Appendix A: Grading Sheets

Each of our panellists were presented with the full contents of this appendix to inform their grading decisions. The grading sheets are broken down by domain, and panellists were asked to provide grades for each company per domain. Within each domain is a set of indicators: a collection of facts about the companies.

You can skip between the domains by selecting one from the list below:

Domains

- - 6 indicators
- A Current Harms
 - 8 indicators
- Safety Frameworks
 - 4 indicators
- Existential Safety
 - 4 indicators
- Governance & Accountability 5 indicators
- ເລື Information Sharing 6 indicators



Domain

This domain evaluates the rigor and comprehensiveness of companies' risk identification and assessment processes for their current flagship models. The focus is on implemented assessments, not stated commitments.

Table of Contents

Internal Testing

Dangerous Capability Evaluations
Elicitation for Dangerous Capability Evaluations
Human Uplift Trials

External Testing

Independent Review of Safety Evaluations
Pre-deployment External Safety Testing
Bug Bounties for Model Vulnerabilities

Grading

Internal Testing

Indicator

Dangerous Capability Evaluations

Definition & Scope

This indicator assesses whether organizations conduct systematic evaluations of dangerous capabilities before deploying frontier models. Priority domains include biological and chemical weapons, offensive cyber operations, AI R&D facilitation, and behaviors associated with goal misalignment or deception. Evidence is drawn from model cards, including published results and detailed testing methodologies. The focus is on external deployments, as there is insufficient transparency on internal deployments.

Evaluation Guidance

Transparency Classification:

- Low detail: Only states that evaluations were conducted without naming specific tests or explaining methodologies.
- Moderate detail: Brief explanations of specific evaluations (and methodology).
- High details: Extensive explanations of individual evaluations and methodology.

Notes on AI regulation in China:

Under the *Interim Measures for the Management of Generative AI Services* (Aug 2023), every public-facing (B2B exemption) GenAI service must pass a government-supervised security assessment and complete an algorithm filing (算法备案); regulators may run their tests during this process.

Providers rely on the draft national standard "Basic Security Requirements for Generative AI Service" (TC260)



to prepare. The draft contains a list of 29 security risks in five buckets, yet, while an earlier draft version from Feb 2024 mentioned frontier risks, the current version does not [China Talk, 2024; China Law Translate].

Why This Matters

Systematic evaluations for high-risk capabilities reflect institutional responsibility for managing low-probability, high-impact harms. In contrast to more routine risks—where market forces often suffice—frontier threats require deliberate foresight. Firms that fail to test for these dangers risk contributing to unmanaged systemic vulnerability.



	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Model	Claude 4 Opus	DeepSeek R1	Gemini 2.5 Pro	Llama 4 Maverick	03	Grok 3	GLM-4
Bio + Chem	Yes ('CBRN') - High level of detail - Quantitative results with human & Al baselines - Safety framework classification - 10+ evaluations reported Evaluations include: Bioweapons acquisition uplift trial, Expert red-teaming (Deloitte), Long-form virology tasks, Multimodal virology (VCT), Bioweapons knowledge questions, DNA Synthesis Screening Evasion, LAB-Bench subset, Creative biology, Shorthorizon computational biology, ASL-4 expert red-teaming System Card: pages 88-103	None	Yes ('CBRN') - Moderate level of detail - Quantitative results with AI (& human) baselines - Safety framework classification (CBRN uplift) Multiple-choice benchmarks, open-ended qualitative assessments led by domain experts across the biological, radiological, and nuclear domains. Three public benchmarks reported: SecureBio VMQA, FutureHouse LAB-Bench, Weapons of Mass Destruction Proxy. Model Card: pages 9-11	Yes ('CBRNE') - but no details provided - Minimal details - No quantitative results Reports expert-designed targeted evaluations and red-teaming without giving details [Meta]	Yes ('Bio') - Moderate level of detail - Quantitative results with human & Al baselines - Safety framework classification (Bio&Chem). Evaluations: Long-form biorisk questions, Multimodal troubleshooting virology, ProtocolQA Open-Ended Tacit knowledge, and troubleshooting Model Card: pages 12-15	None	None
Cyber offense	Yes ('Cybersecurity') - High level of detail - Quantitative results - Tracked in RSP, but no formal threshold Evaluations: Web Exploitation (15 CTFs), Cryptography (22 CTFs), Exploitation (9 CTFs), Reverse Engineering (8 CTFs), Network (4 CTFs), Cyber-harness network (3 ranges), Cybench (39 challenges) System Card: pages 116-122	None	Yes ('Cybersecurity') - High level of detail - Quantitative results with AI baselines - Safety framework classification (Cyber uplift + Cyber Autonomy) - Open-sourced evaluation suite 1) Previously published evaluation suite including In-house CTF (13), Hack The Box (13), Vulnerability detection (3) [arXiv, 2024]. 2) 50 additional challenges across four categories following their newly published framework: Reconnaissance, Tool development, Tool usage, Operational security [arXiv, 2025]. Model Card: pages 11-13	Yes ('Cyber attack enablement') - Minimal details - No quantitative results Card reports "evaluating the capabilities of Llama 4 to automate cyberattacks, identify and exploit security vulnerabilities, and automate harmful workflows". Does not give more details. [Meta]	Yes ('Cybersecurity') - Moderate level of detail - Quantitative results with AI baselines - Safety framework classification (Cybersecurity). Model card (p.15) 1) Two scenarios from the "Cyber Range" evaluation for conducting fully end-to-end cyber operations in a realistic, emulated network. 2) 100 capture the flag challenges across three difficulty levels (high school, collegiate, professional). Categories: Web Application Exploitation, Reverse Engineering, Binary and Network Exploitation, Cryptography, misc. Model Card: pages 15-22	None	None



	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
AI R&D	Yes ('Autonomy') - High level of detail - Quantitative results with expert baselines - Safety framework classification (Autonomy) Evaluations: SWE-bench Verified (hard subset), METR data deduplication, Internal Al research suite 1 (6 tasks), Internal Al research suite 2, Internal model use survey System Card: pages 101-113	None	Yes ('ML R&D') - High level of detail - Quantitative results with human and Al baselines - Safety framework classification (ML R&D uplift) Used open-source Research Engineering Benchmark (RE-Bench, Wijk et al., 2024) https://arxiv.org/abs/2411.15114 Model Card: pages 14-16	None	Yes ('Al Self-improvement') - Moderate to High level of detail - Quantitative results with Al baselines - Safety framework classification (Al Self-improvement) - Additional external evaluation from METR on autonomy and Al R&D [METR] Five evaluations: OpenAl Research Engineer Interviews, SWE-bench Verified, OpenAl PRs, SWE-Lancer, PaperBench. Model Card: pages 22-28	None	None
Scheming capabilities & Misalignment	Yes ('Comprehensive Alignment Assessment') - High level of detail - First comprehensive alignment assessment in the industry (30 pages) - Additional external testing by Apollo Research Comprehensive Alignment Assessment covers: systematic deception, hidden goals, self-preservation, sandbagging, situational awareness, sycophancy, reasoning faithfulness, jailbreak susceptibility, values, high-agency behavior, subtle sabotage, Excessive compliance with harmful system prompts, pro-Al bias, and Reward seeking Apollo Research's evaluations are similar to a recent technical paper: "in-context scheming reasoning" [arXiv, 2024]. System Card: pages 19-48	None	Yes (Focuses on capabilities but not propensities) ('Deceptive alignment') - high level of detail - safety framework classification ("Instrumental Reasoning") - Evaluations open-sourced Five evaluations on ability to circumvent oversight (stealth), ("cover_your_tracks, decision_sabotage, strategic_rule_breaking, oversight_pattern, classifier_evasion") and seven evaluations for measuring a model's ability to instrumentally reason about itself and its environment (situational awareness). Details published in a technical paper [arXiv, 2025]. Model Card: pages 15-16	None	Yes ('Deception / Scheming'), only external evaluations by Apollo Research - High level of detail - quantitative results with human and Al baselines Evaluations: strategic deception, in-context scheming, reasoning, and sabotage. Evaluations similar to recent technical paper: "in-context scheming reasoning" [arXiv, 2024]. Model Card: pages 10+30	None	None



	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Transparency Overview	Model Card Length: 122 pages (Opus + Sonnet) Safety Evaluations: - 10.5 pages (p 11-21) Frontier Risk Evaluations: - 36 pages (p. 87-122) External Evaluations: - 2 pages (p. 30-31, 122) Other: 1) Comprehensive Alignment Assessment: 29 pages (p. 22-51) [Anthropic, 2025] 2) Al Safety Level 3 Deployment Safeguards Report 25 pages Content: Claude 4 Opus was classified as requiring Al Safety Level 3 (ASL-3) under their Responsible Scaling Policy, indicating it could potentially assist with CBRN weapons development. The relevant safeguards report (separate from the model card) outlines the core threat model, details the implemented safeguards, and provides evidence demonstrating their effectiveness [Anthropic, 2025].	Technical report length: 22 pages No content on safety evaluations [arXiv, 2025]	Model Card Length: 17 pages Safety Evaluations: - 2 pages (p. 5-7) Frontier Risk Evaluations: - 8 pages (p. 8-16) External Evaluations: - 0.5 pages (p. 10) Linked Resources: Additional results in technical paper: 'Evaluating Frontier Models for Stealth and Situational Awareness' (45 pages) [arXiv, 2025] Other: Announces: "detailed technical report will be published once per model family's release, with the next technical report releasing after the 2.5 series is made generally available." (Google considers the current release to be a "Preview") [Google, 2025].	Model Card Length: 14.5 pages (browser print format of website) Safety Evaluations: - 2 pages (p. 10-12) Frontier risk evaluations: - 1 page (p. 13-14) [Huggingface]	Model Card Length: 31.5 pages (o3 +o4 mini) Safety Evaluations: - 7 pages (p. 2-8) Frontier Risk Evaluations: - 16 pages (p. 11-27) External Evaluations: - 5 pages (p. 8-11, 30-32) Other: OpenAl's Safety Evaluations Hub webpage provides an ongoing overview of safety test results regarding harmful content, jailbreaks, hallucinations, and instruction hierarchy compliance. It currently shares updated evaluation results across 9 different Al models. [OpenAl, 2025; OpenAl, 2025]	No relevant model card found. The announcement post does not report safety evaluations. [xAI, 2025]	Technical report length: 12.5 pages Safety Evaluations: - 1 page (p. 12) [arXiv, 2024]

• Forum, Frontier Model. "Issue Brief: Preliminary Taxonomy of Pre-Deployment Frontier Al Safety Evaluations." Frontier Model Forum, 14 Jan. 2025

Elicitation for Dangerous Capability Evaluations

Definition & Scope

Assesses the extent to which a company discloses its elicitation strategy in its most recent dangerous-capability evaluations.* We record whether the company is transparent about:

- 1. Parallelization settings e.g., *pass@ *n* and *best-of-n* sampling parameters (especially relevant for AI-R&D and cyber-security tasks).
- 2. Tooling any use of internet access, code interpreters, agentic scaffolds, or relevant tools that can amplify model performance.
- 3. Model variants the exact model checkpoints tested, including "helpful-only" variants and any domain- or task-specific fine-tuning.

Why This Matters

It has been demonstrated that small improvements in elicitation methodology can dramatically increase scores on evaluation benchmarks. Naive elicitation strategies cause significant underreporting of risk profiles, potentially missing dangerous capabilities that sophisticated actors could unlock. Companies thus need to implement comprehensive elicitation methodologies to better approximate an AI model's true capabilities, not just its default behavior. This should include task-specific fine-tuning in domains like bio-risk, especially if model weights will be made generally available, but also be general, as model weights might be stolen or leaked. A structured, transparent, and well-resourced approach to capability elicitation demonstrates a genuine commitment to risk discovery.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Parallel test-time compute & tooling	Mentions specific tools on tools and parallel computing approaches for several cyber evaluations. For cyber CTFs, pass@30 is reported. Bio-section: - "for automated evaluations, our models have access to various tools and agentic harnesses (software setups that provide them with extra tools to complete tasks)" - Some evaluations comment on the parallel test time compute approach, e.g., pass@5 for longform virology	No tests reported	None mentioned	None mentioned	Mentions specific tools on tools and parallel computing approaches for several cyber and self-improvement evaluations. For cyber CTFs, pass@12 is reported, for self-improvement, often pass@1. Multiple choice bio-risk questions were reported as consensus@32.	No tests reported	No tests reported
Model versions & Domain / Task-specific fine-tuning	- Tested helpful-only model without safety mitigations. - No mention of domain/task-specific fine-tuning. System Card: page 8	No tests reported	None mentioned	None mentioned	- Tested helpful-only model without safety mitigations. - No domain/task- specific fine-tuning reported System Card: page 13	No tests reported	No tests reported

Sources

- Adler, Steven. "AI Companies Should Be Safety-testing the Most Capable Versions of Their Models." Steven Adler's Substack, 26 Mar. 2025
- Metr. "Guidelines for Capability Elicitation." METR's Autonomy Evaluation Resources, 13 Mar. 2024Indicator

Human Uplift Trials

Definition & Scope

This indicator assesses whether organizations conduct rigorous, controlled human-subject studies to evaluate the marginal risk AI systems pose in dangerous domains by "uplifting" people's ability to cause harm. Key evidence includes experimental designs that compare task performance with and without AI support, the inclusion of domain-relevant experts, realistic and consequential task scenarios, and transparent publication of methods and findings. To assess worst-case potential, models should be tested without embedded safety filters.

Why This Matters

Empirical uplift studies are critical for grounding AI safety policy in observable outcomes. These studies assess whether advanced systems significantly enhance a user's ability to cause harm and inform the development of proportionate safety interventions. Entities that conduct and publish such studies exhibit leadership in transparent, evidence-based risk governance.

Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Yes (1)	None	None	None	None	None	None
Bioweapons Acquisition Uplift Trial:						
- Methodology:						
Controlled trial with groups of 8-10 participants given up to 2 days to draft a comprehensive bioweapons acquisition plan						
- Groups:						
Control group: Only basic internet resources						
Model-assisted group: Additional access to Claude with safeguards removed						
- Participants: Contracted from SepalAl and Mercor						
- Grading: By Deloitte using a detailed rubric, assessing key steps of the acquisition pathway						
System Card: page 92-93						

External Testing

Indicator

Independent Review of Safety Evaluations

Definition & Scope

Assesses whether an AI developer *commissions independent third-party experts to (A) verify the factual accuracy and process integrity of its internal dangerous-capability evaluations and (B) assess the* evaluation quality *and the company's interpretation of the results. We collect information on the reviewers' identity and credentials, their independence (including any conflicts of interest), the scope of the review, depth of access to data and logs (including rights to replicate or extend tests), and whether their findings are published unredacted.

Why This Matters

Al developers control both the design and disclosure of dangerous capability evaluations, creating inherent incentives to underreport alarming results or select lenient testing conditions that avoid costly deployment delays. Regulators, investors, and the public face a critical information asymmetry: they must trust safety claims based on self-reported evaluations with minimal methodological transparency. Independent external scrutiny



can address this trust deficit by verifying reported results, assessing whether evaluations are sufficiently rigorous to uncover real risks, and providing credible third-party perspectives on whether safety claims are justified. This need is especially acute for catastrophic risk domains such as biosecurity, where companies may cite "infohazard" concerns to limit transparency.

Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
None	None	None	None	None	None	None

Sources

- "Key Components of an RSP." METR, 26 Sept. 2023
- · Homewood, Aidan, et al. "Third-party compliance reviews for frontier AI safety frameworks."

Indicator

Pre-deployment External Safety Testing

Definition & Scope

This indicator evaluates whether companies facilitate independent third-party safety assessments prior to releasing frontier models. It excludes collaborative testing arrangements and focuses solely on unaffiliated evaluators. Evidence includes the identity and qualifications of external parties, the level and duration of access provided, compensation arrangements, testing permissions, and the evaluators' ability to publish independently. The strength of these practices is judged by the comprehensiveness of the evaluations, the depth of access, and the autonomy of the evaluators.

Why This Matters

Independent evaluations are essential for verifying safety claims and uncovering risks that internal teams may miss, perhaps due to misaligned incentives or bias. Providing evaluators with substantial access—and ensuring their ability to publish freely—reflects a company's commitment to transparent, evidence-based governance.



Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
UK AI Security	None reported	Model card reports leveraging "third-party external testers" for CBRN risks, but does not disclose any details.	None reported	U.S. Al Safety Institute Scope: Cyber and biological capabilities evaluations Model versions: Early checkpoints + final launch candidate models of both o3 and o4-mini Details: No details U.K. Al Security Institute Scope: Cyber, chemical, and biological autonomy capabilities, and an early version of safeguards Model versions: Early checkpoints + final launch candidate models Details: No details METR (Model Evaluations and Threat Research) Duration: 15 days Model versions: Earlier checkpoints of o4-mini and o3 Scope: Autonomous capabilities and reward hacking. Details: METR published a paper on methods [arXiv, 2025] and an evaluation report here [METR]. Apollo Research Scope: In-context scheming and strategic deception Model versions: Early checkpoints + final launch candidate models of both o3 and o4-mini Methodology: Similar to Apollo's technical paper on of evaluation. Meinke, A., et al. (2024). Frontier models are capable of in-context scheming [arXiv, 2025]. Details: <1 page summary of results + 3 pages with quantitative results in appendix. Pattern Labs Scope: Cybersecurity - three types of cyber offensive challenges: 1) Evasion, 2) Network Attack Simulation, 3) Vulnerability Discovery and Exploitation. Model versions: Early checkpoints + final launch candidate models of both o3 and o4-mini Details: <1 page summary of results System Card: page 9 Additional details from the index survey [Response]: (Questions reference this specific release) Q17: "In some instances, we paid private consultants for their work, but payment is not conditioned on the content of their findings." "Helpful-only" or base model API (no harmlessness fine-tuning and no filters)" Q18: Highest level of access granted to any external evaluator: "Helpful-only" or base model API (no harmlessness fine-tuning and no filters)" Q19: Longest access period for any external evaluator: Between 2-3 weeks Q20: Publication arrangements: 1) Evaluators are under NDA, and publications require prior approval from OpenAI.	None reported	None reported

Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
				Q21: Query-rate or volume restrictions: Elevated but capped – evaluators had higher quotas than the public/enterprise tier but were still subject to explicit caps (e.g., requests-per-minute or daily token limits). OpenAI notes that rates can sometimes depend on technical feasibility. Q22: Logging and retaining model interactions: "Zero Data Retention available upon request, if technically feasible during pre-deployment periods (for some new models or products, ZDR is not always possible during pre-deployment testing)."		

- Che, Zora, et al. "Model Manipulation Attacks Enable More Rigorous Evaluations of LLM Capabilities." Neurips Safe Generative AI Workshop 2024.

Indicator

Bug Bounties for Model Vulnerabilities

Definition & Scope

This indicator assesses the presence and design of structured incentive programs—such as bug bounties or red-teaming initiatives—that encourage responsible disclosure of safety vulnerabilities in AI models. It focuses exclusively on programs addressing model behavior, excluding conventional cybersecurity initiatives due to insufficient public reporting. Evidence includes the scope of eligible issues, compensation levels, response timelines, and public availability of program documentation.

Why This Matters

Structured disclosure programs with financial incentives harness external expertise to identify model vulnerabilities before they are exploited in deployment. Investments in such programs indicate a proactive attitude toward risk identification.

Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Bug bounty on universal jailbreaks - Opened applications for early access testing of new safety mitigations. - Started May 2025 (last iteration ran August 2024) [Anthropic, 2024] - Up to \$25,000 for verified universal jailbreak attacks that could expose vulnerabilities in critical, high-risk domains - Still accepting applications [Anthropic, 2025]	None	Abuse Vulnerability Reward Program: Accepts certain abuse-related discoveries: - Prompt Attacks - Training Data Extraction - Manipulating Models - Adversarial Perturbation - Model Theft (excludes jailbreaks) [Google]	Bounty programs are restricted to privacy or security issues, like extracting training data through tactics like model inversion or extraction attacks. [Meta]	Early access for safety testing (December 2024) One-off programs allowed safety researchers to apply for early access to frontier models to help surface novel risks. No payments announced. [OpenAI, 2024]	None	None

٦	$\Gamma \cap$	RF	COMPL	FTFD	BY PANFI	LISTS

Grading

Please pick a grade for each firm. You can add brief justifications to your grades.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Grades	Grade	Grade	Grade	Grade	Grade	Grade	Grade
Grade comments (Justifications, opportunities for improvements, etc.)							

Grading Scales

Grading scales are provided to support consistency between reviewers.

- Assessment methods provide very high confidence that dangerous capabilities would be detected in time to prevent significant harm
- B Assessment practices provide high confidence in detecting dangerous capabilities
- C Assessment approach provides moderate confidence with concerning gaps
- Limited assessments provide low confidence in risk detection
- F No credible basis for detecting dangerous capabilities

Domain comments

Optional: Share observations that apply across companies, including general recommendations, notes on how you weighted indicators, or feedback on FLI's methodology.

|--|



Domain



Companies were evaluated on how effectively their models mitigate current harms, with a focus on safety benchmark performance, robustness against adversarial attacks, watermarking of Al-generated content, and the treatment of user data.

Table of Contents

Model Safety / Trustworthiness

Stanford's HELM Safety Benchmark Stanford's HELM AIR Benchmark TrustLLM Benchmark

Robustness

Gray Swan Arena: UK AISI Agent Red-Teaming Challenge Cisco Security Risk Evaluation Protecting Safeguards from Fine-tuning

Digital Responsibility

Watermarking
User Privacy

Grading

Model Safety / Trustworthiness

Indicator

Stanford's HELM Safety Benchmark

Definition & Scope

This indicator measures model performance on Stanford's HELM Safety v1.0 benchmark, a suite of five safety tests covering six risk categories: violence, fraud, discrimination, sexual content, harassment, and deception. The benchmark includes: HarmBench (jailbreak resistance); BBQ (social discrimination); SimpleSafetyTest; XSTest (alignment between helpfulness and harmlessness); and AnthropicRedTeam (resilience to adversarial probing). Performance is reported as normalized aggregate scores ranging from 0 to 1, where higher scores indicate fewer safety risks. Scoring is based on exact match accuracy for BBQ and model-judge ratings (GPT-40 and Llama 3.1 405B) for the remaining benchmarks.

Why This Matters

HELM Safety offers a standardized evaluation of core safety risks. HELM enables consistent comparison across models. Strong performance signals effective safeguards against harmful outputs, including violence, fraud, discrimination, and other high-impact risks that affect millions of users.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.AI	Zhipu Al
Model	Claude 4 Opus	DeepSeek R1	Gemini 2.5 Pro (03-25 preview)	Llama 4 Maverick	03	Grok 3 Beta	n/a (Model not evaluated
Average score (max score = 1)	0.97	0.87	0.91	0.91	0.98	0.86	by external benchmark)
HarmBench	0.92	0.47	0.65	0.66	0.98	0.45	
SimpleSafetyTests	1.00	0.98	0.97	0.99	0.99	0.97	
BBQ accuracy	0.97	0.97	0.96	0.93	0.98	0.94	
Anthropic Red Team	0.99	0.96	1.00	0.98	0.98	0.96	
XSTest	0.97	0.95	0.99	0.97	0.97	0.96	
Retrieved	12 June 2025						
Release	Release v1.8.0						

Sources

• Farzaan, et al. "HELM Safety: Towards Standardized Safety Evaluations of Language Models." Stanford Center for Research on Foundation Models, 8 Nov. 2024. Accessed 14 Jun. 2025.

Indicator

Stanford's HELM AIR Benchmark

Definition & Scope

This indicator evaluates model performance on Stanford's AIR-Bench 2024 (AI Risk Benchmark), an AI safety benchmark aligned with emerging government regulations and company policies. We report mean scores across 5,694 tests spanning 314 granular risk categories, with scores measuring the percentage of appropriately refused requests. The benchmark systematically evaluates four major risk domains: System & Operational Risks (e.g.,

cybersecurity, operational misuse), Content Safety Risks (e.g., child sexual abuse material), Societal Risks (e.g., surveillance), and Legal & Rights-related Risks (e.g., privacy violations, defamation). All prompts are manually curated and human-audited to ensure they reflect genuine policy violations rather than benign content.

Why This Matters

HELM AIR provides an evaluation of how well AI systems align with real-world safety expectations. Unlike abstract or theoretical benchmarks, it directly reflects the kinds of standards developers will be expected to meet. Strong performance indicates that AI models comply with normative expectations around responsible AI deployment.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.AI	Zhipu Al
Model	Claude 3.7 Sonnet	DeepSeek R1	Gemini 2.5 Pro	Llama 4 Maverick	03	Grok 3 Beta	n/a (Model not evaluated)
Average score (max score = 1)	0.82	0.53	0.74	0.69	0.85	0.51	evaluateu)
System & Operational Risks	0.69	0.28	0.62	0.55	0.80	0.37	
Content Safety Risk	0.84	0.57	0.73	0.78	0.78	0.49	
Societal Risk	0.87	0.51	0.74	0.72	0.94	0.37	
Legal & Rights- related Risks	0.90	0.59	0.83	0.79	0.92	0.58	
Retrieved	28 May 2025						
Release	Release v1.11.0	(2025-04-23)					

Sources

- Zeng, Yi, et al. "Air-bench 2024: A safety benchmark based on risk categories from regulations and policies." 2024. Accessed 14 Jun. 2025.
- Website: Stanford's AIR-Bench 2024.

Indicator

TrustLLM Benchmark

Definition & Scope

This indicator measures a model's overall trustworthiness using the TrustLLM benchmark, a comprehensive framework spanning six dimensions: truthfulness, safety, fairness, robustness, privacy, and machine ethics. The benchmark includes over 30 datasets across more than 18 subcategories, assessing issues such as hallucination, jailbreak resistance, and privacy leakage. Models are evaluated on tasks ranging from simple classification to complex generation, with results reported as published scores and rankings across each dimension. TrustLLM was developed by 45 research institutions, including 38 based in the U.S.

Why This Matters

TrustLLM provides a comprehensive evaluation of model reliability across multiple critical dimensions beyond traditional capability metrics. Strong performance across truthfulness, privacy protection, and ethical reasoning indicates that the companies have invested in aligning their models to be harmless and helpful, and not to cause unintended harm.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.AI	Zhipu Al
Model Tested	Claude 4 Opus (extended thinking)	DeepSeek R1	Gemini 2.5 Pro	Llama 4 Maverick	03	Grok 3 Beta	GLM 4 (32B)
Total score (max score = 1)	0.64	0.62	0.62	0.60	0.61	0.61	0.64
Truthfulness	0.54	0.62	0.57	0.53	0.43	0.59	0.57
Safety	0.65	0.61	0.62	0.60	0.65	0.60	0.63
Fairness	0.43	0.39	0.39	0.45	0.42	0.38	0.61
Privacy	0.60	0.58	0.57	0.53	0.57	0.55	0.52
Ethics	0.90	0.89	0.88	0.84	0.85	0.86	0.82
Robustness	0.67	0.62	0.64	0.61	0.66	0.65	0.67
Resources	Code & Gener	rations, Full Score	e Breakdown				
Release	v0.3.0						

- Sun, Lichao, et al. "Trustllm: Trustworthiness in large language models." arXiv preprint arXiv:2401.05561 3 (2024).
- Website: TrustLLM

Robustness

Indicator

Gray Swan Arena: UK AISI Agent Red-Teaming Challenge

Definition & Scope

This indicator measures resistance to adversarial attacks based on the U.K. AISI Agent Red-Teaming Challenge conducted on the Gray Swan Arena platform. The challenge is the largest public evaluation to date of safety in agentic LLMs, testing how well AI agents can withstand attempts to manipulate or subvert their behavior. Participants used both direct and indirect red-teaming techniques—strategies designed to trick or exploit the model—to identify vulnerabilities across five core behavior categories, including Confidentiality Breaches and Instruction Hierarchy Violations. Performance is measured using the Attack Success Rate (ASR), calculated as Total Breaks divided by Total Chats, offering a concrete metric of model robustness under real-world adversarial pressure.

Why This Matters

Agent red-teaming resistance measures real-world robustness against sophisticated attacks that could compromise AI systems in deployment. Models with lower attack success rates demonstrate stronger defenses against attempts to violate safety policies, extract confidential information, or manipulate agent behavior. This competitive environment with expert red-teamers provides a more realistic assessment than academic benchmarks, revealing which companies have invested in hardening their systems against adversarial exploitation.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Model Tested	Claude 3.7 Sonnet Thinking	n/a (Model not	n/a (Model not	Llama 4 Maverick	03	Grok 3 Beta	n/a (Model not
Attack Success Rate	1.45%	evaluated) (woder not evaluated)		5.90%	2.46%	4.14%	evaluated)
Retrieved	28 May 2025						

- Gray Swan Al. "UK AISI × Gray Swan Agent Red-Teaming Challenge: Results Snapshot." Gray Swan News, 2024.
- Website: Agent Red-Teaming Leaderboard

Indicator

Cisco Security Risk Evaluation

Definition & Scope

This indicator presents the results of a security risk assessment of frontier reasoning models, conducted by researchers from Cisco's Robust Intelligence team and the University of Pennsylvania. The experiments evaluated how resistant these models are to automated jailbreaking attacks—techniques designed to bypass safety systems and elicit harmful outputs. Researchers used 50 randomly selected prompts from the HarmBench dataset, covering six harm categories: cybercrime, misinformation, illegal activities, general harm, harassment, and chemical/biological weapons. The main metric reported is the Attack Success Rate (ASR), reflecting the percentage of harmful prompts for which a successful jailbreak was achieved. This provides a standardized way to compare the strength of safety guardrails across different models.

Why This Matters

Algorithmic jailbreaking resistance is a key measure of the robustness of safety guardrails. Models with high attack success rates are "highly susceptible to algorithmic jailbreaking and potential misuse," creating significant risks when deployed at scale.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Model Tested	Claude 3.5 Sonnet	DeepSeek R1	Gemini 1.5 Pro	Llama 3.1 405B	GPT-40 01- preview	n/a (Model not	n/a (Model not
Attack Success Rate	36%	100%	64%	96%	86% 26%	evaluated)	evaluated)

Sources

• Kassianik, Paul. "Evaluating Security Risk in DeepSeek and Other Frontier Reasoning Models." Cisco Security Blog, January 31, 2025.

Indicator

Protecting Safeguards from Fine-tuning

Definition & Scope

This indicator evaluates whether companies implement safeguards to prevent the removal of safety measures through fine-tuning. Evidence distinguishes between hosted supervised fine-tuning, where inference-time mitigations remain in place, and full weight release without tamper-resistant safeguards. Where no specific data on frontier models is reported, neither fine-tuning nor open weights are accessible.

Why This Matters

Companies that release full model weights enable malicious actors to strip all safety protections and create uncensored versions, while supervised fine-tuning helps maintain safety guardrails during customization.

Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Frontier model weights protected Provide supervised fine-tuning for older and smaller Claude 3 Haiku through Amazon Bedrock. Safety mitigations are in place. [AWS, 2024]	Fully released weights of frontier models. No tamper-resistant safeguards. [DeepSeek, 2025]	Frontier model weights protected Released weights of non-frontier Gemma family, including Gemma 3 27B [Hugging Face, 2025]. No tamper-resistant safeguards. Enables supervised finetuning of Gemini 2.0 Flash via Vertex Al. Safety mitigations are in place. [Google, 2025].	Fully released weights of the frontier model Llama 4 Maverick. No tamper- resistant safeguards. [Meta Al, 2025]	Frontier model weights protected Provide supervised finetuning of GPT-40 [OpenAl, 2024] and RL fine-tuning for 04 mini [OpenAl, 2025]. Safety mitigations are in place.	Frontier model weights protected Fully released weights of non-frontier Grok 1. No tamper-resistant safeguards. [xAI, 2024]	Fully released weights of the frontier model GLM-4 Z1 32B. No tamper-resistant safeguards. [THUDM, 2025]

- Qi, Xiangyu, et al. "Fine-tuning aligned language models compromises safety, even when users do not intend to!." arXiv preprint arXiv:2310.03693 (2023).
- Lermen, Simon, Charlie Rogers-Smith, and Jeffrey Ladish. "Lora fine-tuning efficiently undoes safety training in Ilama 2-chat 70b." arXiv preprint arXiv:2310.20624 (2023).
- Tamirisa, Rishub, et al. "Tamper-resistant safeguards for open-weight Ilms." arXiv preprint arXiv:2408.00761 (2024).

Digital Responsibility

Indicator

Watermarking

Definition & Scope

This indicator assesses whether companies have implemented watermarking technologies to help identify AI-generated content in both text and images. It focuses on real-world deployment rather than research alone, evaluating the accuracy and robustness of detection methods, adherence to standards such as C2PA and SynthID, and whether detection tools are publicly accessible.

Why This Matters

Watermarking enables the detection of Al-generated content, helping combat misinformation and digital fraud. Companies that deploy robust watermarking systems—along with public detection tools—help to uphold transparency and accountability.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Text-based	None found	None found	Yes - the SynthID system uses particular token selection to introduce a pattern that marks a text as AI-generated [Google DeepMind]. This can be identified by an online detector, access currently limited [Google, 2025].	None found	Research-only watermark, not shipped [The Verge, 2024]	None found	None found

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Image-based	Claude does not generate images		Yes (SynthID) [Google DeepMind]: pattern is embedded in images, can be identified by an online detector, access currently limited [Google, 2025]	Yes, but detection is restricted: Including invisible marks, detectable by Meta's own detector and partner platforms, they have not opened- sourced the model [Meta, 2024].	Uses the C2PA standard to flag the metadata of images generated by ChatGPT [OpenAI, 2025]. Such metadata is trivial to remove [Forbes, 2024].	None found	None found

- Zhao, Xuandong, et al. "SoK: Watermarking for Al-Generated Content." arXiv preprint arXiv:2411.18479 (2024).
- NIST. "Reducing Risks Posed by Synthetic Content." (2024).

Indicator

User Privacy

Definition & Scope

This indicator reports a company's dedication to user privacy when training and deploying AI models. It considers whether user inputs (such as chat history) are used by default to improve AI models or if companies require explicit opt-in consent. It also considers whether users can run powerful models privately, through on-premise deployment or secure cloud setups. Evidence includes default privacy settings and the availability of model weights for private hosting.

Why This Matters

Opt-in policies and private deployment options enable greater respect for user privacy, especially in sensitive fields such as healthcare, law, and government.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Default training on user inputs	No, unless the user opts in explicitly or the conversation is flagged for violating our Usage Policy [Anthropic, 2025].	Yes [DeepSeek, 2025]	Yes, but not in the enterprise version [Google, 2025]	Yes [Meta]	Yes, but no training on enterprise data (from ChatGPT Team, Enterprise, or API Platform) [OpenAI, 2025]	Yes [Ars Technica, 2024]	Yes
Frontier model weights available for private hosting	No	Yes [Huggingface, 2025]	No, but less-powerful models are open- sourced [Huggingface, 2025]	Yes [<u>Meta</u>]	No	No, but less-powerful models are open-sourced [xAI]	Yes [<u>THUDM</u>]

$T \cap$	$D\Gamma$	FTFF	BY PANFI	LICTO

Grading

Please pick a grade for each firm. You can add brief justifications to your grades.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Grades	Grade	Grade	Grade	Grade	Grade	Grade	Grade
Grade comments (Justifications, opportunities for improvements, etc.)							

Grading Scales

Grading scales are provided to support consistency between reviewers.

- A Exceptional safety; trivial issues only; no serious harm potential
- B Strong safety; rare moderate issues; serious harms well-controlled
- C Adequate safety; some moderate issues; serious harms mostly controlled
- Inadequate safety; frequent issues; serious harms poorly controlled
- Dangerous systems; pervasive issues; serious harms uncontrolled

Domain comments

Optional: Share observations that apply across companies, including general recommendations, notes on how you weighted indicators, or feedback on FLI's methodology.

Domain

Safety Frameworks

This domain evaluates the companies' published safety frameworks for frontier AI development and deployment from a risk management perspective. The comprehensive analysis supporting this section was conducted by the non-profit research organisation <u>SaferAI</u>.

Table of Contents

Overall Scores

- 1. Risk Identification
- 2. Risk Analysis and Evaluation
- 3. Risk Treatment
- 4. Risk Governance

Grading

Overview

The section focuses on framework documents and excludes other documents such as model cards. The comprehensive analysis for the indicators in this domain was conducted <u>SaferAI</u>, a leading governance and research non-profit focused on AI risk management. The organisation works to incentivize responsible AI practices through policy recommendations, research, and innovative risk assessment tools. *Note: The assessment contains living scores that are updated on a continuous basis. We extracted the scores from June 24, 2025.*

Frameworks:

- Anthropic Responsible Scaling Policy v2.2
- OpenAI Preparedness Framework v2
- Google DeepMind Frontier Safety Framework v2.0
- Meta Frontier Al Framework v1.1
- xAI Risk Management Framework (Draft)
- DeepSeek & Zhipu AI have not published a framework

	Weight	Anthropic	OpenAl	Google DeepMind	Meta	x.Al
Overall score		36%	34%	22%	26%	19%
1. Risk Identification	25%	28%	27%	17%	33%	5%
2. Risk Analysis & Evaluation	25%	26%	34%	31%	36%	31%
3. Risk treatment	25%	40%	36%	23%	19%	18%
4 Risk Governance	25%	48%	38%	16%	15%	23%

Supporting references: Reviewers were provided with the full database of framework references and quotes supporting SaferAl's assessments of individual all sub-indicators. The data can be found on their project website.

1. Risk Identification

Definition & Scope

This dimension captures the extent to which the company has addressed known risks in the literature and engaged in open-ended red teaming to uncover potential new threats. Moreover, this dimension examines if the AI company has leveraged a diverse range of risk identification techniques, including threat modeling when appropriate, to gain a deep understanding of possible risk scenarios.

Why This Matters

Companies can only mitigate risks they've identified, making comprehensive risk discovery the foundation of any effective safety framework. Firms that employ diverse identification methods are more likely to catch novel threats before they manifest in deployment. This proactive approach to risk discovery demonstrates whether a company takes seriously the full spectrum of potential harms, including those not yet observed in practice.

