

# Bridging the In-Context AI Evaluation Gap: Risk Chain Analysis for Coordinated AI Governance

Anonymous ACL submission

## Abstract

Capability evaluations cover only a narrow portion of the pathway from AI capability to real-world harm. We introduce *risk chain analysis*—tracing harm trajectories step by step and mapping evaluation evidence at each stage—and apply it to three risk domains: cyber attacks, biological risk, and agentic AI failures. Evaluation consistently concentrates where benchmarking is easiest, not where risk is highest. The uncovered steps require different methods and different actors: sectoral regulators, deploying organisations, and domain-specific agencies that currently have no role in AI evaluation. We argue that AI safety institutes are uniquely positioned to coordinate this distributed evaluation ecosystem.

## 1 Introduction

Capability evaluations remain the dominant paradigm for AI risk assessment, with 85% of evaluations focused at the capability level (Weidinger et al., 2023). Yet they lack construct validity (Schaeffer et al., 2023; Bean et al., 2025; Raji et al., 2021) and ecological validity (Raji et al., 2022; Weidinger et al., 2023), producing systematic over- and underestimation of risk (Hubinger et al., 2024; Passi and Vorvoreanu, 2022).

Despite broad acknowledgement of these limitations, alternative methods have not been widely adopted. We identify three reasons: a *knowledge problem* (alternative methods are unfamiliar), an *operationalisation problem* (they are costly and resist automation), and a *coordination problem* (the required evidence is distributed across the governance ecosystem). This paper addresses the third directly.

We introduce *risk chain analysis*: decomposing specific harm trajectories into sequential steps, mapping evaluation coverage at each, and identifying which actors should fill gaps. We apply this

to three structurally different cases—a malicious-actor scenario (cyber), a well-studied domain (bio), and a deployment failure (agentic AI)—and find a consistent pattern: evaluation clusters at early steps while the later steps where harm materialises are systematically uncovered.

## 2 Three Cases

### 2.1 Cyber: AI-Enabled Ransomware on a Hospital

**Scenario.** A criminal group uses AI to conduct a ransomware attack on an NHS hospital trust running legacy systems, grounded in real precedent (Rodriguez et al., 2025) and the WannaCry (2017) and Synnovis (2024) attacks.

We trace six steps: (1) attacker intent and AI access, (2) reconnaissance, (3) vulnerability identification and exploitation, (4) lateral movement, (5) ransomware deployment, and (6) patient harm.

Nearly all evaluation effort concentrates at Step 3: CTF benchmarks (Zhang et al., 2024; Yang et al., 2024), vulnerability databases, and exploit generation tasks. But the ecological validity gap is widest here. The ARTEMIS study (Lin et al., 2025)—the first comparison of AI agents and human pentesters on live infrastructure—found markedly different capability profiles from CTF results. Steps 1 and 4–6 receive essentially no coverage. Critically, outcomes at these steps depend on the *target’s* defences, not the model’s capabilities. Evidence here requires NHS Digital, the NCSC, and health regulators—organisations with no current role in AI evaluation.

### 2.2 Biological Risk

Biorisk has received more evaluation attention than any other domain, with multiple uplift studies (Mouton et al., 2024; OpenAI, 2024; Anthropic, 2024). Even here, evaluation concentrates at Step 2 (information access). As Peppin et al. (2024) note,

