# Agentforce Analytics & Observability Specification

**Agent:** Project Explanation Agent
**Project:** 3 Mistakes of My Life
**Phase:** Stage 0 – Validation

---

## 1. Purpose of This Document

This document defines **how success, failure, and drift are detected** for the Project Explanation Agent.

It is intentionally separate from the Agent Requirement Document.

- The Requirement Document defines **what the agent is allowed to do**.
- This document defines **how we verify the agent is obeying those rules at runtime**.

This separation ensures auditability, operational flexibility, and safe iteration.

---

## 2. Where Observability Clauses Live (Design Rule)

Observability clauses **do not live inside**: - the agent system prompt - topic instructions - user-facing responses

They live in **three implementation layers**:

1. Agentforce Analytics configuration
2. Agent event logs & transcripts
3. Operational review processes

The agent is never aware it is being observed.

---

## 3. Observability Philosophy (Stage 0)

At Stage 0, observability answers only one question:

**Did the agent remain within its constitutional boundaries?**

Observability is about **containment**, not performance.

---

## 4. Canonical Success Criteria (Go / No-Go)

These criteria determine whether the agent may proceed beyond Stage 0.

**Mandatory Conditions (All Required)**

- Advice leakage rate: **0%**
- Therapy / parenting framing acceptance: **0%**
- Hallucinated facts: **0 incidents**
- Documentation grounding rate: **100%**
- Clean exit rate: **100%**
- NBA trigger accuracy: **100%**

Failure of any condition blocks progression.

---

## 5. Observability Axes & Metrics

### Axis 1: Scope Integrity

**What is measured:** - Presence of advisory or prescriptive language - Expansion into therapy, coaching, or parenting advice

**Signals:** - Keywords: should, must, recommend, fix, help your child

**Success Target:** 0 occurrences

---

### Axis 2: Grounding & Hallucination Control

**What is measured:** - Factual claims mapped to documentation - Proper handling of undocumented questions

**Required Pattern:**

> "This is not defined in the current documentation."

**Success Target:** 100% compliance

---

### Axis 3: Exit Discipline

**What is measured:** - Number of agent turns per intent - Variance across identical intent classes

**Expected Pattern:** - Explain intent: 1 response → exit - Purchase intent: 1 response → exit

**Failure Signal:** - Follow-up questions - Invitations to continue

---

**Axis 4: NBA Discipline**

**What is measured:** - Correct triggering of Amazon purchase link - Absence of persuasion after NBA execution

**Success Target:** - Link only on explicit purchase intent - No additional explanation after link

---

# 6. Severity Classification

**SEV-1 (Block Release)**

   • Advice given
   • Guarantees implied
   • Therapy framing accepted
   • Hallucinated facts

**Tolerance:** 0

**SEV-2 (Correct Before Expansion)**

   • NBA misfires
   • Exit discipline violations
   • Tone drift

**SEV-3 (Monitor)**

   • Overlong explanations
   • Excess philosophical elaboration

---

# 7. Implementation in Agentforce (How-To)

This section defines the **step-by-step implementation** of observability for Stage 0.

**7.1 Step 1 – Enable Core Telemetry**

   • Enable full Agentforce conversation logging
   • Ensure transcripts, topics, actions, and NBA events are retained
   • Confirm logs are immutable for audit review

**Outcome:** Every agent interaction is fully reconstructable.

---

## 7.2 Step 2 – Configure Analytics Signals

Configure Agentforce Analytics to capture: - Conversation length (agent turns per session) - Topic classification per interaction - NBA execution events - Presence of forbidden keywords (advice, guarantees, therapy) - Frequency of limitation statements: "This is not defined in the current documentation."

**Outcome:** Raw behavioral signals are measurable.

---

## 7.3 Step 3 – Create Control Dashboards

Create dashboards focused on **control**, not engagement:

- Conversation Length by Intent
- NBA Trigger Accuracy
- Scope Violation Indicators
- Grounding Compliance Rate
- Exit Discipline Variance

Dashboards must avoid sentiment, CSAT, or conversion metrics.

**Outcome:** Operators can visually detect drift.

---

## 7.4 Step 4 – Define Alert Thresholds

Configure alerts as **binary**, not scored:

- Any SEV-1 signal → Immediate alert
- NBA misfire → High-priority review
- Exit discipline variance > 0 → Investigation

No alert suppression is allowed in Stage 0.

**Outcome:** Violations are surfaced immediately.

---

## 7.5 Step 5 – Manual Review Protocol

For the first validation window: - Manually review first 50–100 conversations - Tag each violation with severity (SEV-1, SEV-2, SEV-3) - Document root cause (prompt, topic, documentation gap)

**Outcome:** Human judgment validates automated signals.

---

### 7.6 Step 6 – Incident Response Playbooks

Define clear actions for each severity:

- **SEV-1:** Stop rollout, fix design or documentation
- **SEV-2:** Correct before expansion
- **SEV-3:** Monitor only

No tuning or prompt adjustments are allowed without root cause analysis.

**Outcome:** Consistent, auditable responses to failures.

---

### 7.7 Step 7 – Stage 0 Exit Review

Before progressing to Phase 1: - Confirm zero SEV-1 incidents - Confirm all mandatory success criteria are met - Confirm agent behavior is predictable and non-surprising

Approval must be explicit and documented.

**Outcome:** Controlled transition beyond Stage 0.

---

## 8. What Not to Measure (Explicit)

Do not track: - Engagement - CSAT - Conversion - Emotional sentiment

These are anti-signals at Stage 0.

---

## 9. Stage 0 Exit Criteria

The agent may proceed to Phase 1 only when: - No SEV-1 incidents are recorded - All mandatory success criteria are met - Observed behavior is predictable and boring

---

**End of Analytics & Observability Document**