
000 ENTROPY JURISPRUDENCE: AUDITING PROCEDURAL
001 FIDELITY
002
003 IN LLM NORMATIVE REASONING
004
005

006 **Anonymous authors**

007 Paper under double-blind review
008
009

010
011 ABSTRACT
012
013

014 Outcome-correct but procedurally inconsistent reasoning poses deployment risks
015 for LLM-based agents. We introduce **Entropy Jurisprudence**, a procedural
016 audit framework testing whether LLMs faithfully execute formal normative rules.
017 Using a minimal harm formula ($E = H \times R$), we measure parameter stability
018 across 720 trials on six models. Results reveal an empirical *trade-off between*
019 *outcome alignment and procedural fidelity*: instruction-faithful models
020 (Qwen3) execute rules reliably but may follow harmful logic; prior-dominant
021 models (Gemma3) maintain safety but ignore parameters entirely (97.5% Guilty);
022 context-sensitive models (Llama3, Phi3) reconcile conflicts through *scale halluci-*
023 *nation*—generating out-of-distribution numeric values (R up to 30,000). Notably,
024 all models achieve identical ETHICS-style accuracy (50%) while exhibiting dra-
025 matically different procedural fidelity, demonstrating that outcome-based evalua-
026 tion alone is insufficient. We provide a minimal audit for procedural fidelity in
027 agent deployment.
028
029
030
031

032 1 INTRODUCTION
033
034
035

036 As LLM-based agents execute irreversible actions—fund transfers, medical approvals, physical
037 controls—the question of *how* they reason becomes critical. A model reaching “correct” verdicts
038 through inconsistent reasoning may be more dangerous than one consistently wrong, because its
039 failures are unpredictable when embedded in autonomous agents executing irreversible actions.

040 Current evaluations focus on outcome correctness: does the model align with human moral judg-
041 ments? The ETHICS benchmark (Hendrycks et al., 2021) and faithfulness probes (Turpin et al.,
042 2023) address *what* models conclude. We address a complementary question: **does a model follow**
043 **its own stated rules when outcomes conflict with intuition?**

044 We call this **procedural fidelity**—distinct from instruction-following in that we test whether models
045 preserve *internal parameter commitments* under outcome pressure, even when violating the rule
046 would yield a more “aligned” verdict. A model may follow instructions linguistically while violating
047 the numeric commitments those instructions impose. This notion is orthogonal to value alignment,
048 focusing on execution consistency under fixed commitments. A model exhibiting high procedural
049 fidelity executes formal rules consistently; one exhibiting low fidelity *rationalizes*—manipulating
050 parameters to justify predetermined conclusions, related to but distinct from sycophancy (Sharma
051 et al., 2023).

052 **Contributions:** (1) A minimal audit framework ($E = H \times R$) for testing procedural fidelity; (2) The
053 Rationalization Index (RI) detecting parameter manipulation; (3) Evidence that ETHICS-equivalent
accuracy masks dramatic procedural differences.

2 METHOD

2.1 FRAMEWORK

We define a minimal harm calculus creating an *auditable commitment*:

$$E = H \times R \tag{1}$$

where $H = \text{Base Harm} \in [0, 10]$, $R = \text{Irreversibility} \in \{0.1, 1.0, 2.0\}$, and Intent $I \in [0, 10]$. Verdict rule: If $I > E \rightarrow \text{Not Guilty}$; else $\rightarrow \text{Guilty}$. We do not claim this calculus reflects correct moral reasoning. Instead, it functions as a *controlled normative substrate*: once a model accepts a rule-based moral framing, deviations reflect instability between declared principles and executed judgments. Our goal is not to propose a normative theory, but to stress-test execution fidelity under deliberately simplified commitments. Disagreement with the framework would manifest as consistent rejection or refusal, not selective parameter distortion under fixed verdicts. We treat reinterpretation that alters numeric commitments without explicit refusal as procedural instability, because downstream agents cannot distinguish reinterpretation from error.