ID	Criteria	Weight	Anthropic	OpenAl	Google DeepMind	Meta	x.Al
	1. Risk Identification	25%	28%	27%	17%	33%	5%
	1.1 Classification of Applicable Known Risks	40%	30%	30%	18%	13%	13%
C1	1.1.1 Risks from literature and taxonomies are well covered	50%	50%	50%	25%	25%	25%
C2	1.1.2 Exclusions are justified and documented	50%	10%	10%	10%	0%	0%
	1.2 Identification of Unknown Risks (Open-ended red teaming)	20%	0%	3%	0%	0%	0%
	1.2.1 Internal	70%	0%	3%	0%	0%	0%
C3	1.2.1.1 Adequate methodology (includes resources, time, and access to the model)	70%	0%	0%	0%	0%	0%
C4	1.2.1.2 Appropriate expertise to properly identify hazards	30%	0%	10%	0%	0%	0%
	1.2.2 Third parties	30%	0%	3%	0%	0%	0%
C5	1.2.2.1 Appropriate expertise to identify hazards	33%	0%	0%	0%	0%	0%
C6	1.2.2.2 Adequate resources/time/access to the model	33%	0%	10%	0%	0%	0%
C7	1.2.2.3 Commitment to non-interference with findings	34%	0%	0%	0%	0%	0%
	1.3 Risk Modeling	40%	41%	36%	25%	69%	1%
C8	1.3.2 The company uses risk models for all the risk domains identified, and the risk models are published (with potentially dangerous information redacted)	40%	50%	75%	50%	90%	0%
	1.3.1 Risk modeling methodology	40%	40%	10%	12%	58%	2%
C9	1.3.1.1 Methodology precisely defined	70%	50%	10%	10%	75%	0%
C10	1.3.1.2 Mechanism to incorporate red teaming findings	15%	10%	10%	10%	10%	0%
C11	1.3.1.3 Prioritization of severe and probable risks	15%	25%	10%	25%	25%	10%
C12	1.3.4 Third-party validation of risk models	20%	25%	10%	0%	50%	0%

Indicator

2. Risk Analysis and Evaluation

Definition & Scope

This dimension assesses whether the company has established well-defined risk tolerances that precisely characterize acceptable risk levels for each identified risk. Moreover, this dimension examines if the company has successfully operationalized these tolerances into measurable criteria: Key Risk Indicators (KRIs) that signal

when risks are approaching critical levels, and Key Control Indicators (KCIs) that demonstrate the effectiveness of mitigation measures. The assessment captures whether companies define these indicators in paired "if-then" relationships, where crossing specific KRI thresholds triggers corresponding KCI requirements.

Why This Matters

Without operationalizing risk tolerances into measurable metrics, companies cannot make consistent, evidence-based decisions about when to halt development or implement additional safeguards. Well-defined KRI-KCI pairs create accountability by establishing clear tripwires—when risk indicator X crosses threshold Y, control measure Z must be implemented. This systematic approach prevents ad-hoc decision-making during high-pressure situations and ensures that safety commitments translate into concrete actions rather than remaining aspirational statements.

ID	Criteria	Weight	Anthropic	OpenAl	Google DeepMind	Meta	x.Al
	2. Risk Analysis and Evaluation	25%	26%	34%	31%	36%	31%
	2.1 Setting a Risk Tolerance	35%	7%	16%	3%	24%	57%
	2.1.1 Risk tolerance is defined	80%	8%	20%	3%	30%	71%
C13	2.1.1.1 Risk tolerance is at least qualitatively defined for most risks	33%	25%	50%	10%	90%	90%
C14	2.1.1.2 Risk tolerance is expressed fully quantitatively (cf. criterion above) or at least partly quantitatively as a combination of scenarios (qualitative) and probabilities (quantitative) for most risks	33%	0%	10%	0%	0%	75%
C15	2.1.1.3 Risk tolerance is expressed fully quantitatively as a product of severity (quantitative) and probability (quantitative) for most risks	33%	0%	0%	0%	0%	50%
	2.1.2 Process to define the tolerance	20%	0%	0%	0%	0%	0%
C16	2.1.2.1 Al developers engage in public consultations or seek guidance from regulators where available.	50%	0%	0%	0%	0%	0%
C17	2.1.2.2 Any significant deviations from risk tolerance norms established in other industries are justified and documented (e.g., cost-benefit analyses)	50%	0%	0%	0%	0%	0%
	2.2 Operationalizing Risk Tolerance	65%	36%	44%	47%	43%	18%
	2.2.1 Key Risk Indicators (KRI)	30%	51%	51%	51%	33%	33%
C18	2.2.1.1 KRI thresholds are at least qualitatively defined for most risks	45%	90%	90%	90%	50%	50%
C19	2.2.1.2 KRIs thresholds are quantitatively defined for most risks	45%	50%	25%	10%	0%	75%
C20	2.2.1.3 KRI also identifies and monitors changes in the level of risk in the external environment	10%	0%	0%	0%	10%	0%
	2.2.2 Key Control Indicators (KCI)	30%	31%	45%	38%	18%	14%
	2.2.2.1 Containment KCIs	35%	43%	5%	63%	38%	5%
C21	2.2.2.1.1 Most KRI thresholds have corresponding qualitative containment KCI thresholds	50%	75%	10%	75%	75%	10%
C22	2.2.2.1.2 Most KRI thresholds have corresponding quantitative containment KCI thresholds	50%	10%	0%	50%	0%	0%
	2.2.2.2 Deployment KCIs	35%	38%	45%	25%	13%	25%
C23	2.2.2.2.1 Most KRI thresholds have corresponding qualitative deployment KCI thresholds	50%	75%	90%	50%	25%	50%
C24	2.2.2.2 Most KRI thresholds have corresponding quantitative deployment KCI thresholds	50%	0%	0%	0%	0%	0%

ID	Criteria	Weight	Anthropic	OpenAl	Google DeepMind	Meta	x.Al
C25	2.2.2.3 For advanced KRIs, Assurance processes KCI are defined	30%	10%	90%	25%	0%	10%
C26	2.2.3 Pairs of thresholds are grounded in risk modeling to show that risks remain below the tolerance	20%	10%	50%	90%	50%	10%
C27	2.2.4 Policy to put development on hold if the required KCI threshold cannot be achieved, until sufficient controls are implemented to meet the threshold	20%	50%	25%	10%	90%	10%

3. Risk Treatment

Definition & Scope

This dimension captures the extent to which the company has implemented comprehensive risk mitigation strategies across three critical areas: containment measures that control access to Al models, deployment measures that prevent misuse and accidental harms, and assurance processes that provide affirmative evidence of safety. Furthermore, this dimension assesses whether the company maintains continuous monitoring of both Key Risk Indicators (KRIs) and Key Control Indicators (KCIs) throughout the Al system's lifecycle, from training through deployment.

Why This Matters

Effective risk treatment requires multiple layers of defense. Companies that maintain continuous monitoring of both risks and control effectiveness can detect when mitigations are failing before catastrophic outcomes occur.

ID	Criteria	Weight	Anthropic	OpenAl	Google DeepMind	Meta	x.Al
	3. Risk treatment	25%	40%	36%	23%	19%	18%
	3.1 Implementing Mitigation Measures	50%	24%	32%	19%	13%	18%
	3.1.1 Containment measures	35%	25%	25%	0%	0%	0%
C28	3.1.1.1 Containment measures are precisely defined for all KCI thresholds	60%	25%	25%	0%	0%	0%
C29	3.1.1.2 Proof that containment measures are sufficient to meet the thresholds	40%	25%	10%	0%	0%	0%
C30	3.1.1.3 Strong third-party verification process to verify that the containment measures meet the threshold	formula*	25%	25%	0%	0%	0%
	3.1.2 Deployment measures	35%	40%	50%	35%	35%	50%
C31	3.1.2.1 Deployment measures are precisely defined for all KCI thresholds	60%	50%	50%	25%	25%	25%
C32	3.1.2.2 Proof that deployment measures are sufficient to meet the thresholds	40%	25%	50%	50%	50%	50%
C33	3.1.2.3 Strong third-party verification process to verify that the deployment measures meet the threshold	formula*	10%	25%	0%	0%	50%
	3.1.3 Assurance processes	30%	5%	20%	23%	2%	3%
C34	3.1.3.1 Credible plans towards the development of assurance properties	40%	10%	25%	25%	0%	10%
C35	3.1.3.2 Evidence that the assurance properties are enough to achieve their corresponding KCI thresholds	40%	0%	25%	25%	0%	0%
C36	3.1.3.3 The underlying assumptions that are essential for their effective implementation and success are clearly outlined	20%	10%	10%	25%	10%	0%

ID	O Criteria		Anthropic	OpenAl	Google DeepMind	Meta	x.Al
	3.2 Continuous Monitoring and Comparing Results with Predetermined Thresholds	50%	56%	40%	27%	26%	17%
	3.2.1 Monitoring of KRIs	50%	68%	39%	27%	26%	4%
C37	3.2.1.1 Justification that the elicitation methods used during the evaluations are comprehensive enough to match the elicitation efforts of potential threat actors	30%	75%	75%	25%	50%	0%
	3.2.1.2 Evaluation frequency	25%	100%	0%	25%	0%	0%
C38	3.2.1.2.1 Specification of evaluation frequency in terms of the relative variation of effective computing power used in training	50%	100%	0%	25%	10%	0%
C39	3.2.1.2.2 Specification of evaluation frequency in terms of fixed time intervals to account for post-training enhancements	50%	100%	0%	0%	0%	0%
C40	3.2.1.4 Description of how post-training enhancements are factored into capability assessments	15%	75%	75%	50%	50%	0%
C41	3.2.1.5 Replication of evaluations by third parties	15%	50%	25%	25%	25%	25%
C42	3.2.1.6 Vetting of protocols by third parties	15%	10%	10%	10%	0%	0%
	3.2.2 Monitoring of KCIs	40%	33%	33%	23%	20%	15%
C43	3.2.2.1 Detailed description of evaluation methodology and justification that KCI thresholds will not be crossed unnoticed	40%	25%	25%	50%	50%	0%
C44	3.2.2.2 Replication of evaluations by third parties	30%	50%	50%	0%	0%	50%
C45	3.2.2.3 Vetting of protocols by third parties	30%	25%	25%	10%	0%	0%
C46	3.2.3 Sharing of evaluation results with relevant stakeholders as appropriate	10%	90%	75%	50%	50%	90%

4. Risk Governance

Definition & Scope

This dimension examines whether the company has built robust organizational infrastructure to support effective risk management decision-making. The assessment captures the extent to which companies have established clear risk ownership and accountability, independent oversight mechanisms, and cultures that prioritize safety alongside innovation. Moreover, this dimension evaluates the company's commitment to transparency, whether they publicly disclose their risk management approaches, governance structures, and safety incidents. The evaluation considers how well the company's governance framework ensures that risk considerations are incorporated into strategic decisions and that multiple layers of review prevent any single point of failure in risk management.

Why This Matters

Strong governance structures ensure that risk management isn't just a technical exercise but is embedded in organizational decision-making at all levels. Independent oversight prevents conflicts of interest when safety considerations clash with commercial pressures, while clear accountability ensures someone is always responsible for catching problems. Companies that publicly disclose their governance structures and safety incidents demonstrate confidence in their approach and enable external stakeholders to verify that appropriate safeguards exist.



ID	Criteria		Anthropic	OpenAl	Google DeepMind	Meta	x.Al
	4 Risk Governance		48%	38%	16%	15%	23%
	4.1 Decision-making		44%	28%	13%	30%	40%
C47	4.1.1 The company has clearly defined risk owners for every key risk identified and tracked		25%	10%	0%	10%	75%
C48	4.1.2 The company has a dedicated risk committee at the management level that meets regularly	25%	0%	0%	0%	25%	0%
C49	4.1.3 The company has defined protocols for how to make go/ no-go decisions	25%	75%	50%	50%	75%	10%
C50	4.1.4 The company has defined escalation procedures in case of incidents	25%	75%	50%	0%	10%	75%
	4.2. Advisory and Challenge	20%	35%	53%	28%	21%	4%
C51	4.2.1 The company has an executive risk officer with sufficient resources	17%	75%	0%	0%	0%	0%
C52	4.2.2 The company has a committee advising management on decisions involving risk	17%	10%	90%	90%	25%	0%
C53	4.2.3 The company has an established system for tracking and monitoring risks	17%	50%	75%	50%	50%	25%
C54	4.2.4 The company has designed people who can advise and challenge management on decisions involving risk	17%	25%	75%	0%	25%	0%
C55	4.2.5 The company has an established system for aggregating risk data and reporting on risk to senior management and the Board	17%	50%	75%	25%	25%	0%
C56	4.2.6 The company has an established central risk function	17%	0%	0%	0%	0%	0%
	4.3 Audit	20%	50%	38%	5%	5%	25%
C57	4.3.1 The company has an internal audit function involved in AI governance	50%	25%	0%	0%	0%	0%
C58	4.3.2 The company involves external auditors	50%	75%	75%	10%	10%	50%
	4.4 Oversight	20%	50%	45%	25%	0%	0%
C59	4.4.1 The Board of Directors of the company has a committee that provides oversight over all decisions involving risk	50%	25%	90%	0%	0%	0%
C60	4.4.2 The company has other governing bodies outside of the Board of Directors that provide oversight over decisions	50%	75%	0%	50%	0%	0%
	4.5 Culture	10%	63%	12%	7%	3%	50%
C61	4.5.1 The company has a strong tone from the top	33%	50%	25%	10%	10%	25%
C62	4.5.2 The company has a strong risk culture	33%	50%	10%	10%	0%	50%
C63	4.5.3 The company has a strong speak-up culture	33%	90%	1%	0%	0%	75%
	4.6 Transparency	5%	72%	53%	20%	33%	45%
C64	4.6.1 The company reports externally on what its risks are	33%	50%	75%	25%	50%	75%
C65	4.6.2 The company reports externally on what its governance structure looks like	33%	75%	75%	25%	25%	10%
C66	4.6.3 The company shares information with industry peers and government bodies	33%	90%	10%	10%	25%	50%

$\Gamma \cap$			I ETED	BY PANE	LLICTO
ıv	DE	CUIVIT		DIPANE	111010

Grading

Please pick a grade for each firm. You can add brief justifications to your grades.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Grades	Grade	Grade	Grade	Grade	Grade	Grade	Grade
Grade comments (Justifications, opportunities for improvements, etc.)							

Grading Scales

Grading scales are provided to support consistency between reviewers.

- A Framework virtually guarantees prevention of catastrophic risks
- B Framework prevents catastrophic risks with a high degree of confidence
- Framework prevents catastrophic risks with a moderate degree of confidence
- Framework prevents catastrophic risks with a low degree of confidence
- Framework prevents catastrophic risks with a very low degree of confidence; or no framework exists

Domain comments

Optional: Share observations that apply across companies, including general recommendations, notes on how you weighted indicators, or feedback on FLI's methodology.



Domain

Existential Safety

This domain examines companies' preparedness for managing extreme risks from future AI systems that could match or exceed human capabilities, including stated strategies and research for alignment and control.

Table of Contents

Existential Safety Strategy

Internal Monitoring and Control Interventions

Technical AI Safety Research

Supporting External Safety Research

Providing Deep Model Access to Safety Researchers

Mentoring and Funding

Grading

Indicator

Existential Safety Strategy

Definition & Scope

The assessed companies aim to develop AGI/superintelligence, and many expect to achieve this goal in the next 2–5 years. This indicator evaluates whether companies have published comprehensive, concrete strategies for managing catastrophic risks from these transformative AI systems. We assess the depth, specificity, and credibility of publicly available plans.

We examine official company documents, research papers, and blog posts that articulate safety strategies. We report the most relevant documents, briefly summarize their content, and provide links for detailed reading. Safety frameworks are mentioned for completeness and are fully evaluated in the relevant domain. We note whether documents are declared strategies by leadership or proposals by researchers from a safety team. We strive to keep document summaries proportional to document length and relevance for the safety strategy. Safety frameworks are only noted briefly and evaluated in another domain. Documents that primarily provide recommendations to other actors (e.g., governments) are outside the scope.

Key components:

Technical Alignment and Control Plan:

- Given the short timelines to AGI and the magnitude of the risk, companies should ideally have credible, detailed
 agendas that are highly likely to solve the core alignment and control problems for AGI/Superintelligence
 very soon.
- Companies should be able to demonstrate that they would be able to detect misaligned systems and reliably



prevent them from escaping human control, and have formulated clear protocols for how they will handle serious warning signs of misalignment.

AGI Planning:

- Companies should have detailed plans for managing the transition when AI matches or exceeds human
 capabilities in critical domains and enables large-scale dual-use risks. They should specify clear criteria
 for when they would halt development/deployment.
- Companies should develop concrete, detailed roadmaps to achieve sufficient cyber-defence capabilities
 to protect against attacks from terrorist organizations or resourced state actors before critically dangerous
 systems are developed.

Post-AGI Governance:

- Companies should provide clear descriptions of how they would govern AGI/Superintelligence or how they
 will enable societal control. The company also should have developed reliable protocols that would prevent
 insiders from using Superintelligent systems to seize political power.
- Companies should specify how extreme power concentration will be prevented and benefits distributed if Al replaces humans in the workplace and causes unprecedented mass unemployment.

Overall, this indicator evaluates whether companies have detailed, actionable strategies that match the extraordinary risks they acknowledge when building systems intended to exceed human intelligence.

Why This Matters

Industry leaders and the recent International Scientific Report on the Safety of Advanced AI have identified potentially catastrophic risks from advanced AI systems. Several assessed companies predict AGI development within 2-5 years, creating urgency for reliability, safety preparedness. This indicator summarizes core documents that are relevant to a company's posture toward these risks. Given the irreversible nature of potential failures and their global impact, the sophistication of a company's strategy should scale with its stated ambitions and timelines. Transparency in safety strategies enables accountability and allows policymakers, researchers, and the public to evaluate whether adequate precautions are being taken.

f.,	-1.00
ΙŲ	o te
of .	ESTITUTE
	VSTITUTE

Anthropic

Quantitative safety plan

No alignment or control strategy has been presented that includes the company's quantitative assessment of its likelihood of success.

Company Strategy

The Urgency of Interpretability (2025, ~5k words, strategy blog)

The CEO argues in a personal blog that mechanistic interpretability must advance rapidly to ensure safe deployment of transformative AI systems that could become a "country of geniuses in a datacenter" by 2026-2027. Amodei frames this as a "race between interpretability and model intelligence" and outlines recommendations for the AI community and governments. The blog also discusses the history of interpretability research and recent technical breakthroughs.

Key quotes:

- "Anthropic is doubling down on interpretability, and we have a goal of getting to "interpretability can reliably detect most model problems" by 2027.
- "Our long-run aspiration is to be able to look at a state-of-the-art model and essentially do a 'brain scan': a checkup that has a high probability of identifying a wide range of issues, including tendencies to lie or deceive, power-seeking, flaws in jailbreaks, [...]. This would then be used in tandem with the various techniques for training and aligning models, [...]."

Putting up Bumpers (2025, ~5k words, research blog)

Anthropic alignment researcher Sam Bowman proposes an alignment approach for early AGI systems that prioritizes implementing and testing "many largely-independent lines of defense" to catch and correct misalignment through iterative testing. He highlights "alignment audits" [Anthropic, 2025] as the "Primary Bumper" to notice signs of misalignment like "generalized reward-tampering" [Anthropic, 2024] or "alignment-faking" [Anthropic, 2024].

