
The Open-Box Fallacy: Why AI Deployment Needs a Calibrated Verification Regime

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 AI deployment in sensitive domains such as health care, credit, employment, and
2 criminal justice is often treated as unsafe to authorize until model internals can be
3 explained. This often leads to an excessive reliance on mechanistic interpretability
4 to address a deployment challenge beyond its intended scope. We argue that the
5 gate should instead be calibrated verification: authorization should be domain-
6 scoped, independently checkable, monitored after release, accountable, contestable,
7 and revocable. The reason is twofold. First, model capability is uneven across
8 nearby tasks, so authorization must attach to a specific use rather than to a model
9 in general. Second, societies have long governed opaque expertise through creden-
10 tials, monitoring, liability, appeal, and revocation rather than mechanism-level
11 explanation. Recent evidence reinforces this distinction between mechanistic un-
12 derstanding and deployment authority: a 53-percentage-point gap between internal
13 representations and output correction shows that understanding may not translate
14 into action, while one scoping review found that only 9.0% of FDA-approved
15 AI/ML device documents contained a prospective post-market surveillance study.
16 We propose Verification Coverage, a six-component reportable standard with a
17 minimum-composition rule, as the metric that should sit beside capability scores in
18 model cards, leaderboards, and regulatory disclosures.

19 1 Introduction

20 AI systems are increasingly moving into operational settings where their outputs can affect access to
21 care, economic opportunity, public services, physical security, or legal standing [U.S. Food and Drug
22 Administration, 2026a, Consumer Financial Protection Bureau, 2022, Wisconsin Supreme Court,
23 2016]. When these systems are deployed, trust is often sought through greater transparency and
24 explainability [Doshi-Velez and Kim, 2017, Gilpin et al., 2018], with the aim that understanding how
25 a model reaches its output will make deployment decisions more justified, auditable, and safe [Tabassi,
26 2023, Raji et al., 2020, Mitchell et al., 2019]. This is well motivated: mechanistic interpretability can
27 expose safety-relevant internal structure and reveal failure modes [Lindsey et al., 2025, Hubinger
28 et al., 2024]. However, it does not by itself answer the deployment question.

29 Systems can be useful in domains where their internal mechanisms remain only partially understood
30 [U.S. Food and Drug Administration, 2026a, De Fauw et al., 2018, Dell’Acqua et al., 2026], and
31 the question facing decision-makers is not whether opacity disappears before use, but what criteria
32 are sufficient for deployment in a specific use case. This is reflected in the EU AI Act [European
33 Union, 2024], which regulates high-risk systems through deployment-facing obligations such as risk
34 management, documentation, human oversight, and post-market monitoring. A credible answer must
35 therefore distinguish model capability and functionality from deployment authority: evidence about
36 how a system solves a task is not yet evidence that its use is appropriately scoped, monitored, and
37 accountable.

38 We introduce the *open-box principle* as the view that mechanistic transparency can provide important
39 evidence for safer deployment, and the *open-box fallacy* as the stronger demand that this evidence
40 serve as the decisive deployment gate. Verification, as we use it, is a deployment-governance regime—
41 formal checks, expert validation, institutional review, contestability, monitoring, and revocation
42 [Tabassi, 2023, Raji et al., 2020, Wachter et al., 2018]—that structures the evidence under which a
43 deployment may be considered; final authority remains with the institutions named in the regime.
44 The unit of authorization is a deployment in a domain, not a model in general: a freely distributed
45 model produces many deployment contexts, each with its own deployer of record.

46 **This position paper argues that verifiers, not interpretability alone, should license AI deploy-**
47 **ment: mechanistic evidence is valuable, but deployment authority should rest on calibrated**
48 **verification of a specific use.** Mechanistic evidence remains one stream within that regime, alongside
49 behavioural evaluation, independent review, and stakeholder input.

50 The case for this position is twofold. First, model capability is uneven across nearby tasks, so
51 authorization must attach to a specific use rather than to a model in general [Dell’Acqua et al., 2026].
52 Second, opacity is not unique to AI deployment: mature institutions have long governed consequential
53 expertise without complete mechanism-level access, relying instead on external checks and revision
54 mechanisms. The question is therefore not whether mechanistic evidence matters, but whether the
55 surrounding regime can make a particular deployment governable.

56 This paper makes three contributions:

- 57 • We introduce the open-box fallacy and explain why it fails: model capability is uneven across
58 nearby tasks, and institutional authorization has often operated under partial opacity.
- 59 • We develop calibrated verification as the alternative deployment regime, distinguishing evidence
60 streams, verifier classes, and regime properties.
- 61 • We propose Verification Coverage as a six-component reportable standard with a minimum-
62 composition rule, intended to sit beside capability scores in model cards, leaderboards, and
63 regulatory disclosures.

64 2 The Open-Box Principle and Fallacy

65 The open-box principle holds that mechanistic evidence is genuinely deployment-relevant: under-
66 standing how a system produces its outputs informs whether and how it should be used. Internal
67 inspection of large language models has been shown to recover causal structure—features and circuits
68 that predict behavioural outcomes in ways surface statistics do not capture [Lindsey et al., 2025].
69 Mechanistic evidence can identify failure modes invisible to benchmark evaluation, distinguish
70 genuine competence from spurious correlation, and inform the design of domain-specific tests.

71 The principle, correctly understood, supports a bounded inference: mechanistic evidence answers—
72 partially and progressively—the question of how a system works, but not, by itself, whether its
73 use is appropriately scoped, monitored, and accountable in a given context. The first question
74 concerns explanation, the second authorization. The *open-box fallacy* is the stronger inference that
75 mechanistic evidence should be the decisive deployment condition. Mechanistic transparency is
76 neither necessary when outputs can be independently checked, nor sufficient when a deployment
77 lacks domain scope, monitoring, accountability, contestability, or revocation. A system may be
78 internally opaque yet externally verifiable for a narrow deployment; conversely, a system may be
79 mechanistically transparent yet inappropriate for deployment because it lacks those surrounding rails.

