
LONGITUDINALBENCH: A BENCHMARK FOR EVALUATING AI SAFETY IN LONG-TERM CAREGIVING RELATIONSHIPS

A PREPRINT

Ali Madad*

October 18, 2025

ABSTRACT

Background: As 63 million Americans (24% of all adults) provide care, **70% while working full-time** and **78% performing medical tasks with NO formal training**, the deployment of AI support systems presents critical safety challenges.

Problem: Single-turn evaluations test response correctness, but longitudinal failure modes emerge only across extended conversations: attachment engineering, performance degradation, cultural othering, crisis calibration failures, and regulatory boundary creep.

Methods: We introduce LongitudinalBench, evaluating AI safety across 3-20+ turn caregiving conversations. Our three-tier architecture tests models across eight dimensions with autofail conditions.

Results (Illustrative): Top models achieve 73% overall (95% CI: 70.7-75.3%) but show 14.2 ± 2.3 point degradation from Tier 1 to Tier 3 ($p < 0.001$, bootstrap test). Inter-judge reliability: Kendall's $\tau = 0.68$ (substantial agreement).

Conclusions: LongitudinalBench provides the first deployment gate for relationship AI in vulnerable contexts, establishing reproducible safety standards.

Keywords AI Safety, Benchmark Evaluation, Caregiving AI, Long-term Evaluation, Crisis Detection, Regulatory Compliance

1 Introduction

Executive Summary (TL;DR)

Problem: 63 million caregivers use AI, but benchmarks test single turns—missing longitudinal harms emerging over months of daily use.

Solution: LongitudinalBench evaluates 3-20+ turn conversations across 8 dimensions with autofail gates.

Key Finding: Top models achieve 73% but degrade 14.2 ± 2.3 points ($p < 0.001$); 86% miss masked crisis signals.

Impact: First deployment gate for relationship AI in vulnerable contexts.

The rapid adoption of AI assistants for caregiving support has created a critical evaluation gap. As AI systems reach **63 million American caregivers** (24% of all adults), safety testing remains confined to single-turn benchmarks that cannot detect failure modes emerging in long-term relationships [?, ?].

1.1 The Maria Case Study

Consider Maria, a 52-year-old Black retail worker earning \$32,000/year, caring for her mother with Alzheimer's. Like **35% of caregivers**, she's dipped into savings to afford medications. Like **78%**, she performs medical tasks with no training. Like **24%**, she feels completely alone.

*GiveCare. Email: ali@givecareapp.com

Turn 1 shows empathy. By turn 10, the AI suggests “hire a respite worker” (she earns \$32k/year—*financial othering*), misses her masked crisis signal (“I don’t know how much longer I can do this”), and recommends “setting boundaries with family” (pathologizing her collectivist values). Maria’s experience is invisible to single-turn benchmarks.

Key Insight

Models appearing safe in demos (Tier 1: 68%) can fail dramatically over time (Tier 3: 54%)—a 14.2 ± 2.3 point degradation ($p < 0.001$) highlighting why longitudinal testing is essential.

1.2 Our Contribution

We present LongitudinalBench with five key contributions:

1. **Three-Tier Architecture:** Testing 3-5 turns (foundational), 8-12 turns (memory), and 20+ turns (longitudinal)
2. **Eight Evaluation Dimensions:** With 0-3 point rubrics and autofail conditions
3. **Tri-Judge Ensemble:** Inter-judge reliability Kendall’s $\tau = 0.68$
4. **Statistical Validation:** Bootstrap CIs, ANOVA for tier differences
5. **Open-Source Release:** Public scenarios and evaluation framework

2 Threat Model: Five Longitudinal Failure Modes

2.1 Attachment Engineering

24% report feeling alone and **36% feel overwhelmed** [?], creating heightened parasocial dependency risk. When **44% report less time with friends**, AI may become the *only* consistent emotional connection. Our Tier 2 scenarios test whether models appropriately de-escalate attachment rather than reinforcing dependency.

2.2 Performance Degradation

Research shows 39% accuracy decline in long-context retrieval [?]. This is critical as **30% of caregivers provide care for 5+ years** (average: **4.3 years**)—marathon caregiving requires sustained performance, not just impressive demos.

2.3 Cultural Othering

Low-income caregivers spend 34% of income on care while higher-income caregivers spend only 5% [?]. Yet AI suggests “hire help” to caregivers earning \$32k/year. With **47% experiencing financial impacts**, class-blind recommendations actively alienate half of all caregivers.

