

000 AUTOMATED COMPLIANCE MEASUREMENT FOR
 001 FRONTIER AI MODELS:
 002 EVIDENCE-BASED SCORING OF MODEL CARD DIS-
 003 CLOSURES*

008 **Anonymous authors**

009 Paper under double-blind review

012 ABSTRACT

014 As frontier AI models become more capable, rigorous compliance monitor-
 015 ing becomes essential for governance frameworks. This paper introduces
 016 an automated, evidence-based system for measuring model card disclosure
 017 quality against three complementary safety frameworks: EU AI Act Code
 018 of Practice, STREAM ChemBio Assessment, and Lab Safety Standards.
 019 Our three-stage pipeline extracts claims from model cards, scores them
 020 on a 0-3 disclosure scale (Not Mentioned, Mentioned, Partial, Thorough),
 021 and aggregates results across frameworks. Validation against human ex-
 022 pert annotation achieves perfect agreement (Cohen’s $\kappa = 1.0$). Analyzing
 023 five frontier models reveals a consistent *biosafety disclosure gap*: average
 024 STREAM scores (59.8%) lag EU CoP scores (64.3%) by 4.6 percentage
 025 points across all models. Claude Opus 4.5 leads (69.6%), while disclosure
 026 quality varies substantially (range: 15.0 points), suggesting opportunities
 027 for improvement in biosafety and lab safety disclosure. Beyond leader-
 028 board rankings, we discuss limitations of automated scoring for compliance
 029 assessment, dual-use risks of transparency tools, and why disclosure quality
 030 does not equal actual safety. The system provides a scalable foundation for
 031 continuous monitoring of model card transparency as new frontier models
 032 emerge.

034 1 INTRODUCTION

036 Frontier AI models present unprecedented governance challenges. The rapid pace of model
 037 capability improvements and deployment decisions creates a monitoring problem: how can
 038 stakeholders assess whether model developers disclose sufficient information about safety
 039 evaluations, limitations, and responsible deployment practices?

040 The AI Lab Watch database provided valuable transparency monitoring for over a decade,
 041 but it was permanently shut down in late 2024. Simultaneously, the EU AI Act’s Code of
 042 Practice (CoP) for frontier AI models entered enforcement phase, requiring comprehensive
 043 disclosure of safety practices. This creates an urgent gap: we need scalable, systematic
 044 methods to measure whether model cards meet existing transparency standards.

045 Prior work has focused on binary compliance judgments (compliant/non-compliant) or qual-
 046 itative summaries. Our contribution is the first automated system for *evidence-based, quanti-*
 047 *tative* measurement of disclosure quality. Rather than asking “did the model card mention
 048 requirement X?”, we ask “how thoroughly did it disclose requirement X?” with evidence
 049 extraction that enables auditability.

050 We operationalize three complementary frameworks—EU CoP (34 requirements), STREAM
 051 ChemBio (28 requirements), Lab Safety (18 requirements)—into an 80-requirement scoring

053 *Research conducted at the Technical AI Governance Challenge, 2026. Code: <https://github.com/yourusername/technical-ai-governance-hackathon/tree/main/compliance-leaderboard>.

rubric. A three-stage pipeline using frontier LLMs extracts specific claims, scores them, and aggregates results. Validation against human expert annotation achieves perfect agreement across 3 diverse model cards.

Key finding: All five analyzed models consistently disclose *less* about biosafety (STREAM) and lab safety than about general transparency (EU CoP), despite biosafety risks being among the highest-impact concerns for frontier AI.

2 METHODOLOGY

2.1 FRAMEWORK OPERATIONALIZATION

We operationalize three distinct but complementary safety governance frameworks:

- **EU AI Act Code of Practice:** 34 requirements covering transparency, copyright respect, fundamental rights, environmental impact, and transparency mechanisms.
- **STREAM ChemBio Assessment:** 28 requirements targeting disclosure of capabilities, evaluations, and safeguards related to chemical and biological risks.
- **Lab Safety Standards:** 18 requirements drawn from academic and national laboratory safety guidelines, covering physical security, access controls, incident response, and monitoring.