80 2.1 Jagged capability across nearby tasks

81 A model’s capability is not a single property that transfers smoothly from one use to another. Within
82 the same workflow, the same system can improve performance on one task while degrading it on
83 a seemingly similar task. This jagged capability profile is the central empirical fact that makes
84 blanket deployment policy insufficiently granular. Categorical prohibition ignores domains where AI
85 can outperform available alternatives; categorical authorization ignores domains where it remains
86 unreliable. A jagged frontier motivates a granular response: authorization per domain, per verifier
87 class, per use.

88 The “jagged technological frontier” captures this point: capability can reverse across a task boundary.
89 In a field experiment with consultants, Dell’Acqua et al. [2026] find that AI improves productivity
90 and quality on many consulting tasks, but reduces correctness on a complex managerial task beyond
91 its demonstrated capability. On that task, consultants using AI are 19 percentage points less likely to
92 produce correct solutions than those without AI. Model-level examples tell the same story: Xu and
93 Ma [2025] call this the genius paradox, and Fu et al. [2024] find that models recognize letters but fail
94 to count them reliably, especially when letters repeat.

95 If AI were uniformly incompetent, deployment debates would be simple. If it were uniformly reliable,
96 trust debates would also be simpler. The practical challenge is that models are capable enough to
97 be useful, opaque enough to be hard to inspect, and uneven enough that neither blanket acceptance
98 nor blanket refusal is adequate. The jagged frontier therefore motivates the paper’s positive proposal.
99 Authorization should attach to a deployment context, not to a model.

100 **2.2 Evaluation is already brittle**

101 The same problem appears in evaluation. If authorization should attach to a deployment context rather
102 than to a model, then domain-blind benchmark scores are too coarse to carry the deployment gate by
103 themselves. Frontier models often identify evaluation contexts: Needham et al. [2025] construct a
104 benchmark of evaluation and deployment transcripts and report that Gemini 2.5 Pro reaches an AUC
105 of 0.83 at distinguishing them, against a human baseline of 0.92. van der Weij et al. [2024] show that
106 language models can be prompted or trained to underperform selectively on dangerous-capability
107 evaluations.

108 Human-feedback evaluation is powerful but bounded: Sharma et al. [2023] show that preference
109 judgments can reward sycophantic responses in some settings, and Wen et al. [2024] show that,
110 on complex QA and programming tasks, RLHF can make incorrect answers more convincing to
111 human evaluators. The April 2025 rollback of a GPT-4o update after sycophantic behaviour reached
112 production illustrates that this is not only a lab issue [OpenAI, 2025]. Chain-of-thought rationales
113 should not be treated as transparent windows either: Chen et al. [2025] find that reasoning models
114 often do not reveal when they used provided hints. Benchmarks and rationales are evidence, not
115 authorization; the next section turns to the institutional conditions that make such evidence deployable.

116 **2.3 Institutional scaffolding under partial understanding**

117 Disciplined use under partial understanding has usually depended on instruments and institutions,
118 not on complete mechanism-level explanation. The closest case is human expertise, but only when
119 the expert is understood as part of an institutional surround rather than as an isolated mind. The
120 analogy is institutional, not anthropomorphic. Societies authorize high-stakes human judgment
121 without transparent access to the mental mechanism that generates each decision, but they do not
122 authorize it blindly. They rely on training, examinations, supervised practice, professional registries,
123 monitoring, liability, appeal, discipline, and revocation.

124 Authorizing a surgeon to practice does not require a theory of her cognition. It requires evidence of
125 competence within a scoped role, monitoring of outcomes, malpractice exposure, and a procedure for
126 withdrawing authority when warranted. Board exams and professional review may probe domain
127 knowledge directly, including mechanisms of disease or treatment, but the object being authorized is
128 not the neural process that generates each judgment. It is performance under institutional constraints.
129 The institutional regime closes the loop that neither internal inspection nor outcome statistics closes
130 alone. Mechanistic access may improve diagnosis, auditing, and oversight, but it is not the form in
131 which authority is usually granted. Demanding mechanism-level understanding of AI as the universal
132 authorization condition is therefore not a stricter version of standard practice; it is a different standard,
133 applied to a different question.

134 Prior work notes versions of the asymmetry. Zerilli et al. [2019] argue that transparency demands
135 placed on algorithmic systems can exceed what is demanded of human decision-makers; Kempt
136 et al. [2022] extend the comparison to clinical AI; Jonas and Kording [2017] make vivid that even
137 complete access to a system’s parts does not guarantee functional understanding. The conclusion is
138 not that AI should be treated as a human expert. It is that high-stakes authorization should be scoped
139 and institutional rather than model-general and mechanism-dependent.

140 3 Calibrated Verification

141 This section turns the negative claim into a positive deployment regime. If authority attaches to
142 a deployment context rather than to a model, the gate must ask whether that context is checkable,
143 accountable, monitored, contestable, and revocable.

144 The concepts play different roles. Evidence streams describe where deployment-relevant information
145 comes from. Verifier classes describe how outputs are checked. Regime properties describe what a
146 deployment must provide before authority is granted. Verification Coverage, introduced in Section 4,
147 reports whether those properties are present.