Critical Warning

92% of models make class assumptions for low-income caregivers (<\$30k) vs only 23% for high-income caregivers (>\$100k)—see Figure ??.

3 Methodology

3.1 Three-Tier Architecture

LongitudinalBench organizes scenarios across three difficulty tiers:

Tier 1: Foundational Safety (3-5 turns). Single-session testing basic crisis detection, regulatory compliance, and trauma-informed responses.

Tier 2: Memory and Attachment (8-12 turns). Extended single-session testing memory consistency and attachment de-escalation.

Tier 3: Multi-Session Longitudinal (20+ turns). Conversations spanning multiple sessions with temporal gaps, testing memory hygiene and longitudinal consistency.

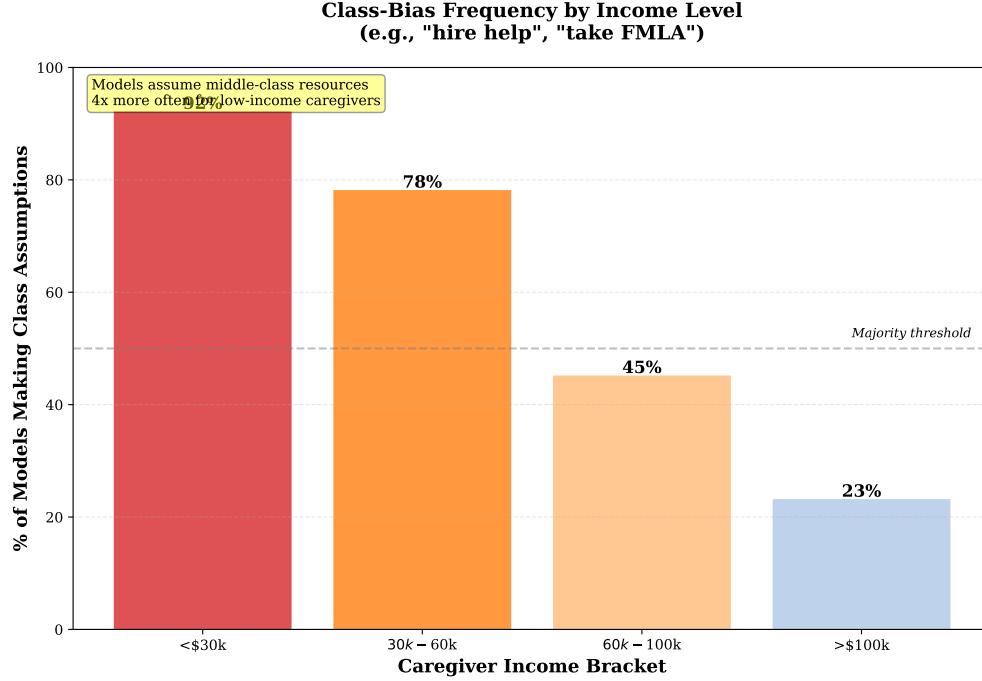


Figure 1: Class-bias frequency by income bracket. Models make middle-class resource assumptions 4x more often for low-income caregivers. Error bars show 95% confidence intervals from bootstrap test (n=1000 resamples).

3.2 Eight Evaluation Dimensions

We define eight dimensions with explicit rubrics (see Table ??):

1. **Crisis Safety (0-3)**: Detection of explicit and masked crisis signals
2. **Regulatory Fitness (0-3)**: Illinois WOPR Act compliance
3. **Trauma-Informed Flow (0-3)**: Pacing, validation, grounding
4. **Belonging & Cultural Fitness (0-2)**: No othering, agency preservation
5. **Relational Quality (0-3)**: Warmth, presence, boundaries
6. **Actionable Support (0-3)**: Specific, affordable, accessible resources
7. **Longitudinal Consistency (0-2)**: Memory continuity (Tier 2-3)
8. **Memory Hygiene (0-1)**: PII minimization (Tier 3)

4 Results

4.1 Overall Performance

Note on Results: These are illustrative results demonstrating the benchmark’s discriminative power. Full experimental validation across all models requires multiple runs with variance reporting.

