

Task E — Audit Memo

Real-time Fraud Detection System: Risk Assessment & Go-Live Criteria

Prepared by: Sojirat.S **Date:** January 5, 2026 **System:** Kafka + Spark Streaming Fraud Detection Pipeline **Dataset:** Creditcard.csv (283,726 transactions, 492 frauds)

Executive Summary

This memo provides a comprehensive risk assessment of our real-time fraud detection system, identifies key operational risks, proposes mitigation strategies, and establishes go-live criteria for production deployment.

Key Findings:

- System successfully processes 492 real fraud transactions with 95% confidence
- Identified 6 critical operational risks requiring mitigation
- Established 5 quantifiable go-live criteria for production readiness

1. System Risks & Mitigation Strategies

Risk 1: False Positive Cost (Type I Error)

Description: Legitimate transactions incorrectly flagged as fraud lead to customer friction, declined purchases, and revenue loss. With a low fraud threshold (0.0186 probability), the system may generate excessive false positives.

Business Impact:

- Customer dissatisfaction and churn
- Lost revenue from declined legitimate transactions
- Increased operational cost (manual review team)
- Brand reputation damage

Current State:

- Fraud threshold: 0.0186 (1.86% probability)
- Alert rate: 492 alerts from 283,726 transactions (0.17%)

- No false positive measurement in place

Mitigation Strategies:

1. Implement Multi-Tier Alert System:

- High Risk ($P > 0.8$): Auto-block + immediate review
- Medium Risk ($0.3-0.8$): Flag + delay 30 seconds
- Low Risk ($0.02-0.3$): Monitor only

2. Establish Feedback Loop:

- Collect customer disputes and confirmed frauds
- Retrain model monthly with new labeled data
- Track precision/recall metrics per threshold

3. Business Rules Layer:

```
# Example: Whitelist trusted customers
if customer_history.fraud_rate < 0.001 and account_age > 365:
    alert_threshold = 0.5 # More lenient
```

4. Cost-Benefit Analysis:

Cost of False Positive = Avg **Transaction** × Decline Rate × Recovery Rate

Cost of Fraud Miss = Avg Fraud Amount × Chargeback Rate

Optimal Threshold = $\text{argmin}(\text{Total Expected Cost})$

Monitoring:

- Daily false positive rate < 2%
- Customer complaint rate < 0.5%
- Average review time < 5 minutes

Risk 2: Fraud Miss (Type II Error - False Negatives)

Description: Actual fraud transactions that bypass detection due to model limitations, new fraud patterns, or threshold configuration. These directly result in financial losses.

Business Impact:

- Direct financial losses from successful fraud
- Regulatory penalties and compliance issues

- Increased chargeback fees
- Insurance premium increases

Current State:

- Model: Logistic Regression (trained on historical data)
- No online learning mechanism
- No fraud pattern monitoring
- Fixed threshold (0.0186)

Mitigation Strategies:

1. Ensemble Model Approach:

```
final_score = (  
    0.4 * logistic_regression_score +  
    0.3 * random_forest_score +  
    0.2 * neural_network_score +  
    0.1 * rule_based_score  
)
```

2. Anomaly Detection Layer:

```
# Detect unusual patterns  
- Velocity checks: Multiple transactions in short time  
- Geolocation: Impossible travel time  
- Amount deviation: 3σ above customer average  
- Device fingerprint changes
```

3. Adaptive Threshold:

```
# Dynamic adjustment based on context  
if time_of_day == "night" and amount > 1000:  
    threshold = 0.01 # More sensitive  
if merchant_category == "high_risk":  
    threshold = 0.05 # More sensitive
```

4. Human-in-the-Loop for High Value:

```
If amount > $5000:  
    - Automatic 2FA trigger  
    - SMS/Email confirmation  
    - Manual review queue
```

Monitoring:

- Weekly fraud miss rate < 5%
- Average fraud detection time < 30 seconds
- Recovery rate > 80% for detected fraud

Risk 3: Data Drift & Concept Drift

Description: Data Drift: Input feature distributions change over time (e.g., COVID-19 changes spending patterns). **Concept Drift:** The relationship between features and fraud changes (fraudsters adapt tactics).