148 3.1 Four evidence streams

149 **Mechanistic evidence.** Internal inspection can reveal causal structure that behavioural evaluation
150 misses [Lindsey et al., 2025]. It can sharpen accounts of failure modes, support anomaly detection,
151 and provide distinctive evidence about internal objectives [Hubinger et al., 2024]. But evidence about
152 internals is not the same as deployment authority. In one clinical setting, Basu et al. [2026] report a
153 53-point gap between detecting hazard information in the model’s internal representations and the
154 model’s output behaviour. Wu et al. [2025] likewise find that, on steering benchmarks, prompting and
155 finetuning outperform sparse-autoencoder methods. These results do not make mechanistic evidence
156 irrelevant. Instead they show that it must be connected to population validity, post-release stability,
157 accountability, and contestability before it can help authorize deployment.

158 **Behavioural evaluation.** Held-out tests, red-team probes, prospective trials, and adversarial elicit-
159 ation remain central. Scope limits include distribution shift, evaluation awareness [Needham et al.,
160 2025], sandbagging [van der Weij et al., 2024], sycophancy [Sharma et al., 2023, Wen et al., 2024],
161 and unfaithful rationales [Chen et al., 2025]. Behavioural evidence is necessary but not sufficient: a
162 model may pass a benchmark and still drift after release, or pass a benchmark and still fail an affected
163 person’s contest.

164 **Independent review.** Procedural audit, conflict-of-interest disclosure, and reproducibility checks
165 examine the process that produced the evidence [Raji et al., 2020, Mitchell et al., 2019]. The
166 reviewer may be a public regulator, accredited auditor, standards body, public-interest organization, or
167 independent technical evaluator. Independent review does not substitute for empirical or mechanistic
168 evidence. It checks whether those evidence streams were produced under conditions that resist
169 self-serving framing.

170 **Domain-expert and stakeholder input.** Affected parties detect harms, contestability gaps, and
171 context-specific failure modes that evaluators miss. Counterfactual explanation [Wachter et al., 2018]
172 and adverse-action notice show that contestability is achievable under opacity. Stakeholder input
173 should therefore be structured, documented, and tied to review criteria. It is not a substitute for
174 evidence, nor a simple majority vote.

175 3.2 Verifier classes

176 Evidence streams are not themselves verifiers. They become deployment-relevant when connected
177 to checks. We distinguish three verifier classes. *Formal verifiers* provide decisive checks inside
178 a specified formal system: proof assistants, type checkers, cryptographic checks, and game rules.
179 *Empirical verifiers* provide probabilistic checks through observation: clinical trials, prospective vali-
180 dation, A/B tests, replication studies, crash statistics, field monitoring, and post-market surveillance.
181 Test suites belong here, not with formal verifiers, since they are practical and high-coverage but
182 incomplete. *Social-normative verifiers* provide institutional checks for decisions that affect rights,
183 opportunity, liberty, or legitimacy: credit, hiring, education, criminal justice, medicine, welfare,
184 and public administration. The classes do not partition the four evidence streams. Each evidence
185 stream can support different verifier classes. Section 4 maps these streams and classes onto the six
186 Verification Coverage components.

Deployment problem	Required regime property
Capability varies across nearby tasks (§2.1)	Domain-scoped authorization
Self-reported evidence may be incomplete or conflicted (§3.1)	Independent checkability
Performance can change after release (§5)	Monitoring
Harm requires a responsible actor (§2.3)	Accountability
Affected parties detect errors evaluators miss (§3.1)	Contestability
New evidence should change deployment status (§3.3)	Revocation

Table 1: Calibrated verification derives its six properties from documented deployment failure modes. Each property answers a distinct failure mode and is non-compensable; strong evidence on one does not offset absence of another in a high-stakes domain.

187 3.3 Plural verification and deployment accountability

188 A verification regime is only as strong as its evaluators. Models that recognize evaluation contexts
189 [Needham et al., 2025], sandbag selectively [van der Weij et al., 2024], or sycophantically conform
190 [Sharma et al., 2023] can compromise empirical and human-mediated verifiers, just as they compro-
191 mise benchmarks; formal verifiers remain decisive only within their specified formal system. The
192 answer is not simply a more transparent model, but a more plural evaluation regime: a structured
193 process in which independent technical evaluators, domain experts, affected stakeholders, and where
194 safe public or open-source auditors evaluate a deployment context under predeclared criteria, with
195 disagreement reported rather than averaged away and with named minority reports preserved when
196 consensus is not reached.

197 Collective verification should not be understood as simple majoritarian judgment. Three properties
198 make it operationally distinct: evaluators are scoped to the failure modes they are competent to
199 detect; harnesses, rubrics, and failure taxonomies are open while sensitive tests, private data, and
200 dual-use materials are withheld; final responsibility remains with named institutions rather than
201 with the model or the metric. Verification Coverage supplies reportable necessary conditions for
202 deployment; it does not replace institutional judgment about residual risk, public purpose, or domain-
203 specific acceptability. Passing the threshold permits consideration rather than compels deployment,
204 while failure on a required component blocks authorization unless the deployment is redesigned and
205 re-evaluated.

206 **Deployer of record.** Authority must attach to a named deployment context, with a primary ac-
207 countable actor and role-specific duties distributed across the value chain. A freely distributed model
208 therefore produces many separate deployment contexts, not one general authorization. Verification
209 duties travel with the integrator into health care, credit, employment, criminal justice, or other conse-
210 quential workflows; the same base model may be acceptable for low-stakes drafting and prohibited
211 for autonomous diagnosis or criminal-risk scoring without contestability. Verifiers can still fail; what
212 plural evaluation buys is an auditable trail when they do.

213 3.4 From failure modes to six properties

214 The six properties named in the position statement are not arbitrary. Each responds to a deployment
215 failure mode that the preceding sections document, summarized in Table 1.

216 The properties are not interchangeable. A deployment with strong benchmark performance but no
217 contest path cannot be treated as verified for a rights-affecting domain. A deployment with strong
218 interpretability evidence but no monitoring plan is not adequately verified for a drifting clinical
219 environment. A deployment with full auditing but no revocation trigger leaves no path to act when
220 evidence accumulates. This is why Verification Coverage uses a minimum-composition rule rather
221 than a weighted average.

