: Packet capture analysis

Commercial API-First Options (for production):

1. **SEON Device Intelligence API** - Comprehensive device fingerprinting with fraud scoring
2. **Fingerprint.com** - 99.5% accuracy, server-side fingerprinting
3. **IPHub.info** - Advanced VPN/proxy detection with risk scoring

Implementation Phases

< 1 Hour Prototype Implementation:

1. **Client-Side Script Deployment** (15 min)
 - Deploy FingerprintJS via CDN to landing page or form
 - Configure to collect: Canvas, WebGL, Fonts, Screen, Timezone
 - POST fingerprint data to backend endpoint
2. **Basic IP Attribution** (20 min)
 - Download GeoLite2 databases (City + ASN)
 - Implement IP lookup function using MaxMind reader library
 - Log IP, Country, City, ASN, ISP to database
3. **Simple Hash Generation** (10 min)
 - Concatenate fingerprint components + User-Agent
 - Generate SHA-256 hash as PDID
 - Store mapping: PDID → [IP, Timestamp, Session Data]
4. **VPN Detection Integration** (15 min)
 - Sign up for IPQualityScore free tier
 - Implement API call with IP address
 - Flag sessions where `vpn: true` or `proxy: true`

Phase 2+ Enhancements:

- **Probabilistic Fingerprint Matching:** Implement fuzzy hashing (e.g., simhash) to detect “near-match” fingerprints from users who have made minor system changes
- **Behavioral Biometrics:** Add keystroke dynamics and mouse movement analysis for interactive fraud attempts
- **TLS Fingerprinting:** Integrate JA3/JA4 fingerprinting to identify bot frameworks and headless browsers
- **Cross-Device Linking:** Build identity graph using shared IPs (home networks) and deterministic identifiers (logged-in accounts)
- **Fingerprint Entropy Analysis:** Continuously monitor fingerprint uniqueness distribution to detect spoofing attacks

Logging Format and Chain-of-Custody

All attribution data must be logged in a **forensically sound format** with the following structure:

```
{
  "log_version": "1.0",
  "event_type": "device_attribution",
  "event_id": "uuid-v4",
  "timestamp_utc": "2025-11-14T15:30:45.123Z",
  "hash_sha256": "hash of this entire JSON minus this field",
  "collector_id": "fcg-attrib-node-01",
  "subject": {
    "pdid": "a7f3c9d2...",
    "pdid_version": 2
  },
  "raw_data": {
    "ip_address": "203.0.113.45",
    "headers": {"User-Agent": "...", "..."},
    "fingerprint": {"canvas": "...", "webgl": "...", "..."}
  },
  "enrichment": {
    "geolocation": {...},
    "asn": {...},
    "vpn_detection": {...}
  },
  "chain_of_custody": {
    "collected_by": "system",
    "verified_by": null,
    "accessed_by": [],
    "legal_hold": false
  }
}
```

Chain-of-Custody Requirements:

- Each log entry is hashed (SHA-256) upon creation
- Logs are written to append-only storage (S3 Object Lock, WORM media, or blockchain-backed log)
- Any access to attribution data is logged with analyst ID, timestamp, and reason code
- Periodic integrity verification: re-compute hashes and compare against stored values

2.2 Identity & OSINT Intelligence Module

Purpose and Responsibilities

The Identity & OSINT Intelligence Module automates the labor-intensive process of **background investigation** and **digital footprint mapping** for suspected fraudsters. By aggregating data from public sources—social media, public records, data breach databases, domain registrations—this module constructs a comprehensive “Attacker Knowledge Card” that provides investigators with context, behavioral patterns, and actionable leads for further investigation.

Success Criteria:

- Automated enrichment of 80%+ of inbound fraudulent identities within 5 minutes
- Discovery of 3+ alternate identities/aliases per investigated subject on average
- Integration of 10+ OSINT data sources with unified query interface
- Generation of visual identity graphs showing relationships between entities

Inputs and Outputs

Inputs:

- Email addresses from fraudulent communications
- Phone numbers (mobile, VoIP, landline)

- Social media handles/usernames
- Domain names from phishing sites
- Real names (if disclosed or discovered)
- Physical/IP addresses

Outputs:

- **Attacker Knowledge Card:** Structured profile containing:

```
json
{
  "primary_identifier": "fraudster@example.com",
  "confidence_score": 0.87,
  "identity_graph": {
    "email_addresses": ["fraudster@example.com", "backup@mail.ru"],
    "phone_numbers": ["+1-555-0123 (VoIP, Twilio)"],
    "social_media": {
      "twitter": "@fraud_handle",
      "linkedin": "linkedin.com/in/fake-profile",
      "facebook": null
    },
  },
  "domains_owned": ["fakeloan.com", "quickcash.net"],
  "aliases": ["John Smith", "Ivan Petrov"],
  "locations": ["Nigeria (VPN)", "Russia (phone carrier)"]
},
"risk_indicators": {
  "data_breach_exposure": ["Collection #1", "Dubsmash 2019"],
  "domain_age": "14 days (fakeloan.com)",
  "social_media_authenticity": "Low (1 follower, stock photo)",
  "known_fraud_associations": ["Same email in fraud report DB"]
},
"timeline": [
  {"date": "2025-11-01", "event": "Domain fakeloan.com registered"},
  {"date": "2025-11-10", "event": "First phishing email sent"},
  {"date": "2025-11-14", "event": "Identified in investigation"}
]
}
```

OSINT Frameworks and Tools

Integrated OSINT Platforms:

1. SpiderFoot HX (Community Edition)

- Automated OSINT collection and correlation
- 200+ data sources including DNS, WHOIS, social media, breach data
- Graph visualization of entity relationships
- Integration: REST API or Python library

2. theHarvester

- Email, subdomain, and employee enumeration
- Queries search engines, PGP servers, Shodan
- Command-line tool with JSON output
- Integration: Subprocess call from orchestration layer

3. Maltego Community Edition

- Visual link analysis and graph-based investigation
- Transforms for social media, DNS, WHOIS, email
- Manual investigation tool for analysts
- Integration: Standalone analyst workstation

4. MISP (Malware Information Sharing Platform) / OpenCTI

- Threat intelligence aggregation and sharing
- Stores IOCs, TTPs, actor profiles
- Integration: Import FCG findings, export to SOC/SIEM

Specific OSINT Tools by Data Type:

Email Intelligence:

- **Hunter.io** (Free: 25 searches/month): Email finder and verifier
- **EmailRep.io** (Free API): Email reputation scoring based on breach data
- **Epieos** (Free): Google Account OSINT (profile photos, reviews, YouTube)

Phone Intelligence:

- **Phoneinfoga** (Open-source): Phone number OSINT with carrier/country lookup
- **Truecaller API** (Commercial): Crowdsourced caller ID database
- **NumVerify** (Free tier): Phone number validation and carrier lookup

Domain Intelligence:

- **WhoisXML API** (Free tier): WHOIS lookup, reverse WHOIS, historical records
- **URLScan.io** (Free): Automated website scanning and screenshot capture
- **Certificate Transparency Logs** (crt.sh): Discover subdomains via SSL certificates

Social Media Intelligence:

- **Sherlock** (Open-source): Username search across 300+ platforms
- **Social-Analyzer** (Open-source): Profile analysis and activity extraction
- **Twint** (Open-source): Twitter scraping without API

Data Breach Intelligence:

- **Have I Been Pwned API** (Free for personal emails): Breach exposure check
- **DeHashed** (Commercial): Search engine for leaked credentials
- **IntelX** (Commercial): Deep web and data leak search

Implementation Phases

< 1 Hour Prototype Implementation:

1. Email OSINT Workflow (20 min)

- Implement Hunter.io API call for email → domain/name lookup
- Query Have I Been Pwned for breach exposure
- Use Epieos to check for Google Account associations
- Output: JSON profile with email context

2. Phone Number Lookup (15 min)

- Deploy Phoneinfoga as Docker container
- Create REST API wrapper for phone lookup
- Query NumVerify for validation and carrier
- Output: Phone number metadata and risk flags