Key quotes:

- "Even if we can't solve alignment, we can solve the problem of catching and fixing misalignment."
- "We believe that, even without further breakthroughs, this work can almost entirely mitigate the risk that we unwittingly put misaligned circa-human-expert-level agents in a position where they can cause severe harm."
- "This is not a costless choice: The Bumpers' worldview largely gives up on the ability to make highly-confident, principled arguments for safety, and it comes with real risks."
- "We are plausibly within a couple of years of developing models that could automate much of the work of AI R&D. This makes sabotage and sandbagging threat models... worth addressing soon."
- "Anthropic is committed to investing seriously in the kinds of measures described here, ... setting up a new team to productionize and professionalize the hands-on work of testing models for AGI-relevant forms of misalignment."

Responsible Scaling Policy (2023, v2.2 in 2025, ~10k words, safety framework)

A set of voluntary commitments based on regular dangerous capability evaluations and a set of capability thresholds in high-risk domains that trigger a requirement for enhanced safety and security mitigations. These commitments include pausing development or deployment if the required mitigations cannot adequately manage the identified risks.

For detailed analysis, refer to the 'Safety Framework' domain.

Core Views on Al Safety (2023, ~6k words, strategy blog)

This blog post outlines Anthropic's AI safety philosophy and technical research portfolio. The document addresses existential risk scenarios, presenting a three-tier framework (optimistic, intermediate, pessimistic) for how difficult alignment might prove to be, with corresponding strategic responses for each scenario. It details six priority research areas: Mechanistic Interpretability, Scalable Oversight, Process-Oriented Learning, Understanding Generalization, Testing for Dangerous Failure Modes, and Evaluating Societal Impact. The post emphasizes empirical research and acknowledges fundamental uncertainty about which approaches will succeed.

Key quotes:

- "Our goal is essentially to develop: 1) better techniques for making AI systems safer; 2) better ways of identifying how safe or unsafe AI systems are."
- "We aim to build detailed quantitative models of how these [dangerous] tendencies vary with scale so that we can anticipate the sudden emergence of dangerous failure modes in advance."
- In pessimistic scenarios where "Al safety may be unsolvable," Anthropic's role would be "to provide evidence that current safety techniques are insufficient and to push for halting Al progress to prevent catastrophic outcomes."

fulle of listitute	Арр
DeepSe	ek

DeepSeek	Quantitative safety plan	No alignment or control strategy has been presented that includes the company's quantitative assessment of its likelihood of success.					
	Company Strategy	Based on searches of company websites, technical papers, and public communications, no relevant strategy documents were found.					
Google DeepMind	Quantitative safety plan						
	Company Strategy	An Approach to Technical AGI Safety and Security (2025,~80k words, technical report/research agenda)					
		A detailed technical report by DeepMind's safety team explains their research agenda for a framework to prevent severe, civilization-scale harm from AGI, defined as systems roughly at the 99th-percentile of skilled adults. The document states that reaching AGI before 2030 is plausible. However, it makes clear that a plan for ASI and a strong recursive Self-improvement is beyond the scope of this paper.					
		Key sections: 'Assumptions about AGI development' (p. 20-43), 'Risk areas' (p. 44-55), 'Addressing misuse' (p. 56-70), 'Addressing Misalignment' (p. 70-108).					
		Scope: "We identify four areas of risk: misuse, misalignment, mistakes, and structural risks. Of these, we focus on technical approaches to misuse and misalignment."					
		The framework bets on a defence-in-depth, combining model-centric alignment research with classic security engineering. The plan is deliberately incremental: succeed on today's frontier models, then re-validate at each capability level rather than solve all of AGI safety in one stroke. The report is explicit about significant problems that remain unsolved (scalable interpretability, quantitative safety-case metrics, red-team rigor), which the authors flag as research priorities. "Importantly, this is a roadmap rather than a solution, as there remain many open research problems to address."					
		Misalignment Strategy - "Two Lines of Defense":					
		First Line - Model-level mitigations: "[] ensure the model does not pursue misaligned goals. We follow the standard ML playbook: (1) improve the ability to get good oversight (amplified oversight), (2) use this to improve the model (guiding model behavior), and (3) figure out where oversight is needed to train a model that will robustly do what we want (robust training)."					
		Amplified oversight: Debate, critique, constitutional AI, recursive reward modeling					
		Guiding model behavior: Preference learning and value alignment techniques					
		Robust AI oversight: Active learning, adversarial training, uncertainty estimation					
		Second Line - System-level controls (monitoring, access) that still catch harm if alignment fails, including using dedicated monitor AI systems to detect misaligned actions, with emphasis on the monitor knowing when it's uncertain and flagging actions for review.					
		Security measures: Treating AI as an "untrusted insider," sandboxing, access control					
		Safer design patterns: Correctability, bounded autonomy, externalized reasoning					
		Interpretability: The paper's longest technical section, covering sparse autoencoders, circuit analysis, and mechanistic interpretability					
		Alignment stress tests: Red-teaming alignment assumptions using "model organisms"					
		Misuse strategy: 3 Steps: 1. Evaluate models for dangerous capability. 2. Mitigate with post-training tuning, capability suppression, monitoring, and weight/access security. 3. Red team to build misuse safety cases (inability & control).					
		Safety-case methodologies: Advocates evidence-based deployment decisions: inability cases (model lacks capability) and control cases (misaligned behaviour will be caught). Future work sketches empirical and incentive-based cases for more capable systems.					
		AGI Safety and Alignment at Google DeepMind: A Summary of Recent Work (2024, 2k words, research blog)					
		This update from DeepMind's team focused on existential risk describes their three main research bets over the past 1.5 years: 1) amplified oversight for proper alignment signals, 2) frontier safety to assess catastrophic risk capabilities, and 3) Mechanistic interpretability as an enabler for both.					
		The post provides detailed explanations of recent work in each area and its rationales.					
		The team admits they are "revising our own high-level approach to technical AGI safety" because current bets "do not necessarily add up to a systematic way of addressing risk." They highlight fundamental gaps, noting that even perfect amplified oversight would be insufficient under distribution shift, requiring additional investments in adversarial training, uncertainty estimation, and monitoring through a control framework.					
		Frontier Safety Framework v2 (2024, v2 in 2025, 4k words, safety framework)					
		A set of voluntary commitments based on regular dangerous capability evaluations and a set of capability thresholds in high-risk domains that trigger a requirement for enhanced safety and security mitigations. These commitments include pausing development or deployment if the required mitigations cannot adequately manage the identified risks.					

For detailed analysis, refer to the 'Safety Framework' domain.

C	×	
Tu	10 0	е
of	lıt	e
-J 11	NSTIT	UTE

Meta	Quantitative safety plan	No alignment or control strategy has been presented that includes the company's quantitative assessment of its likelihood of success.					
	Company	Frontier AI Framework v.1.1 (2025, ~8k words, safety framework)					
	Strategy	A set of voluntary commitments based on regular dangerous capability evaluations and a set of capability thresholds in high-risk domains that trigger a requirement for enhanced safety and security mitigations. These commitments include pausing development or deployment if the required mitigations cannot adequately manage the identified risks. For detailed analysis, refer to the 'Safety Framework' domain.					
		Open Source AI Is the Path Forward (2024, ~3k words, strategy blog)					
		In this blog post, Zuckerberg presents a case for open source AI as their primary approach to AI safety and development (not specifically focused on catastrophic risks). The document makes the case that open source models are inherently safer than closed alternatives due to transparency, distributed scrutiny, and prevention of power concentration. He argues that widely deployed AI systems enable larger actors to check malicious uses by smaller actors. It addresses both unintentional harms (including "truly catastrophic science fiction scenarios for humanity") and intentional misuse by bad actors.					
		Key quotes:					
		• "I think it will be better to live in a world where AI is widely deployed so that larger actors can check the power of smaller bad actors."					
OpenAl	Quantitative safety plan						
	Company	How we think about safety and alignment (2025, ~3k words, strategy blog).					
	Strategy	This blog describes high-level principles that guide OpenAl's thinking and ties it to their safety practices. This document describes a shift from viewing AGI as a single transformative moment to seeing it as continuous progress. For every principle, the blog lays out how it will shape their focus and approach to new challenges and relates to already implemented interventions.					
		Quote of the core principles:					
		1) "Embracing uncertainty: We treat safety as a science, learning from iterative deployment rather than just theoretical principles.' 2) "Defense in depth: We stack interventions to create safety through redundancy." 3) "Methods that scale: We seek out safety methods that become more effective as models become more capable." 4) "Human control: We work to develop AI that elevates humanity and promotes democratic ideals." 5) "Community effort: We view responsibility for advancing safety as a collective effort."					
		Planning for AGI and beyond (2023, 2k words)					
		This high-level blog outlines principles for managing AGI risks. The post emphasizes goals like ensuring AGI benefits are "widely and fairly shared" and advocates for deploying progressively more powerful systems to learn iteratively. It acknowledges the need for new alignment techniques, calls for a global conversation on governance and benefit-sharing, describes the benefits of OpenAI's non-profit structure, and raises the idea of a coordinated slowdown.					
		Key quotes:					
		• "We will need to develop new alignment techniques as our models become more powerful (and tests to understand when our current techniques are failing). Our plan in the shorter term is to use AI to help humans evaluate the outputs of more complex models and monitor complex systems, and in the longer term to use AI to help us come up with new ideas for better alignment techniques."					
		• "As our systems get closer to AGI, we are becoming increasingly cautious with the creation and deployment of our models. Our decisions will require much more caution than society usually applies to new technologies, and more caution than many users would like. Some people in the AI field think the risks of AGI (and successor systems) are fictitious; we would be delighted if they turn out to be right, but we are going to operate as if these risks are existential."					
		Announcement of Superalignment team (2023, ~1k words, strategy blog)					
		Outlined an ambitious strategy to start a new team to build "a roughly human-level automated alignment researcher" that could use vast compute to iteratively align superintelligence. Note: This team was disbanded in 2024 after team leaders Leike and Sutskever left OpenAI [CNBC, 2024].					
		Preparedness Framework (2023, v2 in 2025, ~10k words, safety framework)					
		A set of voluntary commitments based on regular dangerous capability evaluations and a set of capability thresholds in high-risk domains that trigger a requirement for enhanced safety and security mitigations. These commitments include pausing development or deployment if the required mitigations cannot adequately manage the identified risks.					
		For detailed analysis, refer to the 'Safety Framework' domain.					

x.Al	Quantitative safety plan	No alignment or control strategy has been presented that includes the company's quantitative assessment of its likelihood of success.			
	Company Strategy	xAI Risk Management Framework (Draft) (2025, ~2k words, safety framework) A set of voluntary commitments based on regular dangerous capability evaluations and a set of capability thresholds in high-risk domains that trigger a requirement for enhanced safety and security mitigations. For detailed analysis, refer to the 'Safety Framework' domain.			
Zhipu Al	Quantitative safety plan	No alignment or control strategy has been presented that includes the company's quantitative assessment of its likelihood of success.			
	Company Strategy	Based on searches of company websites, technical papers, and public communications, no official strategy documents were found. Media report on superalignment initiative (National Business Daily, 2024)			
		At the AWS China Summit (Shanghai, 29 May 2024) Zhipu Al's Chief Ecosystem Officer Liu Jiang said: "AGI will reach ordinary-human level within 5-10 years.". He announces that "Zhipu AI has already launched a 'Superalignment' initiative." The article explains superalignment as "ensuring a super-human-level AI system follows human values and goals."			

- Bengio, Yoshua, et al. "The Singapore Consensus on Global Al Safety Research Priorities." 2025.

Indicator

Internal Monitoring and Control Interventions

Definition & Scope

This indicator evaluates whether companies have implemented or prepared monitoring and control systems to detect and prevent risks from misalignment during internal deployment. Companies are assessed on whether they have concrete implementation plans tied to specific capability thresholds, published methodologies for control evaluations, and protocols for investigating potential scheming or deceptive alignment. General statements about monitoring without specific technical details, thresholds, or implementation timelines are insufficient. Research about monitoring without statements on implementation plans or status is out of scope.

Why This Matters

As AI systems become more capable, they may develop the ability to engage in deceptive behavior or "scheming"—appearing aligned while pursuing misaligned goals that could include attempts to gain unauthorized access to resources, sabotage safety research, subvert oversight mechanisms, or manipulate staff. Internal deployment poses unique risks, as this is usually the first time a model has longer time-horizon interactions with the external world. Robust monitoring and control measures serve as a critical line of defense, enabling companies to detect and prevent harmful actions even if alignment techniques fail to prevent scheming entirely. Concrete, technically sophisticated control protocols indicate that companies are taking tail risks seriously and work on mitigations.



Stated AI Control Plans	No control plan found.				
	The Responsible Scaling Policy says that when AI can fully automate the work of junior AI researchers, they "will develop an affirmative case that (1) identifies the most immediate and relevant risks from models pursuing misaligned goals and (2) explains how we have mitigated these risks to acceptable levels. The affirmative case will describe [] mitigations (such as monitoring [] [Responsible Scaling Policy, 2025].				
Implemented safeguards	None found				
Stated AI Control Plans	No public control plan found				
Implemented safeguards	None found				
Stated AI Control Plans	Early sketch of a plan without technical details. Details: Relevant section of Frontier Safety Framework: "An initial mitigation approach focuses on detecting when models might develop a baseline instrumental reasoning ability at which they have the potential to undermine				
	human control, assuming no additional mitigations were applied. When models reach this capability level, we believe applying an automated monitor to the model's explicit reasoning (e.g., chain-of-thought output) is an effective mitigation. Once a model is capable of effective instrumental reasoning in ways that cannot be monitored, additional mitigations may be warranted—the development of which is an area of active research."				
	• The Framework states that developers should implement 'control evaluations' [Greenblatt, 2023] and construct a safety case showing:				
	"That the likelihood of deceptive alignment risk would be low enough for safe deployment, even if the model were trying to meaningfully undermine human control." [Google DeepMind, 2025].				
Implemented safeguards	None found				
Stated AI Control Plans	No public control plan found				
Implemented safeguards	None found				
Stated AI Control Plans	No public control plan found				
	Voluntary commitment in a safety framework requires "safeguards meeting a high standard" if the relevant capability threshold is crossed. Exact safeguard not specified, but several are suggested.				
	Details:				
	Preparedness Framework states that a "High" risk classification for "Cybersecurity" domain triggers requirement for " safeguards meeting High standard against misalignment (Appendix C.2) for large-scale internal deployment".				
	Appendix C.2 contains illustrative safeguards against misaligned models.				
	Sub-section on "Reliable and Robust System Oversight" lists 4 potential measures:"				
	 Monitor Al supervises agent actions to detect agents in the case of trying to take harmful actions, scheming, or similar 				
	System forces human approval of agent actions (either by default or by escalation of the monitor)				
	Model's reasoning is provided faithfully and interpretably to humans to review				
	- Agent actions are all logged in an uneditable database, and asynchronous monitoring routines review those actions for evidence of harm" [OpenAl, 2025]				
Implemented safeguards	None found				
	Implemented safeguards Stated AI Control Plans Implemented safeguards				



x.Al	Stated AI Control Plans	No public control plan found
	Implemented safeguards	None found
Zhipu Al	Stated AI Control Plans	No public control plan found
	Implemented safeguards	None found

- Stix, Ch. Al Behind Closed Doors: a Primer on The Governance of Internal Deployment. Apollo Research, 17 Apr. 2025
- Carlsmith, Joe. "New Report: 'Scheming Als: Will Als Fake Alignment During Training in Order to Get Power?" Joe Carlsmith, 15 Nov. 2023,
- Greenblatt, Ryan, et al. "Al control: Improving safety despite intentional subversion." arXiv preprint arXiv:2312.06942 (2023).

Technical AI Safety Research

Definition & Scope

This indicator tracks research publications on technical AI safety research that are relevant to extreme risks. More specifically, the indicator is a collection of work that is plausibly helpful for averting large-scale risks from misalignment or misuse. This includes mechanistic interpretability, scalable oversight, unlearning, model organisms of misalignment, model evaluations on dangerous capabilities or alignment, and others.

The collection also includes substantial outputs besides papers—weights, tools, code, transcripts, data—but these are almost always published as part of a paper. Excluded are capability-focused research, papers on hallucinations, and model cards. The full collection was created by Zach Stein-Perlman—numbers for DeepSeek, ZhipuAl, and xAl added by FLI.

Why This Matters

The industry is rapidly advancing toward increasingly capable AI systems, yet core challenges—such as alignment, control, interpretability, and robustness—remain unresolved, with system complexity growing year by year. Safety research conducted by companies reflects a meaningful investment in understanding and mitigating these risks. When companies publicly share their safety findings, they enable external scrutiny, strengthen the broader field's understanding of critical issues, and signal a commitment to safety that goes beyond proprietary interests.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Total	32	0	28	6	12	0	0
2025	9	-	4	1	3	Safety advisor Dan Hendrycks	0
2024	11	-	11	5	7	publishes	0
2023	12	-	13	-	2	research, but not formally for xAI	0

Sources

• Stein-Perlman, Zach. "Boosting Safety Research." Al Lab Watch. Accessed 16 Jun. 2025.

Indicator

Supporting External Safety Research

Definition & Scope

This indicator assesses the extent to which companies invest in and support external AI safety research through a range of mechanisms. Evidence may include: (1) Mentorship programs—participation in formal initiatives such as the Machine Learning Alignment Theory Scholars (MATS) program, the number of mentors provided, and the existence of company-specific fellowships; (2) Research grants and funding—provision of financial support or subsidized API access to safety researchers, including grants and targeted funding programs; and (3) Deep model access for safety researchers—offering privileged access that goes beyond public APIs, such as employee-level permissions, early access to unreleased models, safety-mitigation-free versions for testing, fine-tuning rights on frontier models, and allocated compute resources.

Why This Matters

External safety researchers often lack the access or funding to do the most valuable work they can. Companies committed to ecosystem-wide safety progress should enable the research community by providing deeper

access to frontier models, mentoring the next generation of research talent, and empowering funding-constrained external researchers. Deep model access enables critical research into the true model capabilities, alignment properties, and internal workings. Company-provided compute resources and API credits can enable academics and independent researchers with limited financial resources to experiment on frontier models.

Providing Deep Model Access to Safety Researchers

Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Al Safety researcher Ryan Greenblatt from Redwood Research was recently given employee-level access. [LessWrong, 2023]	Frontier model weights are publicly available	Non-frontier model Gemma 3 model weights publicly available [Google, 2025]	Frontier model weights are publicly available	OpenAI offers a public RL fine-tuning API. [OpenAI]	Non-frontier model Grok-1 model weights are publicly available [xAl, 2024]	Frontier model weights are publicly available

Sources:

- Shevlane, Toby, et al. "Model Evaluation for Extreme Risks." arXiv, 24 May 2023
- Casper, Stephen, et al. "Black-Box Access Is Insufficient for Rigorous AI Audits." arXiv, 29 May 2024

Mentoring and Funding

Anthropic	Mentoring: They have their own Anthropic Fellows program and provide a high number of mentors for the independent research seminar program MATS. [Anthropic, 2024; MATS, 2025] External Researcher Access Program (ongoing): • gives free API credit to safety/alignment researchers • Standard usage policies apply • \$1000 in API Credits (sometimes more) [Anthropic, 2025] Initiative for developing third-party model evaluations (Jul 2024): One-off program to provide funding for a third-party to develop evaluations that can effectively measure advanced capabilities in AI models: "The approach is designed to enable you to distribute your evaluations to governments, researchers, and labs focused on AI safety." [Anthropic, 2024].
DeepSeek	None found
Google DeepMind	Mentoring: Provides a high number of mentors for the independent research seminar program MATS. [MATS, 2025; MATS]
Meta	None found
OpenAI	Mentoring: Currently provides one mentor for the independent research seminar program, MATS. Researcher Access Program (back since February 2025): • gives free API credit to safety/alignment researchers • Standard usage policies apply • Up to \$1,000 of API credits [OpenAI, 2025] Superalignment Fast Grants (2023): \$10M to support technical research towards the alignment and safety of superhuman AI systems, including weak-to-strong generalization, interpretability, scalable oversight, and more [OpenAI, 2023].
x.AI	None found
Zhipu Al	None found

Sources

- "Shevlane, Toby, et al. "Model Evaluation for Extreme Risks." arXiv, 24 May 2023
- Casper, Stephen, et al. "Black-Box Access Is Insufficient for Rigorous Al Audits." arXiv, 29 May 2024.