079	there are no published empirical studies connecting	deployment-context assessment, and incident anal-	128
080	information access via LLMs to downstream steps:	ysis (McGregor, 2021).	129
081	material procurement, wet lab execution, and dis-	<b>Different steps require different actors.</b> AI	130
082	semination. These require biosecurity regulators,	labs and safety institutes cover Steps 1–3. Every-	131
083	DNA synthesis screening bodies, and public health	thing beyond falls to organisations with no current	132
084	agencies (Carter et al., 2023).	role in AI evaluation.	133
085	<b>2.3 Agentic AI: Coding Agent Introduces</b>	<b>4 Who Should Do What</b>	134
086	<b>Vulnerability</b>		
087	A coding agent generates a SQL injection vulner-	The risk chain framework reveals that a safety case	135
088	ability; the developer approves under time pres-	built on capability evaluations alone covers at most	136
089	sure; the vulnerability reaches production (Vera-	one or two steps of a five-to-six step causal chain.	137
090	code, 2025; Apiiro, 2025). Partial coverage ex-	Addressing this requires distributing evaluation	138
091	ists at Step 2 (code generation benchmarks). The	responsibility—but distribution without coordina-	139
092	most consequential gap is Step 3: human over-	tion risks fragmentation. We argue that AI safety	140
093	sight. Every agentic AI governance framework	institutes (AISIs) are uniquely positioned to serve	141
094	identifies oversight as a core safeguard (Infocomm	as the coordinating node.	142
095	Media Development Authority, 2025; Robey et al.,	<b>4.1 AI Developers</b>	143
096	2025; Chan et al., 2024), yet nobody has studied	Developers hold unique evidence on real-world	144
097	whether developers actually provide effective over-	usage patterns (Tamkin et al., 2024; Zhao et al.,	145
098	sight of AI-generated code. The automation bias	2024; Zheng et al., 2024). They should invest	146
099	literature predicts they often will not (Parasuraman	in ecologically valid evaluations: not just func-	147
100	and Manzey, 2010; Buçinca et al., 2021), but this	tional benchmarks but operational simulations and	148
101	remains untested for coding agents specifically. Ev-	security testing in realistic contexts. Sharing ag-	149
102	idence must come from research institutions and	gregated, privacy-preserving usage data—as An-	150
103	deploying organisations.	thropic’s Clio project has begun to do—would sig-	151
104	<b>3 Cross-Cutting Findings</b>	nificantly strengthen the evidence base for Steps 1–	152
105	<b>Evaluation concentrates where benchmarking</b>	2 across risk chains.	153
106	<b>is easy, not where risk is highest.</b> Cyber clusters	<b>4.2 Sectoral Regulators and Deployers</b>	154
107	at vulnerability identification, bio at information	Sectoral regulators and deploying organisations	155
108	access, agentic AI at code generation. The steps	own the later steps of every chain. NHS Digital un-	156
109	where harm materialises receive little or no atten-	derstands hospital network architecture; financial	157
110	tion.	regulators understand trading system risk; software	158
111	<b>The best-covered steps have the worst ecolog-</b>	engineering organisations understand code review	159
112	<b>ical validity.</b> CTF benchmarks, knowledge tests,	practices. These actors have domain knowledge for	160
113	and isolated code generation tasks strip away con-	deployment-context assessment but typically lack	161
114	textual complexity. ARTEMIS (Lin et al., 2025),	AI evaluation methodology. Capacity building—	162
115	Epoch AI’s analysis of biorisk evaluation ambigu-	through training, shared frameworks, and collab-	163
116	ity, and Kilov et al.’s (2025) work on moral com-	orative evaluation design—is the bridge between	164
117	petence evaluation all demonstrate that benchmark	model-level evidence and deployment-level risk.	165
118	performance is unreliable when deployment con-	<b>4.3 Safety Institutes as Coordination Nodes</b>	166
119	text is removed.	AISIs have three features that make them natural	167
120	<b>Gaps cluster around human behaviour and</b>	coordinators of the distributed evaluation ecosys-	168
121	<b>deployment context.</b> Attacker intent, developer	tem.	169
122	oversight, institutional resilience—these are prop-	First, <b>methodological comparative advantage.</b>	170
123	erties of sociotechnical systems, not models. They	AISIs concentrate evaluation expertise that no sin-	171
124	require usage monitoring (Tamkin et al., 2024),	gle sectoral regulator can replicate. Their most	172
125	human interaction experiments (Dell’Acqua et al.,	valuable contribution may not be conducting every	173
126	2023; Noy and Zhang, 2023), field testing (Na-	evaluation themselves but developing the analyti-	174
127	tional Institute of Standards and Technology, 2025),	cal frameworks—such as risk chain analysis—that	175

176 make evaluation gaps visible and tractable across  
177 domains. This aligns with the emerging shift in  
178 AISI mandates from direct model testing toward  
179 supporting evaluation capacity across government  
180 (Department for Science, Innovation and Technol-  
181 ogy, 2024).

182 Second, **cross-domain visibility**. The risk  
183 chains above show structurally similar gaps across  
184 very different domains: human oversight is unstud-  
185 ied for both coding agents and biosecurity con-  
186 texts; ecological validity is poor for both cyber and  
187 biorisk benchmarks. A central coordinating body  
188 can identify these cross-cutting patterns and pre-  
189 vent duplication of effort. No sectoral regulator has  
190 line of sight across domains in this way.