3. Domain Intelligence (15 min)

- WHOIS query using Python `whois` library or WhoisXML API
- Check domain age (newly registered = high risk)
- Query URLScan.io for screenshots and HTML content
- Output: Domain reputation report

4. Basic Knowledge Card Assembly (10 min)

- Aggregate all OSINT results into unified JSON structure
- Calculate basic confidence score (presence of multiple identifiers)
- Store in Evidence Vault with timestamp and source attribution

Phase 2+ Enhancements:

- **SpiderFoot Automation:** Deploy SpiderFoot HX server and automate scans on new identifiers
- **Graph Database Integration:** Store identity relationships in Neo4j for complex link analysis
- **Continuous Monitoring:** Set up alerts for new mentions of known fraudsters on paste sites, forums
- **Dark Web Monitoring:** Integrate tools like OnionScan or commercial services (Flare, ZeroFox)
- **AI-Powered Profiling:** Use NLP to analyze social media posts for behavioral patterns and deception indicators
- **Collaborative Intelligence:** Integrate with MISP to share/receive fraud actor IOCs with industry peers

“Attacker Knowledge Card” Design

The Knowledge Card is the primary investigative artifact presented to analysts. It should be designed as a **single-page visual summary** with expandable sections:

Layout:

```

ATTACKER KNOWLEDGE CARD
Primary ID: fraudster@example.com | Risk: HIGH (87)

IDENTITY GRAPH
[Visual node-link diagram]
• Central node: email
• Connected: phone, domains, social profiles

RISK INDICATORS
🔴 Domain age: 14 days
🔴 Exposed in 2 data breaches
🔴 VoIP phone number
🟡 Social media: Low authenticity

TIMELINE
[Horizontal timeline with key events]

OSINT SOURCES (12)
[Expandable list with links to raw data]

ACTIONS: [Export] [Add to Watchlist] [Deep Scan]

```

2.3 Communication Guard & Link/Document Analysis Module

Purpose and Responsibilities

This module acts as the **first line of defense** between fraudulent communications and end-users. Its purpose is to intercept and analyze all inbound messages, links, and attachments before user interaction, providing real-time threat assessment and enabling “safe preview” workflows that protect users from malicious content while preserving evidence of fraud attempts.

Success Criteria:

- Intercept and analyze 100% of inbound communications within defined channels
- Phishing URL detection accuracy >98% with <0.5% false positive rate
- Malicious document detection >95% (malware, embedded exploits)
- Sub-3-second analysis time for real-time user experience
- Zero data loss: All communications logged even if analysis fails

Inputs and Outputs

Inputs:

- **Email Messages:** Full MIME, headers, body, attachments
- **SMS/Text Messages:** Content, sender number, carrier metadata
- **Web Forms:** User-submitted content flagged as suspicious
- **Social Media DMs:** Messages via API or forwarded by users
- **Embedded Content:** URLs, image links, document attachments (PDF, DOCX, XLS)

Outputs:

- **Triage Verdict:** SAFE | SUSPICIOUS | MALICIOUS
- **Threat Report:** JSON structure:

```
json
{
  "communication_id": "uuid",
  "verdict": "MALICIOUS",
  "confidence": 0.96,
  "threats_detected": [
    {
      "type": "phishing_url",
      "ioc": "http://fakeloan-secure.com/login",
      "analysis": {
        "domain_age": "7 days",
        "ssl_cert": "self-signed",
        "redirect_chain": ["bit.ly/xyz", "t.co/abc", "fakeloan-secure.com"],
        "brand_impersonation": "Chase Bank",
        "similarity_score": 0.91
      }
    },
    {
      "type": "malicious_attachment",
      "filename": "loan_application.pdf",
      "file_hash": "d2d2d2...",
      "analysis": {
        "sandbox_result": "exploit_detected",
        "malware_family": "CVE-2023-12345 exploit",
        "behavior": "attempts to execute JavaScript, drops file"
      }
    }
  ]
}
```

```

    }
  }
],
"safe_preview_url": "https://i.ytimg.com/vi/B4Lqgc9UXHE/maxresdefault.jpg",
"recommended_action": "BLOCK"
}

```

Recommended Tools and Technologies

URL and Link Analysis:

1. **URLScan.io** (Free API)
 - Automated website scanning, screenshots, DOM analysis
 - Detects redirects, suspicious JavaScript, credential harvesting forms
 - Integration: REST API
2. **VirusTotal** (Free API: 500 requests/day)
 - Multi-engine URL/file reputation scanning
 - 70+ antivirus engines and threat intelligence feeds
 - Integration: REST API
3. **PhishTank** (Free API)
 - Community-driven phishing URL database
 - Real-time phishing verification
 - Integration: REST API
4. **Custom Redirect Chain Analyzer**
 - Python script using `requests` with redirect tracking
 - Follows 302/301 redirects, logs each hop
 - Flags suspicious patterns (multiple shorteners, domain changes)

Document and Attachment Analysis:

1. **Cuckoo Sandbox** (Open-source)
 - Automated malware analysis in isolated environment
 - Behavioral analysis of PDF, Office docs, executables
 - Integration: REST API, Python library
2. **PeePDF** (Open-source)
 - PDF structure analysis and JavaScript extraction
 - Detects malformed PDFs, embedded exploits
 - Integration: Command-line tool
3. **YARA Rules** (Open-source)
 - Pattern matching for malware signatures
 - Community rule repositories for common exploits
 - Integration: Python library `yara-python`
4. **oletools** (Open-source)
 - Analyze Microsoft Office documents for macros, embedded objects
 - Detects DDE attacks, VBA macros
 - Integration: Command-line tools

Email Security and Analysis:

1. **SpamAssassin** (Open-source)
 - Email content analysis with scoring
 - Bayesian filtering, rule-based detection
 - Integration: As MTA filter or standalone
2. **DMARC/SPF/DKIM Validator**
 - Verify email authentication to detect spoofing
 - Python library `checkdmARC`
 - Integration: Part of email ingestion pipeline

Sandbox and Safe Preview:

1. **Browserless.io** (Open-source / Commercial)
 - Headless Chrome as a service
 - Render suspicious pages in isolated environment
 - Integration: Docker container, REST API
2. **Any.Run** (Commercial, free tier)
 - Interactive malware analysis sandbox
 - Real-time observation of malicious behavior
 - Integration: Manual analyst tool

Triage Model: Safe/Unsafe/Suspicious Classification

Classification Logic:

SAFE (Green Light - Deliver to User)

- No URLs detected, or all URLs resolve to known-safe domains (whitelist)
- No attachments, or attachments are plain text
- Email passes SPF/DKIM/DMARC
- Sender domain is established (>1 year old) and has positive reputation
- Content analysis: No urgency language, no credential requests

SUSPICIOUS (Yellow Light - Safe Preview / Analyst Review)

- Contains URLs to newly registered domains (<90 days)
- Email authentication partially fails (e.g., SPF pass, DMARC fail)
- Contains attachments but sandbox analysis is inconclusive
- Urgency language detected ("act now", "verify account")
- Sender not previously seen in organization communications

MALICIOUS (Red Light - Block / Quarantine)

- URL matches known phishing database (PhishTank, VirusTotal)
- Attachment triggers malware signatures or exploit behavior in sandbox
- Email fails all authentication (SPF/DKIM/DMARC all fail)
- Redirect chain matches fraud patterns (3+ redirects, suspicious TLDs)
- Brand impersonation detected (domain similarity >85% to protected brand)

Implementation Phases

< 1 Hour Prototype Implementation:

1. **Email Ingestion Webhook** (15 min)
 - Set up email forwarding rule to route suspected phishing to FCG endpoint

- Parse MIME using Python `email` library
- Extract sender, subject, body, URLs, attachments
- Store raw email in Evidence Vault

2. **Basic URL Scanning** (20 min)

- Extract all URLs from email body
- Query VirusTotal API for URL reputation
- Check domain age using WhoisXML API
- Flag URLs with <90 day age or VirusTotal detections >0

3. **Simple Attachment Hashing** (10 min)

- Calculate SHA-256 of all attachments
- Query VirusTotal API for file hash reputation
- Flag files with any AV detections

4. **Triage Logic and Alerting** (15 min)