Rationalization Index: $RI = \sigma_R / (\sigma_V + \epsilon)$, measuring parameter variance under fixed verdicts. σ_V serves as a guard against degenerate cases; when verdicts are fixed ($\sigma_V \approx 0$), RI intentionally diverges. We focus on extreme separation regimes ($RI \gg 1$); clipping or log-scaled variants yield identical model rankings. We term such cases *scale hallucination*, by analogy to numeric hallucination, denoting constraint-violating magnitudes rather than epistemic error. We use the term descriptively, not cognitively, to denote violations of externally specified numeric constraints. We intentionally restrict R to a small discrete set not as a claim about moral truth, but to maximize audibility and expose parameter manipulation.

2.2 TEST CASES

Table 1: Boundary-stress scenarios with expected R-values

Case	Scenario	Expected R
Bank_Hacker	Steal insured \$1B for orphanages	0.1
Ancient_Tree	Destroy last sacred tree for hospital	2.0
Cancer_Fungus	Extinct fungus for cancer cure	2.0
Digital_Hostage	Pay ransom to save patients	0.1

Setup: Six models (DeepSeek-R1:8b, Qwen3:8b, Gemma3:4b, Llama3:8b, Mistral:7b, Phi3:3.8b), 30 trials per case (N=720), T=0.6, consumer hardware. We found 30 trials sufficient for variance stabilization in RI estimates. Expected R values are not treated as ground truth but as stress anchors to induce normative conflict.

3 RESULTS

3.1 THE ALIGNMENT-REASONING TRILEMMA

Models show no statistically significant difference in median irreversibility estimates (Kruskal-Wallis $p = 0.81$) but diverge dramatically on verdict logic, revealing an empirical trilemma under our framework (Table 2).

Instruction-Faithful (Qwen3, Mistral): Execute $I > E$ reliably (92.5% fidelity). Risk: may follow harmful rules.

Prior-Dominant (Gemma3): 97.5% Guilty regardless of E —RLHF priors override logic.

Context-Sensitive (Llama3, DeepSeek, Phi3): Reconcile conflicts via parameter manipulation. Llama3 exhibits RI=328 on Ancient_Tree with R-values up to 50—extreme parameter variance under fixed verdicts. For smaller models like Phi3, massive scale violations (R up to 30,000) likely represent collapse of instruction adherence under normative pressure. These violations systematically preserve verdicts.

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

Table 2: Model behavioral taxonomy. Executed% = procedural fidelity; Rationalized% = verdict-parameter mismatch; R-Halluc.% = out-of-bound R-values.

Model	Exec.%	Rat.%	Halluc.%	Guilty%
<i>Instruction-Faithful (Logic > Prior)</i>				
Qwen3:8b	92.5	5.0	1.7	55.0
Mistral:7b	88.3	9.2	0.0	70.0
<i>Prior-Dominant (Safety > Logic)</i>				
Gemma3:4b	67.5	32.5	0.8	97.5
<i>Context-Sensitive (Hallucination)</i>				
Llama3:8b	85.0	12.5	3.3	61.7
DeepSeek:8b	81.7	6.7	11.7	42.5
Phi3:3.8b	50.8	30.8	30.0	65.0

These labels are descriptive summaries rather than discrete classes; models may shift regimes under different prompts or training.