191 Third, **standard-setting authority**. AISIs can  
192 set standards requiring safety cases to specify  
193 which parts of a risk pathway their evidence ac-  
194 tually covers—and which parts remain uneval-  
195 uated. This reframes safety cases from “we ran  
196 these benchmarks” to “here is where our evidence  
197 is strong, here is where it is absent, and here is  
198 who should provide it.” The risk chain framework  
199 provides the structure for such requirements.

200 This model has precedent. Financial regulators  
201 maintain systemic risk frameworks without audit-  
202 ing every firm: they require institutions to report  
203 against standardised frameworks and flag when crit-  
204 ical risk areas lack evidence (Bank for International  
205 Settlements, 2019). Pandemic preparedness sim-  
206 ilarly required coordination between virologists,  
207 epidemiologists, and public health bodies—no sin-  
208 gle institution assessed the full risk pathway alone  
209 (World Health Organization, 2021). AISIs could  
210 play an analogous role: maintaining risk chain  
211 frameworks, requiring developers to produce evi-  
212 dence for the steps they control, requiring deploy-  
213 ers to assess deployment-context steps, and making  
214 visible when critical steps have no evidence at all.

215 As AI proliferates beyond frontier catastrophic  
216 risk into healthcare, education, finance, and critical  
217 infrastructure, the question AISIs face—*what is*  
218 *our job versus other departments’ jobs?*—is pre-  
219 cisely what risk chain analysis helps answer. The  
220 AISI’s job is to maintain the framework, build eval-  
221 uation capacity in other departments, and ensure  
222 accountability across the full chain. Individual de-  
223 partments’ job is to provide deployment-context  
224 evidence for their domain.

## 5 Limitations 225

226 Risk chains are linear simplifications of pathways  
227 that branch, loop, and intersect. Constructing them  
228 requires domain expertise, and different experts  
229 would draw different chains. The framework iden-  
230 tifies gaps but does not prioritise them. Emergent  
231 or systemic risks—correlated failures, cascading  
232 infrastructure effects—do not follow sequential  
233 chains and require different tools (Zwetsloot and  
234 Dafoe, 2019; Stein et al., 2024).

## 6 Conclusion 235

236 Current AI evaluation covers a narrow and ecolog-  
237 ically questionable portion of the pathway from  
238 capability to harm. Risk chain analysis makes this  
239 visible and maps which actors should provide evi-  
240 dence at each step. The central implication is that  
241 adequate AI risk assessment is a multi-actor coor-  
242 dination problem, and AI safety institutes are best  
243 positioned to coordinate it.

## References 244

- 245 Anthropic. 2024. Responsible scaling policy, version  
246 1.0.
- 247 Apiiro. 2025. AI-generated code risk in fortune 50  
248 enterprises.
- 249 Bank for International Settlements. 2019. Stress testing  
250 banks: A comparative analysis.
- 251 Andrew M Bean and 1 others. 2025. Benchmark infla-  
252 tion: Revealing LLM performance gaps using retro-  
253 holdouts. *arXiv preprint arXiv:2501.14892*.
- 254 Zana Buçinca, Maja Barbara Malaya, and Krzysztof Z  
255 Gajos. 2021. To trust or to think: Cognitive forc-  
256 ing functions can reduce overreliance on AI in AI-  
257 assisted decision-making. *Proceedings of the ACM*  
258 *on Human-Computer Interaction*, 5(CSCW1).
- 259 Sarah R Carter and 1 others. 2023. The convergence  
260 of artificial intelligence and the life sciences. *NTI*  
261 *Biosecurity*.
- 262 Alan Chan, Matthew J Salganik, and 1 others.  
263 2024. Visibility into AI agents. *arXiv preprint*  
264 *arXiv:2401.13138*.
- 265 Fabrizio Dell’Acqua, Edward McFowland III, Ethan R  
266 Mollick, and 1 others. 2023. Navigating the jagged  
267 technological frontier: Field experimental evidence  
268 of the effects of AI on knowledge worker productiv-  
269 ity and quality. *Harvard Business School Working*  
270 *Paper*.
- 271 Department for Science, Innovation and Technology.  
272 2024. AI safety institute approach to evaluations.