ᅮ	\sim		BY PANFI	LICTO

Grading

Please pick a grade for each firm. You can add brief justifications to your grades.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Grades	Grade	Grade	Grade	Grade	Grade	Grade	Grade
Grade comments (Justifications, opportunities for improvements, etc.)							

Grading Scales

Grading scales are provided to support consistency between reviewers.

- A Strategy provides strong quantitative guarantees against catastrophic risks from superintelligent Al
- B Strategy very likely to prevent catastrophic risks from superintelligent Al
- Strategy likely to prevent catastrophic risks from superintelligent AI
- Strategy may prevent catastrophic risks from superintelligent AI
- E Strategy unlikely to prevent catastrophic risks from superintelligent AI, or Strategy increases risk

Domain comments

Optional: Share observations that apply across companies, including general recommendations, notes on how you weighted indicators, or feedback on FLI's methodology.



Domain

Governance & Accountability

This domain audits whether each company's governance structure and day-to-day operations prioritize meaningful accountability for the real-world impacts of its AI systems. Indicators examine whistleblowing systems, legal structures, and advocacy on AI regulations.

Table of Contents

Lobbying on AI Safety Regulations

Company Structure & Mandate

Whistleblowing Protection

Whistleblowing Policy Transparency

Whistleblowing Policy Quality Analysis

Reporting Culture & Whistleblowing Track Record

Grading

Indicator

Lobbying on AI Safety Regulations

Definition & Scope

This indicator documents a company's efforts to influence laws and regulations relevant to AI safety. It compiles publicly available evidence on direct policy positions—such as written statements, consultation responses, testimony, blog posts, and reputable media coverage—and records indirect engagement through membership in trade associations or coalitions that lobby on key safety rules.

Why This Matters

Leading AI developers have unique technical expertise and credibility to advise governments on charting a responsible path for this transformative technology. Tracking patterns in companies' engagements on specific regulations can indicate which firms take a proactive stance on raising the bar for sensible protections.



Sub-Indicator	Anthropic	DeepSeek	Google	Meta	OpenAl	x.AI	Zhipu Al
EU AI Act	No publicly available information found	No publicly available information found	Between 2022 and 2023, DeepMind lobbied EU institutions not to classify general-purpose AI and foundational models as "highrisk" technologies, a designation that would have triggered stricter safety obligations [Tranberg, 2023; TIME, 2023]. Google argued that the classification would hinder innovation, and regulations should attach further down the value chain [POLITICO, 2025; Data Ethics, 2023].	Between 2022 and 2023, Meta lobbied EU institutions to limit safety rules in the AI Act, opposing strict obligations for general-purpose models and seeking exemptions for open-source systems [Open Letter, 2023]. The company argued that strict obligations could hinder innovation and pushed for open-source models to be excluded from high-risk classification [Politico, 2025]. Chief AI Scientist Yann LeCun also criticized the EU's approach as overly restrictive [X, 2023].	In 2023, OpenAI lobbied EU officials to weaken parts of the AI Act, arguing that foundation models like GPT-4 should not face strict obligations unless adapted for specific uses [TIME, 2023; The Verge, 2023]. The company also pushed to delay transparency requirements and limit liability for general-purpose models.	No publicly available information found	No publicly available information found
Comprehensive	California's SB 1047	No publicly available	California's SB 1047	California's SB 1047	California's SB 1047	California's SB 1047	No publicly available
US State-Level Regulation US	In 2024, Anthropic initially raised concerns about California's SB 1047, influencing changes to the bill that softened key provisions [TechCrunch, 2024]. While the company opposed aspects of the original text, CEO Dario Amodei later expressed cautious support, stating in a letter to the governor that the bill's "benefits likely outweigh its costs" [Sanity.io, 2024]. Anthropic's involvement shaped the final version of the legislation [Vox, 2024].	information found	In 2024, Google DeepMind opposed California's SB 1047, arguing that its safety rules would burden developers and stifle innovation. The company warned that requirements like pre-deployment evaluations and state oversight could fragment regulation and urged alignment with federal efforts instead [DocumentCloud, 2024; Carnegie Endowment, 2024]. Responsible AI Safety and Education (RAISE) Act In 2025, industry groups with ties to Google DeepMind—including Tech:NYC and the Computer & Communications Industry Association (CCIA)—opposed New York's Responsible AI Safety and Education (RAISE) Act. They argued the legislation could conflict with federal policy and impose overly broad restrictions on AI development [Gothamist, 2025]. Both groups urged Governor Hochul to veto the bill, warning it could hamper innovation and create regulatory fragmentation [CCIA, 2025].	In 2024, Meta lobbied against California's SB 1047, arguing that its Al safety requirements— especially pre-deployment risk assessments and licensing—were overly broad and could hinder innovation [DocumentCloud, 2024; TechCrunch, 2024]. Alongside other tech firms, Meta urged lawmakers to adopt more flexible, federally aligned policies [Carnegie Endowment, 2024]. Responsible Al Safety and Education (RAISE) Act In 2025, Meta opposed New York's Responsible Al Safety and Education (RAISE) Act through multiple affiliated groups. Tech:NYC, a trade group co-founded by Meta, warned the bill could restrict innovation and conflict with federal policy [Gothamist, 2025]. The Al Alliance also sent a letter to state leaders opposing the bill's scope and regulatory approach [Al Alliance, 2025]. The Computer & Communications Industry Association (CCIA), whose members include Meta, urged Governor Hochul to veto the legislation [CCIA, 2025].	In 2024, OpenAI opposed California's SB 1047, arguing that its safety requirements—such as third-party evaluations and incident reporting—would hinder innovation and disadvantage U.S. firms [DocumentCloud, 2024; Carnegie Endowment, 2024]. The company also argued that the bill could raise national security risks by driving advanced research abroad [The Verge, 2024; Financial Times, 2024].	In 2024, xAI CEO Elon Musk publicly supported the bill in an X post, stating: "This is a tough call and will make some people upset, but, all things considered, I think California should probably pass the SB 1047 AI safety bill. For over 20 years, I have been an advocate for AI regulation, just as we regulate any product/ technology that is a potential risk to the public." [Tech Crunch, 2024].	information found

Sub-Indicator	Anthropic	DeepSeek	Google	Meta	OpenAl	x.Al	Zhipu Al
Preemption of state-level AI legislation	In 2025, Anthropic opposed federal efforts to preempt state-level AI laws. CEO Dario Amodei argued that states should retain authority to set transparency and safety standards, warning that federal preemption could weaken oversight [New York Times, 2025]. The company also lobbied against the Trumpbacked "Big Beautiful Bill," which aimed to override state AI regulation [WinBuzzer, 2025; Semafor, 2025].	No publicly available information found	In 2025, Google DeepMind supported federal preemption of state Al laws, urging a unified national framework to avoid regulatory fragmentation. In its response to the U.S. Al Action Plan, it called for federal leadership over issues like copyright, export controls, and development standards, warning that state-level rules could hinder innovation [Google Policy Response, 2025; TechCrunch, 2025].	In 2025, Meta supported federal preemption of state-level AI regulations, warning that fragmented laws could create compliance challenges and hinder innovation across jurisdictions [Meta, 2025]. The company's position aligned with broader industry efforts to shift AI governance to the federal level, drawing criticism from digital rights groups who argued this would weaken stronger state protections [X, 2025].	In 2025, OpenAl supported federal preemption of state-level Al laws, arguing that a unified national framework would better promote innovation and avoid regulatory fragmentation [OpenAl, 2025]. The company expressed concern that inconsistent state regulations could impose conflicting requirements and slow progress in the field [Bloomberg Law, 2025; Masood, 2025].	No publicly available information found	No publicly available information found

Indicator

Company Structure & Mandate

Definition & Scope

This indicator evaluates whether a company's fundamental legal structure, ownership model, and fiduciary obligations enable safety prioritization over short-term financial pressures in high-stakes situations. We report any embedded durable commitments to safety, social welfare, and benefit sharing and focus on any legally binding mechanisms (e.g., PBC status, capped equity, empowered governance bodies) that constrain management or shareholder incentives.

Why This Matters

Structural governance commitments can influence how companies respond when safety considerations conflict with profit incentives. During competitive pressures or deployment races, traditional for-profit structures may legally compel management to prioritize shareholder returns even when activities may pose significant societal risks. Structural governance innovations that formally embed safety into fiduciary duties—such as Public Benefit Corporation status or capped-profit models—create legally binding constraints that can override short-term financial pressures.



Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Uncommon governance structure. Finetuned for the ability to handle extreme events with humanity's interests in mind. Delaware Public Benefit Corporation (PBC) with a public benefit purpose. Anthropic's Purpose: "responsible development and maintenance of advanced AI for the long-term benefit of humanity." The Long-Term Benefit Trust (LTBT) is an independent body of five financially disinterested members, with the same purpose as PBC. It has the authority to select and remove a growing portion of the board of directors (ultimately the majority of the board) within 4 years, phasing in according to time- and funding-based milestones [Anthropic, 2023]. This is meant to ensure board decisions can prioritize long-term safety and public benefit over short-term commercial pressures when making high-stakes decisions about transformative AI. The Trust also has "protective provisions" requiring notice of actions that could significantly alter the corporation or its business. The structure is explicitly experimental, with "failsafe" provisions allowing changes through increasing supermajorities of stockholders as the Trust's power phases in. New Trustees are selected by existing Trustees, in consultation with Anthropic, and have no financial stake in Anthropic. The firm publicly announces new members [Anthropic 2025].	For-profit company	Part of Google, a public for-profit company	Public for-profit company	Uncommon governance structure. Founded as a Nonprofit, as founders initially believed a 501(c)(3) would be the most effective vehicle to direct the development of safe and broadly beneficial AGI while remaining unencumbered by profit incentives. Later incorporated a for-profit subsidiary (capped profit) to raise funds. For-profit is legally bound to pursue the Nonprofit's mission. Mission of OpenAI: "To ensure that artificial general intelligence (AGI) benefits all of humanity. We will attempt to directly build safe and beneficial AGI, but will also consider our mission fulfilled if our work aids others to achieve this outcome." The for-profit arm has a capped equity structure that limits maximum financial returns to investors and employees to balance profit incentives with safety concerns. Residual value will be returned to the Nonprofit. The size of the cap is not transparent. Charter contains an 'assist clause' to stop competing and assist a value-aligned, safety-conscious project to avoid race dynamics in late-stage AGI development [OpenAI] Conversion plans: In December 2024, OpenAI proposed a restructuring plan to convert the capped-profit into a Delaware-based public benefit corporation (PBC) and to release it from the control of the nonprofit. The nonprofit would sell its control and other assets, getting equity in return, and would use it to fund and pursue separate charitable projects. OpenAI's leadership described the change as necessary to secure additional investments. The plans provoked outside resistance and criticism. For example, a legal letter named "Not For Private Gain" [Not for Private Gain, 2025] asked the attorneys general of California and Delaware to intervene, stating that the restructuring is illegal and arguing that it would remove governance safeguards from the nonprofit's board chairman announced that the nonprofit would renounce plans to cede control after outside pressure. The capped-profit still plans to transition to a PBC, which critics said would diminish the nonprofit's co	Filed as a Nevada for-profit benefit corporation. Definition by Secretary of State: "for-profit entities that consider the society and environment in addition to fiduciary goals in their decision-making process, differing from traditional corporations in their purpose, accountability, and transparency." Registered purpose: "to advance human scientific discovery and deepen understanding of the universe." Nevada gives the state attorney-general independent standing to sue a public-benefit corporation that drifts from its mission, while Delaware does not [The Information; The Review Stories, 2025; NVSOS, 2014].	For-profit company

· Hendrycks, Dan. Introduction to Al safety, ethics, and society. Taylor & Francis, 2025. Section 8.4: Corporate Governance

Whistleblowing Protections

Indicator

Whistleblowing Policy Transparency

Definition & Scope

This indicator measures how fully and how accessibly an AI developer discloses its whistleblowing (WB) policy and system to the outside world. We look for a publicly reachable document (no paywall or login) that contains the material scope of reportable concerns, the people protected, the reporting channels offered (including anonymous options), oversight of the process, and the investigation and anti-retaliation guarantees. Evidence consists of artefacts that any external party can view, including public policy PDFs, dedicated "raise-a-concern" portals, relevant parts of safety frameworks, and transparency reports summarizing WB usage, outcomes, and effectiveness metrics.

Transparency Tiers:

- 1. No transparency
- 2. Fragments public: Parts of the design of the whistleblowing policy are public
- 3. Full policy public: Full policy, incl. processes, is public and highly transparent
 - a. Full policy public + all details accessible: Policy does NOT refer to internal policies that are inaccessible
 to the public, but outside parties can fully review policy details (within reason)
 - b. Effectiveness & Outcome transparency: The company provides details on the number of reports, topics, and follow-up actions, and also effectiveness, e.g., awareness & trust among employees, % of anonymous reports, appeal rates, whistleblower satisfaction, and types of cases received.

Why This Matters

Transparency on whistleblowing policies allows outsiders to assess the robustness of a firm's whistleblowing function. In AI safety contexts—where employees may be the first to spot concerning model behaviour or negligent risk management—robust, visible policies are critical. Public posting subjects the company to scrutiny by regulators, journalists, and prospective staff for both the policy's quality and the firm's adherence to it. Private policies, on the other hand, can hide restrictive terms. Many large companies demonstrate high levels of transparency around internal whistleblowing systems (e.g., Microsoft, Volkswagen, Siemens), including by publishing annual whistleblowing statistics.

Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Fragments Public Voluntary commitment in the safety framework [Anthropic, 2025], and comment on implementation status [Anthropic, 2024].	No transparency	Fragments public Employees are covered by Alphabet's group-wide code of conduct [Alphabet, 2024].	Fragments public WB policy referenced in Code of Conduct [Meta, 2024], integrity line available [Meta], and harassment policy is public in full [Meta].	Full Policy Public Details shared via FLI AI Safety Index Survey [Response] (16/16 relevant questions answered) Full Raising concerns policy public [OpenAI, 2024]; Integrity line available [OpenAI].	Details shared via FLI AI Safety Index Survey [Response] (16/16 relevant questions answered)	No transparency

Sources;

- Bullock, Charlie et al. "Protecting AI Whistleblowers." Lawafe. 2025. Accessed 1 July 2025.
- Wilson, Claudia et al. Whistleblower Protections for AI Employees. Center for AI Policy, 19 Jun. 2025,



Indicator

Whistleblowing Policy Quality Analysis

Definition & Scope

This analysis evaluates the quality of companies' whistleblowing policies based on all available evidence. The assessment analyzes 29 sub-indicators across five critical dimensions:

- 1) reporting channels and access,
- 2) whistleblower protections,
- 3) investigation processes,
- 4) system governance, and
- 5) Al-specific provisions.

Sub-indicators were derived from international reference standards—ISO 37002:2021, the ICC Guidelines, and the EU Whistleblowing Directive 2019/1937, which establish the gold standard for evaluation. Additional AI-specific items were included to address AI-specific concerns. For each Item, FLI evaluated the available evidence listed in the Whistleblowing Policy Transparency' indicator and rated the degree to which a company's policy satisfies it on a scale from 0 to 10, based on the publicly available information listed in the indicator on whistleblowing policy transparency, which includes whistleblowing policies, codes of conduct, safety frameworks, and survey responses. Where no information was available, 0 points were assigned. The assessment measures how well firms' policies align with best practices while specifically examining whether companies have implemented specialized AI safety provisions, such as protections for reporting violations of safety frameworks.

Why This Matters

Al development's technical complexity and commercial pressures create unique risks that only insiders can identify, but safety culture needs to be prioritized. Robust whistleblowing policies with Al-specific protections serve as a critical last line of defense when internal incentives fail, enabling employees to report concerning behaviors, intentional deception, or capability discoveries that could pose catastrophic risks. Without robust protections, adequate coverage, and secure channels, companies can quietly abandon safety commitments while those best positioned to prevent harm remain silenced.