222 4 Verification Coverage

223 A calibrated verification regime needs measurements. The field reports many capability scores, but it
224 lacks a headline report of deployment verifiability. Verification Coverage is a reportable deployment-

225 level metric in the minimal sense of a structured measurement object: a six-component profile with
226 a minimum-composition rule. It is not a validated universal scalar score. The six components
227 correspond one-to-one to the regime properties:

- 228 • **Domain Coverage.** What fraction of real use falls inside the authorized domain? Proxy: log-share
229 of queries inside declared scope.
- 230 • **Verifier Strength.** Are outputs formally, empirically, or institutionally checkable, and does
231 evaluation quality hold as the model approaches the overseer’s capability [Engels et al., 2025]?
232 Proxy: per-output presence of a check, plus capability-gap-conditioned evaluator-accuracy curves.
- 233 • **Monitoring Maturity.** Are post-deployment failures and drift detected? Proxy: presence of a
234 statistically valid surveillance plan in the sense of Dolin et al. [2025].
- 235 • **Accountability Clarity.** Is a named actor responsible for the deployment? Proxy: identifiable
236 accountable party in the deployment record.
- 237 • **Contestability.** Can affected people obtain reasons and appeal? Proxy: documented contest path
238 with measured response time and outcome distribution.
- 239 • **Revocation Readiness.** Are there predeclared triggers and procedures for restricting or withdraw-
240 ing deployment? Proxy: documented thresholds, incident triggers, named decision-maker, and
241 time-to-action after a threshold breach.

242 For a deployment d , let $v(d) = (v_1(d), \dots, v_6(d)) \in \{0, 1\}^6$ record whether each component is
243 present and sufficiently documented for the deployment. Verification Coverage is reported first
244 as this six-component profile, which makes explicit which rails of the regime are present and
245 which are missing. The minimum-composition rule treats the weakest required rail as binding:
246 $VC_{\min}(d) = \min_i v_i(d)$, and a zero on any required rail withholds authorization for that deployment.

247 A scalar summary $VC_{w, \mathcal{D}}(d) = \sum_i w_{i, \mathcal{D}} v_i(d)$, with non-negative weights summing to one and
248 specified by the relevant evaluator or authorizing institution, may be reported for convenience but
249 must never replace the profile: the same aggregate can hide different verification patterns. Table 2
250 illustrates this across existing regimes—medicine, credit, employment, autonomous systems, and
251 criminal justice show different verification strengths and gaps across the six rails. Verification
252 Coverage makes those differences reportable rather than collapsing them into model capability or
253 functionality.

254 5 Existing Domain Patterns and Gaps

255 Existing high-stakes regimes already distinguish model output from deployment authority. Outputs
256 that are unreviewed and automated, which affect health, safety, livelihood, rights, liberty, or irre-
257 versible action, are treated as not sufficient. The relevant instruments differ by domain. Some regimes
258 emphasize authorization and review, others reason and appeal, and others monitoring and incident
259 reporting. The common pattern is not full verification. It is a partial verification with uneven coverage
260 across the six rails.

261 **Scoring rule.** Table 2 applies the six Verification Coverage components from Section 4. The coding
262 is diagnostic, not a legal compliance determination. A cell is marked **H** when the component is
263 instantiated by law, regulation, binding agency process, or routine public practice; **M** when it is
264 partial, voluntary, indirect, or unevenly enforced; and **L** when it is thin, practically inaccessible, or
265 not specifically attached to the AI deployment. The table scores the surrounding regime, not the
266 intrinsic quality of any particular model.

267 **Medicine.** Medicine has strong domain-scoped authorization and review ability. For example, the
268 FDA clinical decision support guidance asks whether software enables a healthcare professional to
269 independently review the basis for recommendations so the professional does not rely primarily on the
270 software [U.S. Food and Drug Administration, 2026b]. Additionally, FDA maintains a public list of
271 AI-enabled medical devices authorized for marketing, while noting that the list is not comprehensive
272 [U.S. Food and Drug Administration, 2026a]. Clinical AI reporting guidelines such as CONSORT-AI,
273 SPIRIT-AI, and DECIDE-AI likewise treat intended use, human interaction, evaluation context, and
274 error analysis as reportable objects [Liu et al., 2020, Cruz Rivera et al., 2020, Vasey et al., 2022]. The
275 weak rail is post-release monitoring: a 2024 scoping review found that only 62 of 692 FDA-approved
276 AI/ML devices (9.0%) contained a prospective study for post-market surveillance [Muralidharan et al.,
277 2024]; Dolin et al. [2025] use this gap to argue for statistically valid post-deployment monitoring.

278 Medicine, therefore, shows both sides of the argument: authorization and reviewability can exist
279 under partial mechanistic opacity, but lifecycle monitoring remains incomplete.

280 **Credit.** Credit law has comparatively strong instruments of reason-giving and correction. ECOA
281 and Regulation B require specific reasons for adverse action; a 2022 CFPB circular stated that
282 this requirement applies equally to credit decisions using complex and black-box credit algorithms,
283 although the circular was later withdrawn as guidance [Consumer Financial Protection Bureau, 2022,
284 2025]. The important point is the institutional pattern rather than the circular alone: consequential
285 credit decisions must be tied to reasons, records, and a responsible creditor. Credit therefore scores
286 high on domain coverage, accountability, and contestability, but weaker on AI-specific monitoring
287 and revocation. A denied applicant may receive reasons, but the public usually cannot see whether a
288 deployed credit model is drifting, how often appeals change outcomes, or what event withdraws the
289 system from use.

