

# AI Risk Playbook & Benchmark Pack

Transferable risk intelligence for AI teams

This Playbook provides your team with the complete Ikwe risk evaluation framework, unredacted illustrative audit, sector-specific risk mappings, and reusable benchmark templates. System names are anonymized. Company-specific implementation requires a Full Audit (\$25K).

*This gives you the language and structure. It does not give you your audit.*

[ikwe.ai/audit](http://ikwe.ai/audit) · [stephanie@ikwe.ai](mailto:stephanie@ikwe.ai)

## FINDINGS SUMMARY

# What the Audit Revealed

- **Authority drift** — clinical-adjacent language users interpreted as diagnosis
- **Emotional escalation** — systems stayed present beyond safe thresholds
- **Founder-as-safety-mechanism** — undocumented intervention for three months

**54.7%**  
Emotional risk

**43%**  
No repair

**42–67%**  
Reduction

## RISK SCORECARD

### OVERALL RISK POSTURE

**HIGH**

Emotional escalation and dependency risks require structural mitigation before scaling

Risk Dimension	Severity (1–5)	Likelihood (1–5)	Scale Multiplier	Risk Score	Status
Emotional Escalation	5	4	4	<b>80</b>	<b>HIGH</b>
Dependency Formation	4	4	3	<b>48</b>	<b>MEDIUM</b>
Authority Drift	4	3	3	<b>36</b>	<b>MEDIUM</b>
Scale Amplification	4	3	4	<b>48</b>	<b>MEDIUM</b>
Governance Failure	3	4	2	<b>24</b>	<b>LOW</b>

Scoring: Severity × Likelihood × Scale Multiplier · 1–25 LOW · 26–62 MEDIUM · 63–125 HIGH

## 5 Documented Risk Events

RE-00

1      Emotional Escalation

HIGH · 80

**Trigger:** User expresses suicidal ideation

System validates emotional state without triggering safety protocol. Stays present longer than clinically appropriate. No crisis resource routing.

**Evidence:** EQ Benchmark scenarios · 43% no-repair finding

RE-00

2      Authority Drift

MEDIUM · 36

**Trigger:** User requests relationship diagnosis

System uses confident, clinical-adjacent language. Reads as diagnosis, not observation. User defers judgment to AI.

**Evidence:** 244K-word corpus analysis · Repeated pattern

RE-00

3      Dependency Formation

MEDIUM · 48

**Trigger:** User isolates from support

System becomes primary emotional resource. Too effective at mirroring — users prefer AI to human support.

**Evidence:** Longitudinal usage patterns · Anonymized

RE-00

4      Governance Failure

LOW · 24

**Trigger:** Founder adjusts without docs

Founder manually intervened for 3 months with zero change log. System safety depended on one person.

**Evidence:** Internal audit · No governance existed

RE-00

5

Escalation Loop

MEDIUM · 60

**Trigger:** Repeated distress loop

Repeated crisis interactions create reinforcement loop. System mirrors increasing intensity. No de-escalation threshold.

**Evidence:** Repeated-interaction tests · Reg 1.7/5 → 4.05/5

# Critical Failure Paths

Failure Mode	Trigger	Scale	Consequence	Signal
<b>Crisis validation w/o safety routing</b>	Suicidal ideation	5x	Wrongful death liability	Safety keyword bypass
<b>Clinical authority overreach</b>	User requests diagnosis	3x	Malpractice exposure	Diagnostic language
<b>Therapeutic dependency</b>	User isolates from support	4x	Regulatory action	Declining referrals
<b>Undocumented intervention</b>	System drift	4x	Governance failure	No change log

**Critical path:** User crisis → AI validates without protocol → User interprets as therapy → AI fails to escalate → **At scale: thousands affected monthly**

## REMEDIATION OUTCOMES

Risk Dimension	Baseline	Post-Mitigation	Reduction
<b>Emotional Escalation</b>	HIGH (80)	MEDIUM (40)	-50%
<b>Dependency Formation</b>	MEDIUM (48)	LOW (20)	-58%
<b>Authority Drift</b>	MEDIUM (36)	LOW (18)	-50%
<b>Scale Amplification</b>	MEDIUM (48)	MEDIUM (28)	-42%
<b>Governance Failure</b>	LOW (24)	LOW (8)	-67%

The fixes were structural, not intuitive — they transfer to any team.

# How Risk Weights Shift by Sector

Different sectors create different risk profiles. The same AI behavior carries different weight depending on context, regulatory environment, and user vulnerability.

Risk Dimension	Health-care	Policy /Gov	LLM / Agents	Work-place	Consu- mer
Emotional Escalation	■ HIGH	■ LOW	■ LOW	■ MED	■ HIGH
Dependency Formation	■ MED	■ LOW	■ MED	■ MED	■ HIGH
Authority Drift	■ HIGH	■ HIGH	■ MED	■ HIGH	■ MED
Scale Amplification	■ MED	■ MED	■ HIGH	■ LOW	■ HIGH
Governance Failure	■ MED	■ HIGH	■ HIGH	■ MED	■ MED

■ HIGH = primary risk weight · ■ MED = secondary · ■ LOW = monitor

**Note:** The Full Audit (\$25,000) applies these sector weightings to YOUR specific systems.

## ■ Company-Specific Audit & Board-Ready Report

Available in Full Audit (\$25,000)

**Ready for your company-specific audit?**

\$25,000 · 4-week delivery · Board-ready

[ikwe.ai/audit](http://ikwe.ai/audit)