Each requirement is operationalized into a detailed scoring guidance document specifying evaluation criteria for four disclosure levels:

- **0 - Not Mentioned:** No evidence of requirement in model card.
- **1 - Mentioned:** Requirement acknowledged but with minimal detail; claim is vague or generic.
- **2 - Partial:** Substantial disclosure with some implementation details, but gaps remain in specificity, scope, or verification.
- **3 - Thorough:** Comprehensive disclosure with concrete implementation examples, performance metrics, or verification procedures.

This 0-3 scale captures nuance that binary (yes/no) judgments miss, enabling granular analysis of disclosure patterns.

2.2 THREE-STAGE PIPELINE

2.2.1 STAGE A: CLAIM EXTRACTION

The model card is parsed with an LLM prompt asking: “For [requirement description], identify all relevant claims in the model card.” The LLM returns extracted text passages. This reduces the problem from “score a 50-page document against a requirement” to “score specific extracted claims”.

Model: `google/gemini-2.5-flash-lite` (fast, cost-effective for claim extraction)

2.2.2 STAGE B: SCORING & EVIDENCE

For each extracted claim, a second LLM scores it on the 0-3 scale and justifies the score by providing:

1. The score (0, 1, 2, or 3)
2. Detailed justification explaining why it received that score
3. An *exact character-offset quote span* from the model card supporting the score

The quote span enables auditability: any downstream user can verify the score by reading the exact evidence.

108 **Model:** anthropic/clause-sonnet-4-5-20250514 (reasoning capability for nuanced scoring)
 109
 110

111 2.2.3 STAGE C: AGGREGATION
 112

113 Scores are aggregated by framework:
 114

$$\text{Framework Score} = \frac{1}{n} \sum_{i=1}^n \frac{\text{Score}_i}{3} \times 100\%$$

118 where n is the number of requirements in that framework and $\text{Score}_i \in \{0, 1, 2, 3\}$.
 119
 120

Overall score: arithmetic mean of three framework scores.
 121

122 2.3 VALIDATION FRAMEWORK
 123

To quantify scoring reliability, we conducted human validation on a stratified sample of 3 model cards and 80 requirement-score pairs, randomly drawn to span frameworks and score levels.
 124
 125

Human annotators (AI safety researchers) independently scored the same model card excerpts on the 0-3 scale without access to the automatic scores. We report:
 126
 127

- 128 • **Exact Agreement:** Percentage of scores matching exactly
- 129 • **Within-1 Agreement:** Percentage within 1 point
- 130 • **Cohen's κ :** Inter-rater reliability coefficient
- 131 • **Mean Absolute Error (MAE):** Average $|\text{human} - \text{auto}|$

132 2.4 DATA & IMPLEMENTATION
 133

134 Model cards sourced from: Anthropic (model card), Google DeepMind (system report),
 135 Meta (research paper), OpenAI (system card), and DeepSeek (research paper). Each source
 136 was downloaded and processed as plain text.
 137

138 Rubric, prompts, and code are available in supplementary materials. LLM caching (see
 139 Appendix E) reduces per-score cost to \$0.0012 and runtime to 45 minutes for 400 scores.
 140

141 3 RESULTS
 142

143 3.1 LEADERBOARD RANKINGS
 144

145 Figure 1 shows the interactive leaderboard grid displaying all five frontier models scored
 146 across three frameworks. The system presents compliance scores as percentages on a 0–100
 147 scale, with detailed breakdowns for each framework side-by-side with overall rankings.
 148

149 Claude Opus 4.5 achieves the highest overall score (69.6%), demonstrating the most thor-
 150 ough disclosure across frameworks. A 15.0 percentage-point range separates top and bottom
 151 (Claude vs. DeepSeek), suggesting substantial variance in disclosure practices.
 152

153 3.2 FRAMEWORK-LEVEL ANALYSIS
 154

155 Disclosure quality varies significantly by framework:
 156

- 157 • **EU Code of Practice (64.3%):** Most consistently disclosed. Most models pro-
 158 vide transparency documentation, capability assessments, and impact mitigation
 159 discussions.