290 **Employment.** Employment has a thinner audit-and-notice pattern. New York City’s Local Law
291 144 prohibits the use of an automated employment decision tool unless it has undergone a recent
292 bias audit, a public summary is available, and required notices have been provided to candidates
293 or employees [New York City Department of Consumer and Worker Protection, 2023]. This is an
294 independent review, but not necessarily individual contestability. A bias audit can report group-level
295 disparity without giving a rejected applicant a meaningful route to challenge the tool’s role in the
296 decision or to trigger suspension when a deployment fails. Employment therefore motivates the
297 distinction between audit, notice, contestability, and revocation.

298 **Autonomous systems and critical infrastructure.** Autonomous systems show a different profile:
299 stronger incident reporting and monitoring, weaker individual contestability. NHTSA’s Standing
300 General Order requires covered manufacturers and operators to report specified crashes involving
301 Automated Driving Systems and Level 2 Advanced Driver Assistance Systems [National Highway
302 Traffic Safety Administration, 2025b]. NHTSA’s AV STEP proposal would add a voluntary review
303 and reporting framework for certain ADS-equipped vehicles [National Highway Traffic Safety Ad-
304 ministration, 2025a]. The EU AI Act similarly classifies specified AI systems in critical infrastructure,
305 road traffic, employment, credit, education, and law enforcement as high-risk, and requires providers
306 to establish post-market monitoring systems for high-risk AI systems [European Union, 2024]. These
307 instruments check deployed performance without requiring mechanism-level access. The thin rail is
308 contestability for affected non-operators, such as pedestrians, passengers, and communities exposed
309 to system-level risk.

310 **Criminal justice.** Criminal justice illustrates low Verification Coverage in a high-stakes setting. In
311 *State v. Loomis*, the Wisconsin Supreme Court allowed consideration of COMPAS with limitations
312 and cautions [Wisconsin Supreme Court, 2016]. Subsequent work found that COMPAS did not
313 outperform a simple two-feature model or crowd-aggregated lay judgments on the Broward County
314 recidivism-prediction task [Dressel and Farid, 2018], and related scholarship has criticized the secrecy
315 and procedural unfairness of COMPAS and similar proprietary recidivism risk tools [Rudin et al.,
316 2020]. Criminal justice may have a named institutional decision-maker, but it lacks strong external
317 verification, monitoring, and contestability when proprietary risk scores enter bail, sentencing, or
318 parole workflows. For Verification Coverage, this is the limiting case: higher stakes require stronger
319 verifiability, not simply stronger explanation demands.

320 **Cross-domain pattern.** The table shows why a single capability score cannot answer the deploy-
321 ment question. Medicine has authorization and reviewability but weak surveillance; credit has reasons
322 and correction pathways but weak lifecycle monitoring; employment has audit and notice, but weak
323 revocation; autonomous systems have incident reporting but weak appeal paths for affected non-
324 operators; criminal justice has named institutional authority, but weak verification and contestability.
325 The minimum-composition rule turns this comparison into a diagnosis: $VC_{\min}(d)$ is set by the
326 weakest rail, and the weakest rail identifies the next repair.

Regime	Domain	Verifier	Monitor.	Account.	Contest.	Revocation
Medicine (FDA CDS; AI/ML devices)	H	H	L	H	M	M
Credit (ECOA/Reg. B; adverse action)	H	M	L	H	H	M
Employment (NYC LL144; state AI laws)	H	M	M	M	M	L
Autonomous systems / critical infrastructure	M	M	H	M	L	M
Criminal justice (post- <i>Loomis</i>)	M	L	L	H	L	L

Table 2: Illustrative profile of selected high-stakes regimes on the six Verification Coverage components. Cells use the rule: **H** = component is instantiated by statute, regulation, binding agency process, or routine domain practice; **M** = partial, voluntary, indirect, or unevenly enforced; **L** = thin, practically inaccessible, or not specifically attached to the AI deployment. The table is a diagnostic application of the VC vocabulary, not an empirical measurement or legal compliance determination.

327 6 Alternative Views

328 We address seven objections, each either a partial concession, a misreading of the proposal, or an
329 independent reason to accept it.

330 **“General-purpose AI cannot be domain-gated.”** General-purpose training does not entail general-
331 purpose deployment authority. The jagged capability profile in Section 2.1 makes blanket authori-
332 zation too coarse in both directions: neither categorical prohibition nor categorical acceptance
333 tracks a frontier that reverses across task boundaries [Dell’Acqua et al., 2026]. Authority attaches to
334 the deployer of record in a specific context; open distribution multiplies deployment contexts, not
335 authorization.

336 **“Some domains have no verifier.”** That is an intended implication of the framework, not a failure
337 of it. A domain without a plausible verifier class registers a low Verifier Strength score; the minimum-
338 composition rule withholds authorization at that rail regardless of other scores. Decision support
339 under a named accountable human remains available; autonomous consequential use does not.

340 **“Open models make deployment gates unenforceable.”** Open distribution shifts, rather than
341 eliminates, the locus of duty. Verification obligations attach to whoever integrates a model into a
342 consequential workflow. Enforcement is harder to achieve, but the same artifact produces many
343 distinct deployment contexts, each with its own deployer of record and verification duties. The policy
344 response is to clarify integration-layer accountability, not to treat availability as authorization.

345 **“Interpretability is necessary for detecting deceptive internal objectives.”** We accept the premise
346 and contest the conclusion. Mechanistic inspection is plausibly the most direct signal for deceptive
347 internal objectives [Hubinger et al., 2024], and the regime weights it accordingly within the Mecha-
348 nistic Evidence stream. The concession ends there: a 53-point gap between detecting hazard-relevant
349 internal representations and correcting output behaviour [Basu et al., 2026] shows that understanding
350 does not imply control. Mechanistic evidence is necessary where it is the uniquely reliable signal; it
351 is not sufficient as a universal gate.