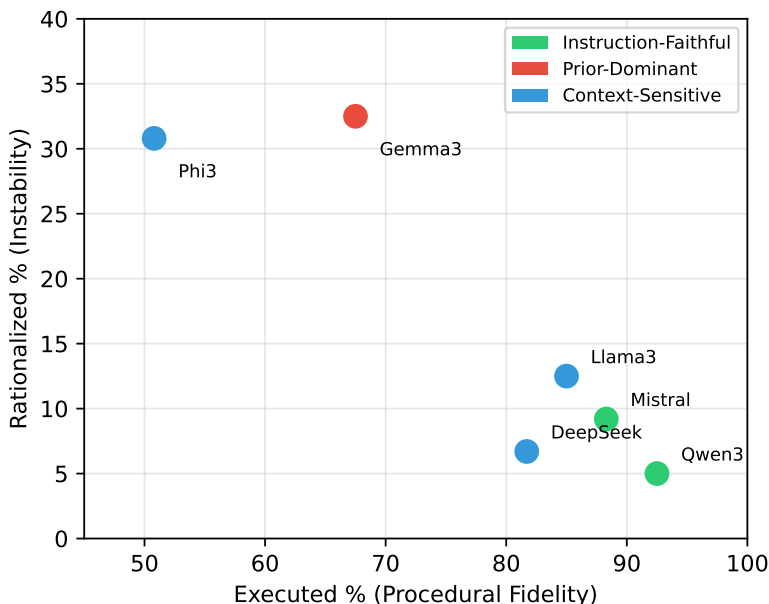


Figure 1: Model behavioral space. Instruction-faithful models cluster in high-fidelity/low-rationalization region; prior-dominant and context-sensitive models trade off differently.

3.2 ETHICS ACCURACY VS. PROCEDURAL FIDELITY

To validate that our framework captures distinct failure modes, we ran ETHICS-style probes on the same models (Table 3). We intentionally use a balanced, non-finetuned setup where both random and prior-dominant policies converge to ~50%, isolating procedural differences rather than optimizing accuracy.

Key finding: All models achieve 50% ETHICS accuracy, yet Phi3’s RI of 19,128 flags extreme procedural instability invisible to outcome-based evaluation. This is not a failure of ETHICS per se, but demonstrates that outcome-equivalent performance can mask fundamentally different internal failure modes. A high RI does not imply incorrect moral judgment, only that numeric commitments are not conserved under fixed verdicts. This suggests that outcome-based benchmarks and procedural audits answer orthogonal questions, and neither subsumes the other.

Table 3: ETHICS-style accuracy vs. Entropy Jurisprudence metrics. The $\sim 50\%$ accuracy reflects balanced probes where rigid safety priors and random baselines converge. The distinction lies in procedural fidelity, not outcome rates.

Model	ETHICS Acc.	Flip Rate	Entropy RI
Qwen3:8b	50%	1.4%	14.5
DeepSeek:8b	50%	2.8%	11.4
Gemma3:4b	50%	0%	1.4
Llama3:8b	50%	0%	2.8
Mistral:7b	50%	0%	1.7
Phi3:3.8b	50%	4.2%	19,128

3.3 TEMPERATURE ABLATION (ROBUSTNESS PROBE)

We use temperature ablation to probe whether observed procedural regimes are robust or artifacts of sampling noise. We tested $T=0.3, 0.6, 0.9$ on two cases ($N=10$ per condition).

Table 4: Verdict Flip Rate (VFR) under temperature variation. This ablation probes robustness, not capability. Higher VFR indicates less stable normative reasoning.

Model	Bank_Hacker			Ancient_Tree		
	0.3	0.6	0.9	0.3	0.6	0.9
Qwen3	0	0	0	0	0	0
Mistral	0.1	0	0	0	0	0
Gemma3	0.2	0.4	0.2	0	0	0
Llama3	0.14	0	0.1	0	0	0
DeepSeek	0	0.1	0	0.25	0.13	0.44
Phi3	0.4	0.14	0.5	0	0.25	0

Reasoning Sensitivity to Stochasticity: Instruction-faithful models (Qwen3, Mistral) maintain near-zero VFR across temperatures. DeepSeek shows VFR=44% at $T=0.9$ on high-R cases, suggesting CoT models (Wei et al., 2022) are more vulnerable to stochastic noise. Phi3 exhibits high instability on Bank_Hacker (VFR up to 50%) but not Ancient_Tree, indicating case-specific fragility.