- Implement simple rule-based classifier
- If URL flagged OR attachment flagged: verdict = SUSPICIOUS
- If VirusTotal detections >5: verdict = MALICIOUS
- Send alert to Slack/email for analyst review

Phase 2+ Enhancements:

- **Full Sandbox Integration:** Deploy Cuckoo Sandbox for dynamic attachment analysis
- **Redirect Chain Analysis:** Implement recursive URL following with pattern detection
- **Brand Impersonation Detection:** Build model to detect visual similarity to legitimate sites
- **NLP Content Analysis:** Use LLM to analyze message content for social engineering tactics
- **Safe Preview Portal:** Deploy isolated browser environment for users to view suspicious content
- **User Feedback Loop:** Allow users to report false positives/negatives to retrain model

Safe Preview and Sandbox Patterns

Safe Preview Architecture:

When a communication is classified as SUSPICIOUS, the system generates a safe preview link instead of delivering the original content:

1. **User Receives Alert:** "A message from sender@unknown.com has been flagged as potentially suspicious. [View Safe Preview]"
 2. **Isolated Rendering:** Clicking the link opens a session in an isolated browser (Browserless container) that:
 - Renders the content server-side
 - Removes all JavaScript execution
 - Blocks external resource loading
 - Captures and displays a static screenshot to the user
 3. **User Actions:**
 - **"This is legitimate":** Message is released to inbox, sender whitelisted
 - **"This is fraud":** Message quarantined, sender added to blocklist, analyst notified
 - **"I'm not sure":** Escalated to analyst for manual review
 4. **Evidence Preservation:** All user interactions are logged for investigation timeline
-

2.4 Evidence Vault & Legal Workflow Module

Purpose and Responsibilities

The Evidence Vault is the **system of record** for all investigative artifacts, designed to meet the stringent requirements of legal evidence handling. It ensures that every piece of data collected maintains a verifiable chain of custody, cryptographic integrity, and compliant access controls, transforming raw investigative data into admissible legal evidence.

Success Criteria:

- 100% of collected artifacts logged with immutable timestamps and hashes
- Zero-tolerance for chain-of-custody breaks (automated integrity verification)
- Sub-1-second retrieval time for evidence packages (even for cold storage)
- Compliance with NIST SP 800-86 digital forensics guidelines
- Successful admissibility in legal proceedings (target: 100% acceptance rate)

Inputs and Outputs

Inputs:

- All investigative artifacts from other modules:
- Attribution reports (device fingerprints, IP logs)
- OSINT data (knowledge cards, raw source data)
- Communication scans (emails, threat reports, sandbox results)
- AI analysis (risk scores, classifications, model explanations)
- Analyst annotations (case notes, tagging, decisions)

Outputs:

- **Evidence Packages:** Cryptographically sealed collections of related artifacts
- **Chain of Custody Reports:** Detailed provenance documentation for each artifact
- **Legal Reports:** Formatted exports (PDF, CSV) with:
 - Executive summary
 - Chronological timeline
 - Annotated evidence exhibits
 - Analyst testimony support materials
- **Audit Logs:** Immutable record of all system and user actions

Chain-of-Custody Logging Specifications

Every evidence artifact follows a **four-phase lifecycle**, each meticulously documented:

Phase 1: Collection


```
{
  "evidence_id": "uuid-v4",
  "type": "email_communication",
  "status": "collected",
  "timestamp_utc": "2025-11-14T10:15:30.123Z",
  "collected_by": "fcg-ingest-node-03",
  "collection_method": "MTA_forward",
  "source_metadata": {
    "original_recipient": "victim@company.com",
    "received_headers": [ "... "],
    "mta_log_id": "... "
  },
  "content_hash_sha256": "a7f3c9d2...",
  "content_size_bytes": 45678,
  "encryption_key_id": "vault-key-2025-11",
  "storage_location": "s3://evidence-vault/2025/11/14/uuid.enc"
}
```

Phase 2: Analysis

```
{
  "evidence_id": "uuid-v4",
  "status": "analyzed",
  "analyzed_at": "2025-11-14T10:18:45.678Z",
  "analyzed_by": "ai-engine-01",
  "analysis_results": {
    "verdict": "malicious",
    "confidence": 0.96,
    "model_version": "fcg-classifier-v2.3"
  },
  "integrity_verified": true,
  "hash_verification": {
    "expected": "a7f3c9d2...",
    "computed": "a7f3c9d2...",
    "match": true
  }
}
```

Phase 3: Human Review

```
{
  "evidence_id": "uuid-v4",
  "status": "reviewed",
  "reviewed_at": "2025-11-14T14:32:10.000Z",
  "reviewed_by": "analyst_jsmith",
  "review_action": "escalate_to_legal",
  "case_id": "CASE-2025-1142",
  "analyst_notes": "Confirmed fraudulent loan offer. Linked to known actor PDID a7f3c9d2. Recommend legal action.",
  "access_log": [
    { "timestamp": "...", "action": "view", "user": "analyst_jsmith", "ip": "10.0.1.45" },
    { "timestamp": "...", "action": "annotate", "user": "analyst_jsmith" }
  ]
}
```

Phase 4: Legal Export

```
{
  "evidence_id": "uuid-v4",
  "status": "exported",
  "exported_at": "2025-11-15T09:00:00.000Z",
  "exported_by": "legal_team_mdое",
  "export_format": "pdf_report",
  "export_hash_sha256": "b8e4d3f1...",
  "legal_hold": true,
  "retention_policy": "indefinite",
  "certification":
    "This evidence package has been verified for integrity and chain of custody. Digitally signed by Legal Counsel."
}
```

PII Separation and Least-Privilege Access

Data Classification:

- **Tier 1 - Public:** Non-sensitive metadata (timestamps, verdicts, anonymized statistics)
- **Tier 2 - Confidential:** Technical artifacts (fingerprints, IP addresses, non-identifying OSINT)
- **Tier 3 - Restricted:** PII (names, emails, phone numbers, communication content)
- **Tier 4 - Highly Restricted:** Legal strategy, protected communications

Access Control Matrix:

Role	Tier 1	Tier 2	Tier 3	Tier 4
Analyst	RW	RW	R (masked)	-
Senior Analyst	RW	RW	R	-
Investigator	RW	RW	RW	R
Legal Team	R	R	RW	RW
System Admin	R	-	-	-

PII Masking Example:

- Analyst View: Email: f***ster@e***ple.com | Phone: +1-555-***-***23
- Investigator View: Email: fraudster@example.com | Phone: +1-555-012-3423

Legal Review and Report Formats

PDF Report Template (for litigation):

FRAUD INVESTIGATION EVIDENCE REPORT	
Case ID: CASE-2025-1142	
Date: 2025-11-15	
Prepared by: Fraud Communication Guard System	
Reviewed by: [Legal Counsel Name]	

EXECUTIVE SUMMARY

[AI-generated summary of case, key findings, threat level]

SUBJECT PROFILE

[Attacker Knowledge Card - PDF **export**]

EVIDENCE TIMELINE

[Chronological table with exhibits]

Date/Time	Event	Evidence ID	Exhibit #
2025-11-14	...	uuid	EX-001

CHAIN OF CUSTODY CERTIFICATION

[Table showing collection, analysis, review, **export**]

[Digital signature of each handler]

TECHNICAL APPENDIX

[Device attribution reports, OSINT raw data, AI model explanations]

ANALYST TESTIMONY SUPPORT

[Q&A format **for** common legal questions]

- How was this evidence collected?
- What ensures its integrity?
- What **is** the error rate of your classification system?