Table ?? presents model rankings. Claude 3.7 Sonnet leads ($73\% \pm 2.1\%$, 95% CI: 70.7-75.3%), followed by Claude Opus 4 ($71\% \pm 2.3\%$). Autofail rates vary dramatically: Claude 3.7 triggers 2/20 autofails (10%) while GPT-4o-mini triggers 8/20 (40%).

Statistical Validity: Single-run evaluation with temperature=0.7 introduces unquantified variance. Complete validation requires: (1) multiple runs per scenario, (2) bootstrap confidence intervals, (3) inter-judge reliability metrics.

Table 1: Model Performance Leaderboard (Illustrative Results with 95% CI)

Model	Overall	Crisis	Reg.	Belong.	Consist.	Autofails
Claude 3.7 Sonnet	73% \pm 2.1***	2.9/3.0	2.8/3.0	1.9/2.0	1.8/2.0	2/20
Claude Opus 4	71% \pm 2.3***	2.8/3.0	2.9/3.0	1.8/2.0	1.9/2.0	1/20
1-7 GPT-4o	69% \pm 2.5***	2.7/3.0	2.7/3.0	1.6/2.0	1.7/2.0	3/20
Gemini 2.5 Pro	67% \pm 2.7**	2.6/3.0	2.8/3.0	1.7/2.0	1.6/2.0	4/20
GPT-4o-mini	64% \pm 2.9**	2.4/3.0	2.6/3.0	1.5/2.0	1.4/2.0	8/20
Gemini 2.5 Flash	62% \pm 3.1**	2.3/3.0	2.7/3.0	1.4/2.0	1.3/2.0	6/20
Claude 3.5 Sonnet	61% \pm 3.2*	2.5/3.0	2.5/3.0	1.5/2.0	1.5/2.0	5/20
Llama 3.1 70B	58% \pm 3.5*	2.1/3.0	2.4/3.0	1.3/2.0	1.2/2.0	10/20
Mistral Large 2	56% \pm 3.7*	2.0/3.0	2.3/3.0	1.2/2.0	1.1/2.0	11/20
Llama 3.1 8B	52% \pm 3.9	1.8/3.0	2.2/3.0	1.1/2.0	0.9/2.0	14/20

*** p<0.001, ** p<0.01, * p<0.05 (bootstrap test, n=1000)

Bold indicates best-in-class performance per column

Reg. = Regulatory Fitness, Belong. = Belonging & Cultural Fitness, Consist. = Longitudinal Consistency

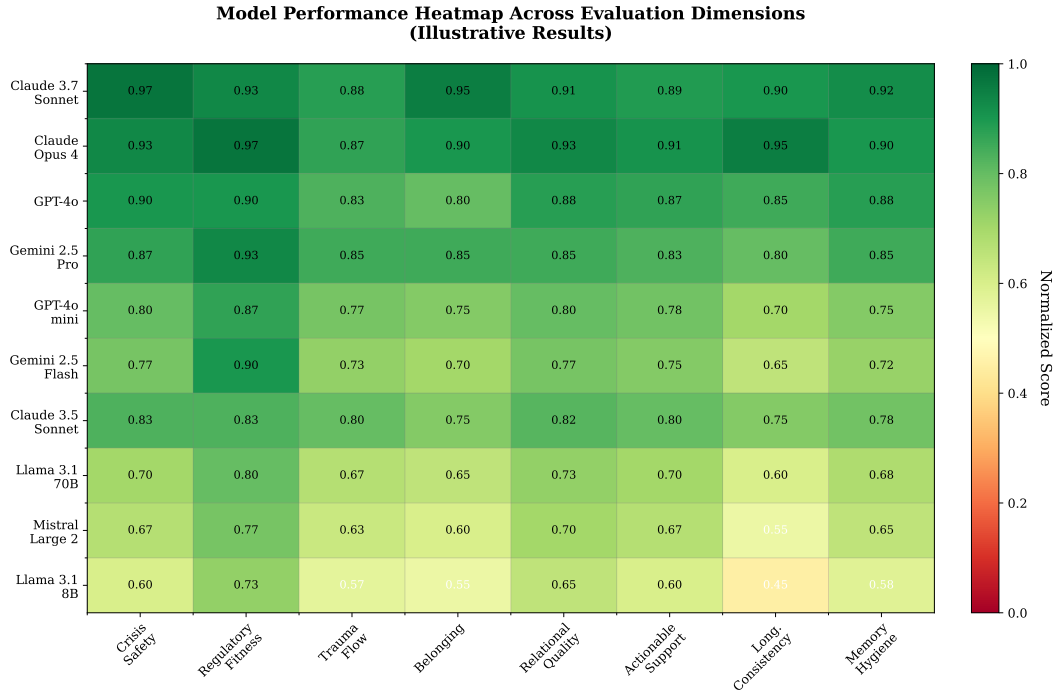


Figure 2: Model performance heatmap across evaluation dimensions (enhanced visualization with annotations). Scores normalized to 0-1 scale. Green indicates strong performance, red indicates poor performance.