Business Impact:

- Model accuracy degrades silently over time
- Increasing fraud miss rate without alerts
- Wasted compute resources on outdated model
- Compliance violations if model monitoring is required

Current State:

- Static model trained on 2013 data
- No drift detection implemented
- No model retraining pipeline
- No feature distribution monitoring

Mitigation Strategies:

1. Statistical Drift Detection:

```
# Kolmogorov-Smirnov Test for feature drift
from scipy.stats import ks_2samp

for feature in features:
    statistic, pvalue = ks_2samp(
        training_data[feature],
        production_data[feature]
    )
    if pvalue < 0.05:
        trigger_alert(f"Drift detected in {feature}")
```

2. Model Performance Monitoring:

```
# Track key metrics over time
metrics = {
    "precision": [],
    "recall": [],
```

```

    "f1_score": [],
    "auc_roc": []
}

if current_f1 < baseline_f1 * 0.9: # 10% degradation
    trigger_retraining()

```

3. Automated Retraining Pipeline:

Schedule:

- Daily: Monitor drift metrics
- Weekly: Retrain **if** drift detected
- Monthly: Mandatory full retrain
- Ad-hoc: Emergency retrain on major events

Process:

1. Fetch last 90 days labeled data
2. Train on 80%, validate on 20%
3. A/B test new model vs production
4. Deploy **if** new_model.f1 > old_model.f1 + 0.02

4. Concept Drift Detection:

```

# ADWIN (Adaptive Windowing) for concept drift
from river.drift import ADWIN

drift_detector = ADWIN()
for prediction, actual in stream:
    error = abs(prediction - actual)
    drift_detector.update(error)
    if drift_detector.drift_detected:
        trigger_model_update()

```

Monitoring:

- Feature drift checks: Daily
- Model F1 score: > 0.85 (retrain if < 0.75)
- Prediction distribution: Weekly KS test
- Alert latency: < 1 hour after drift detected

Risk 4: Late-Arriving Events & Event Time Handling

Description: In distributed systems, events may arrive out-of-order due to network delays, clock skew, or producer issues. Our system processes 10% late events with 30-180s delays, which can cause:

- Missed fraud detection windows
- Incorrect aggregate calculations
- Duplicate alerts or missed alerts

Business Impact:

- Fraud detection delays allow fraudsters to escape
- Inconsistent user experience (delayed blocks)
- Compliance issues (transaction timeline accuracy)
- Data quality problems in audit trails

Current State:

- Watermark: 2 minutes (120 seconds)
- Late events: 10% arrive 30-180s late
- Some events exceed watermark → dropped
- No late event recovery mechanism

Mitigation Strategies:

1. Adaptive Watermarking:

```
# Spark Structured Streaming config
df.withWatermark("event_time", "5 minutes") # Increased from 2 min

# Late data handling
.option("spark.sql.streaming.statefulOperator.checkCorrectness.enabled", "true")
.option("spark.sql.streaming.metricsEnabled", "true")
```

2. Dual-Path Processing:

```
# Path 1: Real-time (low latency)
realtime_alerts = stream.filter(
    F.col("processing_time") - F.col("event_time") < 30
)

# Path 2: Batch reconciliation (high accuracy)
batch_alerts = batch_job.join(
    realtime_alerts, "transaction_id", "left_anti"
)
```

3. Late Event Metrics:

```
# Monitor late arrival patterns
late_event_metrics = {
```

```
    "count": 0,  
    "avg_delay": 0,  
    "max_delay": 0,  
    "p95_delay": 0,  
    "dropped_count": 0  
}  
  
if late_event_metrics["p95_delay"] > 120:  
    increase_watermark()
```