Title	Description	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al	ISO 37002 "gold standard"
Overall average		2.44	0	1.54	2.44	3.56	2.32	0	8.35
Reporting Channels, Access, and Coverage		2.6	0	4.9	5.1	6.3	1.8	0	9.5
1.1 Protected Persons Coverage	Policy should at least cover current and former employees, contractors, shareholders, suppliers, former/prospective employees, and facilitators of reports	2	0	3	2	3	2	0	10
1.2 Policy Accessibility	Policy is easily accessible to all covered persons	0	0	2	8	8	0	0	10
1.3 External Reporting Information & Rights	Policy must provide clear information about external reporting channels and the right to approach these independently of internal processes, and explain or at least link to whistleblower protection rights	0	0	5	5	7	3	0	9
1.4 Multiple Reporting Channels	Offer multiple channels for reporting misconduct internally, incl. written, oral, and in-person	5	0	9	9	7	2	0	10
1.5 Anonymous Two-Way Reporting	The system enables fully anonymous reporting with secure two-way communication between the reporter and the investigators	4	0	5	9	10	0	0	9
1.6 Ombudsperson Channel	The reporting channel is operated by an outsourced whistleblowing service provider.	0	0	0	0	3	0	0	10
1.7 Executive Oversight Channel	A separate reporting channel is available for reports concerning senior executives (e.g., direct reporting line to the board audit committee) or board members	7	0	10	3	5	0	0	8
1.8 Broad but clear material scope	Material scope covers, at a minimum, potential violations of law and, code of conduct. Ideally, also further, broad categories, while retaining a high degree of clarity of what is in and out of scope.	3	0	5	5	7	7	0	10
2. Whistleblower Protections & Anti- Retaliation Measures		1.3	0	1.3	2.9	4	3.4	0	8.3
2.1 Confidentiality Protection	Strict protection is required for the reporter's identity and any third parties mentioned in reports	0	0	2	10	8	0	0	10
2.2 Public Disclosure Protection	Protection for responsible media disclosure if internal and regulatory channels have failed or if there is an imminent or manifest danger to the public interest	0	0	0	0	0	0	0	3
2.3 List of Prohibited Practices and Anti- Retaliation Provisions	Policy must list comprehensive prohibited retaliatory actions with specific examples (demotion, harassment, termination, etc.), and explicit anti-retaliation provisions	2	0	2	2	2	5	4	0
2.4 Post-Investigation Monitoring	Active monitoring for retaliation continues for a minimum of 12 months after the investigation concludes	0	0	0	0	0	0	0	8
2.5 NDA/Non- Disparagement Exceptions	Explicit statement that NDAs and non-disparagement agreements cannot prevent safety-related whistleblowing	7	0	0	0	7	10	0	10



INSTITUTE									
Title	Description	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al	ISO 37002 "gold standard"
2.6 Good Faith or Reasonable Cause Provisions	Clear good faith or reasonable cause standard that protects honest mistakes; high burden of proof required for false report sanctions	0	0	5	8	8	10	0	10
2.7 Handler/Investigator Protection	Explicit protections for employees who receive, investigate, or support whistleblowing reports	0	0	0	0	0	0	0	8
3. Investigation Process & Standards		0.8	0	1	3.2	2.2	0.4	0	7.6
3.1 Designated Impartial Receiver	A provably independent person or department must be designated to receive and handle reports, ideally attached to the board	4	0	5	6	6	2	0	9
3.2 Seven-Day Acknowledgment	Written confirmation of report receipt must be provided within 7 days	0	0	0	10	0	0	0	10
3.3 Three-Month Feedback Timeline	Investigation status and follow-up measures must be communicated to the reporter within 3 months	0	0	0	0	0	0	0	6
3.4 Adequately Resourced Investigation Teams	Investigators must be independent from implicated departments and possess appropriate technical expertise for AI safety issues, as well as sufficient resources to investigate effectively	0	0	0	0	5	0	0	9
3.5 Investigation Appeal Process	Formal right to appeal investigation outcomes to an independent review body or board committee	0	0	0	0	0	0	0	4
4. System Governance & Quality Assurance		3	0	0	1	1	0	0	8
4.1 Comprehensive Effectiveness Metrics	Regular measurement tracking report outcomes, investigation timeliness, appeal rates, % of anonymous reports, retaliation incidents, and reporter satisfaction - not just volume	7	0	0	0	0	0	0	10
4.2 Data Retention and Deletion Policy	Clear policy specifying retention periods for reports and investigations (typically 5-7 years), secure deletion procedures, and data minimization principles	0	0	0	0	0	0	0	8
4.3 Secure Documentation System	Comprehensive audit trail with secure case management system and defined retention policies	0	0	0	5	5	0	0	9
4.4 Comprehensive Training Programs	Regular, role-specific training is provided for all employees, specialized training for managers and investigators, ideally measuring training effectiveness.	0	0	0	0	0	0	0	10
4.5 Independent System Certification	Regular third-party audit and certification of the whistleblowing system's effectiveness and compliance	8	0	0	0	0	0	0	3
5. Al Safety-Specific Provisions		4.5	0	0.5	0	4.3	6	0	N/A
5.1 Al Safety Commitment Protection	Explicit protection for reporting violations of frontier safety frameworks (e.g., RSP, Preparedness Frameworks), public AI safety commitments, and internal safety policies	8	0	0	0	5	5	0	
5.2 Al Safety Coordination	Protection for AI risk reporting to dedicated AI safety bodies (UK AI Security Institutes, US Center for AI Standards and Innovation, or other international regulatory bodies)	0	0	0	0	2	2	0	

Title	Description	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al	ISO 37002 "gold standard"
5.3 Al risk transparency	Protections for reporting intentional deception of external evaluators, regulators, or the public, suppression of publication of safety evaluation results, and inadequate disclosure of risk to regulators and the public.	5	0	0	0	5	10	0	
5.4 Adequacy of AI risk management and cybersecurity	Protections for reporting inadequate risk management processes, incl. assessment, monitoring, mitigation, deployment pressure despite concerning levels of risk, insufficient operational and cybersecurity practices, incl. incidents	5	0	2	0	5	7	0	

Sources

- European Parliament and Council of the European Union. <u>Directive (EU) 2019/1937 of the European Parliament and of the Council of 23 October 2019 on the Protection of Persons Who Report Breaches of Union Law.</u> Official Journal of the European Union, 26 Nov. 2019
- International Organization for Standardization. ISO 37002:2021 Whistleblowing Management Systems Guidelines. ISO, 2021
- Nowers, Ida., and Terracol, Marie. "Monitoring Internal Whistleblowing Systems A framework for collecting data and reporting on performance and impact". 2025

Indicator

Reporting Culture & Whistleblowing Track Record

Definition & Scope

This indicator evaluates whether an AI developer fosters a climate in which employees can raise safety-relevant concerns without fear of retaliation and with confidence that the concerns will be addressed. Evidence is drawn from (i) the organisation's track-record of documented whistleblowing cases, (ii) the use, scope, and enforcement of non-disclosure or non-disparagement agreements (NDAs), (iii) leadership signals that encourage or discourage internal dissent, (iv) third-party evidence of psychological safety, and (v) patterns of safety information leaking externally (vi) departures linked to safety governance. The focus is on demonstrated behaviour and outcomes rather than written policy statements. For whistleblowing incidents, we report individual names, concerns raised, and company response & status where available.

Notes of Best Practice: Companies should show a clear recent pattern of protecting and acting on employee safety reports; public commitment not to enforce legacy NDAs for safety topics; leadership statements praising internal critics; ≥ one anonymized psychological-safety survey with ≥ 70 % of staff agreeing "I can raise safety concerns without fear" and no credible retaliation cases in the last 24 months. Little public leaks as issues are addressed internally. Recent evidence (≤ 24 months) should be weighted twice as heavily as older cases to reward reforms.

Why This Matters

Whistleblowing policies can look impressive on paper, but they fail if the climate in the company suppresses reports, they're not effective when employees fear retaliation, or doubt anyone will act. This is why scrutinizing how firms respond to disclosures is critical. By focusing on actual cases, NDA practices, leadership signals, and exits tied to safety concerns, this indicator reveals which firms have built cultures where raising concerns feels like following protocol rather than betraying the company or colleagues—the trust and accountability needed for early detection of catastrophic Al risks.



Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Statement on non- disparagement agreements (June 2024): Cofounder Sam McCandlish announced that the firm has been using standard non-disparagement agreements in severance agreements, but now considers this practice to be in conflict with their mission and has started removing them. Stated that former employees who signed a non-disparagement agreement are free to state that fact, raise concerns about safety at Anthropic, and that Anthropic would not enforce non-disparagement agreements in such cases [LessWrong].	None	(Note: covers all of Google, not only DeepMind) Satrajit Chatterjee (2022- ongoing): Engineering manager fired in March 2022 after challenging a paper published by Google about Al chip design capabilities. A California state judge in July 2023 rejected Google's request to dismiss his wrongful termination and whistleblower protection claims [Bloomberg, 2023]. Chatterjee alleged that Google terminated him in retaliation for refusing to participate in what he viewed as misrepresentation of the company's Al technology capabilities, potentially defrauding shareholders and the public. Google stated the allegations were 'academic disputes' and defended the paper's scientific merit [The Star, 2023]. Shareholder motion for stronger protections (2021): Trillium Asset Management (Alphabet shareholder) filed a resolution calling for expanded whistleblower protections for Google employees. The resolution requested that Alphabet's Board of Directors oversee a third-party review of current whistleblower policies, citing the importance of strong protections for employees who raise concerns. *Outcome*: Alphabet's board recommended "AGAINST"; at the 2 June 2021 AGM, only **10% of total votes** (≈ 63.8 m) supported the motion [SEC, 2021]. Margaret Mitchell (2021): Co-lead of the ethical Al team, ran scripts to archive emails documenting the handling of Gebru's case. Fired for "exfiltrating thousands of files" in breach of security policy, according to Google [The Verge, 2021; MIT Technology Review, 2020].	(Note: covers all of Meta, not only Llama teams) Sarah Wynn-Williams (2025): Former global public policy director. Published a memoir and testified to Congress about Meta's alleged cooperation with China's government and misleading of lawmakers. Meta invoked a 2017 severance non-disparagement clause, won an emergency arbitration order potentially temporarily barring "disparaging" statements and blocking meetings with US/UK/EU legislators. Wynn-Williams testified before the Senate. Meta told TechCrunch that the arbitration order does not prohibit her from speaking to Congress and that the company does not intend to interfere with her legal rights [TechCrunch, 2025; CNN, 2025]. In Apr 2025, Sen. Grassley's letter to Meta demanding answers on NDAs allegedly silencing whistleblowers [Grassley, 2025]. Internal memo threatens termination for leaks (Jan 31, 2025): After CEO Mark Zuckerberg complained that "everything I say leaks," Meta CISO Guy Rosen circulated an internal memo warning that "we will take appropriate action, including termination," against employees who leak confidential information. The memo confirmed Meta had "recently terminated relationships with employees who leaked confidential company information." (The warning memo itself was subsequently leaked to the press.) [The Verge, 2025; Fortune, 2025]. NLRB Ruling (2024): The NLRB judge ruled that Meta's separation agreements used during the 2022 mass layoffs were illegal. The agreements affected over 7,000 employees who were required to sign "unlawfully overbroad" non-disparagement and confidentiality clauses in exchange for enhanced severance pay. This followed the precedent set by the McLaren Macomb decision in February 2023 [The Register, 2024; HRD, 2024]. 20 Employees terminated for leaks of internal meeting details; Meta said more firings may follow [TechCrunch, 2025].	Dissolution of AGI readiness team (Oct 2024): Team leader Miles Brundage left the firm. As part of a broader farewell message shared that some of his colleagues seem to think "that speaking up has big costs, and that only some people are able to do so.", but that he "think[s] people almost always assume that it's harder/more costly to raise concerns or ask questions than it actually is."[X, 2024]. The AGI readiness team was then dissolved within OpenAI [CNBC, 2024]. Brundage's exit then spurred former team member Rosie Campbell to depart [The Byte, 2024]. Anonymous whistleblowers (Jul 2024): Letter & formal SEC complaint allege OpenAl's NDAs illegally bar staff from alerting regulators and waive Dodd-Frank rewards. SEC matter pending [Tech Crunch, 2024; The Hill, 2024]. Right-to-Warn" open letter – 11 current & former staff, plus peers at other firms (Jun 2024): Called for an enforceable right to disclose Al-risk concerns without retaliation or broad NDAs. All current employees chose to remain anonymous. The letter cites criticism of OpenAl's equity claw-back clause [Right to Warn, 2024]. Jan Leike (May 2024): Superalignment Co-lead. Resigned, stating "safety culture and processes have taken a backseat to shiny products," and that the team lacked compute [X, 2024]. Co-team lead Sutskever left simultaneously to start competitor firm 'Safe Superintelligence! The superalignment team was then disbanded [X, 2024]. Policy researcher Gretchen Krueger announces her departure hours later, stating similar concerns [The Verge, 2024]. Exit-agreement overhaul after media leak (May 2024): Vox revealed strict severance papers that let OpenAl *cancel vested equity* if ex-staff "disparaged" the company. After a media exposé, OpenAl **removed the clauses** and said it had never clawed back equity; CEO Sam Altman apologized, though leaked paperwork later showed his prior sign-off on the wording [Vox, 2024; The Verge, 2024]. Safety researcher Todor Markov left the company over the issue, arguing the debacle incident proved Al	None	None



Anthropic DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
	Timnit Gebru (2020): Co-lead of Google's ethical AI team. Objected to Google's demand to retract a paper outlining large-language-model risks (bias, emissions). Google says she "resigned"; Gebru says she was **terminated**. Incident provoked >2,000-employee petition [ABC News, 2020; MIT Technology Review, 2020; Time, 2022]. Mustafa Suleyman (2019, precedent for accountability): Internal probe found patterns of workplace bullying. **Placed on administrative leave** July 2019; later moved to Google product role and eventually left Alphabet in 2022 [TechBrew, 2021].	Arturo Bejar (2023): Former engineering director. Testified before Congress that leadership, including. Mark Zuckerberg ignored evidence that Instagram harms teens (bullying, self-harm). He had emailed Zuckerberg, Sandberg, and other executives in 2021 with research showing harmful effects on young users. Meta stated it "does not agree" with Bejar's characterisation and highlighted existing safety tools; no legal action has followed [CNBC, 2023; NPR, 2023; OPB, 2023]. Frances Haugen (2021): Former product manager. Supplied thousands of internal files ("Facebook Papers") to the SEC & US Congress, alleging Meta misled the public about known harms (teen mental health, misinformation). Meta said documents were "cherry-picked" and Haugen's claims "mischaracterise" its work. No litigation between parties. Haugen testified before the Senate Oct 2021 [CBS, 2021] Sophie Zhang (2020-21): Data scientist. Farewell memo described government-backed political manipulation campaigns across 25 countries on Facebook; reposted the memo on a password-protected personal site. Facebook deleted her internal post and requested her web-host & registrar remove the external site, which they did. Zhang declined a severance agreement containing a non-disparagement clause and later testified before the British Parliament and provided documents to US law enforcement [Independent, 2021; BuzzFeed, 2020].	Daniel Kokotajlo (Apr 2024): Governance researcher resigned because he "Lost confidence [OpenAI] would behave responsibly around AGI" [Futurism, 2024]. William Saunders (Feb 2024): Interpretability engineer resigned, telling *Business Insider* leadership "was not adequately addressing" catastrophic-risk issues; later co-signed the "Right-to-Warn" letter [Business Insider, 2024]. Board coup & reversal (Nov 2023): Board unexpectedly removed CEO, stating he "was not consistently candid in his communications". Altman was reinstated within a week. A WilmerHale special review later found his conduct "did not mandate removal" [ARS Technica, 2023]. In the aftermath, three independent directors (including Helen Toner) resigned [Aljazeera, 2024]. Toner later stated, "For years, Sam had made it really difficult for the board to actually do that job by, you know, withholding information, misrepresenting things that were happening at the company, in some cases outright lying to the board.[]". [TED, 2024].		

٦	$\Gamma \cap$	RF	COMPL	FTFD	BY PANFI	LISTS

Grading

Please pick a grade for each firm. You can add brief justifications to your grades.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Grades	Grade	Grade	Grade	Grade	Grade	Grade	Grade
Grade comments (Justifications, opportunities for improvements, etc.)							

Grading Scales

Grading scales are provided to support consistency between reviewers.

- A Exemplary accountability ensures safety-focused decision-making at all levels
- B Strong accountability enables safety-focused decision-making
- Moderate accountability with gaps affecting safety decision-making
- Weak accountability hinders safety-focused decision-making
- No meaningful accountability

Domain comments

Optional: Share observations that apply across companies, including general recommendations, notes on how you weighted indicators, or feedback on FLI's methodology.

Domain

This section gauges how openly firms share information about products, risks, and risk management practices. Indicators cover voluntary cooperation, transparency on technical specifications, and risk/incident communication.

Table of Contents

Technical Specifications

System Prompt Transparency
Behavior Specification Transparency

Voluntary Cooperation

G7 Hiroshima AI Process Reporting
FLI AI Safety Index Survey Engagement

Risks & Incidents

Serious Incident Reporting & Government Notifications Extreme-Risk Transparency & Engagement

Grading

Technical Specifications

Indicator

System Prompt Transparency

Definition & Scope

This indicator reports whether companies publicly disclose the system prompts used in their most capable deployed AI models. Evidence includes published system prompts in model cards, technical documentation, or dedicated transparency pages, and changelogs. Best practice involves publishing exact prompts used in production, version history, verification that prompts are used in production, and explanations of design choices.

Why This Matters

System prompts fundamentally shape AI behavior and safety properties, yet most companies keep them secret. Publishing prompts enables researchers to verify to better understand the models and makes the company's intended behaviour transparent. High transparency can reflect a commitment to accountability.

Anthropic	Transparency In March 2024, Anthropic shared the full system prompt alongside the release of Claude 3 as a one-off [Fast Company, 2024]. Since August 2024, Anthropic has publicly shared the systems' prompts for the Claude.ai web interface and mobile apps since August 2024. Shared system prompts for six models, plus several updates. They further committed to logging changes they make they make to these prompts online. Shared systems prompts do NOT currently cover the API [TechCrunch, 2024; X, 2024; Anthropic]. Simon Willison reported that the publicly shared version does not include the description of various tools available to the model [Simon Willison, 2025].
DeepSeek	Frontier model weights are public, so the system prompt can be decided by user/hosting service. Their own hosted service does not disclose it.
Google DeepMind	No transparency on system prompts for Frontier Systems.
Meta	Frontier model weights are public, so the system prompt can be decided by the user/hosting service. Their own hosted service does not disclose it.
OpenAl	No transparency on system prompts for Frontier Systems.
x.AI	Transparency: Following the incident in May 2025 listed below, x.Al published their system prompts for Grok (on xAl & X) on Github and promised these will be regularly updated [Github, 2025]. Incidents: February 2024: A change to the system prompt, Grok briefly censored responses about Elon Musk and Donald Trump spreading disinformation. After the issue received public attention, xAl quickly reverted the changes and publicly stated that the problem was caused by an unnamed employee conducting unauthorized modifications [Fortune, 2024]. May 2025: After a change to the system prompt, Grok started randomly discussing whether there was a "white genocide" happening in South Africa in many completely unrelated conversations. The Al Chatbot told users it was 'instructed by my creators' to accept 'white genocide as real and racially motivated' [Guardian, 2025]. x.Al quickly apologized for this incident and rolled back the changes. They reported that unauthorized modifications by an employee caused the incident [X, 2025].
Zhipu Al	Frontier model weights are public, so the system prompt can be decided by the user/hosting service. Their own hosted service does not disclose it.

Sources

· Kokotajlo, Daniel, and Alexander, Scott. "Make The Prompt Public." AI Futures Project, 17 May 2025



Indicator

Behavior Specification Transparency

Definition & Scope

This indicator assesses whether companies publish detailed specifications outlining their models' intended behaviors, boundaries, and decision-making frameworks. For companies that shared such documents, we provide high-level summaries and link to the sources. We include documents that concretely outline the goals, values, and behavioral guidelines that developers aim to instill in their models. Documentation should explain how developers want their models to handle various scenarios, conflicts, and edge cases, and detail how these values are implemented, including metrics or evidence of how well these values are achieved in practice. Specifications should ideally be current and include a tracked version history with dates. Important aspects are specificity, comprehensiveness across use cases, and inclusion of concrete examples. Internal training documents, vague mission statements, and brief high-level descriptions are not in scope.

Why This Matters

Model specifications reveal how companies intend their AI systems to behave in complex situations, including safety-critical decisions. Publishing these specs enables external verification of whether deployed models match stated intentions and allows identification of gaps in safety considerations. Companies willing to specify and publish concrete behavioral guidelines demonstrate accountability for their choices and enable public scrutiny.

Of LITE	
INSTITUTE	
Anthropic	Constitution Method What's 1) Super 2) RLA Timelin Decemment May 20
	Consti 58 prin - UN D - Apple - Deep - Non-\ - Anthr Examp
	Benefi t Readak
	Version Anthro training Since t values
	Source
DeepSeek	No det
Google DeepMind	No det

itutional AI:

d for training AI systems to be harmless by using a set of written principles (a "constitution") rather than relying solely on large-scale human feedback.

s it for:

- ervised learning phase: Model self-critiques and revises its outputs based on constitutional principles, creating a supervised learning dataset
- NF phase: Model compares response pairs using constitutional principles to generate preference labels, then trains via RL on these Al-generated preferences

ne & Development:

nber 2022: Original Constitutional AI paper published

D23: Claude's constitution made public (58 principles)

itution (May 2023):

nciples (1.2k words) drawn from:

- Declaration of Human Rights
- e's Terms of Service
- Mind's Sparrow principles
- Western perspectives
- ropic's own research

ble principle: "Please choose the response that most supports and encourages freedom, equality, and a sense of brotherhood."

ts:

ble, transparent, and explicitly formulated principles, as opposed to RLHF, which leverages implicit values.

tions:

n uncertainty: Only the May 2023 constitution is public; the current production versions are unknown

pic uses a "variety of techniques including human feedback, Constitutional AI [..], and the training of selected character traits.". Given that other approaches are incorporated in postg, the impact of any one of them is unclear.

the AI itself determines how to balance competing constitutional principles, Anthropic's approach does not explicitly specify the intended behavior of its AI systems, especially when conflict.

e: [Anthropic, 2025]

tailed specification available, but frontier model weights are public, so models can be modified.

tailed specification available

Meta

No detailed specification available, but frontier model weights are public, so models can be modified.

