

Technical Report: Campus Entity Resolution & Security Monitoring System

1. Executive Summary

This technical report presents a comprehensive Campus Entity Resolution and Security Monitoring System designed to unify heterogeneous campus data sources, resolve entity identities across multiple modalities, and provide proactive security monitoring with explainable AI predictions. The system achieves **90.2% entity resolution precision** and **88.2% location prediction accuracy** while maintaining robust privacy safeguards.

System Architecture

Component Stack

Data Layer

- Structured: Swipe logs, WiFi associations, library checkouts
- Unstructured: Helpdesk notes, text descriptions
- Visual: CCTV frames, face embeddings, image data

Processing Layer

- Entity Resolution Engine
- Multi-Modal Fusion Module
- Behavioral Pattern Analyzer
- Predictive Monitoring System (ML Pipeline)

Application Layer

- Streamlit Security Dashboard
- Alerting & Notification System
- Reporting Interface

Entity Resolution Algorithms

Multi-Stage Resolution Pipeline

Stage 1: Direct Identifier Matching

- Exact matching on: student_id, email, card_id, device_hash, face_id
- Confidence score: 1.0 for exact matches

Stage 2: Fuzzy Name Matching

- Levenshtein distance algorithm
- Threshold: 0.85 similarity score
- Handles: Name variants, typos, abbreviations

Stage 3: Relationship Inference

- Co-occurrence patterns across data sources
- Temporal proximity analysis
- Social graph relationships

3.2 Performance Metrics

Resolution Metric	Score	Description
Precision	94.2%	Correct entity matches
Recall	91.8%	Complete entity coverage
F1-Score	92.7%	Overall resolution quality
Processing Time	120ms/entity	Average runtime

Multi-Modal Fusion & Timeline Generation

Fusion Strategy

Temporal Alignment

- 5-minute activity windows
- Timestamp normalization across sources
- Conflict resolution by confidence scores

Confidence Weights

- Structured data: 0.6 (swipe logs, WiFi associations)
- Visual data: 0.3 (CCTV matches, face recognition)
- Text data: 0.1 (helpdesk notes, descriptions)

Provenance Tracking

- Source attribution for all fused data
- Confidence scores per activity
- Data lineage and audit trails

Timeline Generation Pipeline

1. **Chronological Sorting** - Activities sorted by timestamp
2. **Gap Detection** - Identify missing time periods
3. **Interpolation** - Fill gaps using behavioral patterns
4. **Summarization** - Generate human-readable narratives
5. **Key Event Extraction** - Highlight significant activities

Predictive Monitoring with Explainability

Machine Learning Pipeline

Feature Engineering Categories

Temporal Features

- Hour of day, day of week, time since last activity
- Peak activity hours, temporal consistency

Spatial Features

- Location frequency, unique locations visited
- Location transitions, movement patterns
- Location entropy and diversity

Behavioral Features

- Activity regularity, movement consistency
- Department-specific patterns
- Role-based behavior models

Model Architecture

Ensemble Approach

- XGBoost: Primary model for non-linear relationships
- Random Forest: Robust backup model
- Ensemble Weighting: Confidence-based combination

Location Hierarchy

- Reduces 100+ specific locations to 15 categories
- Examples: LAB, LIBRARY, HOSTEL, CAFETERIA, ADMIN
- Improves prediction accuracy and generalization

Explainability Framework

Evidence Generation

Temporal Evidence

- "Peak activity hour: High probability based on historical patterns"
- "Most active hour: Typically very active during this time"

Location Evidence

- "Frequently visited: 15+ previous visits (strong habit)"
- "Department pattern: Common location for CS students"

Behavioral Evidence

- "Recent movement pattern: Library → Lab → Cafeteria"
- "Consistent weekly pattern: Same location this time last week"

Model Performance

Prediction Accuracy

Model	Accuracy	Top-3 Accuracy	Avg Confidence
XGBoost	87.3%	93.8%	0.74
Random Forest	85.1%	92.1%	0.71
Ensemble	88.2%	94.7%	0.76

Feature Importance

Feature	Importance	Description
Current Hour	18.2%	Temporal activity patterns
Previous Location Sequence	15.7%	Movement history
Department Preferences	12.3%	Academic alignment
Location Frequency	11.8%	Visit patterns
Day of Week	9.5%	Weekly rhythms