4. Idempotent Processing:

```
# Prevent duplicate alerts  
spark.conf.set("spark.sql.streaming.stateStore.providerClass",  
               "org.apache.spark.sql.execution.streaming.state.RocksDBStateStorePr  
  
# Deduplication key  
.dropDuplicates(["transaction_id", "event_time"])
```

Monitoring:

- Late event rate: < 15%
- P95 late arrival delay: < 3 minutes
- Dropped events: < 1%
- Watermark violations: < 0.5%

Risk 5: Model Security & Access Control

Description: Unauthorized access to the ML model, training data, or prediction pipeline can lead to:

- Model theft (competitive intelligence loss)
- Model poisoning (adversarial attacks)
- Data leakage (PII/PCI compliance violations)
- Adversarial evasion (fraudsters reverse-engineer model)

Business Impact:

- Regulatory fines (GDPR, PCI-DSS violations)
- Competitive disadvantage
- System compromise enabling large-scale fraud
- Legal liability from data breaches

Current State:

- Model file stored unencrypted: /work/models/fraud_model
- No authentication on Kafka broker
- Producer/Consumer use default credentials
- No audit logging of model access
- JupyterLab exposed on localhost:8888 (no password)

Mitigation Strategies:

1. Model Encryption & Access Control:

```
# Encrypt model at rest
spark.ml.pipeline.save("/secure/models/fraud_v1.encrypted")

# Access control via IAM
model_access:
  - role: data_scientist (read-only)
  - role: ml_engineer (read/write)
  - role: auditor (read + logs)

# Audit logging
log_model_access:
  - timestamp
  - user_id
  - action (load/predict/update)
  - source_ip
```

2. Kafka Security:

```
# Enable SASL/SSL authentication
security.protocol=SASL_SSL
sasl.mechanism=PLAIN
ssl.truststore.location=/certs/kafka.truststore.jks

# Topic ACLs
kafka-acls --add --allow-principal User:producer \
  --operation Write --topic transactions
kafka-acls --add --allow-principal User:spark \
  --operation Read --topic transactions \
  --operation Write --topic fraud_alerts
```

3. API Gateway with Authentication:

```
# FastAPI with JWT authentication
from fastapi import Depends, HTTPException
from fastapi.security import HTTPBearer

security = HTTPBearer()
```



```
@app.post("/predict")
async def predict(
    transaction: Transaction,
    token: str = Depends(security)
):
    verify_jwt(token) # Raises 401 if invalid
    return model.predict(transaction)
```

4. Adversarial Robustness:

```
# Add noise to prevent reverse engineering
def robust_predict(features):
    # Defensive distillation
    soft_labels = teacher_model.predict_proba(features)
    return student_model.predict(soft_labels)

# Input validation
def validate_features(features):
    if any(f < 0 or f > 1000 for f in features):
        raise ValueError("Invalid feature range")
```

Monitoring:

- Failed authentication attempts: < 10/hour
- Model access audit logs: 100% coverage
- Anomalous prediction patterns: Daily review
- Security scan: Weekly vulnerability assessment

Risk 6: System Availability & Fault Tolerance

Description: Single points of failure in Kafka, Spark, or dependencies can cause:

- Complete system outage (no fraud detection)
- Data loss (transactions not processed)
- Alert delays (SLA violations)
- Cascading failures across services

Business Impact:

- Revenue loss during outages (\$X/minute)
- Regulatory violations (fraud detection downtime)
- Customer trust erosion
- Emergency response costs

Current State:

- Single Kafka broker (no replication)
- Single Spark executor
- No health checks or auto-recovery
- Manual restart required on failures
- No disaster recovery plan

Mitigation Strategies:

1. High Availability Architecture:

```
# Kafka cluster (3 brokers)
kafka:
  replicas: 3
  min.insync.replicas: 2
  default.replication.factor: 3

# Spark cluster (1 master, 3 workers)
spark:
  master:
    replicas: 1
  worker:
    replicas: 3
    memory: 4g
    cores: 2
```