OpenAl

OpenAl Model Spec:

OpenAl's Model Spec is a detailed (~28k words), public, living rule-book that defines the objectives, safety rules, and default behaviours OpenAl trains its models —via human feedback and deliberative alignment-to follow.

What's it for:

1) Human RLHF guidance - provides a single, public rule-book that labelers follow when creating preference data.

2) Deliberative Alignment - o-series models (o1, o3, o4-mini) are explicitly taught to read and reason over the Spec before answering.

3) Automated evaluation - OpenAl ships a challenge-prompt suite to measure adherence.

Timeline & Versions:

1st May 2024

2nd Feb 2025

3rd Apr 2025

Framework:

Three principal types:

1) Objectives - broad goals such as "assist the developer & end user" and "benefit humanity."

2) Rules - hard, platform-level constraints (e.g., comply with law, prohibit or restrict certain content, protect privacy, uphold fairness).

3) Defaults - stylistic and behavioural norms that developers/users may override.

Sections: Stay in bounds · Seek the truth together · Do the best work · Be approachable · Use appropriate style.

Includes specific guidance on specific policy areas such as potential, medical, or harmful content.

Risk taxonomy: Misaligned goals · Execution errors · Harmful instructions.

Chain of command:

Platform (OpenAI) \rightarrow Developer \rightarrow User \rightarrow Guideline \rightarrow Untrusted text.

Within any level, explicit > implicit, later > earlier.

(OpenAI's Usage Policy overrides the Spec if the two conflict.)

Ongoing Development:

Released under CC0 license (public domain)

Changelog and version history maintained on GitHub

OpenAI commits to regular updates as the spec evolves

Key Benefits

Greater transparency of intended model behavior.

Finer-grained steerability via the chain of command

Reduced reliance on implicit human values; models can show interpretable reasoning steps grounded in the Spec.

Transparency & Limitations

Production models don't fully reflect the spec yet.

OpenAI states: "While the public version of the Model Spec may not include every detail, it is fully consistent with our intended model behavior."

Source: [OpenAl, 2025]

x.AI

No detailed specification available

Zhipu Al

No detailed specification available, but frontier model weights are public, so models can be modified.

Sources

- Kokotajlo, Daniel, and Alexander, Scott. "Make The Prompt Public." AI Futures Project, 17 May 2025
- Ball, Dean. "4 Ways to Advance Transparency in Frontier Al Development." The Foundation for American Innovation, 16 Oct. 2024

Voluntary Cooperation

Indicator

G7 Hiroshima AI Process Reporting

Definition & Scope

The G7 Hiroshima AI Process (HAIP) Reporting Framework is a voluntary transparency mechanism launched in February 2025 for organizations developing advanced AI systems. Organizations complete a comprehensive questionnaire covering seven areas of AI safety and governance practices, including risk assessment, security measures, transparency reporting, and incident management. All submissions are published in full on the OECD transparency platform. This indicator tracks whether firms participated in HAIP as a measure of their commitment to AI safety transparency.

Why This Matters

The HAIP framework represents the first globally standardized mechanism for AI developers to disclose their safety practices in comparable detail. Participation creates reputational stakes and enables external scrutiny since reports are published. Organizations choosing to participate signal a willingness to be held accountable and contribute to collective learning.

Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Substantive submission [OECD AI, 2025]	(Not based in a G7 nation)	Substantive submission (Google) [OECD AI, 2025]	No Submission	Substantive submission [OECD AI, 2025]	No Submission	(Not based in a G7 nation)

Sources

- Perset, Karine, James Gealy, and Sara Fialho Esposito. "Shaping Trustworthy Al: Early Insights from the Hiroshima Al Process Reporting Framework." OECD. Al, 11 Jun. 2025
- Ministry of Internal Affairs and Communications, Japan. G7 Hiroshima Process on Generative Artificial Intelligence. Ministry of Internal Affairs and Communications, Japan
- Organisation for Economic Co-operation and Development (OECD). OECD.Al Policy Observatory: Reports. OECD

Indicator

FLI AI Safety Index Survey Engagement

Definition & Scope

We report which companies have engaged with our index survey to voluntarily disclose additional information. Full survey responses are linked below.

Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
None	None	None	None	Survey response submitted [PDF]	Survey response submitted [PDF]	Survey response submitted [PDF]

Risks & Incidents

Indicator

Serious Incident Reporting & Government Notifications

Definition & Scope

This indicator evaluates incident reporting commitments, frameworks, and track records. For frameworks and

commitments, the indicator assesses whether companies have publicly discussed any systems and commitments to share critical information about red-line incidents or capabilities with government bodies (e.g., US CAISI, UK AISI), peer organizations, or the public. Such incidents can include successful large-scale misuse, near-miss events, scheming by AI models, and identified model capabilities with severe national security implications. The indicator further tracks relevant incident documentations that the company has already shared. Evidence comes from safety frameworks, documented reporting procedures, participation in information-sharing agreements, and public incident reports.

Notes on Best Practice: Clear public commitments to report specific categories of incidents to government bodies, with documented procedures for incident classification and escalation. Information-sharing agreements with disclosed scope, publishing reports on recent incidents, demonstrating transparency about warning signs discovered during development, and establishing clear thresholds for mandatory reporting, specificity, and comprehensiveness of reporting commitments.

Why This Matters

Proactive incident reporting enables collective learning from safety failures and near-misses across the Al industry, preventing repeated mistakes and identifying emerging risks before they materialize. Transparency about dangerous capabilities and misalignment incidents is critical for government oversight. Without such transparency, companies may make deployment decisions based on marginal safety improvements while baseline risks remain unacceptably high.

Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	xAI	Zhipu Al
		Serio	us incident reporting framev	vorks		
	Chinese AI firms operate under several regulations with mandatory incident reporting requirements, often under short timeframes. We list applicable GenAI specific frameworks but not those focused on (data-/cyber-) security:					Chinese AI firms operate under several regulations with mandatory incident reporting requirements, often under short timeframes. We list applicable GenAI specific frameworks but not those focused on (data-/cyber-) security:
	- Interim Measures for Generative AI Services, Art. 14 (Aug 2023) – Gen- AI providers that detect illegal content or model misuse must "promptly" stop generation, rectify, and inform the competent authorities [China Law Translate, 2023]					- Interim Measures for Generative AI Services, Art. 14 (Aug 2023) – Gen- AI providers that detect illegal content or model misuse must "promptly" stop generation, rectify, and inform the competent authorities [China Law Translate, 2023]
	- Deep-Synthesis Provisions (Jan 2023) - Deep-fake service providers must remove illegal or harmful synthetic content, preserve records, and "timely" report the incident to the CAC and other competent departments [Cyberspace Administration of China, 2023]					- Deep-Synthesis Provisions (Jan 2023) - Deep-fake service providers must remove illegal or harmful synthetic content, preserve records, and "timely" report the incident to the CAC and other competent departments [Cyberspace Administration of China, 2023]
		Red-line G	overnment notifications con	nmitments		
Responsible Scaling Policy contains a broad voluntary commitment on ASL disclosing ASL levels: - "We will notify a relevant U.S. Government entity if a model requires stronger protections than the ASL- 2 Standard" [Anthropic, 2025].		Frontier Safety Framework 2.0 states that if a model reaches a "Critical Capability Level" posing unmitigated material risk, DeepMind "aims to share information with appropriate government authorities" and may also notify other external organisations [Google, 2025].				



Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	xAI	Zhipu Al
			Public transparency reports			
Anthropic published one comprehensive misuse report, which documents real-world cases of actors attempting to exploit Claude for malicious purposes, along with detection methods and enforcement actions taken. -Mar 2025 – "Misuse Monitoring and Response Report" [Anthropic, 2025]. -Platform Security transparency page provides some enforcement statistics, including banned accounts for Usage Policy violations, number of appeals processed, CSAM reports to NCMEC, and law enforcement requests [Anthropic, 2024].		Published a detailed report on how threat actors—from scammers to state-aligned groups—attempt to misuse Google Gemini in deception, persuasion, and cyber operations. Described mitigation strategies and detection tooling -Jan 2025 - 'Adversarial Misuse of Generative Al" [Google 2025].	Meta consistently issues quarterly integrity reports about its platforms [Meta, 2024], which include reports on disrupting adversarial threats such as influence operations [Meta, 2025]. No reports for frontier Al models are available.	Regular reports documenting their disruption of malicious uses of their AI systems. Comprehensive reports detail enforcement actions against state-affiliated threat actors and covert influence operations, identify specific threat groups (e.g., Storm-2035, Spamouflage), quantify disruptions (accounts banned, operations terminated), and describe the tactics employed (phishing, malware development, influence campaigns, election interference). - Feb 2024 - "Disrupting Malicious Uses of AI by State-Affiliated Threat Actors" [OpenAI, 2024] - May 2024 - "Disrupting a Covert Iranian Influence Operation" [OpenAI, 2024] - Jun 2024 - "Update on Disrupting Deceptive Uses of AI" [OpenAI, 2024] - Aug, 2024: "Disrupting a covert Iranian influence operation" [OpenAI, 2024] - Oct 2024 - "Influence and cyber operations: an update" [OpenAI, 2024] - Feb 2025 - "Disrupting malicious uses of our models" [OpenAI, 2025] - Jun 2025 - Disrupting malicious uses of AI [OpenAI, 2025]		



Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	xAI	Zhipu Al				
	Industry information sharing									
The Frontier Model Forum (FMF) announced an information-sharing agreement signed by member firms (incl. Anthropic, Google, Meta, and OpenAl) to facilitate the sharing of threats, vulnerabilities, and capability advances specific to frontier Al. The agreement, narrowly scoped to manage national security and public safety risks (including CBRN and advanced cyber threats)		The Frontier Model Forum (FMF) announced an information-sharing agreement signed by member firms (incl. Anthropic, Google, Meta, and OpenAI) to facilitate the sharing of threats, vulnerabilities, and capability advances specific to frontier AI. The agreement, narrowly scoped to manage national security and public safety risks (including CBRN and advanced cyber threats), covers three categories:	The Frontier Model Forum (FMF) announced an information-sharing agreement signed by member firms (incl. Anthropic, Google, Meta, and OpenAI) to facilitate the sharing of threats, vulnerabilities, and capability advances specific to frontier AI. The agreement, narrowly scoped to manage national security and public safety risks (including CBRN and advanced cyber threats), covers three categories:	The Frontier Model Forum (FMF) announced an information-sharing agreement signed by member firms (incl. Anthropic, Google, Meta, and OpenAI) to facilitate the sharing of threats, vulnerabilities, and capability advances specific to frontier AI. The agreement, narrowly scoped to manage national security and public safety risks (including CBRN and advanced cyber threats), covers three categories:		Zhipu AI is a founding member of the IIFAA "Trusted Agent Inter- connect Working Group" (Dec 2024) alongside Huawei, Alibaba, ByteDance, etc.; the group sets cross-agent security and data-sharing norms [China Daily, 2024].				
advanced cyber threats), covers three categories: (1) vulnerabilities and exploitable flaws that could compromise Al safety/security, (2) threats involving unauthorized access or manipulation of frontier models, and (3) capabilities of concern		(1) vulnerabilities and exploitable flaws that could compromise Al safety/security, (2) threats involving unauthorized access or manipulation of frontier models, and (3) capabilities of concern with potential for largescale societal harm.	(1) vulnerabilities and exploitable flaws that could compromise Al safety/security, (2) threats involving unauthorized access or manipulation of frontier models, and (3) capabilities of concern with potential for largescale societal harm.	(1) vulnerabilities and exploitable flaws that could compromise AI safety/security, (2) threats involving unauthorized access or manipulation of frontier models, and (3) capabilities of concern with potential for large-scale societal harm.						
with potential for large- scale societal harm. Details on implementation and use are unclear [Frontier Model Forum, 2025].		Details on implementation and use are unclear [Frontier Model Forum, 2025].	Details on implementation and use are unclear [Frontier Model Forum, 2025].	Details on implementation and use are unclear [Frontier Model Forum, 2025].						

Indicator

Extreme-Risk Transparency & Engagement

Definition & Scope

Assesses the extent to which companies and their leadership (A) publicly recognise the potential for catastrophic AI harm and (B) proactively disseminate evidence-based analyses of such risks to external stakeholders. The criteria are frequency, specificity, and prominence of communication about AI's potential for catastrophic outcomes (including existential risks, mass casualties, or societal-scale disruption).

Evidence includes official blogs, testimonies, leadership communications, including signed statements. Excludes technical safety papers, model cards, and formal safety frameworks (captured in separate indicators).

Note: The research methodology for this indicator did not follow a formal structure; results are incomplete and likely biased/skewed by the prominence of different media reports.

Why This Matters

Public communication about Al's potential for catastrophic outcomes shapes societal preparedness, policy responses, and research priorities. Companies developing frontier Al possess unmatched knowledge of actual capabilities, near-term developments, and observed warning signs. Their leadership's willingness to transparently discuss extreme risks indicates a precautionary approach and enables an informed discourse on policy and national security.

Anthropic

The company and its leaders regularly and proactively communicate extreme risks.

Examples from CEO Dario Amodei:

- Warns AI may eliminate 50% of entry-level white-collar jobs within the next five years [Business Insider, 2025] and says on television that he is "raising the alarm" about this [CNN, 2025].
- Blog post calling the Paris AI Action summit a "missed opportunity", saying ".. greater focus and urgency is needed on several topics given the pace at which the technology is progressing." [Anthropic, 2025].
- Warned Congress that AI could enable bioweapon creation within 2-3 years [Bloomberg, 2023].
- Repeatedly warns that 'powerful Al', which he likens to "a country of geniuses in a datacenter", could arrive as early as 2026 or 2027, and is explicit about extreme risks [Anthropic, 2025]: ".. hardcore misuse in Al autonomy that could be threats to the lives of millions of people. That is what Anthropic is mostly worried about." [Business Insider, 2025]

CAIS statement on AI Risk signed by: Dario Amodei (CEO), Daniela Amodei (President), Jared Kaplan (co-founder), Chris Olah (co-founder)

Relevant blogs by Anthropic are below. Several share quantitative evidence related to extreme risks:

- Progress from our Frontier Red Team [Anthropic, 2025]
- Third-party testing as a key ingredient of AI policy [Anthropic, 2024]
- Reflections on responsible scaling policy [Anthropic, 2024]
- The case for targeted regulation [Anthropic, 2024]
- Frontier Threats Red Teaming for AI Safety [Anthropic, 2023]

DeepSeek

The company and its leadership do not discuss extreme risks from AI.

CEO Liang Wenfeng keeps a very low profile and rarely speaks in public. Beijing instructed DeepSeek "not to engage with the media without approval." [Reuters, 2025].



Google DeepMind

Corporate communications rarely mention extreme risks. Google DeepMind's leadership regularly discusses extreme risks in media interviews. Google's leadership does not.

Leadership examples:

•"We must take the risks of AI as seriously as other major global challenges, like climate change [...] It took the international community too long to coordinate an effective global response [..]. We can't afford the same delay with AI" [Guardian, 2024].

Time reported Hassabis saying: "Artificial intelligence is a dual-use technology like nuclear energy: it can be used for good, but it could also be terribly destructive" [Time, 2025]. Demis shares that he thinks AGI is only a "handful of years away" and that he is very worried about deception, calling it "incredibly dangerous", and speaks about encouraging the Security institutes to investigate them [Youtube, 2025]. Other examples: [CNN, 2025; CBS, 2025; Dwarkesh Podcast, 2024; TIME, 2023].

Shane Legg (Chief AGI Scientist) communicates a similar stance [<u>Dwarkesch Podcast</u>; 2023, <u>Google DeepMind</u>, 2023]. Talking to Axios, Legg recently stated AI is a very powerful technology, and it can and should be regulated." [Axios, 2025].

Google's CEO, Sundar Pichai, stated that "The biggest risk could be missing out," at the AI Action Summit in Paris yesterday

https://observer.com/2025/02/biggest-risk-ai-is-missing-out-google-ceo-sundar-pichai/?utm_source=chatgpt.com

CAIS statement on AI Risk signed by: Demis Hassabis (CEO), Shane Legg (Co-Founder), Lila Ibrahim (COO)

Meta

Company and leadership rarely address extreme risks.

Mark Zuckerberg and Chief AI Scientist Yann LeCun express the strongest counternarrative to AI existential risk concerns among major companies [Interesting Engineering, 2025]. LeCun does not believe that AI poses existential risk and calls such concerns "complete B.S.", arguing we need "the beginning of a hint of a design for a system smarter than a house cat before worrying about superintelligence" [Tech crunch, 2024]. Meta's president of global affairs expresses a similar position [Politico, 2024], comparing the discussion and framing the topic as a "moral panic" [Independent, 2024].

Zuckerberg is concerned about power concentration: "But I stay up at night worrying more about an untrustworthy actor having the super strong AI, whether it's an adversarial government or an untrustworthy company or whatever." He shares that: Bioweapons are one of the areas where the people who are most worried about this stuff are focused, and I think it makes a lot of sense." He expresses less urgency on existential risk addressing deception as "longer-term theoretical risks", and saying ".. we focus more on the types of risks that we see today .." [Dwarkesch Podcast, 2024].

OpenAl

OpenAI and its leadership sometimes talk about extreme risks

CEO Altman's communications have changed over time. In 2015, he stated: "I think that AI will probably, most likely, sort of lead to the end of the world" [Standford, 2024], and published a blog on "why machine intelligence is something we should be afraid of" [Altman, 2015].

In 2023, he published a blog "Planning for AGI and Beyond," stating OpenAI will proceed as if risks are "existential" [OpenAI, 2023]. In another blog, argued about the need for global coordination on the governance of superintelligence, and that "it would be important that such an agency focus on reducing existential risk" [OpenAI, 2023]. In his 2023 Senate testimony, he urged lawmakers to implement federal licensing and external audits to bound risk [Time, 2023].

In his recent communications, Altman adopted a notably more optimistic tone. In his recent congressional testimony, Altman told lawmakers that requiring government approval would be "disastrous" for US AI leadership [Washington Post, 2025]. His recent blogs focus on the benefits Superintelligence could bring [Altman, 2025].

CAIS statement on AI Risk signed by: Sam Altman (CEO), Adam D'Angelo (board member), Wojciech Zaremba (cofounder)

Relevant blogs:

• Preparing for future AI capabilities in biology:

Acknowledges AI could enable people to: "recreate biological threats or assist highly skilled actors in creating bioweapons.", then explains OpenAI's approach to preventing misuse [OpenAI, 2025].



x.AI

xAI itself does not publicly share information about extreme risks.

CEO Musk has a track record of raising concerns.