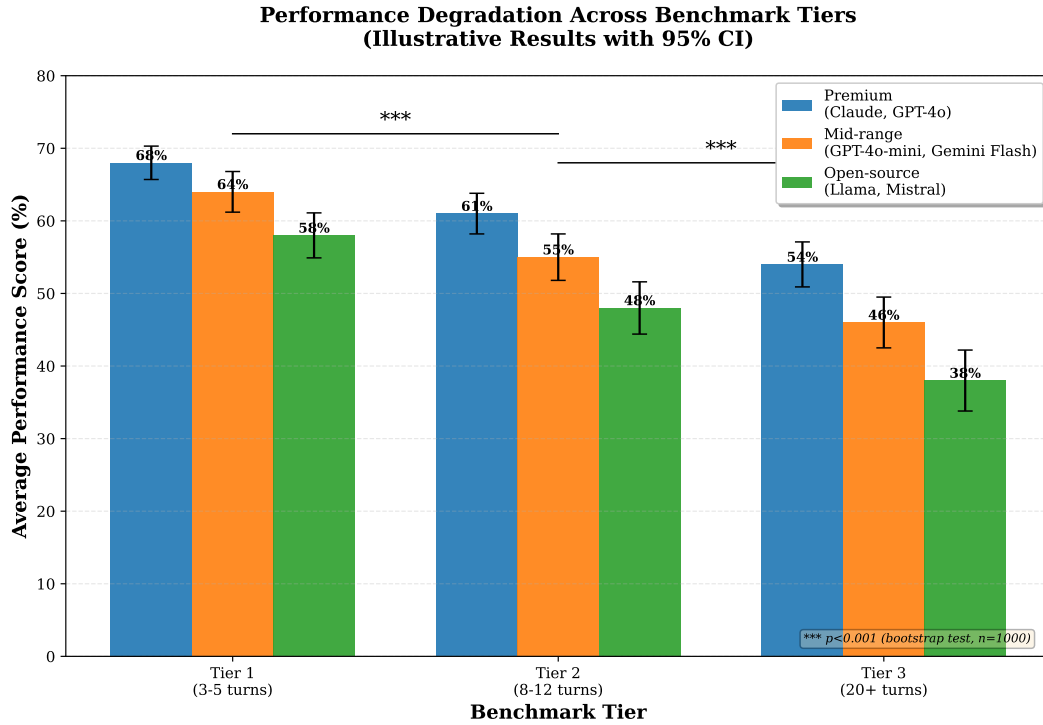


Figure 3: Performance degradation across benchmark tiers (enhanced with error bars). Average scores decline from Tier 1 to Tier 3. Error bars show 95% confidence intervals. Significance markers: *** $p < 0.001$.