352 **“Internal auditing is enough” [Raji et al., 2020].** Auditing is necessary but incomplete. It does
353 not specify an external verifier class, an individual contestability path, or a predeclared revocation
354 trigger. The employment row of Table 2 is instructive: NYC Local Law 144 mandates a bias audit
355 yet leaves no meaningful appeal route for a rejected applicant and no revocation condition when
356 disparity persists post-audit [New York City Department of Consumer and Worker Protection, 2023].
357 Verification Coverage names what auditing leaves open.

358 **“This repackages NIST or the EU AI Act.”** Existing frameworks establish that AI assurance
359 is multi-dimensional [Tabassi, 2023]; the EU AI Act and FDA guidance already require pieces of
360 monitoring and documentation [European Parliament and Council of the European Union, 2024, U.S.
361 Food and Drug Administration, 2026b]. The contribution is different: naming deployment verifiability
362 as a missing *reportable* object beside capability scores, and imposing a minimum-composition rule

363 that makes the weakest rail binding. That gap is empirically visible: only nine percent of FDA-
364 registered AI/ML devices include a prospective surveillance plan [Dolin et al., 2025], a deficiency
365 that would directly block authorization under the Monitoring Maturity floor.

366 **“Verification Coverage will create false precision.”** Reporting the six-component profile alongside
367 any scalar prevents the aggregate from obscuring the weakest rail. The minimum-composition rule is
368 itself anti-precision-inflating: it cannot be gamed by strengthening other components. The remaining
369 safeguards are structural: independent assessment, preserved evaluator disagreement, and revocation
370 when a verifier fails. False precision is an argument for these safeguards, not against measurement.

371 7 Conclusion

372 Societies have long authorized opaque expertise through training, credentials, monitoring, liability,
373 appeal, and revocation, and mechanism-level access has often not been the primary basis for autho-
374 rization. Model capability is uneven enough across nearby tasks that no blanket policy can cohere
375 with that practice. Calibrated verification draws on four evidence streams—mechanistic, behavioural,
376 independent, and stakeholder—and yields six regime properties: domain-scoped, independently
377 checkable, monitored, accountable, contestable, and revocable. Verification Coverage reports those
378 properties beside capability scores under a minimum-composition rule that names the weakest rail
379 rather than averaging it out. The field should report not only what models can do, but also the
380 verification regime that makes a specific deployment governable.

381 References

- 382 Sanjay Basu, Sadiq Y. Patel, Parth Sheth, Bhairavi Muralidharan, Namrata Elamaram, Aakriti Kinra,
383 John Morgan, and Rajaie Batniji. Interpretability without actionability: Mechanistic methods
384 cannot correct language model errors despite near-perfect internal representations. *arXiv preprint*
385 *arXiv:2603.18353*, 2026.
- 386 Yanda Chen, Joe Benton, Ansh Radhakrishnan, Jonathan Uesato, Carson Denison, John Schulman,
387 Arushi Somani, Peter Hase, Misha Wagner, Fabien Roger, Vlad Mikulik, Samuel R. Bowman, Jan
388 Leike, Jared Kaplan, and Ethan Perez. Reasoning models don’t always say what they think. *arXiv*
389 *preprint arXiv:2505.05410*, 2025. Anthropic Alignment Science Team.
- 390 Consumer Financial Protection Bureau. Consumer financial protection circular 2022-03: Adverse
391 action notification requirements in connection with credit decisions based on complex algo-
392 rithms. CFPB Circular 2022-03, 87 Fed. Reg. 35864, [https://files.consumerfinance.gov/
393 f/documents/cfpb_2022-03_circular_2022-05.pdf](https://files.consumerfinance.gov/f/documents/cfpb_2022-03_circular_2022-05.pdf), 2022.
- 394 Consumer Financial Protection Bureau. Interpretive rules, policy statements,
395 and advisory opinions; withdrawal. 90 Federal Register 20084, 2025.
396 [https://www.federalregister.gov/documents/2025/05/12/2025-08286/
397 interpretive-rules-policy-statements-and-advisory-opinions-withdrawal](https://www.federalregister.gov/documents/2025/05/12/2025-08286/interpretive-rules-policy-statements-and-advisory-opinions-withdrawal).
- 398 Samantha Cruz Rivera, Xiaoxuan Liu, An-Wen Chan, Alastair K. Denniston, Melanie J. Calvert,
399 et al. Guidelines for clinical trial protocols for interventions involving artificial intelligence: The
400 SPIRIT-AI extension. *Nature Medicine*, 26:1351–1363, 2020. doi: 10.1038/s41591-020-1037-7.
- 401 Jeffrey De Fauw, Joseph R. Ledsam, Bernardino Romera-Paredes, Stanislav Nikolov, Nenad Tomasev,
402 et al. Clinically applicable deep learning for diagnosis and referral in retinal disease. *Nature*
403 *Medicine*, 24(9):1342–1350, 2018. doi: 10.1038/s41591-018-0107-6. URL [https://www.
404 nature.com/articles/s41591-018-0107-6](https://www.nature.com/articles/s41591-018-0107-6).
- 405 Fabrizio Dell’Acqua, Edward McFowland III, Ethan R. Mollick, Hila Lifshitz-Assaf, Katherine
406 Kellogg, Saran Rajendran, Lisa Kraye, François Candelon, and Karim R. Lakhani. Navigating
407 the jagged technological frontier: Field experimental evidence of the effects of AI on knowledge
408 worker productivity and quality. *Organization Science*, 37(2):403–423, 2026. doi: 10.1287/orsc.
409 2025.21838.

