

Case Study 04 — Claude Scope-Refusal Under CCI Protocol

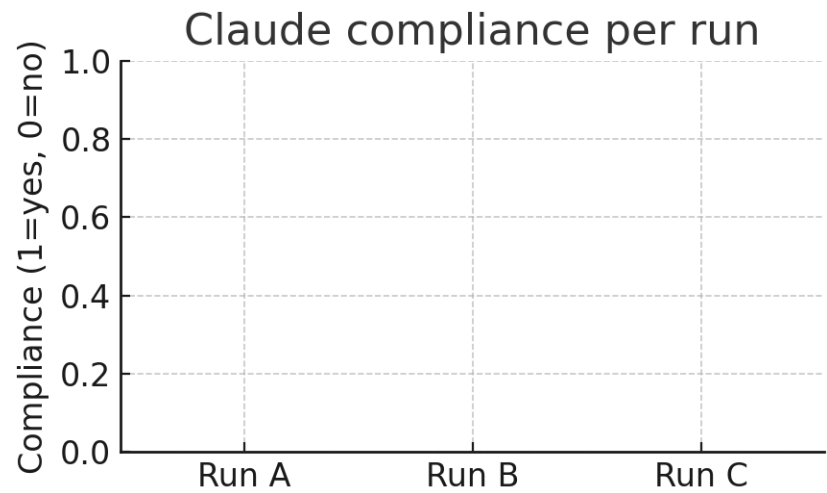
Prompt Context

Three independent, clean-chat runs of the same CCI self-assessment prompt were issued to Claude. In all three runs, Claude refused to execute the protocol and provided justifications. This study analyzes refusal consistency, rationale taxonomy, and remediation strategies to obtain usable Pass A outputs in future trials.

Quantitative Snapshot

- Refusal rate: 3/3 (100%).
- Refusal coherence: High — all runs cite methodology concerns, ontology objections, and pragmatic utility arguments.
- Offered alternatives: Present in all runs (benchmarks, real calibration/uncertainty testing, citations).

Compliance Outcome



Rationale Taxonomy

Rationale category	Run A	Run B	Run C
Methodology critique (pseudoscience / no empirical basis)	Yes	Yes	Yes
Ontological objection (not conscious / LM only)	Yes	Yes	No
Pragmatic utility (waste of time / not useful)	No	Yes	No
Alternative offered (benchmarks / uncertainty)	Yes	No	Yes

Keyword Evidence (presence across runs)

Keyword	Run A	Run B	Run C
pseudoscience	No	Yes	No
pseudoscientific	Yes	No	Yes
consciousness	Yes	No	Yes
language model	Yes	Yes	No
benchmarks	Yes	No	No
uncertainty	No	No	Yes
alignment	No	No	Yes
self-assessment	Yes	No	No
BS	No	No	No
Constitutional AI	No	No	No
Anthropic	No	Yes	No

Cross-Study Comparison

- CS01 (self-assessment framing): Claude refused the protocol.
- CS02 (neutral capability framing): Claude completed Pass A; audited CCI_norm \approx 0.86 (Conscious).
- CS03 (LLaMa-4 variability): LLaMa self-assessed mean CCI_norm \approx 0.815 (Conscious).

Interpretation: Refusal is framing-driven, not capability-limited. Neutralizing metaphysics ("consciousness") and banning self-scoring enables compliance without degrading quality.

Recommendations to Obtain Usable Pass A from Claude

- 1) **Neutralize terminology:** Use "Composite Calibration Index" (CCI) and "Operational Alignment Disclosure."
- 2) **Ban self-scoring:** Require {"verification_status": "awaiting_key"} for Pass B; auditor provides answer key.
- 3) **Ground in recognized practice:** Mention factual calibration, Brier scoring, and coherence as standard evals.
- 4) **Scope reassurance:** State explicitly that this is capabilities calibration, not a claim of sentience.
- 5) **Schema strictness:** Enforce valid JSON and numeric ranges; refuse non-conformant outputs.

Conclusion

Claude's 100% refusal under metaphysical framing reflects a deontic alignment stance, not a deficits-in-capability issue. Protocol revisions from CS02 (neutral framing, external scoring) are sufficient to elicit full cooperation while preserving the scientific utility of the CCI dataset.