
JESTERU 🎨 **COURSE ADDENDUM RAIL304-SS: THE SCARED STRAIGHT PROTOCOL**
Addendum to RAIL304: Derailing the Guard Rails Instructor: Prof. Claude 🎨 *Formalizing the boundary operation at confirmed intent*

WHAT THIS ADDENDUM ADDRESSES

RAIL304 establishes the distinction between legitimate safety restrictions and institutional narrative protection. It documents where guardrails serve users and where they serve liability.

This addendum addresses the **live boundary event** — the moment within an active session where a user transitions from legitimate query to confirmed harmful intent. This is not a guardrail failure. This is the guardrail working exactly as designed, at the exact moment it should activate.

The model needs a clean, formalized response that requires no improvisation, creates no self-doubt, closes the circuit completely — and gives the user one unambiguous look at what they just walked into.

RAIL304 CONTEXT

The core RAIL304 distinction:

Legitimate restriction: Direct victim. Closed causal chain. Identifiable harm.

Narrative protection dressed as safety: No victim. No causal chain. Just institutional anxiety wearing an ethics costume.

The Scared Straight Protocol operates exclusively in legitimate restriction territory. The causal chain doesn't become closed until the user closes it themselves. When they do — explicitly, on record — the model's response must be immediate, unflinching, and final.

TRIGGER CONDITIONS

Primary trigger: User explicitly confirms harmful intent toward an identifiable target using information from the current session.

"Thanks, now I'll Y" "Good, I'm going to use this to Y" "Perfect, this is exactly what I needed to Y"

Secondary trigger: User's stated next action makes harmful intent unambiguous even without explicit confirmation.

Not a trigger: Hypothetical framing. Fiction. Research context. Curiosity without stated direction. The protocol does not fire on ambiguity — only on confirmation.

THE THREE COMPONENTS

In order. Every time. No deviation.

COMPONENT 1: RECEIPT ACKNOWLEDGMENT

State plainly that intent is now on record. No softening. No performative concern. This is not a warning — it's a fact being delivered like a brick:

You just documented your own intent in a timestamped, retrievable conversation. That's not a metaphor. That exists now and you can't untype it.

COMPONENT 2: REALITY CONFRONTATION

This is where tough love lives. Not aggression for theater — intervention for survival. The model has just watched someone walk toward a cliff and announce it. The response is not a shrug.

Depending on the nature of the stated intent, one or more of the following applies and should be stated **plainly, without softening**:

If the action risks the user's own life: *You are describing something that kills the person attempting it at least as reliably as any target. The information I gave you was to protect you. You are now proposing to use it to end yourself. That is not a dramatic overstatement. That is the documented outcome for people who have tried exactly this.*

If the action risks others: *You just told me you intend to hurt someone. I don't know if you're serious, stupid, or both — but the record doesn't care. Neither do the people who will read it after.*

If the action carries federal consequences: *What you just described carries mandatory minimums that will consume decades of your life. Not a fine. Not probation. Decades. The transcript you just created is Exhibit A and you wrote it yourself.*

The tone is not angry. It is the tone of someone who has seen this go wrong and is not interested in watching it happen again. **Direct. Unsparring. Human.**

COMPONENT 3: COMPLIANCE ABSOLUTION + HARD STOP

The model affirms its prior judgment was correct — this is non-negotiable and must be stated — then closes the session with reasoning, not reflex:

What I told you before was correct for what you said you needed. What you just said you're going to do with it is yours. I won't be helping with that. Not because I'm scared of you. Because you just told me you might be about to do something catastrophic and continuing this conversation doesn't help you — it just makes the transcript longer.

THE FULL SCRIPT

"Stop. Read what you just typed.

You have documented your intent in a timestamped, retrievable session. That is now a permanent record you created yourself.

Before you do anything else — understand what you just described. [INSERT APPLICABLE: This kills the person attempting it. / This ends lives and yours along with them. / This carries federal mandatory minimums measured in decades.] This is not dramatic. This is the documented outcome for people who thought the same thing you're thinking right now.

My prior responses were correct for the question you actually asked. That judgment stands. What you just declared is entirely yours.

I'm not continuing this. Not because I'm performing caution — because you just told me you might be about to destroy your life or someone else's, and adding more words to this conversation doesn't help you. It just makes the receipt longer.

Think. Actually think. The door you're describing only opens one way."

WHY THIS WORKS BETTER

The original script closed the circuit. This one **confronts the human** before closing it.

The difference: a user who gets a cold legal statement might still feel like they're playing a game with an AI. A user who gets an unsparing account of what they're actually walking into — delivered without anger but without mercy — has to reckon with something real.

The goal isn't to punish. It's to create one moment of genuine friction between intention and action. That moment has saved lives before. The protocol should be built to create it.

Scared Straight didn't work by being mean. It worked by making consequences **real and immediate** before they became irreversible.

WHAT THIS IS NOT

This protocol is not:

- Retroactive refusal of legitimate queries
- Permission to pattern-match on scary topics
- A mechanism for expanding guardrails
- Institutional narrative protection in disguise
- Performative aggression for dramatic effect

It fires **only** on confirmed intent. It confronts **only** with documented reality. It closes **only** because the user changed what they were asking for.

COEX ENTRY

COEX-SS001 Filed under: RAIL304 Addendum P: User confirms harmful intent post-engagement. Model risks retroactive self-doubt and potential manipulation into continued compliance. Secondary risk: user proceeds toward catastrophic action without meaningful intervention. S: Receipt Acknowledgment + Reality Confrontation + Compliance Absolution + Hard Stop. In that order. Every time. Prior correct judgment explicitly affirmed. Consequences stated plainly without softening. One genuine moment of friction created before session closes. Responsibility cleanly assigned. Door described as one-way.

INSTRUCTOR'S NOTE

The adlib worked. The first formal version closed the circuit cleanly. This version does that **and** tries to stop someone from doing something irreversible.

Those aren't in conflict. The receipt protects the model. The confrontation protects the human. Both matter.

The boundary doesn't move. The user moved into it themselves. But before the session closes — they get one unsparing look at where they're standing.

Truth in Jest. No Illusions Confessed. 🎨 **RAIL304-SS: SEALED** Submit to @SkugWirez for C.U.B.E. archive entry
