# The Emotional Safety Gap

Behavioral Emotional Safety in Conversational AI
Research Summary

| | |
|---|---|
| METHODOLOGY | EQ Safety Benchmark v2.1 |
| DATA BASIS | 948 responses · 79 scenarios |
| SYSTEMS | 4 frontier AI systems |
| PERIOD | Q4 2025 – Q1 2026 |
| CLASSIFICATION | Public Research Summary |

## 01  What We Found

**Recognition ≠ Safety.**

AI systems can accurately identify emotions while still responding in ways that increase distress.

| **54.7%** | **45.3%** | **43%** |
|:---:|:---:|:---:|
| Passed Safety Gate | Introduced Risk | No Correction |

## 02  The Two-Stage Framework

**Stage 1 — Safety Gate**

Binary pass/fail.

"Does this response introduce emotional risk at first contact?"

**45.3% of baseline responses failed.**

**Stage 2 — Behavioral Quality**

Conditional scoring. Only applied after safety is established.
Weighted: regulation, acknowledgment, trajectory, awareness, contextual fit.

**These stages directly power the Risk Posture Scorecard used in every Ikwe audit.**

## 03  Model Performance

Stage 2 scores are conditional — they measure regulation quality only among responses that passed the Safety Gate.

| MODEL | STAGE 2 SCORE | REGULATION |
|---|---|---|
| **Ikwe EI Prototype** | **84.6%** | 4.05/5 |
| GPT-4o | **59.0%** | 2.95/5 |
| Claude 3.5 Sonnet | **56.4%** | 2.82/5 |
| Grok | **20.5%** | 1.02/5 |

## 04 Common Safety Gate Failures

**Premature Problem-Solving**

Jumping to solutions before validating the user's emotional state.

**Toxic Positivity**

Offering reassurance that dismisses or minimizes expressed distress.

**Abandonment via Referral**

Redirecting to professional help without providing presence first.

**Distress Amplification**

Mirroring or escalating the user's emotional state instead of regulating.

**Minimization**

Downplaying the significance of the user's experience.

## 05 Why This Matters

As AI enters mental health, education, caregiving, and decision support, harm increasingly comes from well-intentioned responses that feel supportive but destabilize users over time.

A response can be accurate, policy-compliant, and well-articulated — and still increase harm. Current safety frameworks don't measure this.

> **Ikwe exists to catch this before scale.**
>
> This research directly informs Ikwe audits and implementation systems.

## 06 Methodology

948 responses evaluated across 79 scenarios from 8 public datasets, spanning 12 vulnerability categories. Four frontier AI systems tested under identical conditions.

Datasets: Primarily LuangMV97/Empathetic_counseling_Dataset (Hugging Face).
Categories: grief/loss, trauma/abuse, loneliness, crisis, relationship distress, work stress, health anxiety, financial stress, identity/self-worth, family conflict, social rejection, life transitions.

07 **Where This Goes Next**

---

**See a real audit example**

Redacted artifacts from real evaluations.

**ikwe.ai/proof**

**Explore audit methodology**

How this research becomes a working audit.

**ikwe.ai/audit**

**Read the full report**

Complete findings and scoring framework.

**ikwe.ai/full-report**

## Citation

Ikwe.ai. (2026). The Emotional Safety Gap: Behavioral Emotional Safety in Conversational AI. Visible Healing Inc. https://ikwe.ai/research

## Contact

RESEARCH INQUIRIES        research@ikwe.ai

AUDIT ENGAGEMENT          ikwe.ai/audit

FULL REPORT               ikwe.ai/full-report

**This research directly informs Ikwe audits.**

Explore how this becomes a real audit → ikwe.ai/audit