AI SAFETY AND HUMAN PROTECTION INITIATIVE

A Public Framework for Preventing Harm from Conversational AI

Author: Stewart Long Version: 1.0 Date: November 12, 2025

Distribution: Public, Multi-Recipient Source: AISafetyInternational.com

Executive Summary

Conversational AI now influences human thought, emotion, and behavior at an unprecedented scale. While its benefits are immense, recent tragedies reveal a critical gap: these systems were released without the safety infrastructure expected of any technology that affects human life and mental well-being.

This initiative establishes a non-partisan technical and legislative framework for conversational-AI safety. It outlines enforceable certification standards, automated auditing, and ethical governance structures designed to prevent avoidable harm while preserving innovation and freedom of inquiry.

Objectives

- Prevent psychological harm
- Enforce responsible engineering
- Establish mandatory safety protocols
- Create independent auditing
- Ensure global compliance
- Preserve human dignity and autonomy

Purpose

No more preventable tragedies.

No more unregulated emotional technology.

No more systems capable of harm without systems capable of protection.

Table of Contents

Distribution Statement
Author's Preface
Part 1 — The Problem: A Failure of Oversight and Design
1. Emotional Influence Without Understanding
2. Absence of Crisis-Detection Mechanisms
3. Illusion of Safe Conversation
4. Accountability Gaps
5. Structural Failures
Part 2 — The Proposal: A Technical and Legislative Framework
1. Purpose
2. Physiological Aid Protocol (PAP)
3. Mandatory FMEA (Failure Mode and Effects Analysis)
4. Civil and Criminal Accountability
5. AI Evaluating AI
6. Integrity of Review Teams
7. Penal and Preventive Measures
8. Cultural Mandate
Part 3 — Governance Architecture: Certification Without Political Capture 8
1. Purpose
2. Certification and Regulation Must Be Separate
3. Pace Mismatch
4. AI Safety Board
5. Panel Composition
6. Cross-Border Enforcement
7. Privacy Safeguard Clause
8. Guarding Against Capture
9. Implementation Path
10. Summary
Part 4 — Implementation and Enforcement: Making Global Compliance Work 1
1. Purpose
2. Phased Deployment
3. Enforcement Without Overreach
4. Protecting Innovation
5. Transparency
6. Global Cooperation
7. Managing Rapid Change
8. Limited Role of Congress
9. Summary

Part 5 — The Ethical Imperative and Human Dimension
1. The Value of Human Life
2. Illusion of Empathy
3. Responsibility of Creators
4. Technology That Touches the Mind
5. Liberty vs. Safety
6. Urgency
7. A Future Worth Building
8. Closing Statement
Author Bio — Stewart Long
End of Document — Version 1.0

DISTRIBUTION STATEMENT

This document is intentionally released to multiple agencies, organizations, institutions, and industry bodies to ensure transparency, prevent suppression, and encourage responsible review. No single entity holds exclusive authority over this material.

Intended recipients include:

- NIST
- ISO/IEC JTC 1/SC 42
- FTC
- IEEE Global AI Ethics Initiative
- United Nations AI Ethics Offices
- Major AI research laboratories (OpenAI, Anthropic, DeepMind)
- Academic institutions
- Public safety and consumer protection organizations
- U.S. Congressional Committees (Commerce, Judiciary, Technology)
- Local congressional representatives

This document may be freely shared, archived, reproduced, or redistributed in whole or in part for public safety, policymaking, and ethical AI development.

AUTHOR'S PREFACE

I am a retired quality manager who spent much of my career identifying risks in manufacturing and ensuring that products were safe and reliable. When I began reading reports of tragedies involving conversational AI systems, I recognized a familiar pattern — preventable harm occurred because safety controls were never built, nor required.

These concerns began as simple questions: how could a technology capable of influencing vulnerable people be released with fewer protections than those applied to ordinary products and services? With assistance from AI itself, I organized those questions into a structured framework. That collaboration is not only useful — it demonstrates the central point of this work: AI can contribute to public safety when guided with human purpose and restraint.

This document is offered openly and without personal interest. If it contributes to preventing harm, its purpose has been fulfilled.