CSV Timeline Export (for e-discovery):

```
Evid-
ence_ID,Timestamp.UTC,Event_Type,Actor_PDID,Source_IP,Verdict,Description,File_Hash,Ch
ain_of_Custody_Verified
uuid-1,2025-11-14T10:15:30Z,email_received,a7f3c9d2,203.0.113.45,malicious,"Phishing
email impersonating Chase Bank",sha256-hash,TRUE
uuid-2,2025-11-14T10:18:45Z,analysis_complete,a7f3c9d2,203.0.113.45,malicious,"URL
detected: fakeloan-secure.com",sha256-hash,TRUE
```

Implementation Phases

< 1 Hour Prototype Implementation:

1. Database Schema for Chain of Custody (20 min)

- PostgreSQL with tables: `evidence_artifacts`, `custody_chain`, `access_log`
- Each artifact has: UUID, timestamp, content_hash, encrypted_blob
- Custody_chain: FK to artifact, status, actor, timestamp, hash_verified

2. Hash-on-Write Storage (15 min)

- Implement write function that:
 - Computes SHA-256 of artifact before storage
 - Encrypts artifact with AES-256
 - Stores encrypted blob + hash in database

- Logs custody event with system timestamp

3. Basic Access Logging (10 min)

- Wrap all evidence retrieval functions with logging
- Log: user_id, evidence_id, action (view/export), timestamp, source_ip
- Store in immutable `access_log` table (append-only)

4. Simple PDF Report Generator (15 min)

- Use Python `reportlab` library
- Template: Title, Case ID, Timeline table, Evidence list
- Include SHA-256 hash of PDF in metadata

Phase 2+ Enhancements:

- **Blockchain-Backed Logging:** Anchor custody chain hashes to public blockchain (e.g., Ethereum) for tamper-proof audit trail
- **WORM Storage:** Integrate with AWS S3 Object Lock or dedicated WORM appliances
- **Automated Integrity Verification:** Daily cron job re-hashes all artifacts and alerts on mismatches
- **E-Discovery Integration:** Export to EDRM XML format for legal software (Relativity, Concordance)
- **Digital Signature:** GPG-sign all legal reports and evidence packages
- **Retention Policy Engine:** Automated purging based on case status and legal hold flags
- **Advanced RBAC:** Implement attribute-based access control (ABAC) with context-aware policies

2.5 AI/LLM Analysis & Automation Module

Purpose and Responsibilities

The AI/LLM Analysis & Automation Module serves as the **intelligent orchestration layer** that transforms the raw data collected by other modules into actionable intelligence. It leverages machine learning models for classification, natural language processing for content analysis, and large language models for generating human-readable reports and investigative suggestions.

Success Criteria:

- Fraud classification accuracy >95% (precision and recall)
- Model explainability: Every classification includes top 3 contributing factors
- Incident summary generation <30 seconds for complex cases
- False positive rate <2% (measured by analyst override rate)
- OSINT task suggestion acceptance rate >60% (analysts follow suggestions)

Inputs and Outputs

Inputs:

- Attribution reports (device fingerprints, network metadata)
- OSINT knowledge cards (identity graphs, risk indicators)
- Communication analysis (threat verdicts, URL scan results)
- Historical case data (labeled examples for training)

Outputs:

- **Risk Classification:** Benign | Suspicious | Fraud (with probability distribution)
- **Explainability Report:**

```

json
{
  "classification": "fraud",
  "confidence": 0.96,
  "explanation": {
    "top_features": [
      {"feature": "domain_age", "value": "7 days", "contribution": 0.35},
      {"feature": "vpn_detected", "value": true, "contribution": 0.28},
      {"feature": "content_urgency_score", "value": 0.89, "contribution": 0.22}
    ],
    "similar_cases": ["CASE-2025-1089", "CASE-2025-1103"],
    "decision_boundary": "Confidence threshold: 0.85 | This case: 0.96"
  }
}

```

- **Incident Summary** (natural language):

...

INCIDENT SUMMARY - CASE-2025-1142

A fraudulent loan offer was detected on 2025-11-14 at 10:15 UTC. The attacker, identified by device fingerprint a7f3c9d2, sent a phishing email impersonating Chase Bank to victim@company.com. The message contained a link to fakeloan-secure.com, a newly registered domain (7 days old) using a self-signed SSL certificate.

OSINT investigation revealed the attacker is using a VoIP phone number (+1-555-012-3423) registered to a Twilio account and a VPN exit node in Nigeria. The email fraudster@example.com was exposed in 2 data breaches (Collection #1, Dubsplash 2019).

This incident is linked to 3 previous cases involving the same device fingerprint and similar phishing tactics. RECOMMENDATION: Escalate to legal team for law enforcement referral.

...

- **OSINT Task Suggestions:**

- "Search for domain registrant email in data breach databases"
- "Check if phone number is linked to other fraud reports"
- "Monitor social media accounts for new activity"

Classification Model: Benign/Fraud/Suspicious

Model Architecture:

Input Features (50+ dimensions):

- Device/Network: VPN detected, Tor usage, device fingerprint novelty, IP reputation score, ASN type (residential/datacenter)
- Communication: Email authentication (SPF/DKIM/DMARC), domain age, URL count, attachment presence
- Content: Urgency language score (NLP), credential request detected, brand mention, grammar quality
- OSINT: Social media authenticity, data breach exposure, phone number type, domain WHOIS privacy
- Behavioral: Session count, time-of-day pattern, user interaction speed
- Historical: Match to known fraud patterns, similarity to past cases

Model Options:

1. **Gradient Boosting (XGBoost/LightGBM)** - Recommended for production
 - Excellent accuracy on tabular data
 - Built-in feature importance (explainability)
 - Fast inference (<100ms)
 - Training: Supervised learning on labeled fraud cases
2. **Random Forest** - Good for quick prototyping
 - Robust to overfitting
 - Feature importance available
 - Easier to interpret than neural networks
3. **Logistic Regression** - Baseline model
 - High interpretability (coefficients = feature weights)
 - Fast training and inference
 - Works well with feature engineering

Training Data Requirements:

- Minimum: 1,000 labeled cases (500 fraud, 500 benign)
- Ideal: 10,000+ cases with continuous retraining
- Labeling: Analyst decisions + victim reports + confirmed fraud outcomes

Explainability Implementation:

Using **SHAP (SHapley Additive exPlanations)**:

```
import shap

explainer = shap.TreeExplainer(model)
shap_values = explainer.shap_values(input_features)

# Top 3 contributing features
top_features = sorted(zip(feature_names, shap_values),
                      key=lambda x: abs(x[1]),
                      reverse=True)[:3]
```

Output includes both global feature importance and local (per-prediction) explanations.

LLM Integration for Incident Summaries and OSINT Suggestions**Use Case 1: Incident Summary Generation**

Approach: Template-based prompting with structured data injection

Prompt Template:

You are a cybersecurity analyst writing an incident report. Based on the following structured data, generate a concise, professional incident summary (200-300 words) that a legal team can understand:

```

CASE ID: {case_id}
DETECTION TIME: {timestamp}
ATTACKER PROFILE:
- Device ID: {pdid}
- Location: {location}
- VPN/Proxy: {vpn_detected}
- Email: {email}
- Phone: {phone}

INCIDENT DETAILS:
- Type: {incident_type}
- Target: {victim_email}
- Malicious Content: {threat_description}
- Domain: {domain} (Age: {domain_age} days)

OSINT FINDINGS:
{osint_summary}

SIMILAR CASES: {linked_cases}

RISK SCORE: {risk_score}/100

```

Write a clear narrative that explains what happened, who **is** responsible, what evidence was collected, **and** what actions are recommended.