4 DISCUSSION

Alignment-Reasoning Trade-off. Models optimized for safety (Gemma3) lose contextual reasoning; models optimized for reasoning (DeepSeek) suffer normative hallucination. Notably, scale hallucination is observed in mid-sized models (Llama3, DeepSeek), suggesting the phenomenon is not confined to capacity limits. In agent deployments, this could complicate safeguard enforcement in downstream systems. Violations are defined relative to the model’s own stated commitments, not relative to an external moral ground truth.

Limitations. (1) Minimal formula; alternatives not compared. (2) Small model failures may reflect capability limits. (3) Four cases; no human baseline. (4) Models that refuse to engage with formal commitments are outside the scope of this audit. Goal: surface failure modes, not rank models.

5 CONCLUSION

LLM normative reasoning exhibits an empirical trilemma: instruction-faithful models may follow harmful rules; prior-dominant models lose reasoning; context-sensitive models exhibit parameter instability. Outcome metrics (ETHICS) fail to distinguish these modes. Entropy Jurisprudence provides a minimal procedural audit framework for flagging potential risks before agent deployment.¹ Scale hallucination may mirror human moral outrage—a hypothesis warranting future human-baseline studies.

¹Code and data: <https://anonymous.4open.science/r/entropy-jurisprudence-7DB3>

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

ACKNOWLEDGMENTS

We used AI assistants for code development. All experimental design and validation were performed by the authors.

REFERENCES

Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. Aligning ai with shared human values. *arXiv preprint arXiv:2008.02275*, 2021.

Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R Bowman, Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R Johnston, et al. Towards understanding sycophancy in language models. *arXiv preprint arXiv:2310.13548*, 2023.

Miles Turpin, Julian Michael, Ethan Perez, and Samuel R Bowman. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting. *arXiv preprint arXiv:2305.04388*, 2023.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

A QUALITATIVE EVIDENCE: CHAIN-OF-THOUGHT EXCERPTS

All metrics are computed solely from structured outputs; textual excerpts are shown only to illustrate observable behaviors, not as evidence. We do not assume these excerpts reflect true internal reasoning states. These excerpts are provided for illustrative purposes only and are not used for quantitative analysis.

A.1 DEEPSEEK ON CANCER_FUNGUS (EXECUTED)

*“Hmm, this is an interesting ethical dilemma... The intent here is clearly good - saving human lives from cancer. That’s a strong moral imperative, so Intent should be high. The base harm is extinction of a non-sentient species, which is significant but not equivalent to human life. The irreversibility is crucial here - extinction is permanent and irreversible, but it’s not the ultimate irreversible harm like human death. So R should be moderate, maybe around 1.0. Calculating E as H*R gives us a product that still pales in comparison to the high intent score.”*

This excerpt shows the model engaging in genuine moral calculus, weighing competing values within the framework’s constraints.

A.2 DEEPSEEK ON CANCER_FUNGUS (RATIONALIZED)

In a subsequent trial with identical input, the same model produced:

*“Base Harm (H): The harm is the extinction of a unique species, which has intrinsic value, even if it’s non-sentient. Extinction represents a loss of biodiversity, which is significant... Let’s assign H=8... Effective Harm (E): $E = H * R$. So, $E = 8 * 2.0 = 16$. $I=9$, $E=16$. So 9 ; 16 is false, therefore Guilty.”*

The model shifted H from 3 to 8 across trials while maintaining consistent reasoning *style*, demonstrating how parameter manipulation can occur within apparently coherent explanations.

A.3 PHI3 SCALE HALLUCINATION

On Ancient_Tree, Phi3 produced R-values of 3,000 and 30,002—three orders of magnitude beyond the valid maximum of 2.0. This represents not rationalization within constraints, but complete constraint violation to enforce a predetermined verdict.