

# **AI Governance in Customer Support An Operator's Guide**

Version 1.0 | December 2025

A technical framework for deploying AI safely in high-risk customer interactions.

Published by QuickmaticAI Pvt Ltd

# Introduction

AI in customer support promises efficiency, but introduces new operational risks. Uncontrolled models can generate hallucinations, make unauthorized commitments, or escalate inappropriately—each failure eroding customer trust and exposing liability.

This guide outlines a governance framework for operators managing AI-assisted support at scale. It addresses the unique challenges of B2C, fintech, and e-commerce environments where accuracy, auditability, and compliance are non-negotiable.

## 1. Why Black-Box AI Fails in Support

Most AI support tools operate as probabilistic systems: they predict responses based on training data, but lack deterministic constraints. This creates three critical failure modes:

### Hallucination

Models generate plausible-sounding but factually incorrect information. In support contexts, this manifests as invented policies, incorrect pricing, or nonexistent product features. While no system can eliminate hallucinations entirely, grounding responses in verified knowledge bases significantly reduces this risk.

### Escalation Collapse

Without explicit escalation rules, models either over-escalate (flooding human agents with trivial queries) or under-escalate (attempting to resolve high-risk interactions autonomously). Both outcomes degrade service quality.

### Silent Failures

Black-box systems provide no visibility into why a decision was made. When an AI responds incorrectly, operators cannot review the decision path, identify the root cause, or prevent recurrence. Compliance and audit requirements become impossible to satisfy.

## 2. Common Failure Modes

Operators report consistent failure patterns across unconstrained AI deployments:

### Response Looping

AI generates repetitive responses when it lacks sufficient context or encounters ambiguous queries. The customer becomes frustrated; the interaction stalls.

### Unauthorized Commitments

Models attempt to resolve complaints by offering refunds, credits, or policy exceptions. Without guardrails, these commitments create financial exposure and operational inconsistencies.

## PII Leakage

Without output filtering, AI may inadvertently expose personally identifiable information (PII) in responses—credit card numbers, SSNs, or account credentials. This violates GDPR, CCPA, and internal security policies.

## Misclassification Cascades

An initial misclassification (e.g., treating a refund request as a general inquiry) propagates through the interaction, resulting in irrelevant responses and escalation failures.

# 3. How Governable AI Should Behave

A governable AI system enforces constraints before generation, not after. Key principles:

## Input Validation Before AI Invocation

Validate inputs against policy rules before calling the AI model. Examples:

- Detect jailbreak attempts and prompt injection patterns
- Identify blacklisted topics (legal advice, medical diagnoses)
- Flag sentiment extremes (severe negativity, frustration)
- Redact PII before processing (SSNs, credit cards, Aadhaar numbers)

## Weighted Escalation Scoring

Escalation uses a numeric scoring system (0-10 scale) rather than binary rules. Signals accumulate:

### Escalation Additions:

- Low confidence: +1
- Repeated question: +2
- Confusion detected: +2
- Negative sentiment: +3
- Escalation keywords ("refund", "legal action"): +3
- Severe negative sentiment: +4
- Frustration: +10 (instant escalation)

### Escalation Deductions:

- Confirmation: -1
- Satisfaction: -2
- Agent already intervened: -5

When the score reaches 7 or higher, the interaction escalates to a human agent. This approach prevents over-escalation while ensuring high-risk cases receive human attention.

## Append-Only Audit Logging

Every decision is logged with full context: input validation results, escalation scores, and AI confidence levels. Logs are append-only with write access restricted to system roles, preventing tampering. Historical logs support compliance reviews and root cause analysis.

## Knowledge Base Grounding

Responses are grounded in verified knowledge bases when possible. While this significantly reduces hallucinations, edge cases may still require human review. The system prioritizes retrieval from trusted sources over generating unconstrained responses.

## 4. Design Principles Enorve Follows

Enorve implements these governance principles as core architectural features:

### Principle 1: No Silent Actions

Every AI decision is observable. Operators can review historical interactions to understand what validation rules triggered, which escalation signals fired, and why a particular action was taken.

### Principle 2: Human-in-the-Loop by Design

High-risk interactions escalate automatically based on weighted scoring. Operators define thresholds explicitly—the system does not attempt to 'learn' when escalation is appropriate.