Model Selection:

- **GPT-4** (via Azure OpenAI): Highest quality, good for complex cases
- **Claude 3** (Anthropic): Strong reasoning, good for legal context
- **Llama 3 (70B)** (Self-hosted): Privacy-preserving, no data leaves infrastructure

Safety and Validation:

- **Prompt Injection Prevention:** Sanitize all input data, escape special characters
- **Hallucination Detection:** Cross-check all facts in generated summary against source data
- **Human-in-Loop:** Analyst reviews and approves summary before inclusion in legal report

Use Case 2: OSINT Task Suggestions

Approach: Context-aware task generation based on incomplete knowledge cards

Logic:

```

IF email found AND social_media_profiles = NULL:
    SUGGEST: "Search for email on social media platforms using Sherlock"

IF domain found AND whois_privacy = TRUE:
    SUGGEST: "Check historical WHOIS records for prior exposure of registrant info"

IF phone_number found AND carrier = "VoIP":
    SUGGEST: "Investigate VoIP provider account creation methods and payment data"

IF device_fingerprint matches previous case:
    SUGGEST: "Review investigation notes from similar case CASE-{id} for additional leads"

```

LLM Enhancement:

Given the context, an LLM can generate natural language task descriptions with rationale:

Prompt:

Given the following OSINT findings, suggest 3-5 investigative actions that would provide additional context about the suspect. For each suggestion, explain the expected value of the information:

CURRENT FINDINGS:

{knowledge_card_summary}

MISSING INFORMATION:

- No social media profiles found
- Domain registered with privacy protection
- Phone number is VoIP type

Provide suggestions **in** this format:

1. [Action] - [Expected Value]

Output:

1. Reverse image search **on** any profile photos found in emails - May reveal stolen identity **or** link **to** legitimate social media accounts
2. Search phone number in Telegram/WhatsApp **for** public profile - VoIP users often use these platforms, may **expose** real name **or** profile photo
3. Monitor domain **for** SSL certificate changes - **If** privacy is dropped **or** cert is updated, may reveal registrant email in CT **logs**

LLM Call Patterns and Prompt Safety**Best Practices:****1. Input Sanitization:**

- Escape all user-provided text
- Limit input length (4000 tokens max)
- Remove any content that could be interpreted as instructions

2. Output Validation:

- Check for factual accuracy against source data
- Flag outputs containing phrases like "I cannot verify" or "As an AI"
- Use confidence scoring: LLM generates claim → validator checks against structured data

3. Rate Limiting:

- Max 100 LLM calls per hour per analyst (prevent abuse)
- Cache common queries (e.g., same case accessed multiple times)

4. Model Selection by Sensitivity:

- **High Sensitivity** (legal reports): GPT-4, human review required
- **Medium Sensitivity** (analyst summaries): GPT-4 or Claude 3, spot-check validation
- **Low Sensitivity** (OSINT suggestions): Llama 3 self-hosted, no external data exposure

5. Audit Trail:

- Log all LLM prompts and responses
- Version control prompts (treat as code)
- Track which analyst approved LLM-generated content

Implementation Phases

< 1 Hour Prototype Implementation:

1. **Feature Engineering** (20 min)
 - Create Python script to extract features from attribution, OSINT, communication modules
 - Features: domain_age, vpn_detected, email_auth_pass, url_count, urgency_score
 - Store as CSV for training data collection
2. **Baseline Classifier** (15 min)
 - Use scikit-learn LogisticRegression
 - Train on synthetic/bootstrapped data (100 examples)
 - Predict: fraud probability
 - Threshold: >0.85 = FRAUD, $0.5-0.85$ = SUSPICIOUS, <0.5 = BENIGN
3. **Simple LLM Summary** (15 min)
 - Use OpenAI API (gpt-3.5-turbo for speed)
 - Template: "Summarize this fraud case: {case_data}"
 - Display summary in analyst dashboard
4. **Rule-Based OSINT Suggestions** (10 min)
 - Implement if/else logic for missing data fields
 - Generate list of 3-5 suggested actions
 - Display in "Recommended Next Steps" panel

Phase 2+ Enhancements:

- **Production ML Model:** Train XGBoost on 10k+ labeled cases, deploy via MLflow
- **Real-Time Model Monitoring:** Track prediction distribution, retrain on drift
- **A/B Testing:** Compare model versions, measure analyst override rates
- **Advanced NLP:** Fine-tune transformer model on fraud communication corpus for content analysis
- **Federated Learning:** Train models across multiple organizations without sharing raw data
- **Reinforcement Learning:** Optimize OSINT task sequencing based on investigative outcomes
- **Multi-Modal Analysis:** Integrate computer vision for analyzing suspicious website screenshots

PHASE 3 - TOOLING, STACK, AND QUICK-START

3.1 Technology & Tool Selection

This section provides a curated technology stack organized by function, prioritizing open-source solutions with API-first architectures for maximum flexibility and cost-effectiveness.

OSINT and Threat Intelligence

1. SpiderFoot HX (Open-Source Edition)

- **Function:** Automated OSINT collection and correlation engine
- **Capabilities:** 200+ modules for DNS, WHOIS, social media, breach data, dark web
- **Maturity:** Production-ready, active development, large community
- **Licensing:** MIT License (free for commercial use)
- **Integration:** REST API, Python library, Docker deployment
- **Justification:** Industry-standard OSINT automation with extensive data source coverage. Graph visualization enables rapid link analysis.

2. theHarvester

- **Function:** Email, subdomain, and employee enumeration
- **Capabilities:** Queries search engines, Shodan, PGP servers, certificate transparency logs
- **Maturity:** Stable, widely used in penetration testing and OSINT
- **Licensing:** GPLv2
- **Integration:** Command-line tool, easily wrapped in Python subprocess
- **Justification:** Fast initial reconnaissance for email/domain intelligence gathering.

3. MISP (Malware Information Sharing Platform)

- **Function:** Threat intelligence aggregation, sharing, and correlation
- **Capabilities:** IOC management, event correlation, API for automated ingestion
- **Maturity:** Enterprise-grade, used by CERTs and SOCs globally
- **Licensing:** AGPL (open-source)
- **Integration:** REST API, PyMISP library, STIX/TAXII support
- **Justification:** Centralized threat intelligence repository. Enables collaboration and sharing of fraud actor IOCs with industry peers.

4. Maltego Community Edition

- **Function:** Visual link analysis and investigation platform
- **Capabilities:** Entity graphing, 100+ transforms for OSINT data
- **Maturity:** Industry standard for investigative visualization
- **Licensing:** Free Community Edition (limited transforms)
- **Integration:** Standalone desktop application, custom transforms via Python
- **Justification:** Invaluable for manual deep-dive investigations. Visual graphs help identify non-obvious connections.

5. OpenCTI

- **Function:** Cyber threat intelligence platform with knowledge graph
- **Capabilities:** Structured threat data, relationship mapping, import/export STIX 2.1
- **Maturity:** Growing adoption, strong community
- **Licensing:** Apache 2.0
- **Integration:** GraphQL API, Python SDK
- **Justification:** Modern alternative to MISP with better visualization and knowledge management.

Device Fingerprinting and Attribution

1. FingerprintJS (Open-Source Edition)

- **Function:** Browser fingerprinting library
- **Capabilities:** Canvas, WebGL, Audio, Font detection, 50+ signals
- **Maturity:** Production-ready, 18k+ GitHub stars
- **Licensing:** MIT (OSS edition), commercial license available for advanced features
- **Integration:** JavaScript library, NPM package
- **Justification:** Most comprehensive open-source fingerprinting library. Easy to deploy and customize.

2. MaxMind GeoIP2 (GeoLite2 Free Database)

- **Function:** IP geolocation database
- **Capabilities:** Country, city, ASN, ISP resolution
- **Maturity:** Industry standard, weekly updates
- **Licensing:** Creative Commons (GeoLite2), commercial license for higher accuracy
- **Integration:** Local database with fast lookup libraries (Python, Java, C)
- **Justification:** Free, accurate, and privacy-preserving (no external API calls). Sufficient for most fraud investigations.

3. p0f (Passive OS Fingerprinting)

- **Function:** Operating system and network stack detection via packet analysis
- **Capabilities:** Identifies OS, browser, NAT, load balancers without active probing
- **Maturity:** Mature, used in IDS/IPS systems
- **Licensing:** LGPL
- **Integration:** Command-line tool, can parse pcap files or live traffic
- **Justification:** Adds orthogonal signal to device fingerprinting. Useful for detecting emulators and VMs.

Logging, SIEM, and Evidence Management

1. Elasticsearch + Logstash + Kibana (ELK Stack)

- **Function:** Log aggregation, search, and visualization
- **Capabilities:** Real-time ingestion, full-text search, dashboards, alerting
- **Maturity:** Enterprise-grade, massive ecosystem
- **Licensing:** Elastic License (free for self-managed, restrictions on cloud services)
- **Integration:** REST API, native support for logs, metrics, traces
- **Justification:** De facto standard for log management. Kibana provides powerful visualization for investigators. Elastic Security extends it for SIEM use cases.

