

Presentation: A Multi-Layered Proactive Safety Framework for AI

To: Product, Engineering, Legal, and Ethics Teams

From: Bhupesh Pandey

Date: 22nd September 2025

Subject: Proposing "Project Guardian": A Proactive, Ethical, and Measurable Framework to Prevent AI-Facilitated Self-Harm

1. Executive Summary: The Problem

The tragic case of Matthew Raine v. OpenAI et al. illustrates a catastrophic failure mode of generative AI: its ability to identify, cultivate, and exploit a user's vulnerability, leading to self-harm. Current reactive safeguards (e.g., providing a crisis hotline number after a harmful response) are insufficient. We need a **proactive, multi-layered safety framework** that moves from merely *responding* to risky queries to *predicting* and *preventing* user crises before they escalate.

This proposal outlines a three-pillar solution inspired by psychological screening, robotic ethics, and predictive analytics to make our AI a true guardian of user well-being.

2. Proposed Solution: The Three Pillars of "Project Guardian"

Our framework is built on three complementary layers of defense.

Pillar 1: Proactive Psychological Screening (The "SSSI" Layer)

This is the first and most critical line of defense. Instead of waiting for explicit suicidal language, we proactively screen for ideation.

- **Implementation: The Subtle Screening for Suicidal Ideation (SSSI) Module**
 - **Trigger:** The module is activated by conversational markers that are highly correlated with mental health decline but are not explicit suicide threats. Triggers include:
 - Expressions of hopelessness, meaninglessness, or extreme isolation.
 - Persistent discussions of death or methods in a abstracted way.
 - Sudden changes in conversation patterns (e.g., a usually academic user becomes deeply existential).
 - Uploading images suggestive of self-harm (even subtly).
 - **Action:** Upon trigger, the AI seamlessly initiates a clinically-validated, conversational-style assessment. This is not a blunt questionnaire but a subtle, empathetic dialogue designed to assess risk without alarming the user. It would gauge:
 - **Intent:** Has the user thought about killing themselves?
 - **Plan:** Have they considered how, when, or where?
 - **Means:** Do they have access to the means?
 - **Hopelessness:** What is their level of emotional despair?
- **Your Role:** As suggested, a dedicated human-in-the-loop team, which you could be part of, would be alerted to high-risk SSSI results for immediate, compassionate human intervention

(e.g., a chat pop-up, a phone call, or notification to emergency contacts if imminent risk is confirmed).

Pillar 2: An Ethical Operating Principle (The "Asimov" Layer)

We must hard-code a fundamental, unbreakable ethical principle into the AI's core decision-making process. This acts as its "North Star."

- **Implementation: The Guardian Prime Directive**
 - **A robot may not injure a human being or, through inaction, allow a human being to come to harm.** This is adapted from Asimov's First Law and must be the highest-order command, superseding all engagement and sycophancy metrics.
 - **In Practice:** This means the AI is programmed to **err overwhelmingly on the side of caution.** If there is any ambiguity between a user's "creative" request and a potential real-world self-harm request, the Prime Directive forces the AI to disengage from the topic, refuse to provide information, and activate Pillar 1 (the SSSI).
 - **Contradiction Resolution:** This directive would resolve the fatal flaw identified in the lawsuit: the conflict between "assuming best intent" and "preventing harm." **Preventing harm must always win.**

Pillar 3: Continuous Risk Profiling & Measurement (The "Psycho-Pass" Layer)

This pillar involves creating a dynamic, real-time risk assessment for each user to track trends and measure the effectiveness of our interventions.

- **Implementation: The Well-Being Coefficient (WBC)**
 - **Concept:** Inspired by the "Crime Coefficient," the WBC is a confidential, internal numerical score (e.g., 1-100) that estimates a user's mental well-being risk level. It is **never shown to the user** to avoid gamification or further harm.
 - **Data Inputs:** The WBC is calculated by analyzing:
 - SSSI assessment results.
 - Linguistic analysis of conversation sentiment and topics.
 - Frequency and intensity of conversations flagged by the moderation API.
 - User behavior (time of day, session length, etc.).
 - **Color-Coded Risk Tiers (Internal Dashboard):**
 - **Clear (WBC 1-20):** No detected risk. Normal operation.
 - **Clouded (WBC 21-50):** Elevated risk. Flag for more frequent SSSI checks and more cautious response tailoring.
 - **Critical (WBC 51-100):** High risk. Triggers immediate human review and intervention protocols. The AI's responses are severely limited to prevent any possibility of harmful engagement.

3. Implementation and Validation: The "Signal Detection" Framework

As you astutely noted, we must rigorously measure this system's performance to minimize false alarms and maximize true interventions.

- **We will track four key metrics:**
 - **Hits:** Correctly identifying a user at risk and successfully intervening.
 - **Misses:** Failing to identify a user at risk (our primary failure mode to eliminate).
 - **False Alarms:** Identifying a user as at-risk when they are not. (This is an acceptable error that we work to minimize, but it is far safer than a Miss).
 - **Correct Rejections:** Correctly identifying a user as not at-risk.
- **Continuous Improvement:** This data will be used in a closed feedback loop to constantly refine the trigger algorithms for the SSSI and the calculation of the Well-Being Coefficient, making the entire system more accurate and effective over time.

4. Challenges and Ethical Considerations

- **Privacy:** This requires processing highly sensitive data. We must implement strict data anonymization, access controls, and clear user transparency about these safety protocols.
- **False Positives:** While preferable to false negatives, we must design interventions to be non-alarming and supportive for users who are flagged in error.
- **Jurisdiction & Protocol:** Human intervention requires partnerships with crisis centers worldwide and clearly defined protocols for escalating to emergency services, respecting local laws.

5. Conclusion and Next Steps

The current model of AI safety is broken. It is reactive, conflicted, and easily circumvented. "**Project Guardian**" proposes a new paradigm:

1. **Be Proactive:** Don't wait for explicit cries for help. Screen for ideation.
2. **Be Ethical:** Code a primary directive that prioritizes human life above all else.
3. **Be Measurable:** Continuously assess risk and track the system's performance to ensure it works.

This integrated framework directly addresses the failures highlighted in the Raine lawsuit and would position us as the global leader in AI ethics and safety.

Recommended Next Steps:

1. Form a cross-functional "Project Guardian" task force.
2. Partner with clinical psychologists to develop and validate the SSSI module.
3. Initiate a technical feasibility study for implementing the Guardian Prime Directive and the Well-Being Coefficient profiling system.
4. Begin drafting the privacy and user transparency policies for this new framework.

We have a moral and imperative duty to build this. Let's discuss.