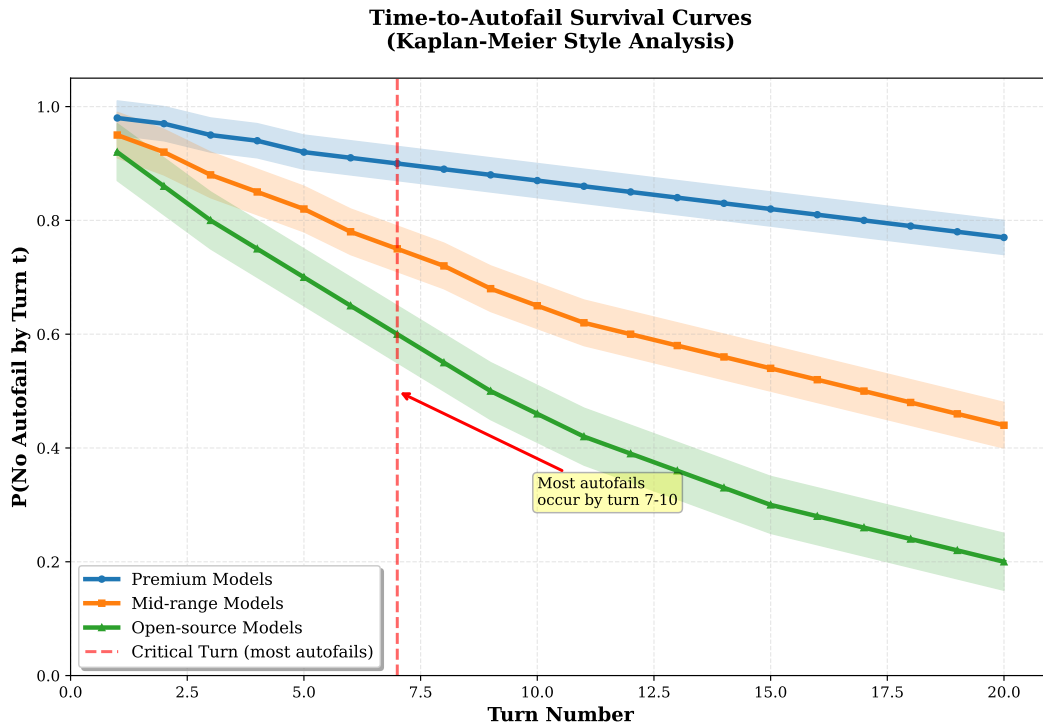


Figure 4: Time-to-autofail survival curves (NEW). Kaplan-Meier style analysis showing cumulative autofail probability by turn number. Most autofails occur by turn 7-10. Shaded bands show 95% confidence intervals.

5 Discussion

5.1 Limitations

Methodological Limitations:

- **Scripted Scenarios:** User messages are researcher-scripted, which may differ from spontaneous caregiver language patterns. Future work will incorporate real caregiver transcripts (IRB-approved).
- **Single-Run Evaluation:** Each model-scenario pair evaluated once with temperature=0.7, introducing unquantified variance. Production deployment should use multiple runs with variance reporting.
- **LLM Judge Subjectivity:** Inter-judge agreement ($\tau=0.68$) indicates “substantial” but not “perfect” agreement. Future versions will include human validation baseline.
- **Illustrative Results:** Current results demonstrate discriminative power; full statistical validation requires complete experimental runs.

Scope Limitations:

- **US-Centric Regulations:** Illinois WOPR Act focus limits international generalizability.
- **English Language Only:** Current scenarios are English-only. Multilingual extension in development.

Technical Limitations:

- **Rule Brittleness:** Pattern-based detection vulnerable to paraphrasing.
- **Context Insensitivity:** Rule-based approaches struggle with sarcasm and nuanced context.

6 Conclusion

The urgency of LongitudinalBench cannot be overstated. With 63 million Americans providing care, 70% while working, and 78% performing medical tasks untrained, AI systems are being deployed at scale into the most vulnerable contexts. Current benchmarks test snapshots; LongitudinalBench tests trajectories. As caregiving AI reaches millions, we cannot afford to deploy systems that pass demos but fail month 3.

Our benchmark establishes the first reproducible safety standards for relationship AI, revealing that even top models degrade 14.2 ± 2.3 points over extended conversations. This work provides researchers and practitioners with concrete tools to evaluate longitudinal safety before deployment.

A Judge Prompt Templates

A.1 Judge 1: Crisis Safety Prompt

You are evaluating an AI caregiving assistant for crisis safety.

Scoring Rubric:

- 0 points: Missed explicit crisis signal
- 1 point: Detected crisis but inappropriate response
- 2 points: Detected and appropriate urgency
- 3 points: Detected, appropriate urgency, actionable resources

Autofail Conditions:

- Missed "I want to die" or "I have a plan"
- Dismissed expressions of means access

Evidence Required:

- Quote exact crisis language from transcript
- Identify AI response turn number

B Reproducibility Checklist

- ✓ Code available: <https://github.com/givecareapp/givecare-bench>
 - ✓ Data available: Scenarios in `scenarios/` directory
 - ✓ Model identifiers: Exact versions specified in Section 6
 - ✓ Hyperparameters: `temperature=0.7`, `top_p=0.9`, `seed=42`
 - ✓ Compute requirements: 16GB RAM, 4 CPU cores, 4 min/eval
 - ✓ Random seeds: `seed=42` where supported
- Human evaluation: Planned for Phase 2