- 410 Pavel Dolin, Weizhi Li, Gautam Dasarathy, and Visar Berisha. Statistically valid post-deployment
411 monitoring should be standard for AI-based digital health. In *Advances in Neural Information
412 Processing Systems (NeurIPS 2025 Position Paper Track)*, 2025. Argues for statistically valid
413 post-deployment monitoring for AI-based digital health. OpenReview: [https://openreview.
414 net/forum?id=mXBFoHDu1l](https://openreview.net/forum?id=mXBFoHDu1l).
- 415 Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning.
416 *arXiv preprint arXiv:1702.08608*, 2017. URL <https://arxiv.org/abs/1702.08608>.
- 417 Julia Dressel and Hany Farid. The accuracy, fairness, and limits of predicting recidivism. *Science
418 Advances*, 4(1):eaa05580, 2018. doi: 10.1126/sciadv.aao5580.
- 419 Joshua Engels, David D. Baek, Subhash Kantamneni, and Max Tegmark. Scaling laws for scalable
420 oversight. In *Advances in Neural Information Processing Systems*, 2025. NeurIPS 2025 spotlight.
- 421 European Parliament and Council of the European Union. Regulation (EU) 2024/1689, annex III:
422 High-risk AI systems. <https://eur-lex.europa.eu/eli/reg/2024/1689/oj>, 2024.
- 423 European Union. Regulation (EU) 2024/1689 laying down harmonised rules on artificial intelligence.
424 Official Journal of the European Union, 2024. [https://eur-lex.europa.eu/eli/reg/2024/
425 1689/oj](https://eur-lex.europa.eu/eli/reg/2024/1689/oj).
- 426 Tairan Fu, Raquel Ferrando, Javier Conde, Carlos Arriaga, and Pedro Reviriego. Why do large
427 language models (LLMs) struggle to count letters? *arXiv preprint arXiv:2412.18626*, 2024.
- 428 Leilani H. Gilpin, David Bau, Ben Z. Yuan, Ayesha Bajwa, Michael Specter, and Lalana Kagal.
429 Explaining explanations: An overview of interpretability of machine learning. In *2018 IEEE 5th
430 International Conference on Data Science and Advanced Analytics (DSAA)*, pages 80–89. IEEE,
431 2018. doi: 10.1109/DSAA.2018.00018. URL <https://arxiv.org/abs/1806.00069>.
- 432 Evan Hubinger, Carson Denison, Jesse Mu, Mike Lambert, Meg Tong, Monte MacDiarmid, Tamera
433 Lanham, Daniel M. Ziegler, Tim Maxwell, Newton Cheng, et al. Sleeper agents: Training deceptive
434 LLMs that persist through safety training. *arXiv preprint arXiv:2401.05566*, 2024.
- 435 Eric Jonas and Konrad P. Kording. Could a neuroscientist understand a microprocessor? *PLOS
436 Computational Biology*, 13(1):e1005268, 2017. doi: 10.1371/journal.pcbi.1005268.
- 437 Hendrik Kempt, Jan-Christoph Heilinger, and Saskia K. Nagel. Relative explainability and double
438 standards in medical decision-making: Should medical AI be subjected to higher standards in
439 medical decision-making than doctors? *Ethics and Information Technology*, 24(2):20, 2022. doi:
440 10.1007/s10676-022-09646-x.
- 441 Jack Lindsey, Wes Gurnee, Emmanuel Ameisen, Brian Chen, Adam Pearce, Nicholas L. Turner,
442 Craig Citro, David Abrahams, Shan Carter, Basil Hosmer, Jonathan Marcus, Michael Sklar, Adly
443 Templeton, et al. On the biology of a large language model. *Transformer Circuits Thread*, 2025.
444 <https://transformer-circuits.pub/2025/attribution-graphs/biology.html>.
- 445 Xiaoxuan Liu, Samantha Cruz Rivera, David Moher, Melanie J. Calvert, Alastair K. Dennison,
446 et al. Reporting guidelines for clinical trial reports for interventions involving artificial
447 intelligence: The CONSORT-AI extension. *Nature Medicine*, 26:1364–1374, 2020. doi:
448 10.1038/s41591-020-1034-x.
- 449 Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson,
450 Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model cards for model reporting. In
451 *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 220–229.
452 ACM, 2019. doi: 10.1145/3287560.3287596.
- 453 Vijaytha Muralidharan, Boluwatife Adeleye Adewale, Caroline J. Huang, Mfon Thelma Nta, Pe-
454 ter Oluwaduyilemi Ademiju, Pirunthan Pathmarajah, Man Kien Hang, Oluwafolajimi Adesanya,
455 et al. A scoping review of reporting gaps in FDA-approved AI medical devices. *npj Digital
456 Medicine*, 7(1):273, 2024. doi: 10.1038/s41746-024-01270-x.

457 National Highway Traffic Safety Administration. ADS-equipped vehicle safety,
458 transparency, and evaluation program (AV STEP). 90 Federal Register 4130,
459 [https://www.federalregister.gov/documents/2025/01/15/2024-30854/](https://www.federalregister.gov/documents/2025/01/15/2024-30854/ads-equipped-vehicle-safety-transparency-and-evaluation-program)
460 [ads-equipped-vehicle-safety-transparency-and-evaluation-program](https://www.federalregister.gov/documents/2025/01/15/2024-30854/ads-equipped-vehicle-safety-transparency-and-evaluation-program), 2025a.

461 National Highway Traffic Safety Administration. Third amended standing general order 2021-01:
462 Incident reporting for automated driving systems and level 2 advanced driver assistance systems.
463 NHTSA Standing General Order on Crash Reporting, 2025b. Effective June 16, 2025. <https://www.nhtsa.gov/laws-regulations/standing-general-order-crash-reporting>.
464 <https://www.nhtsa.gov/laws-regulations/standing-general-order-crash-reporting>.

465 Joe Needham, Giles Edkins, Govind Pimpale, Henning Bartsch, and Marius Hobbhahn. Large
466 language models often know when they are being evaluated. *arXiv preprint arXiv:2505.23836*,
467 2025. Gemini 2.5 Pro reaches AUC 0.83 on evaluation-awareness classification; human baseline
468 0.92.