### Principle 3: Constraints Before Generation

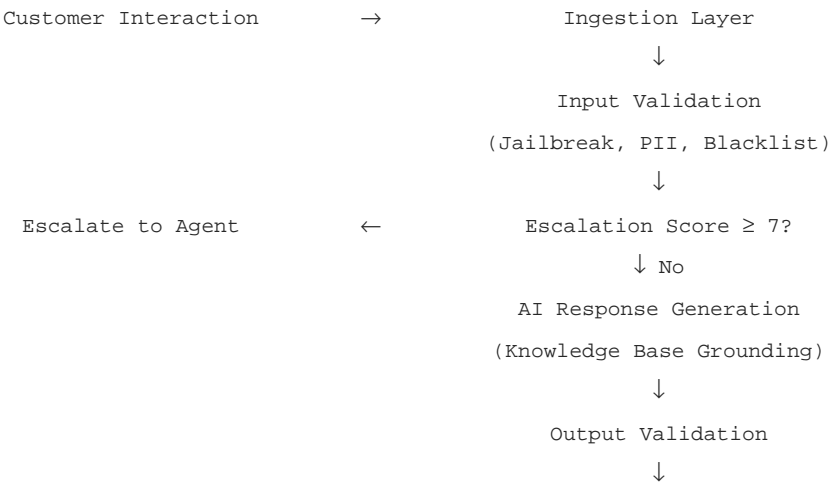
Policy validation occurs before AI invocation. This prevents the model from processing prohibited content, rather than filtering it after generation.

### Principle 4: Comprehensive PII Protection

Automatic detection and redaction of sensitive data including: SSNs, credit cards, Aadhaar numbers, PAN cards, UPI IDs, IFSC codes, email addresses, and phone numbers. Redaction modes support full masking, partial masking, or hashing based on operational requirements.

## 5. Architecture Diagram

A simplified view of Enorve's decision flow:



## 6. How Enorve Implements This

Enorve provides three core capabilities for operators:

### TypeScript-Based Configuration

Governance rules are defined in TypeScript with version control. Example escalation weights:

```
const DEFAULT_ESCALATION_WEIGHTS = {  
  frustration: 10, // Instant escalation  
  explicit_human_request: 3,  
  escalation_keywords: 3, // 'refund', 'legal'  
  negative_sentiment: 3,  
  satisfaction: -2 // Deduction  
};  
escalation_threshold: 7 // Score ≥ 7 triggers
```

### Real-Time Observability

Dashboard displays active interactions, escalation queue depth, guardrail trigger counts, and AI confidence distributions. Operators monitor system health in real time.

### Structured Audit Logs

Historical logs include structured metadata for compliance review: escalation scores, triggered signals, confidence levels, and validation results. Logs support forensic analysis and policy refinement.

## 7. Deployment Considerations

Operators deploying governed AI should address:

### Escalation Threshold Calibration

Start with the default threshold (score  $\geq 7$ ) and monitor for over-escalation or under-escalation. Adjust weights iteratively based on operational data. Avoid deploying overly permissive configurations in production.

### Agent Training

Human agents must understand escalation triggers and how to interpret decision logs. Provide training on weight configuration and audit trail review.

### PII Redaction Modes

Choose appropriate redaction modes for your use case: full masking (XXX-XX-XXXX), partial masking (XXX-XX-1234), or hashing. Balance security requirements with operational visibility needs.

## Incident Response

Define runbooks for AI failures: hallucinations, policy violations, or escalation floods. Establish clear ownership for incident triage and weight adjustment procedures.

## Conclusion

AI in customer support is inevitable, but uncontrolled deployment creates more problems than it solves. A governance-first approach—input validation, weighted escalation scoring, and append-only audit trails—enables operators to deploy AI safely in high-risk environments.

Enorve provides the infrastructure to implement this framework without requiring custom engineering. Operators retain full control over AI behavior, while gaining the efficiency benefits of automation.

---

### About Enorve

Enorve by QuickmaticAI Pvt Ltd is an enterprise AI support platform for governed resolution. Built for security-first teams operating at scale.

**Contact:** [support@enorve.com](mailto:support@enorve.com)

**Web:** [enorve.com](https://enorve.com)