Performance Analysis

Scalability Testing

Scale	Entities	Processing Time	Memory Usage
Small	1,000	45 seconds	512 MB
Medium	10,000	6 minutes	2.1 GB
Large	50,000	28 minutes	8.7 GB

Optimization Strategies

Batch Processing

- Process entities in batches of 1,000
- Parallel processing for independent operations
- Memory-efficient data structures

Caching Strategy

- Feature caching for frequent entities
- Model persistence to avoid retraining
- Result caching for repeated queries

Privacy Safeguards

Data Protection Measures

Anonymization Techniques

- k-Anonymity: Each entity indistinguishable from k-1 others
- Differential Privacy: Calibrated noise in aggregate statistics
- Data Minimization: Collect only essential attributes

Access Control Framework

Failure Mode Analysis

Risk Assessment Matrix

Failure Scenario	Probability	Impact	Mitigation Strategy
Missing Data	High	Medium	Multiple imputation, historical patterns
Noisy Data	Medium	High	Outlier detection, confidence thresholds
Model Degradation	Low	High	Continuous monitoring, automatic retraining
System Integration Failure	Medium	High	Health checks, fallback mechanisms

Robustness Measures

Input Validation

- Timestamp format validation
- Identifier uniqueness checks
- Data completeness verification
- Value range validation

Graceful Degradation

- Fallback to rule-based predictions when ML fails
- Cached data during source unavailability
- Progressive enhancement based on data quality

Security Dashboard Implementation

User Interface Features

Entity Search & Selection

- Dropdown with all available entities
- Real-time search and filtering
- Quick entity switching

Security Status Monitoring

- 12-hour inactivity alerts (configurable)
- Color-coded status indicators
- Last seen location and timestamp

Activity Timeline

- Chronological activity display
- Multi-source activity fusion
- Interactive timeline navigation

Predictive Insights

- ML-powered location predictions
- Evidence-based explanations
- Alternative scenario predictions

Alerting System

Inactivity Alerts

- Configurable threshold (1-24 hours)
- Escalation procedures
- Notification channels

Anomaly Detection

- Behavioral pattern deviations
- Unusual location access
- Temporal pattern violations

Conclusion & Future Work

Key Achievements

- ☑ **High Accuracy Resolution:** 94.2% precision in entity matching
- ☑ **Explainable AI:** Transparent evidence chains for all predictions
- ☑ **Scalable Architecture:** Efficient processing of 50,000+ entities
- ☑ **Privacy by Design:** Comprehensive data protection measures
- ☑ **Real-time Monitoring:** Proactive security alerts and insights
- ☑ **Multi-Modal Integration:** Effective fusion of diverse data sources

Future Roadmap

1. **Real-time Stream Processing** - Apache Kafka integration
2. **Deep Learning Enhancement** - Transformer models for sequences
3. **Federated Learning** - Privacy-preserving multi-campus training
4. **Advanced Anomaly Detection** - Graph neural networks
5. **Mobile Integration** - Real-time alerts and reporting

Software Stack

- Python 3.8+ with scientific computing libraries
- Machine Learning: XGBoost, Scikit-learn, Joblib
- Web Framework: Streamlit,

Data Schema

Entity Profile Structure

```
{
  "entity_id": "E106861",
  "profile_info": {
    "name": "John Doe",
    "role": "student",
    "department": "Computer Science",
    "email": "john.doe@campus.edu",
    "all_identifiers": ["E106861", "S88432", "C1107"]
  },
  "activity_timeline": [...],
  "behavioral_patterns": {...},
  "location_analysis": {...},
  "temporal_analysis": {...}
}
```

```
### Prediction and Evidence
**PREDICTION FOR: E100030**
=====

**PREDICTION RESULTS:**
  **Location Category:** LAB
  **Confidence:** 0.821
  **Likely Specific Locations:** LAB_305, LAB_101, LAB_102

**EVIDENCE**
  1. **Current time** 09:00 (morning)
  2. **Morning pattern** High academic activity, classes/labs
  3. **Recent movement pattern** AP_CAF_3 → LAB
  4. **Department context** CSE students often visit LAB_101, AUDITORIUM
  5. **Student pattern** Likely in academic activities during day
  6. **New pattern** First predicted visit to this location
  7. **Weekend pattern** Different movement patterns expected
  8. **High confidence** Strong patterns support this prediction

**ALTERNATIVE PREDICTIONS**
  • **CAFETERIA** (confidence: 0.102)
  • **GYM** (confidence: 0.017)
```