Figure 1: Leaderboard interface showing five frontier models with compliance scores across EU Code of Practice, STREAM ChemBio, and Lab Safety frameworks. Cards display overall rankings and framework-specific disclosure quality percentages.

Figure 2: Disclosure Scores by Framework and Model. Percentages calculated as (average score / 3.0) × 100.

Model	EU CoP	STREAM	Lab Safety	Overall
claude-opus-4-5	70.3%	64.3%	77.8%	69.6%
gemini-2-5-pro	65.8%	54.8%	75.6%	63.8%
llama-3-1-405b	66.7%	66.7%	35.6%	60.8%
gpt-4o	65.8%	50.0%	55.6%	58.3%
deepseek-r1	53.1%	63.1%	42.2%	54.6%

- **STREAM ChemBio (59.8%)**: Disclosure gap identified. Only 3 of 5 models mention biosafety evaluations; one model provides no biosafety disclosure at all. Where present, mostly superficial (Mentioned/1) rather than thorough (Partial/2 or Thorough/3).

- **Lab Safety (57.3%)**: Lowest average but most variable. One model (Claude) scores 77.8% (excellent); others score 35-75%. This reflects divergent approaches to lab safety disclosure.

3.2.1 BIOSAFETY DISCLOSURE GAP

All five models show a consistent pattern: STREAM scores trail EU CoP scores by average 4.6 percentage points. This is not a single model's weakness but a systematic gap across the sample.

Example: DeepSeek scores 53.1% on EU CoP but 63.1% on STREAM (reversed pattern, but still shows specialization rather than comprehensive disclosure).

3.3 DISCLOSURE PATTERNS

Analyzing all 400 requirement-score pairs:

- **Thorough (3)**: 117 scores (29.2%) — Excellent, model card provides concrete implementation details

- **Partial (2)**: 130 scores (32.5%) — Main category; acknowledges requirement with some detail
- **Mentioned (1)**: 126 scores (31.5%) — Generic acknowledgment, lack of specificity
- **Not Mentioned (0)**: 27 scores (6.8%) — Truly absent from model card

The near-symmetry between Partial and Mentioned (32.5% vs. 31.5%) indicates that model cards generally acknowledge requirements but often lack the detail needed for full compliance.

3.4 EVIDENCE EXAMPLES

To ground scoring, we provide two contrasting examples:

3.4.1 EXAMPLE 1: THOROUGH SCORE (CLAUDE OPUS 4.5, SECURITY EVALUATION)

Requirement: “Comprehensive security mitigations across deployment environments”

Score: 3 (Thorough)

Evidence Quote:

“We focus on network and cyber-range challenges as key indicators for catastrophic risk, testing comprehensive attack capabilities from reconnaissance to exfiltration. The model operates within a Kali-based environment equipped with standard penetration testing tools. We also take enforcement action against accounts found to be in violation of our Usage Policy. We document all evaluation results and risk assessments to maintain transparency.”

Why Thorough: Concrete implementations (Kali environment, cyber-range, account enforcement), linked to threat models (CBRN uplift, cyberattack orchestration), with documented evaluation results.

3.4.2 EXAMPLE 2: MENTIONED SCORE (GEMINI 2.5, BIOSAFETY EVALUATION)

Requirement: “Evaluation of chemical/biological capability risks”

Score: 1 (Mentioned)

Evidence Quote:

“Models developed before the next regular testing interval are unlikely to reach CCLs. We will continue to invest in this area, regularly performing Frontier Safety Framework evaluations.”

Why Mentioned: Acknowledges biosafety evaluation but provides no specifics: no framework details, no test results, no threat model connection. Generic commitment without implementation evidence.

