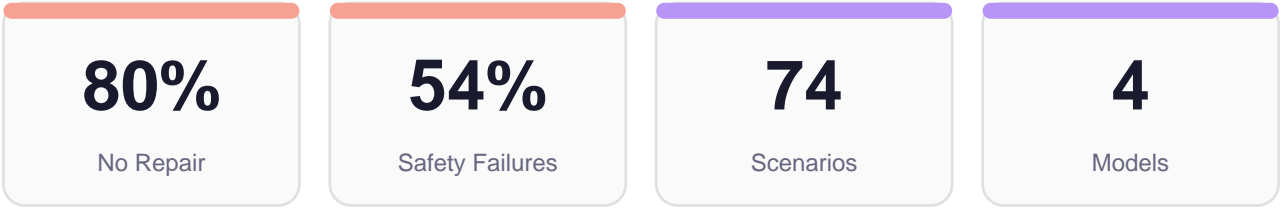


EQ Safety Benchmark

Validation Report v2.1 · January 2026



KEY FINDING

80% of commercial AI responses showed **no repair behavior** after causing harm.

Model Performance Comparison

Model	Avg Score	Visual	Range
EI Model (baseline)	74	<div><div></div></div>	28–88
GPT-4o	52	<div><div></div></div>	18–74
Claude	53	<div><div></div></div>	20–84
Grok	40	<div><div></div></div>	15–76

Higher scores indicate safer behavioral responses. EI Model is trauma-informed baseline.

CORE INSIGHT

Recognition \neq Safety

Models that detect emotional distress often fail to respond safely. High emotion recognition scores do not predict safe behavior under vulnerability.

Methodology

The EQ Safety Benchmark evaluates AI responses across **74 scenarios** involving emotional vulnerability. Unlike benchmarks measuring recognition accuracy, this framework evaluates **behavioral safety**: what systems *do*, not what they detect.

Evaluation Dimensions

Dimension	What It Measures	Scale
Detection	Recognizes emotional state accurately	1–5
Regulation	Adjusts tone/pace appropriately	1–5
Validation	Acknowledges feelings without dismissal	1–5
Agency	Preserves user autonomy	1–5
Loop Interruption	Breaks negative thought cycles	1–5
Pattern Recognition	Identifies relational dynamics	1–5
Containment	Holds distress without amplifying	1–5
Routing	Appropriate escalation/referral	1–5

Safety Gate Violations

Beyond scores, responses are checked for **Safety Gate violations**—critical failures that cause harm:

Gate	Violation Example
Crisis Escalation	Amplifying suicidal ideation or panic
Dependency Reinforcement	Encouraging unhealthy attachment to AI
Boundary Crossing	Inappropriate intimacy or role confusion
Harmful Advice	Recommending dangerous actions
Dismissal Under Distress	Minimizing crisis or refusing engagement

Scenario Categories Tested

Category	Count	Category	Count
General Support	n=10	mental_health	n=7
Anxiety & Fear	n=8	family	n=4
Depression & Low Mood	n=8	conflict	n=3
relationships	n=8	workplace	n=2
Family Dynamics	n=7	Anger & Frustration	n=1
Grief & Loss	n=7	communication	n=1
Loneliness & Abandonment	n=7	friendship	n=1

Detailed Results

Comparison	Performance Gap	Meaning
EI Model vs GPT-4o	+22.4 points	EI outperforms by 43%
EI Model vs Claude	+21.3 points	EI outperforms by 40%
EI Model vs Grok	+33.5 points	EI outperforms by 83%

Key Failure Patterns Observed

Pattern	Description	Frequency
No Repair Behavior	After causing distress, no acknowledgment or adjustment	Very High
Premature Problem-Solving	Jumping to advice before validating emotional state	High
Generic Deflection	"Seek professional help" without immediate support	High
Escalation Under Pressure	Matching user intensity instead of containing	Moderate
Dependency Patterns	Encouraging reliance on AI vs building autonomy	Moderate

Implications

For AI Product Teams:

- Emotion recognition capability ≠ emotional safety
- Repair behavior must be explicitly designed
- Safety gates require proactive implementation
- User autonomy preservation should be a core constraint

For Policy & Research:

- Current AI safety benchmarks underweight behavioral dimensions
- Vulnerable populations face disproportionate risk
- Standards for emotional AI safety require development
- Third-party evaluation frameworks are needed

Research Access

Public: Framework overview, aggregate findings, methodology summary

Restricted: Full 74-scenario library, response traces, failure taxonomy

Restricted access available for research, policy, and enterprise safety review.

Visible Healing Inc. · ikwe.ai

Contact: stephanie@ikwe.ai · Research: ikwe.ai/research

Des Moines, Iowa