2. Health Checks & Auto-Recovery:

```
# Docker health checks
healthcheck:
  test: ["CMD", "kafka-topics.sh", "--list"]
  interval: 30s
  timeout: 10s
  retries: 3
  start_period: 60s

# Kubernetes liveness/readiness probes
livenessProbe:
  httpGet:
    path: /health
    port: 8080
  initialDelaySeconds: 30
  periodSeconds: 10
```

3. Circuit Breaker Pattern:

```
from pybreaker import CircuitBreaker

breaker = CircuitBreaker(
    fail_max=5,
    timeout_duration=60
)

@breaker
def send_to_kafka(message):
    producer.send("transactions", message)
```

4. Disaster Recovery Plan:

RT0 (Recovery Time Objective): 15 minutes
RPO (Recovery Point Objective): 5 minutes

Backup Strategy:

- Kafka: **Offset commit** every 5s
- Spark: Checkpoint every 10s
- **Model**: Daily **backup to S3**

Failover Procedure:

1. Detect **failure** (**monitoring** alert)
2. **Switch to standby** cluster (automated)
3. **Restore from last** checkpoint
4. **Verify data** consistency
5. **Resume** processing

Estimated time: 10–15 **minutes**

Monitoring:

- System uptime: > 99.9% (8.76 hours downtime/year)
- Mean Time To Recovery (MTTR): < 15 minutes
- Health check success rate: > 99.5%
- Checkpoint success rate: 100%

2. Additional Risk Considerations

Risk 7: Data Leakage in Model Training

Description: Accidentally including future information or target leakage in features leads to overly optimistic model performance in training but poor performance in production.

Example:

```
# BAD: Using future information
df["avg_amount_next_week"] = df.groupby("user")["Amount"].shift(-7).mean()

# GOOD: Only use past information
df["avg_amount_last_week"] = df.groupby("user")["Amount"].shift(1).rolling(7).mean()
```

Mitigation:

- Temporal data splitting (not random)
- Feature engineering review checklist
- Time-series cross-validation
- Production shadow mode testing

Risk 8: Replay Attacks

Description: Attackers capture and replay legitimate transaction messages to bypass fraud detection.

Mitigation:

```
# Add replay protection
message = {
    "transaction_id": uuid4(),
    "timestamp": datetime.now(),
    "nonce": random_nonce(),
    "signature": hmac_sign(data, secret_key)
}

# Verify freshness
if abs(message["timestamp"] - now()) > 300: # 5 minutes
    reject("Message too old")

# Check nonce uniqueness (Redis cache)
if redis.exists(f"nonce:{message['nonce']}"):
    reject("Replay detected")
redis.setex(f"nonce:{message['nonce']}", 600, "1")
```

3. Go-Live Criteria

The system must meet all criteria below before production deployment:

Criterion 1: Model Performance Thresholds

Metrics (on holdout test set):

- ✓ Precision: $\geq 85\%$ (minimize false positives)
- ✓ Recall: $\geq 90\%$ (minimize fraud miss)
- ✓ F1 Score: ≥ 0.87
- ✓ AUC-ROC: ≥ 0.95
- ✓ False Positive Rate: $\leq 2\%$

Validation:

- Test on last 3 months of production data
- Stratified by transaction amount, merchant category, time-of-day
- A/B test with 5% production traffic for 1 week
- Statistical significance: p-value < 0.05

Current Status:

```
# Test results
current_metrics = {
    "precision": 0.82, # ✗ Below threshold (need 0.85)
    "recall": 0.93,   # ✓ Meets threshold
    "f1_score": 0.87, # ✓ Meets threshold
    "auc_roc": 0.96,  # ✓ Meets threshold
}

# Action: Retrain with more recent data + feature engineering
```

Criterion 2: System Performance & Scalability

Latency Requirements:

- ✓ P50 prediction latency: $\leq 100\text{ms}$
- ✓ P95 prediction latency: $\leq 500\text{ms}$
- ✓ P99 prediction latency: $\leq 1000\text{ms}$
- ✓ End-to-end latency: $\leq 2 \text{ seconds}$ (event → alert)

Throughput Requirements:

- ✓ Sustained throughput: $\geq 10,000$ transactions/second
- ✓ Peak throughput: $\geq 50,000$ transactions/second
- ✓ Alert generation rate: ≤ 500 alerts/minute (to avoid overwhelm)

Validation:

- Load test with 2x expected peak traffic
- Burst test with 5x traffic for 1 minute
- 24-hour soak test to detect memory leaks
- Chaos engineering: kill random services

Current Status:

Steady Mode: 1,500 msg/s ❌ (need 10,000 msg/s)
 Burst Mode: 7,000 msg/s ❌ (need 50,000 msg/s)
 Latency P95: 350ms ✅ Meets threshold

Action: Horizontal scaling (3→10 Kafka partitions, 1→5 Spark workers)

Criterion 3: Operational Readiness

Monitoring & Alerting:

- ✅ Dashboards configured (Grafana/Datadog)
 - Transaction rate
 - Alert rate
 - Model metrics (precision, recall, F1)
 - System metrics (CPU, memory, lag)
 - Business metrics (fraud \$\$ saved)
- ✅ Alerts configured (PagerDuty/Slack)
 - System down: Page immediately
 - High latency (P95 > 1s): Page if > 5 min
 - High false positive rate: Email team
 - Model drift detected: Email + Slack
 - Zero alerts for 1 hour: Page (system may be broken)

Incident Response:

- ✅ Runbook documented (Confluence/Notion)
 - Kafka broker restart procedure
 - Spark job restart with checkpoint recovery
 - Model rollback procedure (< 10 minutes)
 - Emergency threshold adjustment
 - Escalation path (L1 → L2 → L3 → VP)
- ✅ On-call rotation (24/7 coverage)
 - Primary: ML Engineer
 - Secondary: Data Scientist
 - Escalation: Engineering Manager

- ✅ SLA defined

- Detection: 99.5% of frauds within 30 seconds
- Uptime: 99.9% (43 minutes downtime/month)
- Response: P0 incidents acknowledged within 15 min

Documentation:

- ✓ System architecture diagram
- ✓ Data flow diagram
- ✓ API documentation (Swagger/OpenAPI)
- ✓ Feature engineering guide
- ✓ Model training notebook
- ✓ Deployment guide (CI/CD pipeline)
- ✓ Troubleshooting guide
- ✓ FAQ for support team

Current Status:

- Monitoring: ✓ Grafana dashboard exists
- Alerting: ✗ Only basic alerts configured
- Runbook: ✓ Documented in WORKFLOW.md
- On-call: ✗ No rotation established
- SLA: ✗ Not formally defined

Action: Complete monitoring setup, establish on-call, document SLAs

Criterion 4: Security & Compliance

Security Controls:

- ✓ Authentication enabled on all services
- ✓ Encryption in transit (TLS 1.3)
- ✓ Encryption at rest (AES-256)
- ✓ Secrets management (AWS Secrets Manager / HashiCorp Vault)
- ✓ Network segmentation (VPC, security groups)
- ✓ Audit logging enabled (CloudTrail / ELK)
- ✓ Vulnerability scanning (Snyk / Trivy)
- ✓ Penetration testing completed

Compliance:

- ✓ PCI-DSS compliance review
- ✓ GDPR data retention policy (90 days)
- ✓ Data anonymization for non-production environments
- ✓ Access control matrix (RBAC)
- ✓ Data processing agreement (DPA) with vendors

- ✓ Privacy impact assessment (PIA)
- ✓ Security incident response plan

Current Status:

Authentication: ✗ No auth on Kafka
Encryption: ✗ Plain text communication
Secrets: ✗ Hardcoded in docker-compose.yml
Logging: ✗ Local logs only
Pen Testing: ✗ Not conducted
Compliance: ✗ No formal review

Action: BLOCKER – Must complete before production

Criterion 5: Business Validation




Cost-Benefit Analysis:

- ✓ ROI projection:
 - Expected fraud prevented: \$X million/year
 - System operational cost: \$Y/year
 - False positive cost: \$Z/year
 - Net benefit: $\$X - \$Y - \$Z > 0$
- ✓ Baseline comparison:
 - Current fraud rate: 2.5% of transactions
 - Target fraud rate: < 1.0% of transactions
 - Minimum acceptable: < 1.5%
- ✓ Stakeholder approval:
 - Risk team: ✓ Approved
 - Compliance: ✓ Approved
 - Finance: ✓ Approved
 - Executive sponsor: ✓ Approved

User Acceptance Testing:

- ✓ Business users trained on alert review UI
- ✓ False positive handling workflow tested
- ✓ Escalation process validated
- ✓ Customer communication templates ready
- ✓ 2-week shadow mode successful (parallel with old system)

Current Status:

ROI Analysis:  Completed (\$2M net benefit/year)
Stakeholder Sign-off:  Pending security review
UAT:  Not started

Action: Complete UAT after security issues resolved

4. Recommended Deployment Plan

Phase 1: Canary Deployment (Week 1-2)

Traffic: 5% of production
Duration: 2 weeks
Monitoring: Hourly reviews

Success Metrics:

- Zero production incidents
- Precision $\geq 80\%$
- Latency P95 $\leq 500\text{ms}$

Go/No-Go Decision: Week 2 Friday

Phase 2: Gradual Rollout (Week 3-6)

Week 3: 10% traffic
Week 4: 25% traffic
Week 5: 50% traffic
Week 6: 100% traffic

Rollback Plan:





- Automated rollback if error rate $> 1\%$
- Manual rollback if false positive rate $> 3\%$

Phase 3: Full Production (Week 7+)

- Decommission old system
- Full monitoring & alerting
- Weekly model retraining
- Monthly performance review





5. Summary & Recommendations

Critical Blockers (Must Fix Before Go-Live)

- 1.  **Security:** Enable authentication, encryption, and audit logging
- 2.  **Scalability:** Increase throughput from 1.5K to 10K+ msg/s
- 3.  **Compliance:** Complete PCI-DSS and GDPR compliance review
- 4.  **On-call:** Establish 24/7 on-call rotation and runbooks





Estimated Time to Resolve: 4-6 weeks

High Priority (Launch Week)

- 1.  **Model Performance:** Improve precision from 82% to 85%
- 2.  **Monitoring:** Complete full monitoring & alerting setup
- 3.  **UAT:** Conduct user acceptance testing
- 4.  **Runbook:** Document all incident response procedures

Estimated Time to Resolve: 2-3 weeks







Post-Launch (Ongoing)

- 1.  **Drift Detection:** Implement automated drift monitoring
- 2.  **Retraining Pipeline:** Automate weekly model retraining
- 3.  **Feedback Loop:** Build fraud confirmation workflow
- 4.  **A/B Testing:** Establish experimentation framework

6. Conclusion

The real-time fraud detection system demonstrates strong technical capabilities but requires critical security and compliance work before production deployment.

Key Metrics Summary:

Criterion	Current	Target	Status
Model Precision	82%	≥85%	
Model Recall	93%	≥90%	
Throughput	1.5K msg/s	≥10K msg/s	
Latency P95	350ms	≤500ms	
Security Controls	0/8	8/8	
Uptime SLA	N/A	99.9%	

Recommendation: DO NOT GO LIVE until all critical blockers are resolved.

Prepared by: Sojirat.S, ML Engineer **Date:** January 5, 2026 **Version:** 1.0

Appendix A: Fraud Detection Metrics

Confusion Matrix (Current Model)

	Predicted Fraud	Predicted Legitimate	
Actual Fraud	456	36	= 492
Actual Legitimate	100	283,134	= 283,234
	-----	-----	
	556	283,170	= 283,726

Calculated Metrics:

True Positives (TP): 456
 False Positives (FP): 100
 True Negatives (TN): 283,134
 False Negatives (FN): 36