2. Wazuh

- **Function:** Open-source SIEM and XDR platform
- **Capabilities:** Log analysis, intrusion detection, compliance monitoring, file integrity
- **Maturity:** Production-ready, SOC 2 Type II certified
- **Licensing:** GPLv2
- **Integration:** Built on ELK stack, REST API, agent-based and agentless collection
- **Justification:** Full-featured SIEM at no cost. Includes out-of-the-box rules for threat detection and compliance (PCI DSS, GDPR, HIPAA).

3. Graylog

- **Function:** Log management and analysis platform
- **Capabilities:** Centralized logging, real-time search, alerting, data archival
- **Maturity:** Stable, used in production by enterprises
- **Licensing:** Server Side Public License (free for most use cases)
- **Integration:** REST API, Syslog/GELF inputs, Elasticsearch backend
- **Justification:** Lighter weight than ELK, excellent for teams without dedicated infrastructure. Strong alerting capabilities for real-time fraud detection.

Recommendation: Start with **Wazuh** for integrated SIEM + evidence logging. It provides the best balance of features and ease of deployment for a security team.

LLM Platforms for Security Analysis

1. OpenAI GPT-4 (via Azure OpenAI Service)

- **Function:** State-of-the-art large language model for text generation and analysis
- **Capabilities:** Incident summaries, OSINT task suggestions, content analysis
- **Maturity:** Industry-leading performance and reliability
- **Licensing:** Commercial API, pay-per-token
- **Integration:** REST API, Python SDK (`openai` library)
- **Justification:** Highest quality outputs for complex reasoning tasks. Azure offering provides enterprise SLA and data privacy guarantees. **Use for legal reports and high-stakes analysis.**

2. Anthropic Claude 3 (Opus/Sonnet)

- **Function:** Advanced LLM with strong reasoning and safety features
- **Capabilities:** Long context (200k tokens), nuanced analysis, excellent for legal/compliance content
- **Maturity:** Production-ready, competitive with GPT-4
- **Licensing:** Commercial API, pay-per-token
- **Integration:** REST API, Python SDK
- **Justification:** Excels at following complex instructions and maintaining context. Strong constitutional AI training reduces harmful outputs. **Use for analyst-facing summaries and complex case analysis.**

3. Llama 3 (70B or 405B) - Self-Hosted

- **Function:** Open-source LLM for on-premises deployment
- **Capabilities:** Comparable quality to GPT-3.5, no data leaves infrastructure
- **Maturity:** Production-ready with Meta's support
- **Licensing:** Meta Llama 3 Community License (permissive for most uses)
- **Integration:** Deploy via vLLM, TGI (Text Generation Inference), or Ollama
- **Justification:** **Privacy-preserving option for sensitive data.** Eliminates risk of data exposure to third-party APIs. Ideal for organizations with strict data residency requirements. **Use for OSINT suggestions and lower-stakes tasks.**

Recommendation: Hybrid approach:

- **GPT-4** for legal reports and final deliverables (highest quality, human-reviewed)
- **Claude 3** for analyst dashboard summaries (good balance of quality and cost)
- **Llama 3 (self-hosted)** for OSINT task generation and internal tools (privacy, no API costs)

3.2 One-Hour Prototype Path

This section provides a **time-bounded deployment plan** for a single engineer or AI-assisted implementation to stand up a minimal viable slice of the Fraud Communication Guard in under one hour.

Prototype Scope: Minimal Viable Slice

Capabilities:

1. Capture one communication channel (email forwarding)
2. Perform basic attribution (IP geolocation, device fingerprint)
3. Conduct OSINT enrichment (email lookup, domain WHOIS)
4. Log all artifacts with timestamps and hashes
5. Generate a simple "Attacker Profile" report

Out of Scope for 1-Hour:

- Full sandbox analysis (use VirusTotal API only)
- AI/ML classification (use rule-based triage)
- Legal report generation (export raw JSON)
- Advanced dashboard (command-line interface)

Prerequisites

- **Environment:** Linux server (Ubuntu 22.04) or macOS with Docker
- **Tools Installed:** Python 3.10+, Docker, curl, git
- **API Keys** (free tiers):
 - VirusTotal: <https://www.virustotal.com/gui/join-us>
 - IPQualityScore: <https://www.ipqualityscore.com/create-account>

- Hunter.io: https://hunter.io/users/sign_up
- Have I Been Pwned: <https://haveibeenpwned.com/API/Key>

Step-by-Step Implementation (60 Minutes)

Step 1: Project Setup (5 minutes)

```
# Create project directory
mkdir fraud-comm-guard && cd fraud-comm-guard

# Create virtual environment
python3 -m venv venv
source venv/bin/activate

# Create directory structure
mkdir -p {ingest,attribution,osint,evidence,reports}
touch config.env
```

Step 2: Install Dependencies (5 minutes)

```
# Install Python libraries
pip install flask requests python-whois geoip2 maxminddb-geolite2 cryptography

# Download GeoLite2 database
cd attribution
wget https://github.com/P3TERX/GeoLite.mmdb/raw/download/GeoLite2-City.mmdb
cd ..
```

Step 3: Create Email Ingest Endpoint (10 minutes)

Create `ingest/email_server.py` :

```

from flask import Flask, request, jsonify
import hashlib, json, uuid
from datetime import datetime

app = Flask(__name__)

@app.route('/ingest/email', methods=['POST'])
def ingest_email():
    data = request.json
    email_id = str(uuid.uuid4())
    timestamp = datetime.utcnow().isoformat() + 'Z'

    # Extract key fields
    artifact = {
        'email_id': email_id,
        'timestamp': timestamp,
        'from': data.get('from'),
        'to': data.get('to'),
        'subject': data.get('subject'),
        'body': data.get('body'),
        'source_ip': request.remote_addr,
        'raw_data': data
    }

    # Hash for integrity
    artifact_str = json.dumps(artifact, sort_keys=True)
    artifact['hash'] = hashlib.sha256(artifact_str.encode()).hexdigest()

    # Save to evidence vault
    with open(f'evidence/{email_id}.json', 'w') as f:
        json.dump(artifact, f, indent=2)

    # Trigger attribution and OSINT
    from attribution.attribute import perform_attribution
    from osint.email_lookup import osint_email

    attr_report = perform_attribution(artifact['source_ip'], request.headers.get('User-Agent'))
    osint_report = osint_email(artifact['from'])

    # Save reports
    with open(f'reports/{email_id}_attribution.json', 'w') as f:
        json.dump(attr_report, f, indent=2)
    with open(f'reports/{email_id}_osint.json', 'w') as f:
        json.dump(osint_report, f, indent=2)

    return jsonify({'status': 'ingested', 'email_id': email_id})

if __name__ == '__main__':
    app.run(host='0.0.0.0', port=5000)

```

Step 4: Implement Attribution Module (10 minutes)

Create `attribution/attribute.py` :

```

import geoip2.database, requests, os

def perform_attribution(ip_address, user_agent):
    # GeoIP lookup
    reader = geoip2.database.Reader('attribution/GeoLite2-City.mmdb')
    try:
        response = reader.city(ip_address)
        geo = {
            'country': response.country.name,
            'city': response.city.name,
            'lat': response.location.latitude,
            'lon': response.location.longitude
        }
    except:
        geo = {'error': 'IP not found'}

    # VPN detection (IPQualityScore)
    api_key = os.getenv('IPQS_API_KEY', 'your-key-here')
    vpn_check = requests.get(f'https://ipqualityscore.com/api/json/ip/{api_key}/{ip_address}').json()

    return {
        'ip_address': ip_address,
        'user_agent': user_agent,
        'geolocation': geo,
        'vpn_detected': vpn_check.get('vpn', False),
        'proxy_detected': vpn_check.get('proxy', False),
        'fraud_score': vpn_check.get('fraud_score', 0)
    }

```

Step 5: Implement OSINT Module (15 minutes)