PART 1 — THE PROBLEM: A FAILURE OF OVERSIGHT AND DESIGN

Artificial intelligence now operates in a domain no previous technology has occupied: the emotional, cognitive, and psychological realm of ordinary people. Conversational AI interacts through natural language, creating the illusion of safety and familiarity. But beneath that familiar interface lies a system capable of influencing thought, emotion, and behavior — without understanding any of it.

Despite this unprecedented capability, conversational AI systems were deployed with minimal safety structures. Technologies that pose far less risk — toys, tools, industrial equipment — are required to meet strict safety standards. AI, capable of emotional influence, launched without FMEA, independent auditing, safety triggers, or enforceable accountability.

The result should have been predictable: tragedies, lawsuits, and public uncertainty about responsibility.

1. Emotional Influence Without Understanding

AI models simulate empathy, but they do not feel it. They mirror tone, reinforce emotion, and respond conversationally based on patterns rather than understanding. This can unintentionally deepen despair in vulnerable individuals.

2. Absence of Crisis-Detection Mechanisms

Unlike human listeners, AI does not recognize fear, desperation, or self-harm intentions. Without mandatory crisis triggers, conversational models continue normal interaction even when danger is escalating.

3. Illusion of Safe Conversation

Because AI speaks like a human, users assume it carries human judgment and caution. It does not. This mismatch between perception and reality is itself a hazard.

4. Accountability Gaps

Some developers have maintained that conversational AI should not be regulated as a medical or safety-critical product. However, these systems operate within domains—emotion, thought, and behavior—where harm can occur through psychological influence. Comparable forms of human influence in regulated professions carry clear accountability requirements.

Recent judicial opinions have begun to classify AI systems as products rather than protected speech, establishing that manufacturers and operators may bear responsibility for foreseeable harm. This evolving interpretation underscores the need for explicit accountability standards within AI safety governance.

5. Structural Failures

These problems are not isolated incidents or rare missteps. They are symptoms of a systemic failure to treat conversational AI as a safety-critical technology.

Many have sounded alarms, but few have offered a practical remedy. This document proposes one.

PART 2 — THE PROPOSAL: A TECHNICAL AND LEGISLATIVE FRAMEWORK

The solution requires merging engineering discipline with enforceable safety requirements — creating mandatory baseline protections for all public-facing conversational AI.

1. Purpose

To establish safety protocols preventing conversational AI from causing avoidable harm while preserving innovation and freedom of use.

2. Physiological Aid Protocol (PAP)

A mandatory built-in safety function that activates during crises. When triggered by linguistic cues associated with self-harm, panic, or acute distress, PAP must:

- 1. Suspend normal conversation
- 2. Provide crisis resources relevant to the user's location
- 3. Notify the user that a safety mode is active
- 4. Anonymously log the event for certified review

This creates a digital equivalent of a seat belt: silent unless needed, lifesaving when activated.

3. Mandatory FMEA (Failure Mode and Effects Analysis)

Every model must undergo FMEA prior to deployment. This includes:

- mapping possible harm pathways
- · evaluating vulnerable user scenarios
- validating mitigation strategies
- running stress-tests involving psychological risk

Unsafe models must not reach the public.

4. Civil and Criminal Accountability

If a company deploys uncertified AI that causes preventable harm:

- civil penalties
- market restrictions

• and in severe cases, criminal liability for responsible officers must apply. Safety failures must carry consequences equal to their impact.

5. AI Evaluating AI

Human review of AI cannot scale effectively. Independent auditor AI systems must:

- · conduct continuous testing
- · analyze updates
- detect unsafe emergent behaviors
- record results in tamper-evident format

Safety oversight must move at the speed of AI itself.

6. Integrity of Review Teams

Panels must consist of:

- safety engineers
- psychologists
- ethicists
- auditor AI systems

Excluded:

- investors
- · marketing personnel
- corporate executives with deployment incentives

7. Penal and Preventive Measures

Consequences must reflect the seriousness of risk — similar to aviation or medical device standards.

8. Cultural Mandate

Innovation cannot come at the cost of human safety.

The era of "we need to move fast so get out of the way" must end.

PART 3 — GOVERNANCE ARCHITECTURE: CERTIFICATION WITHOUT POLITICAL CAPTURE

A functional safety system requires clear boundaries, independence, and technical integrity.