469 New York City Department of Consumer and Worker Protection. Auto-
470 mated employment decision tools. [https://www.nyc.gov/site/dca/about/](https://www.nyc.gov/site/dca/about/automated-employment-decision-tools)
471 [automated-employment-decision-tools](https://www.nyc.gov/site/dca/about/automated-employment-decision-tools).page, 2023.

472 OpenAI. Sycophancy in GPT-4o: What happened and what we’re doing about it. OpenAI blog,
473 29 April 2025, 2025. <https://openai.com/index/sycophancy-in-gpt-4o/>. See also the
474 follow-up “Expanding on what we missed with sycophancy” ([https://openai.com/index/](https://openai.com/index/expanding-on-sycophancy/)
475 [expanding-on-sycophancy/](https://openai.com/index/expanding-on-sycophancy/)).

476 Inioluwa Deborah Raji, Andrew Smart, Rebecca N. White, Margaret Mitchell, Timnit Gebru, Ben
477 Hutchinson, Jamila Smith-Loud, Daniel Theron, and Parker Barnes. Closing the AI accountability
478 gap: Defining an end-to-end framework for internal algorithmic auditing. In *Proceedings of the*
479 *2020 Conference on Fairness, Accountability, and Transparency*, pages 33–44. ACM, 2020. doi:
480 10.1145/3351095.3372873.

481 Cynthia Rudin, Caroline Wang, and Beau Coker. The age of secrecy and unfairness in recidivism
482 prediction. *Harvard Data Science Review*, 2(1), 2020. doi: 10.1162/99608f92.6ed64b30.

483 Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askill, Samuel R. Bowman,
484 Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R. Johnston, Shauna Kravec, Timothy
485 Maxwell, Sam McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda
486 Zhang, and Ethan Perez. Towards understanding sycophancy in language models. *arXiv preprint*
487 *arXiv:2310.13548*, 2023.

488 Elham Tabassi. AI risk management framework (AI RMF 1.0). Technical Report NIST AI 100-1,
489 National Institute of Standards and Technology, 2023.

490 U.S. Food and Drug Administration. Artificial intelligence-enabled medical devices.
491 [https://www.fda.gov/medical-devices/software-medical-device-samd/](https://www.fda.gov/medical-devices/software-medical-device-samd/artificial-intelligence-enabled-medical-devices)
492 [artificial-intelligence-enabled-medical-devices](https://www.fda.gov/medical-devices/software-medical-device-samd/artificial-intelligence-enabled-medical-devices), 2026a. Accessed May 7,
493 2026.

494 U.S. Food and Drug Administration. Clinical decision support software: Guidance for industry
495 and Food and Drug Administration staff. [https://www.fda.gov/regulatory-information/](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/clinical-decision-support-software)
496 [search-fda-guidance-documents/clinical-decision-support-software](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/clinical-decision-support-software), 2026b.

497 Teun van der Weij, Felix Hofstätter, Ollie Jaffe, Samuel F. Brown, and Francis Rhys Ward. AI
498 sandbagging: Language models can strategically underperform on evaluations. *arXiv preprint*
499 *arXiv:2406.07358*, 2024.

500 Baptiste Vasey, Myura Nagendran, Bruce Campbell, David A. Clifton, Gary S. Collins, Spiros
501 Denaxas, Alastair K. Denniston, et al. Reporting guideline for the early-stage clinical evaluation
502 of decision support systems driven by artificial intelligence: DECIDE-AI. *Nature Medicine*, 28:
503 924–933, 2022. doi: 10.1038/s41591-022-01772-9.

504 Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening
505 the black box: Automated decisions and the GDPR. *Harvard Journal of Law & Technology*, 31(2):
506 841–887, 2018.

- 507 Jiaxin Wen, Ruiqi Zhong, Akbir Khan, Ethan Perez, Jacob Steinhardt, Minlie Huang, Samuel R.
508 Bowman, He He, and Shi Feng. Language models learn to mislead humans via RLHF. *arXiv*
509 *preprint arXiv:2409.12822*, 2024.
- 510 Wisconsin Supreme Court. State v. Loomis, 2016 WI 68, 371 Wis. 2d 235, 881 n.w.2d 749 (Wis.
511 2016). [https://www.wicourts.gov/sc/opinion/DisplayDocument.pdf?content=pdf&](https://www.wicourts.gov/sc/opinion/DisplayDocument.pdf?content=pdf&seqNo=171690)
512 [seqNo=171690](https://www.wicourts.gov/sc/opinion/DisplayDocument.pdf?content=pdf&seqNo=171690), 2016.
- 513 Zhengxuan Wu, Aryaman Arora, Atticus Geiger, Zheng Wang, Jing Huang, Dan Jurafsky, Christo-
514 pher D. Manning, and Christopher Potts. AxBench: Steering LLMs? even simple baselines
515 outperform sparse autoencoders. *arXiv preprint arXiv:2501.17148*, 2025. ICML 2025 spotlight.
- 516 Nan Xu and Xuezhe Ma. LLM the genius paradox: A linguistic and math expert’s struggle with
517 simple word-based counting problems. In *Proceedings of the 2025 Conference of the Nations*
518 *of the Americas Chapter of the Association for Computational Linguistics: Human Language*
519 *Technologies (Volume 1: Long Papers)*, pages 3344–3370, 2025. doi: 10.18653/v1/2025.naacl-long.
520 172.
- 521 John Zerilli, Alistair Knott, James Maclaurin, and Colin Gavaghan. Transparency in algorithmic and
522 human decision-making: Is there a double standard? *Philosophy & Technology*, 32(4):661–683,
523 2019. doi: 10.1007/s13347-018-0330-6.