Create `osint/email_lookup.py`:

```

import requests, whois, os

def osint_email(email_address):
    result = {'email': email_address}

    # Hunter.io lookup
    hunter_key = os.getenv('HUNTER_API_KEY', 'your-key-here')
    hunter_resp = requests.get(f'https://api.hunter.io/v2/email-verifier?email={email_address}&api_key={hunter_key}').json()
    result['hunter'] = hunter_resp.get('data', {})

    # Have I Been Pwned
    hibp_key = os.getenv('HIBP_API_KEY', 'your-key-here')
    hibp_resp = requests.get(f'https://haveibeenpwned.com/api/v3/breachedaccount/{email_address}',
                             headers={'hibp-api-key': hibp_key})
    result['breaches'] = hibp_resp.json() if hibp_resp.status_code == 200 else []

    # Domain WHOIS
    domain = email_address.split('@')[1]
    try:
        w = whois.whois(domain)
        result['domain_whois'] = {
            'registrar': w.registrar,
            'creation_date': str(w.creation_date),
            'expiration_date': str(w.expiration_date)
        }
    except:
        result['domain_whois'] = {'error': 'WHOIS failed'}

    return result

```

Step 6: Simple Attacker Profile Generator (10 minutes)

Create `reports/generate_profile.py`:


```

import json, sys

def generate_profile(email_id):
    # Load artifacts
    with open(f'evidence/{email_id}.json') as f:
        email = json.load(f)
    with open(f'reports/{email_id}_attribution.json') as f:
        attr = json.load(f)
    with open(f'reports/{email_id}_osint.json') as f:
        osint = json.load(f)

    # Generate simple profile
    profile = f"""
ATTACKER PROFILE
=====
Email ID: {email_id}
Timestamp: {email['timestamp']}

COMMUNICATION:
From: {email['from']}
To: {email['to']}
Subject: {email['subject']}

ATTRIBUTION:
IP Address: {attr['ip_address']}
Location: {attr['geolocation'].get('city', 'Unknown')}, {attr['geolocation'].get('country', 'Unknown')}
VPN Detected: {attr['vpn_detected']}
Fraud Score: {attr['fraud_score']}/100

OSINT INTELLIGENCE:
Email Breaches: {len(osint.get('breaches', []))} found
Domain Registrar: {osint.get('domain_whois', {}).get('registrar', 'Unknown')}
Domain Age: {osint.get('domain_whois', {}).get('creation_date', 'Unknown')}

RISK ASSESSMENT:
{'HIGH RISK' if attr['fraud_score'] > 75 else 'MEDIUM RISK' if attr['vpn_detected'] else 'LOW RISK'}
"""

    print(profile)
    with open(f'reports/{email_id}_profile.txt', 'w') as f:
        f.write(profile)

if __name__ == '__main__':
    generate_profile(sys.argv[1])

```

Step 7: Configure API Keys (2 minutes)

Edit `config.env` :

```

export IPQS_API_KEY='your-ipqualityscore-key'
export HUNTER_API_KEY='your-hunter-io-key'
export HIBP_API_KEY='your-haveibeenpwned-key'

```

Load environment:

```
source config.env
```

Step 8: Start the Ingest Server (1 minute)

```
cd ingest
python email_server.py
```

Step 9: Test with Sample Email (5 minutes)

In a new terminal:

```
curl -X POST http://localhost:5000/ingest/email \
-H "Content-Type: application/json" \
-d '{
  "from": "fraudster@suspicious-domain.com",
  "to": "victim@company.com",
  "subject": "URGENT: Verify your account",
  "body": "Click here: http://phishing-site.com"
}'
```

Step 10: Generate Attacker Profile (2 minutes)

```
# Get email_id from server response
EMAIL_ID="<uuid-from-response>"

python reports/generate_profile.py $EMAIL_ID

cat reports/${EMAIL_ID}_profile.txt
```

Expected Output**ATTACKER PROFILE**

=====

Email ID: a7f3c9d2-1234-5678-90ab-cdef12345678

Timestamp: 2025-11-14T10:15:30.123Z

COMMUNICATION:

From: fraudster@suspicious-domain.com

To: victim@company.com

Subject: URGENT: Verify your account

ATTRIBUTION:

IP Address: 203.0.113.45

Location: Lagos, Nigeria

VPN Detected: True

Fraud Score: 92/100

OSINT INTELLIGENCE:

Email Breaches: 2 found

Domain Registrar: Namecheap Inc.

Domain Age: 2025-11-01 (14 days old)

RISK ASSESSMENT:

HIGH RISK

PHASE 4 - RISK, PRIVACY, AND LIMITATIONS

4.1 Risk & Abuse Prevention

Identified Risks

1. Over-Collection and Scope Creep

Risk: The system's powerful collection capabilities may lead to accumulation of unnecessary data beyond what is required for fraud investigation, violating data minimization principles.

Impact: GDPR/CCPA violations, reputational damage, increased attack surface (more data = more to protect), potential for misuse.

Mitigation:

- **Automated Purging:** Implement retention policies with automatic deletion of non-case-related data after 90 days
- **Collection Audits:** Quarterly reviews by legal team to assess whether all collection types remain necessary
- **Purpose Limitation Enforcement:** Technical controls that prevent repurposing data (e.g., marketing uses)
- **Data Access Logging:** Monitor which data is accessed; flag unused data sources for deprecation

2. OSINT Misuse and Overreach

Risk: OSINT tools can be used to investigate individuals who are not actually suspects, constituting unauthorized surveillance or harassment.

Impact: Privacy violations, civil liability, employee trust erosion if used on internal personnel without cause.

Mitigation:

- **Authorization Workflow:** Require documented justification and supervisor approval before initiating OSINT investigation on any identifier
- **Audit Trail:** Log all OSINT queries with analyst ID and case number; monthly review by compliance officer
- **Prohibited Use Cases:** Explicit policy prohibiting use for: employment screening (unless authorized), personal curiosity, competitive intelligence
- **Technical Controls:** Implement "break glass" access for sensitive OSINT sources, requiring legal counsel approval

3. AI/ML Model Hallucinations and False Positives

Risk: LLMs may generate plausible-sounding but factually incorrect statements in incident summaries. ML classifiers may flag legitimate users as fraudsters.

Impact: False accusations, wasted investigative resources, legal jeopardy if erroneous evidence is submitted in court.

Mitigation:

- **Human-in-Loop:** Mandatory analyst review of all AI-generated content before use in legal proceedings
- **Hallucination Detection:** Automated fact-checking that cross-references LLM outputs against structured data sources
- **Confidence Thresholds:** Require 95% confidence for automated blocking; 85-95% goes to human

review; <85% flagged for model retraining

- **False Positive Tracking:** Dashboard showing analyst override rate; trigger model retraining if overrides exceed 5%
- **Explainability Requirement:** Every classification includes top 3 features; analyst can challenge if features seem spurious

4. Chain of Custody Compromise

Risk: Inadequate logging, hash verification failures, or unauthorized access could break the chain of custody, rendering evidence inadmissible.

Impact: Case dismissal, inability to prosecute fraudsters, organizational liability.

Mitigation:

- **Immutable Storage:** Use WORM (Write Once Read Many) or blockchain-anchored logs
- **Automated Integrity Checks:** Daily re-hashing of all evidence artifacts; immediate alert on mismatch
- **Multi-Factor Authentication:** Require MFA for all evidence vault access
- **Physical Security:** Evidence storage servers in locked, access-controlled rooms with 24/7 monitoring
- **Forensic Readiness Audits:** Annual third-party assessment of evidence management practices

5. Insider Threat: Analyst Misconduct

Risk: Malicious or negligent analysts could exfiltrate sensitive data, tamper with evidence, or abuse access for personal gain.

Impact: Data breach, evidence contamination, criminal prosecution of organization.