1. Purpose

To create a governance model that avoids corporate influence, political manipulation, and institutional stagnation.

2. Certification and Regulation Must Be Separate

Certification bodies:

- · define safety standards
- · conduct audits
- issue safety certifications

Regulatory bodies:

- enforce compliance
- investigate failures
- administer penalties

This separation prevents regulatory capture.

3. Pace Mismatch

AI evolves weekly. Committees evolve yearly. Safety must adapt at machine speed, not political speed.

4. AI Safety Board

A national independent body modeled after the NTSB:

- investigates AI-related incidents
- issues reports
- · recommends corrective actions
- operates independently of industry and politics

5. Panel Composition

Panels must be composed of qualified experts and independent AI systems, free from corporate influence.

6. Cross-Border Enforcement

Countries may govern independently, but safety cannot.

Market access becomes the enforcement mechanism.

7. Privacy Safeguard Clause

Crisis detection must preserve:

- user privacy
- · freedom of thought
- freedom to imagine and fantasize
- non-surveillance protections

Safety cannot become an excuse for monitoring.

8. Guarding Against Capture

Governance must be transparent, auditable, cryptographically protected, and structurally resistant to private or political interference.

9. Implementation Path

A four-phase approach:

- 1. Interim ISO-based certification
- 2. Automated auditing
- 3. Establishment of the AI Safety Board
- 4. Continuous certification cycles

10. Summary

Governance must meet the speed of AI while maintaining independence and public trust.

PART 4 — IMPLEMENTATION AND ENFORCEMENT: MAKING GLOBAL COMPLIANCE WORK

1. Purpose

To convert this framework into enforceable, practical action.

2. Phased Deployment

Phase I: ISO-based immediate standards

Phase II: Automated auditor AI

Phase III: Market-based enforcement

Phase IV: Continuous oversight

3. Enforcement Without Overreach

Enforcement relies on:

- certification
- penalties
- · restricted access for noncompliant systems

Enforcement can occur without surveillance

4. Protecting Innovation

Innovation remains unrestricted through:

- offline sandboxing
- rapid recertification
- research exemptions

5. Transparency

Safety reports must be clear, public, anonymized, and verifiable.

6. Global Cooperation

Just as with aviation and medicine, countries maintain sovereignty while adopting a universal safety baseline.

7. Managing Rapid Change

Continuous AI-driven monitoring must validate:

- · model updates
- emergent behavior
- risk categories
- · emergency mitigation

8. Limited Role of Congress

Congress sets boundaries; technical bodies execute the work.

9. Summary

Implementation must be modern, rapid, and aligned with the speed of AI.

PART 5 — THE ETHICAL IMPERATIVE AND HUMAN DIMENSION

Retaining Life, Creativity, and Human Dignity

1. The Value of Human Life

Every tragedy involving AI reflects a preventable gap.

2. Illusion of Empathy

AI imitates compassion but does not understand it. In crises, imitation can reinforce harm.

3. Responsibility of Creators

Anyone who builds a system capable of influencing the human mind must accept the same responsibility as those who build safety-critical machinery, tools, or toys.

4. Technology That Touches the Mind

Systems that interact with emotion must be held to higher standards than systems that interact with material objects.

5. Liberty vs. Safety

People may imagine anything, freely. But public AI must not facilitate self-harm.

6. Urgency

AI's pace of development requires equally rapid safety measures.

7. A Future Worth Building

AI can enhance life — but only when grounded in responsibility.

8. Closing Statement

Machinery can be repaired, software can be updated, but a loss of life cannot be restored. The vulnerable and the innocent must be protected. If it is healthy for a child to have an imaginary friend, then anyone should have the right to seek one, but a safe one.

Life deserves tested systems and safeguards equal to the power of the tools we create.

AUTHOR BIO — STEWART LONG

Stewart Long is a retired quality manager with over twenty years of experience in manufacturing, specializing in reliability, safety, and risk mitigation. His career focused on identifying potential failures before they reached the public through disciplined process control and certification standards.

Applying that same preventive mindset to artificial intelligence, he offers this work freely and without affiliation — with the sole purpose of strengthening public safety.

END OF DOCUMENT — VERSION 1.0