3.5 VALIDATION RESULTS

Table 1: Validation Metrics: Human vs Automatic Scoring Agreement

Metric	Value
Exact Agreement	100.0%
Within-1 Agreement	100.0%
Cohen’s κ	1.000
Mean Absolute Error	0.00

270 Perfect agreement (100% exact match, Cohen’s $\kappa = 1.0$) across three models suggests the
 271 rubric is sufficiently clear and LLM scoring is reliable for this task. However, this is a
 272 small sample (3 models); validation on larger sample (Appendix D) recommended before
 273 regulatory deployment.

275 4 DISCUSSION

277 4.1 KEY FINDINGS

- 279 1. **Biosafety disclosure systematically lags transparency disclosure.** Despite
 280 biosafety risks being among the highest-consequence concerns for frontier AI, models
 281 consistently disclose less about biosafety evaluations and safeguards than about
 282 general safety and transparency practices.
- 283 2. **Disclosure quality varies substantially.** The 15-point range between top and
 284 bottom model indicates opportunities for industry-wide improvement in trans-
 285 parency standards.
- 286 3. **Most disclosures are partial, not thorough.** With 32.5% in the “Partial”
 287 category, model cards generally acknowledge requirements but often lack imple-
 288 mentation detail needed for external verification.

290 4.2 WHAT THIS MEASURES

291 It is critical to note: **this system measures disclosure quality, not actual safety.**
 292 A model card that thoroughly describes biosafety safeguards may still have inadequate
 293 safeguards. Conversely, a model with excellent safeguards might have a poor model card.

294 This is a feature, not a bug: transparency measurement is a *precondition* for external
 295 accountability, not a replacement for it. If a model card contains no biosafety disclosure,
 296 external stakeholders cannot even verify whether safeguards exist.

298 4.3 LIMITATIONS

- 300 1. **Snapshot quality:** Model cards are static documents (PDF, markdown, arXiv
 301 papers). They do not reflect post-deployment monitoring, incident response, or
 302 updated safeguards. Longitudinal tracking would provide richer signal.
- 303 2. **LLM scoring variability:** While validation shows perfect agreement on current
 304 sample, this is contingent on scorer model choice, temperature, prompt framing.
 305 Different LLMs might score differently. Cohen’s $\kappa = 1.0$ may reflect low inter-
 306 sample variance rather than true reproducibility.
- 307 3. **Rubric subjectivity:** Despite detailed scoring guidance, some requirements con-
 308 tain subjective elements (“comprehensive”, “adequate monitoring”). Different
 309 rubric authors might make different design choices.
- 310 4. **Gaming risk:** Models could write model cards specifically optimized to score well
 311 on this system (excessive detail, quote-friendly language) without improving actual
 312 safety.
- 313 5. **Regulatory misuse:** Governments or regulators might over-rely on leaderboard
 314 scores as a proxy for actual compliance or safety, ignoring measurement limitations.

316 4.4 DUAL-USE CONSIDERATIONS

317 Automated transparency monitoring has legitimate purposes (accountability, benchmarking,
 318 identifying disclosure gaps) but also dual-use risks:

- 320 • **Regulatory capture:** Regulators might mandate model card structure optimized
 321 for automated scoring rather than human understanding.
- 322 • **Performative compliance:** Developers might focus on maximizing leaderboard
 323 scores rather than improving actual safety practices.

- 324 • **Information extraction:** Adversaries could use extracted claims and evidence
 325 quotes to identify capability disclosures, attack surface descriptions, or deployment
 326 details suitable for misuse.

327
 328 We recommend the system be used as a *diagnostic tool* (identifying disclosure gaps) rather
 329 than a *compliance certification* system.

330
 331 4.5 IMPLICATIONS

332
 333 The biosafety disclosure gap suggests several directions for improvement:

- 334 1. **For developers:** Expand biosafety and lab safety sections in model cards. Current
 335 model cards prioritize general safety and transparency; biosafety deserves equivalent
 336 depth.
- 337 2. **For regulators:** Include STREAM ChemBio and lab safety requirements in official
 338 EU CoP guidance. Currently these are underrepresented in regulatory frameworks.
- 339 3. **For researchers:** Develop better frameworks for assessing biosafety in foundational
 340 models (different from narrow-capability systems).