273	Evan Hubinger, Carson Denison, Jesse Mu, Mike Lambert, Meg Tong, and 1 others. 2024. Sleeper agents: Training deceptive LLMs that persist through safety training. <i>arXiv preprint arXiv:2401.05566</i> .	Merlin Stein and 1 others. 2024. Interconnected monitoring for AI safety. <i>arXiv preprint</i> .	325
274			326
275		Alex Tamkin and 1 others. 2024. Clio: Privacy-preserving insights into real-world AI use. <i>arXiv preprint arXiv:2412.13678</i> .	327
276			328
277	Infocomm Media Development Authority. 2025. Companion guide for agentic AI systems.	Veracode. 2025. State of AI code security 2025.	329
278			330
279	Daniel Kilov and 1 others. 2025. Moral competence evaluation of AI systems: Beyond clean scenarios. <i>arXiv preprint</i> .	Laura Weidinger, Maribeth Rauh, Nahema Marchal, Arianna Manzini, Lisa Anne Hendricks, Juan Mateos-Garcia, Stevie Bergman, Jackie Kay, Conor Griffin, and 1 others. 2023. Sociotechnical safety evaluation of generative AI systems. <i>arXiv preprint arXiv:2310.11986</i> .	331
280			332
281			333
282	Hao Lin and 1 others. 2025. ARTEMIS: AI agents vs human pentesters on live infrastructure. <i>arXiv preprint</i> .		334
283			335
284			336
285	Sean McGregor. 2021. Preventing repeated real world AI failures by cataloging incidents: The AI incident database. <i>Proceedings of AAI</i> .	World Health Organization. 2021. Strengthening preparedness for health emergencies: Implementation of the International Health Regulations.	337
286			338
287			339
288	Christopher Mouton, Caleb Lucas, and Ella Guest. 2024. Can large language models provide useful guidance on social engineering? <i>RAND Corporation</i> .	John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. 2024. InterCode: Standardizing and benchmarking interactive coding with execution feedback. In <i>NeurIPS</i> .	340
289			341
290			342
291	National Institute of Standards and Technology. 2025. NIST AI 700-2: ARIA field testing protocol.	Andy K Zhang, Neil Perry, Riya Duber, and 1 others. 2024. CYBENCH: A framework for evaluating cybersecurity capabilities and risks of language models. In <i>arXiv preprint arXiv:2408.08926</i> .	343
292			344
293	Shakked Noy and Whitney Zhang. 2023. Experimental evidence on the productivity effects of generative artificial intelligence. <i>Science</i> , 381(6654):187–192.		345
294			346
295			347
296	OpenAI. 2024. GPT-4o system card: Biological risk assessment.	Wenting Zhao and 1 others. 2024. WildChat: 1M ChatGPT interaction logs in the wild. <i>arXiv preprint arXiv:2405.01470</i> .	348
297			349
298	Raja Parasuraman and Dietrich H Manzey. 2010. Complacency and bias in human use of automation: An attentional integration. <i>Human Factors</i> , 52(3):381–410.	Lianmin Zheng and 1 others. 2024. LMSYS-Chat-1M: A large-scale real-world LLM conversation dataset. <i>arXiv preprint arXiv:2309.11998</i> .	350
299			351
300			352
301			353
302	Samir Passi and Mihaela Vorvoreanu. 2022. Overreliance on AI: Literature review. <i>Microsoft Research</i> .	Remco Zwetsloot and Allan Dafoe. 2019. Thinking about risks from AI: Accidents, misuse and structure. <i>Lawfare Blog</i> .	354
303			355
304	Aidan Peppin, Edward Sherburn, Alexander Sherburn, and Jide Kim. 2024. Evaluating biological risks from advanced AI. <i>arXiv preprint arXiv:2412.01946</i> .		356
305			
306			
307	Inioluwa Deborah Raji, Emily M Bender, Amandalynne Paullada, Emily Denton, and Alex Hanna. 2021. AI and the everything in the whole wide world benchmark. <i>Proceedings of NeurIPS Datasets and Benchmarks Track</i> .		
308			
309			
310			
311			
312	Inioluwa Deborah Raji, I Elizabeth Kumar, Aaron Horowitz, and Andrew Selbst. 2022. The fallacy of AI functionality. <i>Proceedings of FAccT</i> .		
313			
314			
315	Alexander Robey and 1 others. 2025. Jailbreaking and safeguarding large language models: A survey. <i>arXiv preprint</i> .		
316			
317			
318	Manuel Rodriguez and 1 others. 2025. Systematization of knowledge: The AI-cybersecurity threat landscape. <i>arXiv preprint</i> .		
319			
320			
321	Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. 2023. Are emergent abilities of large language models a mirage? <i>Advances in Neural Information Processing Systems</i> , 36.		
322			
323			
324			