# Human-in-the-Loop is Not Enough

## The Persistent and Evolving Limits of Human Oversight in High-Stakes AI

# About me

**Ankur Vatsa**

**Experience, Expertise and Education:**

- Built SaaS-based regulatory reporting systems with ML-powered validation
- Led AI/ML adoption across 180+ applications at global investment banks
- Specialist in AI-powered regulatory reporting, risk analytics, and compliance automation
- Expert in Human-AI collaboration patterns in high-stakes financial environments
- AWS Solutions Architect – Professional, with 17+ years in global investment banking technology
- M.Tech. in Artificial Intelligence and Machine Learning – BITS Pilani

# Another look at

- **HITL** "Human in the loop" as a remedy we commonly bank upon

- **Diagnose** HITL failure modes in AI systems

- **Evaluate** the insufficiencies of HITL and safeguards to **Consider**

- **Implement** discriminator agents, monitoring, and governance controls

- **Prioritize** remediation steps using a risk-based roadmap

# Human-in-the-Loop (HITL)

- **Common remedy** where intelligent systems hand over control to humans in several key scenarios:

| Key scenario | Description |
|---|---|
| **Errors and anomalies observed** | AI detects anomalies or unexpected outputs |
| | System confidence drops below threshold |
| | Error conditions trigger human intervention |
| **Low risk appetite** | High-stakes decisions require human approval |
| | Regulatory compliance mandates human oversight |
| | Business-critical processes need human validation |

- not just "humans reviewing AI outputs" — it's a structured handover mechanism

# HITL Triggers (Continued)

| Key scenario | Description |
| --- | --- |
| **Need for domain expertise** | Complex edge cases beyond AI training scope |
| | Nuanced judgment calls requiring contextual understanding |
| | Ethical considerations and value-based decisions |
| **Legal/regulatory mandates** | Healthcare diagnostics requiring physician approval |
| | Financial lending decisions with fair lending requirements |
| | Criminal justice applications with due process rights |

- HITL offers defence to AI limitations; but faces issues — *speed, scale, bias, skillsgap*

# The Case for Enhanced AI Oversight

- **$460M in 45 minutes**: Knight Capital flash crash (2012)

- **10,000 families wrongfully penalized**: Dutch childcare scandal

- **30% error miss rate**: Hospital audits of AI diagnostic reviews

- **95% false positive rate**: AML alert systems overwhelming investigators

HITL alone is insufficient for high-stakes AI deployment.

# Five Failure Modes of HITL

1. **Speed Mismatch**: Microsecond AI vs. second-scale humans *(Knight Capital: 45 minutes, $460M)*

2. **Scale Mismatch**: 10,000+ daily alerts vs. limited reviewers *(AML: 95% false positives)*

3. **Overtrust**: Automation bias hides critical errors *(Hospital study: 30% miss rate)*

4. **Skill Erosion**: "Out-of-the-loop" degraded expertise *(Tesla Autopilot incidents)*

5. **Coordination Gaps**: Unclear handovers and protocols *(Dutch childcare: role confusion)*

# Financial Services: When Structure Deceives

**The Problem:**

- AI extracts regulatory requirements from EMIR, MiFID, CFTC rules

- **Hallucination risk**: LLMs invent non-existent field requirements

- **False confidence**: Well-formatted JSON output appears trustworthy

- **Scale challenge**: Hundreds of fields vs. limited compliance staff

**Real Impact:**

- Potential fines: >$50M for insufficient swap reporting

- Deutsche Bank: $150M AML penalty for ongoing failures

# Healthcare: Life-or-Death Automation Bias

**The Evidence:**

- Hospital audits: Clinicians miss **30% of AI diagnostic errors**

- Emergency departments: Time pressure prevents thorough AI output validation

- Training gaps: Physicians lack awareness of AI system limitations

**Why HITL Fails:**

- **Automation bias**: Trusting confident AI assessments of "normal" results

- **Mode confusion**: Unclear guidance on when to override AI recommendations

- **Workload pressure**: No time for careful review during patient surges

# Transportation: Split-Second Decisions

**NHTSA Findings (2023):**

- Multiple fatal Tesla Autopilot incidents
- Drivers had **seconds** to react, weren't ready to take control
- "Out-of-the-loop" problem: Skill erosion from over-reliance

**The Core Issue:**

- AI operates in **milliseconds**
- Human reaction time: **1-3 seconds**
- Handover complexity: Mode awareness, situational context, skill maintenance

# Discriminator Agents: Automated Validation

**What they are:**

Specialized AI systems that detect errors, hallucinations, or inconsistencies in other AI outputs

**Architecture Pattern:**

```
Source → Primary AI → Output → Discriminator AI → Flagged Items → Human Review
```

**Types:**

- **Binary classifiers**: Error/No-error detection

- **Consistency checkers**: Cross-reference with source material

- **Adversarial critics**: Challenge model outputs

- **Provenance verifiers**: Trace data lineage

# Concrete Monitoring Metrics

**Performance Thresholds:**

- **Model drift**: Alert if accuracy drops >5% from baseline

- **Human interception**: Warn if <80% error catch rate

- **Confidence escalation**: Review if AI confidence <70% on critical outputs

- **Anomaly rate**: Flag if >10% of outputs marked unusual in 24h

**Real Examples:**

- **GLUE/SuperGLUE**: NLP benchmarks, target >90% accuracy

- **ImageNet**: Computer vision, target >95% top-5 accuracy

- **Financial datasets**: <1% hallucination rate in regulatory extraction

# 12-Month Implementation Roadmap

**Phase 1: Foundation (0-3 months)**

✓ Establish baselines with benchmark datasets

✓ Implement structured logging and audit trails

✓ Define concrete alerting thresholds

**Phase 2: Detection (3-6 months)**

✓ Deploy continuous drift monitoring

✓ Add discriminator agents for high-risk outputs

✓ Start measuring human interception rates

**Phase 3: Prevention (6-12 months)**

✓ Implement automated circuit breakers

✓ Add ensemble validation for critical decisions

✓ Establish incident response playbooks

# Audit Trail Schema

```json
{
  "event_id": "uuid",
  "timestamp": "2025-09-02T10:00:00Z",
  "input": {"source": "doc123"},
  "ai_output": {"model": "v1.2", "prediction": "...", "confidence": 0.85},
  "discriminator": {"model": "discA", "flag": true, "reason": "inconsistency"},
  "human_review": {"reviewer": "user456", "action": "approve", "notes": "..."},
  "outcome": {"final": "approved", "error_detected": false}
}
```

**Key Requirements:**

- Immutable storage for regulatory compliance

- Accessible to auditors with proper retention policies

- Captures full decision lineage for post-incident analysis

# Key Implementation Patterns

**Automated Circuit Breakers:**

- Financial: Trading halts when anomalies detected (<1 second response)

- Healthcare: Confidence-based escalation to senior clinicians

- Regulatory: Stop processing when hallucination rate exceeds threshold

**UI Design for Trust Calibration:**

- Show confidence scores and uncertainty ranges

- Highlight source material for AI conclusions

- Force deliberation on high-risk overrides

- Provide "second opinion" views from ensemble models

# Summary and Conclusions

**HITL is necessary but not sufficient**

**The evidence:**

- Knight Capital: $460M in 45 minutes (speed mismatch)

- Dutch childcare: 10,000 families wrongfully penalized (automation bias)

- Hospital study: 30% diagnostic error miss rate (overtrust)

**The solution:**

Layered defenses with discriminator agents, automated safeguards, and measurable monitoring

**Next steps:**

Start with benchmarks and baseline measurements—systems cannot be improved without measurement