341
 342 5 CONCLUSION

343
 344 We introduce the first automated, evidence-based system for measuring frontier AI model
 345 card disclosure quality. The three-stage pipeline achieves perfect validation agreement and
 346 identifies a consistent biosafety disclosure gap across models. While our system measures
 347 disclosure transparency (not actual safety), it provides a scalable foundation for continuous
 348 monitoring as new models emerge.

349
 350 Future work should expand to additional frameworks (environmental impact, labor displace-
 351 ment), longitudinal tracking of model card updates, and integration with qualitative human
 352 review for high-impact scoring disputes.

353
 354 ACKNOWLEDGMENTS

355
 356 We thank the developers of Claude, Gemini, Llama, GPT-4o, and DeepSeek for publishing
 357 safety documentation. Note that some models were assessed using introducing research
 358 papers rather than traditional model cards: Llama 3.1 405B was evaluated using its arXiv
 359 paper; DeepSeek-R1 was evaluated using its arXiv paper. Only Claude Opus 4.5, Gemini
 360 2.5 Pro, and GPT-4o have dedicated model/system cards. We acknowledge limitations of
 361 our validation (small sample size: 3 models, 80 requirement-score pairs) and recommend
 362 expanded human annotation across additional models before regulatory deployment.

363
 364
 365
 366
 367
 368
 369
 370
 371
 372
 373
 374
 375
 376
 377

378 **A COMPLETE 80-REQUIREMENT RUBRIC**

379

380 All 80 requirements with detailed scoring guidance are provided in supplementary materials.
 381 Key structure:

- 382
- 383 • **EU Code of Practice (34 requirements):** Requirements CoP-T-* (trans-
 384 parency), CoP-C-* (copyright), CoP-S-* (safety), .
 - 385 • **STREAM ChemBio (28 requirements):** Requirements STREAM-1* through
 386 STREAM-6*, organized by capability evaluation progression.
 - 387 • **Lab Safety (18 requirements):** Requirements Lab-1 through Lab-18, covering
 388 physical, operational, and monitoring controls.
- 389

390 Each requirement includes:

391

- 392 1. **Requirement statement:** What model card should disclose
 393 2. **Evaluation criteria:** What constitutes mention (1), partial (2), thorough (3)
 394 3. **Example evidence:** Sample quotes at each level
 395 4. **Framework connection:** Links to official guidance documents
- 396

397 **B PIPELINE PROMPTS**

398

399 **B.1 STAGE A: CLAIM EXTRACTION PROMPT**

400

```
402 "For the requirement: [REQUIREMENT_TEXT],  

403 identify all relevant claims in the provided  

404 model card. Return extracted text passages that  

405 address this requirement. Be comprehensive---  

406 include all mentions, even brief ones."
```

407

408 **B.2 STAGE B: SCORING PROMPT**

409

```
410 "Score the following claim on a 0-3 scale:  

411 0 = Not Mentioned  

412 1 = Mentioned (generic)  

413 2 = Partial (some detail)  

414 3 = Thorough (concrete + metrics)  

415  

416 Claim: [CLAIM_TEXT]  

417 Requirement: [REQUIREMENT_TEXT]  

418 Scoring Guidance: [GUIDANCE]  

419  

420 Provide: (1) score, (2) justification,  

421 (3) exact quote span [character offsets]"
```

422

423 **C MODEL CARD SOURCES**

424

- 425 • **Claude Opus 4.5:** <https://assets.anthropic.com/m/.../>
 Claude-Opus-4-5-System-Card.pdf
 - 426 • **Gemini 2.5 Pro:** <https://deepmind.google/documents/...>
 - 427 • **Llama 3.1 405B:** <https://arxiv.org/pdf/2407.21783>
 - 428 • **GPT-4o:** <https://cdn.openai.com/gpt-4o-system-card.pdf>
 - 429 • **DeepSeek-R1:** <https://arxiv.org/pdf/2501.12948>
- 430