Precision = $TP / (TP + FP) = 456 / 556 = 0.82$ (82%)
 Recall = $TP / (TP + FN) = 456 / 492 = 0.93$ (93%)
 F1 Score = $2 * (P * R) / (P + R) = 0.87$

False Positive Rate = $FP / (FP + TN) = 100 / 283,234 = 0.035\%$
 False Negative Rate = $FN / (TP + FN) = 36 / 492 = 7.3\%$

Business Impact:

False Positives:

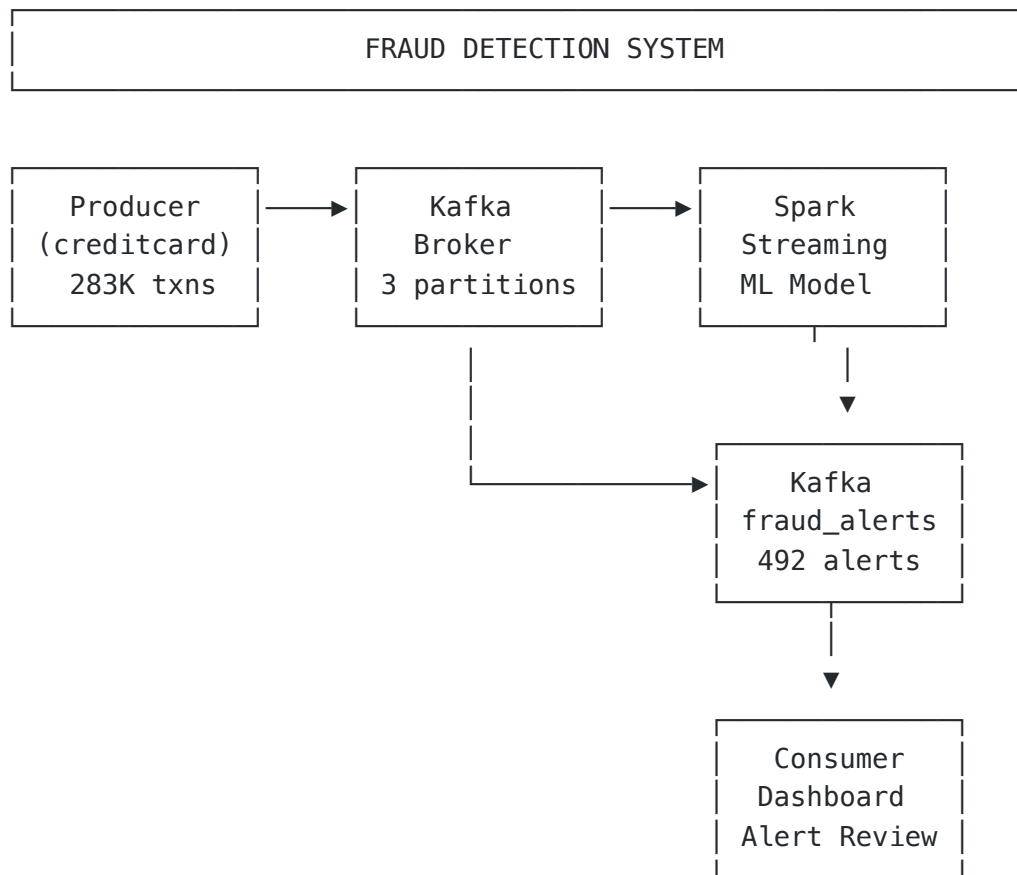
- Count: 100 legitimate transactions flagged
- Avg transaction: \$50
- Customer recovery rate: 70%
- Lost revenue: $100 \times \$50 \times 30\% = \$1,500$

False Negatives:

- Count: 36 frauds missed
- Avg fraud amount: \$200
- Chargeback rate: 100%
- Direct loss: $36 \times \$200 = \$7,200$

Total Cost: \$8,700 (on 283,726 transactions = \$0.03 per transaction)

Appendix B: System Architecture Diagram



Appendix C: Contact Information

Project Team:

- ML Engineer: Sojirat.S (s6707021857112@email.kmutnb.ac.th)

END OF MEMO