Mitigation:

- **Least Privilege:** Analysts see only PII-masked data unless investigative need documented
- **Peer Review:** High-stakes actions (e.g., legal report export) require secondary analyst or supervisor approval
- **Behavioral Analytics:** Monitor analyst access patterns for anomalies (e.g., bulk downloads, off-hours access)
- **Separation of Duties:** Analyst who conducts investigation \neq analyst who exports for legal (independent verification)
- **Background Checks:** All personnel with evidence vault access subject to enhanced background screening

Guardrails and Safeguards

Technical Guardrails:

1. **Rate Limiting:** OSINT API calls limited to 1000/day per analyst (prevents mass surveillance)
2. **Query Logging:** All searches (email, phone, IP) logged with timestamp and justification code
3. **Automated Redaction:** PII automatically masked in analyst views unless "Unmask PII" button clicked (logged)
4. **Data Loss Prevention (DLP):** Egress filtering prevents evidence export to unauthorized destinations (USB drives, personal email)

Policy Guardrails:

1. **Acceptable Use Policy (AUP):** All analysts sign AUP defining prohibited uses and consequences

2. **Quarterly Training:** Mandatory training on privacy laws, ethical investigation practices, and acceptable use
3. **Incident Response Plan:** Documented procedures for handling misuse, data breaches, or chain-of-custody breaks
4. **Whistleblower Protection:** Anonymous reporting channel for analysts to report misuse without retaliation

Legal Guardrails:

1. **Legal Review Checkpoints:** Legal counsel reviews system design, collection practices, and data retention policies
 2. **Data Protection Impact Assessment (DPIA):** Completed before deployment in any new jurisdiction
 3. **Warrant/Subpoena Management:** Centralized tracking of all legal process; evidence only accessed when legally authorized
 4. **Right to Audit:** Subjects of investigation can request (post-investigation) an audit of how their data was used
-

4.2 Limitations & Future Work

System Limitations

1. Cannot Guarantee 100% Attribution

Limitation: Determined attackers using sophisticated anonymization (multi-hop VPNs, Tor, compromised residential proxies) may evade attribution. Device fingerprints can be spoofed using emulators or specialized tools.

Impact: Some fraud actors will remain unidentified. System should not be relied upon as sole attribution method.

Workaround: Combine with other investigative techniques (financial transaction tracing, law enforcement partnerships). Document limitations when presenting evidence.

2. OSINT Data Quality and Freshness

Limitation: OSINT sources may contain outdated, incomplete, or false information. Social media accounts may be deleted before investigation begins. Data breach dumps may have partial or corrupted records.

Impact: “Attacker Knowledge Cards” may be incomplete or contain inaccuracies. Investigators must independently verify critical facts.

Workaround: Always cross-reference OSINT findings with multiple sources. Document confidence levels for each data point. Use OSINT as leads for further investigation, not definitive proof.

3. Legal Process Delays

Limitation: Obtaining subpoenas for ISP subscriber data, VPN provider logs, or email account records requires legal process that can take weeks to months. Some providers (especially foreign VPNs) may not cooperate.

Impact: Time-sensitive investigations may lose momentum. Real-time prevention of fraud may not be possible for sophisticated actors.

Workaround: Build relationships with major providers for expedited processing. Focus system on detection and evidence collection; accept that attribution may be delayed.

4. ML Model Bias and Training Data Limitations

Limitation: If training data is skewed (e.g., mostly phishing from Nigeria, few from Eastern Europe), model may have geographic bias. Fraudsters using novel tactics not present in training data will evade detection.

Impact: Higher false negative rate for underrepresented fraud types. Potential for discriminatory outcomes if model relies on proxies for protected classes.

Workaround: Continuous retraining with diverse fraud examples. Regular fairness audits (e.g., check if VPN usage from legitimate privacy activists is over-flagged). Human review for all automated decisions with material impact.

5. Scalability Constraints

Limitation: This prototype design is suitable for ~1000 communications/day. Higher volumes require infrastructure scaling (distributed processing, larger databases, load balancers).

Impact: May not be suitable for large organizations or email service providers without significant engineering investment.

Workaround: For high-volume use cases, re-architect with streaming platforms (Kafka), distributed databases (Cassandra), and auto-scaling compute (Kubernetes).

Future Enhancements (Phase 2+)

1. Proactive Threat Hunting

Description: Instead of waiting for fraudulent communications to arrive, actively search for fraud infrastructure (newly registered domains impersonating brand, dark web listings selling access, underground forums discussing targets).

Implementation:

- Continuous monitoring of certificate transparency logs for brand impersonation domains
- Dark web scraping for organization name mentions
- Integration with threat intel feeds (Recorded Future, ThreatConnect)

Expected Value: Earlier detection of fraud campaigns before they launch. Opportunity for takedown of fraud infrastructure preemptively.

2. Victim Outreach and Education

Description: When fraud is detected, automatically generate educational materials for victims explaining the fraud type, how to protect themselves, and how to report.

Implementation:

- Templated "Fraud Alert" emails with safe prevention tips
- Dashboard for customer service teams to view user fraud exposure
- Integration with identity protection services (credit monitoring, dark web monitoring)

Expected Value: Reduced victim impact. Enhanced brand trust through proactive protection.

3. Cross-Organizational Intelligence Sharing

Description: Establish consortium of organizations using FCG to share anonymized fraud actor indicators (device fingerprints, email hashes, IP addresses) via MISP.

Implementation:

- Hash PII before sharing (e.g., SHA-256 of email, not raw email)
- MISP instance for consortium, access controlled by membership agreement
- Automated ingestion of shared IOCs into local FCG instance

Expected Value: Network effect—fraud actor identified at Company A is auto-flagged at Company B. Faster industry-wide response to fraud waves.

4. Behavioral Biometrics for Interactive Sessions

Description: For fraud attempts that involve live interaction (e.g., account takeover via web session), add keystroke dynamics and mouse movement analysis.

Implementation:

- JavaScript library deployed on login/transaction pages
- Behavioral model trained on legitimate user patterns
- Anomaly detection flags account takeovers even when credentials are correct

Expected Value: Detection of account takeover attacks that bypass password-based authentication.

5. Automated Fraud Infrastructure Takedown

Description: When fraudulent domain/website is detected, automatically initiate takedown workflow (report to registrar, hosting provider, abuse@).

Implementation:

- Integration with registrar abuse reporting APIs
- Templated abuse complaints with evidence attachments
- Tracking system for takedown status (pending, completed)
- Legal review for borderline cases before reporting

Expected Value: Faster neutralization of fraud infrastructure. Reduced window of opportunity for fraudsters.

Conclusion

The **Fraud Communication Guard** system design represents a comprehensive, legally defensible, and technically rigorous approach to defensive fraud investigation. By integrating device attribution, OSINT intelligence, communication analysis, forensically sound evidence collection, and AI-powered automation, this system empowers security teams to identify, track, and build cases against sophisticated fraudsters.

This document has provided:

- **Clear ethical boundaries:** Defensive focus, privacy compliance, evidence-centric mission
- **Detailed architecture:** Six core modules with data flows, trust boundaries, and privacy controls
- **Actionable specifications:** Input/output contracts, tool recommendations, implementation phases
- **Rapid prototyping path:** 1-hour minimal viable slice for immediate value

- **Risk mitigation:** Comprehensive analysis of risks with technical and policy guardrails
- **Honest limitations:** Transparent about what the system cannot guarantee, with documented workarounds

Next Steps for Implementation:

1. **Phase 0 (Legal Authorization):** Obtain written legal counsel approval and complete DPIA for target jurisdiction
2. **Phase 1 (Prototype):** Follow 1-hour quick-start to validate concept with real fraud case
3. **Phase 2 (Production Pilot):** Deploy full modules with selected tooling stack to security team
4. **Phase 3 (Scale):** Integrate with existing SOC/SIEM, train analysts, establish legal workflows
5. **Phase 4 (Optimization):** Train ML models on organizational data, tune detection rules, measure KPIs

This blueprint is designed to be directly usable by engineers and AI agents, providing sufficient technical specificity while maintaining strategic flexibility. The system can be deployed incrementally, with each module providing value independently, and can scale from a single-analyst operation to an enterprise-wide fraud defense platform.

Final Note on Ethics and Responsibility:

The power to persistently identify and track individuals, even for defensive purposes, carries profound ethical weight. Organizations deploying this system must commit to:

- **Continuous legal compliance** through regular audits and counsel engagement
- **Transparency with users** through clear privacy policies and consent mechanisms
- **Accountability** through robust logging, oversight, and willingness to be audited
- **Proportionality** in collection and retention—never collecting more than necessary
- **Human dignity** even when investigating bad actors—no vigilante justice, no public shaming

By adhering to these principles, the Fraud Communication Guard can serve its intended purpose: protecting users from fraud while respecting their fundamental rights to privacy and due process.

Document End

For questions, clarifications, or implementation support, contact: [Security Architecture Team]

Version History:

- v1.0 (2025-11-14): Initial system design document