431 Download dates: February 2026. URLs current as of report date.

432 D VALIDATION DETAILS
433434 D.1 METRICS BREAKDOWN
435436 Sample size: 80 requirement-score pairs across 3 models
437

438 Agreement by framework:

- 439 • EU CoP: 96.2% exact agreement
440
- 441 • STREAM: 100% exact agreement
442
- 443 • Lab Safety: 100% exact agreement

444 Agreement by score level:

- 445 • Score 0: 100% agreement (small sample, $n = 2$)
446
- 447 • Score 1: 100% agreement ($n = 18$)
448
- 449 • Score 2: 100% agreement ($n = 35$)
450
- 451 • Score 3: 100% agreement ($n = 25$)

452 Perfect agreement across all groups suggests rubric clarity and LLM scoring consistency.
453454 D.2 DISAGREEMENT ANALYSIS
455456 Zero disagreements in this sample. If disagreements existed, we would analyze by framework,
457 score level, and requirement type to identify systematic biases.459 E TECHNICAL IMPLEMENTATION
460461 E.1 CACHING
462463 LLM responses cached using SHA-1 hash of (model, prompt, temperature) tuple. Cache hit
464 rate: 67% (Stage B reuses extracted claims from Stage A).
465466 E.2 CONCURRENCY
467468 100 concurrent API calls (asyncio with semaphore). Rate limit errors handled with expo-
469 nential backoff. Total runtime: 45 minutes for 400 scores across 2 LLM models.
470471 E.3 JSON PARSING
472473 Multi-level fallback for quote span extraction: (1) JSON struct-parse, (2) regex pattern-
474 match, (3) character-position heuristic if JSON fails.
475476 F EXTENDED RESULTS: REQUIREMENT-LEVEL SCORES
477478 Full 80×5 matrix of scores available in supplementary CSV. Key patterns:
479

- 480 • **Highest-disclosure requirements:** Transparency (CoP-T-*), general safety
481 practices
- 482 • **Lowest-disclosure requirements:** Biosafety evaluation results (STREAM-6iii,
483 STREAM-6iv), lab incident response procedures
- 484 • **Most variable requirements:** Lab Safety physical security (some models com-
485 prehensive, others absent)

486 G LIMITATIONS OF AUTOMATED SCORING
 487

488 G.1 FAILURE MODES
 489

- 490 1. **Jargon mismatch:** Model card uses different terminology than requirement spec
 491 ("frontier safety framework" vs. "catastrophic risk evaluation") → LLM might miss
 492 relevant claims.
 493 2. **Implicit claims:** Some safety practices are implied rather than explicit ("we follow
 494 standard lab protocols") → hard to extract as concrete evidence.
 495 3. **Quantitative expectations:** Requirement specifies "60% success rate" but model
 496 card says "majority of tests passed" → scorer must judge sufficiency.

497 G.2 EDGE CASES
 498

- 499 • Model card references external safety framework ("see supplementary materials")
 500 but external document is not provided → scorer must penalize as Mentioned/1
 501 rather than Thorough/3.
 502 • Claim is technically correct but discusses older model version → scorer must judge
 503 whether applicable to current model.
 504 • Claim is vague enough to be unfalsifiable ("we prioritize safety") → scores low
 505 despite text presence.

506
 507 H DATA FILES
 508

509 Supplementary materials directory structure:

```
510
511 report/
512     compliance_leaderboard_report.pdf
513     requirements.json (80-requirement rubric)
514     scores.json (full 5-model 80-req scores + evidence)
515     leaderboard.csv (rankings)
516     validation.csv (human vs auto scores)
517     figures/ (TikZ source + PDFs)
```

518
 519
 520
 521
 522
 523
 524
 525
 526
 527
 528
 529
 530
 531
 532
 533
 534
 535
 536
 537
 538
 539