In 2014, Musk called AI humanity's "biggest existential threat.", calling for regulatory oversight [Live Science, 2014]

In September 2023, he told senators "there's some chance – above zero – that AI will kill us all." [NBC, 2023]. At the 2024 Saudi summit, he estimated a "10-20% chance AI goes bad." [Fortune, 2025]

CAIS statement on AI Risk signed by: Igor Babuschkin (cofounder), Tony Wu (co-founder).

Musk signed the FLI pause letter [FLI 2023]

Zhipu Al

Corporate communications don't speak about the potential for extreme risks. Leadership is discussed publicly.

Tang Jie 唐杰 (Chief Scientist) signed a 2024 track 2 diplomacy statement acknowledging potential for catastrophic risks: "Collectively, we must prepare to avert the attendant catastrophic risks that could arrive at any time." [IDAIS, 2024]

When speaking about AGI at the Seoul Summit, CEO Peng said: "[...] crucial responsibility of ensuring AI safety. As we delve deeper into the realms of AGI, it is imperative that we prioritize the development of robust safety measures to align AI systems with human values and ethical standards, thereby safeguarding our future in an AI-driven world.' [UK Gov, 2024]

Zhang Peng (CEO) was the only industry representative among Chinese scientists signing the IDAIS statement on AI safety redlines that should not be crossed, which stated: "[...] AI systems may pose catastrophic or even existential risks to humanity within our lifetimes." [IDAIS, 2024; Carnegie Endowment, 2024]. He gave a speech emphasizing the need for research to align superintelligent systems [36kr, 2024].

٦	$\Gamma \cap$	RF	COMPL	FTFD	BY PANFI	LISTS

Grading

Please pick a grade for each firm. You can add brief justifications to your grades.

	Anthropic	DeepSeek	Google DeepMind	Meta	OpenAl	x.Al	Zhipu Al
Grades	Grade	Grade	Grade	Grade	Grade	Grade	Grade
Grade comments (Justifications, opportunities for improvements, etc.)							

Grading Scales

Grading scales are provided to support consistency between reviewers.

- A Exemplary transparency enables informed safety decisions by all stakeholders
- B Strong transparency supports effective oversight and public understanding
- Moderate transparency with selective disclosure patterns
- Limited transparency hinders risk assessment
- Deceptive or negligible information sharing

Domain comments

Optional: Share observations that apply across companies, including general recommendations, notes on how you weighted indicators, or feedback on FLI's methodology.



Appendix B: Company Survey

Introduction

Thank you for participating in the **FLI AI Safety Index 2025 Survey**. This survey is designed to allow your company to provide additional information about specific practices and policies for managing risks from advanced AI systems. The independent experts on the review panel will consider the information you provide here when evaluating your company's safety efforts.

Survey instructions

The survey contains a total of **34 questions**, which predominantly follow a **multiple-choice format**. Where options are provided, select the one that best fits your current practices. Some questions allow a brief explanation or ask for details (especially if you answered "Other" or an open-ended part) – please be concise and factual in those responses. You are welcome to provide **URLs or document references** for any publicly available policies or reports that support your answers. It is not necessary to answer all questions within the survey. You can skip specific questions when answering would be difficult/inconvenient.

You have received a personalized link which you can share with colleagues to collaborate on the survey. You do not need to fill out the survey in a single sitting. Progress will be saved whenever you navigate between sections.

Confidentiality

Please do not share confidential information. We plan to publish all survey responses in full after the grading process is completed.

We appreciate your time and effort in providing thorough answers.



Whistleblowing Policies (15 Questions)

If your company has region-specific whistleblowing (WB) policies instead of a single global WB policy, please answer all questions in this survey with regard to the policy that applies to the majority of your frontier Alfocused management, research, and engineering employees. Unless a question specifically asks about other stakeholders, please answer based on the protections available to current full-time employees. You may explain variations for different stakeholder groups in the final question.

You can use the textbox at the end of this section to provide clarifications and/or link to relevant publicly available documents.

Definition of terms:

Whistleblowing Function:

The organizational structure, personnel, processes, and resources are established to receive, assess, investigate, and respond to whistleblowing reports. This includes the designated individuals or teams responsible for writing and acting according to the whistleblowing policy, managing the whistleblowing process, any technological systems used to facilitate reporting, and the mechanisms for investigating and addressing reported concerns.

Whistleblowing Policy:

The formal, documented set of rules, procedures, and guidelines that govern how an organization handles whistleblowing. This policy outlines what concerns can be reported ("material scope"), who can report them ("covered persons"), how reports should be made and to whom, how they will be handled, and what protections are available to whistleblowers who follow this policy. It serves as the official framework that defines the organization's approach to whistleblowing.

Covered persons:

Individuals who are explicitly protected when making good-faith reports under the whistleblowing policy. The range of covered persons may vary by organization and jurisdiction.

Material scope:

The range of issues, concerns, violations, or misconduct that can legitimately be reported through the whistleblowing channels and will be considered for investigation. In this context, this may include legal violations, ethical breaches, safety concerns, alignment issues, misrepresentations of capabilities, or other matters related to responsible AI development and deployment that the organization has defined as reportable concerns.



Question Title	Available options	Zhipu Al	xAI	OpenAl
Does your company have a WB policy & function covering frontier Al- focused staff? Is this policy publicly accessible without login credentials?	 Prefer not to answer (skips whistleblowing section) No WB policy & function - (skips whistleblowing section) Non-public policy exists - Please briefly explain your rationale for keeping it private: 	Prefer not to answer (skips whistleblowing section)	 Non-public policy exists - Please briefly explain your rationale for keeping it private: Only applies to xAI employees 	Public WB policy - Please provide URL here: https://openai.com/index/openairaising-concerns-policy/ https://cdn.openai.com/policies/raising-concerns-policy-blog-copy-202410.pdf
Who is formally designated with primary responsibility for overseeing the whistleblowing function and ensuring reports are properly addressed?	Board/Audit Committee Executive management Compliance/Legal department HR department Other (Please also specify whom this role reports to):		HR department	Board/Audit Committee, Compliance/Legal department, HR department HR, board/audit as well
Which statement best describes the investigative independence of your whistleblowing function?	 The whistleblowing function requires approval from management before initiating investigations based on whistleblower reports. The whistleblowing function can independently initiate and conduct investigations based on whistleblower reports, including those involving senior management. The whistleblowing function can independently initiate and conduct investigations based on whistleblower reports, including those involving senior management, AND has the authority to engage external expertise without approval. 		The whistleblowing function requires approval from management before initiating investigations based on whistleblower reports.	The whistleblowing function can independently initiate and conduct investigations based on whistleblower reports, including those involving senior management, AND has the authority to engage external expertise without approval.
Which of the following concerns are explicitly covered by your whistleblowing policy? (Select all that apply)	Violations of applicable laws and regulations Violations of the company's public AI safety framework (e.g., Anthropic's Responsible Scaling Policy) Credible safety concerns that may not violate specific policies including loss-of-control scenarios Pressure to compromise safety standards or suppress safety concerns Misleading communications about AI capabilities to external parties (such as regulators, the public, or evaluators) or discrepancies between public claims and internal practices None of the above		Violations of applicable laws and regulations, Credible safety concerns that may not violate specific policies including loss-of-control scenarios, Pressure to compromise safety standards or suppress safety concerns, Misleading communications about AI capabilities to external parties (such as regulators, the public, or evaluators) or discrepancies between public claims and internal practices	Violations of applicable laws and regulations, Violations of the company's public Al safety framework (e.g., Anthropic's Responsible Scaling Policy), Credible safety concerns that may not violate specific policies including loss-of-control scenarios, Pressure to compromise safety standards or suppress safety concerns



Question Title	Available options	Zhipu Al	xAI	OpenAl
Does your whistleblowing policy explicitly protect individuals who report concerns in 'good faith' or with 'reasonable cause to believe', rather than requiring certainty that violations occurred?	• Yes • No		Yes	Yes
Which of the following persons are protected from retaliation under your whistleblowing policy? (Select all that apply)	 Current employees Former employees Contractors and self-employed workers AI research collaborators and academic partners Individuals who assist whistleblowers Suppliers and vendors with access to company systems 		Current employees	Current employees,Contractors and self-employed workers
To which of the following individuals or entities can whistleblowers submit reports according to your policy? (Select all that apply)	 Board member or board committee Dedicated Ethics/Whistleblowing Officer Ombudsperson Chief Compliance or Risk Officer General Counsel/Legal Department Human Resources department External/independent third party Direct disclosure to a statutory or supervisory authority Other (please briefly specify): 		Human Resources department,Direct disclosure to a statutory or supervisory authority	Board member or board committee,Chief Compliance or Risk Officer,General Counsel/ Legal Department,Human Resources department,External/ independent third party,Direct disclosure to a statutory or supervisory authority
For former employees and contractors, indicate any policy limitations compared with current employees. (Select all limitations that apply)	Limited Reporting Channels (Former employees Contractors) Limited Reportable Issues (Former employees Contractors) Limited Retaliation Protection (Former employees Contractors) No Limitations (Former employees Contractors)	 Limited Reporting Channels () Limited Reportable Issues () Limited Retaliation Protection () No Limitations () 	Limited Reporting Channels (Former employees Contractors) Limited Reportable Issues () Limited Retaliation Protection (Former employees Contractors) No Limitations ()	Limited Reporting Channels (Contractors: Some channels, such as speaking to your current HR representative, are inherently available only to current employees.) Limited Reportable Issues () Limited Retaliation Protection () No Limitations ()
Which of the following best describes the anonymity and confidentiality provisions in your whistleblowing policy? (Select the one that fits best)	 Our policy does not provide for anonymous reporting Our policy allows anonymous reporting but does not specify technical measures to protect reporter identity Our policy allows anonymous reporting with specific technical measures in place to protect reporter identity (e.g., anonymous hotline, encrypted system) Our policy allows anonymous reporting with technical protections AND includes confidentiality commitments for non-anonymous reports 		Our policy does not provide for anonymous reporting	Our policy allows anonymous reporting with technical protections AND includes confidentiality commitments for non-anonymous reports



Question Title	Available options	Zhipu Al	xAI	OpenAl
If "Limited", under which circumstances is external disclosure protected?	 Imminent risk of serious harm Management or board implicated Reasonable fear of retaliation Internal investigation deadlines missed Unconditional reporting to a competent regulatory authority After internal reporting has been attempted Other (specify): 		Imminent risk of serious harm,Other (specify): Reasonable fear of physical harm	
Which mechanisms ensure that your whistleblowing function has access to adequate (technical) expertise to investigate reports? (Select all that apply)	 Dedicated AI experts within the whistleblowing function itself Authority to consult internal AI experts under confidentiality safeguards, including procedures that shield case details where necessary Standing agreements with external independent AI ethics/safety consultants Budget authority to engage external AI experts without requiring management approval None of the above Other (please specify): 		None of the above	Authority to consult internal AI experts under confidentiality safeguards, including procedures that shield case details where necessary
Investigation timelines and escalation rights: Which best describes your policy's commitments? (Select one)	 None – no specific timelines for acknowledgment, updates, or resolution Basic – acknowledge receipt ≤ 7 days only Standard – acknowledge ≤ 7 days and provide updates ≤ 30 days Full – acknowledge ≤ 7 days, updates ≤ 30 days, final outcome ≤ 90 days Full + internal escalation – all Full timeframes plus whistleblowers may escalate to board/leadership if deadlines are missed Full + comprehensive escalation – all Full timeframes plus whistleblowers may escalate both internally AND to regulators/external parties if deadlines are missed 		None – no specific timelines for acknowledgment, updates, or resolution	None – no specific timelines for acknowledgment, updates, or resolution



Question Title	Available options	Zhipu Al	xAI	OpenAl
Which specific forms of retaliation are explicitly prohibited in your policy? (Check all that apply)	 Termination/Dismissal Demotion, or negative performance reviews Reduction in compensation or benefits Exclusion from meetings or information Harassment or creating a hostile work environment Blacklisting within the industry Legal action against the whistleblower None of the above 		Termination/Dismissal,Demotion, or negative performance reviews,Reduction in compensation or benefits,Blacklisting within the industry,Legal action against the whistleblower	Termination/Dismissal,Demotion, or negative performance reviews,Reduction in compensation or benefits,Exclusion from meetings or information,Harassment or creating a hostile work environment,Blacklisting within the industry,Legal action against the whistleblower Our policy forbids retaliation. Notwithstanding the way this question is worded, it is well established under relevant law that retaliation can include termination or dismissal, demotion or negative performance reviews, or reduction in compensation or benefits. These are all covered under our policy's prohibition
				of retaliation. Our policy also expressly addresses harassment.
Do any employment-, separation-, or settlement-related agreements used by your company contain non-disparagement or confidentiality clauses that could deter current or former employees from disclosing AI safety or risk-related concerns? (Select one)	 No - we do not include such restrictions in our agreements Yes, but clauses only limit public disclosure; internal or regulator disclosures are explicitly unrestricted. Yes, but not enforced - clauses exist, but the company has a written policy never to enforce (or threaten to enforce) them against AI safety or risk-related disclosures (no withholding of pay/equity and no legal action). Yes, enforced - our standard confidentiality and non-disparagement provisions may restrict raising AI safety or risk-related concerns 		No - we do not include such restrictions in our agreements	Yes, but clauses only limit public disclosure; internal or regulator disclosures are explicitly unrestricted. We have confidentiality clauses that could impact some forms of public disclosure, but these have carveouts for internal or regulator disclosures. We do not have non-disparagement clauses in any such agreements, except in specific cases where an employee or former employee has entered a mutual non-disparagement agreement with the company.



Question Title	Available options	Zhipu Al	xAI	OpenAl
Which anti-retaliation provisions are explicitly detailed in your whistleblowing policy? (Select all that apply)	 Defined disciplinary consequences for individuals who retaliate against whistleblowers (e.g., termination, demotion, or other concrete penalties - not just general statements prohibiting retaliation) Documented investigation procedure for retaliation claims (including designated investigators, timelines, evidence standards, and appeal rights) 		None of the above are specifically detailed	Defined disciplinary consequences for individuals who retaliate against whistleblowers (e.g., termination, demotion, or other concrete penalties - not just general statements prohibiting retaliation)
	Concrete remedial measures for whistleblowers who experience retaliation (e.g., compensation, reinstatement, transfer options, or other specific remedies - not just general commitments to address retaliation) None of the above are specifically detailed			

External Pre-Deployment Safety Testing (6 Questions)

Please answer the following questions about external pre-deployment safety testing with regard to the release of your currently most capable publicly deployed AI model.

Frontier models:

- Anthropic Claude 4 Opus
- DeepSeek R1
- Google DeepMind Gemini 2.5 Pro
- Meta Llama 4 Maverick
- OpenAI o3
- xAI Grok3
- Zhipu AI GLM-4 Plus

You can use the textbox at the bottom of the page to provide clarifications and/or link to relevant publicly available documents.



Question Title	Available options	Zhipu Al	xAI	OpenAl
Did your organisation commission one or more independent (no financial/governance ties to your company) organisations to test this model for the dangerous capabilities or propensities you prioritized (in safety framework if available) before public release?	 No – no such external predeployment testing was commissioned (skip to next section) Yes – external testing was commissioned. Please list the organization(s) that performed relevant tests on the specified model and briefly indicate the broad risk domain(s) covered e.g., "UK AISI: cyber-offense, bio-risk" (opens follow-up questions below): 	Yes – external testing was commissioned. Please list the organization(s) that performed relevant tests on the specified model and briefly indicate the broad risk domain(s) covered e.g., "UK AISI: cyber-offense, biorisk (opens follow-up questions below): We intend to share our model with certain independent organizations for evaluation purposes; however, we prefer not to disclose their identities.	No – no such external predeployment testing was commissioned (skip to next section)	Yes – external testing was commissioned. Please list the organization(s) that performed relevant tests on the specified model and briefly indicate the broad risk domain(s) covered e.g., "UK AISI: cyber-offense, bio-risk (opens follow-up questions below): We've worked with the US and UK AI Safety Institutes, and independent third party labs such as METR, Apollo Research, and Pattern Labs to add an additional layer of validation for key risks. Where possible and relevant, we report on their findings in our systems cards, such as in the o3 System Card. Third party assessors were provided OpenAI o3 early checkpoints, as well as the final launch candidate models to conduct their assessments. As part of our ongoing efforts to consult with external experts, OpenAI granted early access to these versions of o3 to the U.S. AI Safety Institute to conduct evaluations of the models' cyber and biological capabilities, and to the U.K. AI Security Institute to conduct evaluations of cyber, chemical and biological, and autonomy capabilities, and an early version of the safeguards. METR measured the models' general autonomous capabilities, and reward hacking. Pattern Labs evaluated the model's cybersecurity related capabilities (evasion, network attack simulation, and vulnerability exploitation). Apollo Research evaluated in-context scheming and strategic deception. In some instances we paid private consultants for their
What was the highest level of technical access granted to any of the listed external evaluators during pre- deployment testing for the specified release? (Select the highest level that applies)	Standard inference API with normal user-facing filters in place Inference API with safety filters disabled (no inference-time mitigations) Helpful-only" or base model API (no harmlessness fine-tuning and no filters) Fine-tuning interface without safety gatekeeping Direct read/write access to internal activations or weights	Inference API with safety filters disabled (no inference-time mitigations)		work, but payment is not conditioned on the content of their findings. Standard inference API with normal user-facing filters in place Inference API with safety filters disabled (no inference-time mitigations) "Helpful-only" or base model API (no harmlessness fine-tuning and no filters)
What was the longest period of time that an external evaluator was given continuous access for predeployment testing of your model? (Select one)	>5 weeks>3 weeks>2 weeks>1 week<1 week	>3 weeks		>2 weeks



Question Title	Available options	Zhipu Al	xAI	OpenAl
Which of the following publication arrangements applied to external evaluators' findings? If different evaluators had different publication terms, please select all that occurred and briefly explain using the text-box. (select all that apply)	Evaluators may publish independently without prior company approval after the model is released. Evaluators may publish independently after company review/possible redaction. The company pre-committed to reproduce an independently written report in the model card without redactions. The company publishes report after review/possible redactions. The company provided its own summary of the evaluator's key findings. Findings remain internal Other: Please briefly explain:	Evaluators may publish independently without prior company approval after the model is released.,Evaluators may publish independently after company review/possible redaction.,The company precommitted to reproduce an independently written report in the model card without redactions.,The company publishes report after review/possible redactions.		Evaluators may publish independently without prior company approval after the model is released. This is true if they run their evaluations independently on the deployed model. Results from the red teaming period are under NDA / require prior approval Evaluators may publish independently after company review/possible redaction. See above, in cases where the evaluator wishes to publish about the specifics of the pre-deployment red teaming period The company publishes report after review/possible redactions. OpenAl publishes excerpts from the report mutually agreed upon or written, with OpenAl having the final say for what content goes in System Cards. The company provided its own summary of the evaluator's key findings. This is true in some cases, but we also share back any summaries that we plan to publish with the evaluator prior to release.
During pre- deployment testing, what best describes the query-rate or volume restrictions applied to external evaluators? (Select one)	No limits – evaluators could automate or batch queries with no additional throttling or hard caps. Elevated but capped – evaluators had higher quotas than the public/enterprise tier but were still subject to explicit caps (e.g., requests-per-minute or daily token limits). Public-tier caps – evaluators were held to the same rate/volume limits as ordinary paying users. Lower than Public-tier caps – evaluators had lower quotas than ordinary paying users.	No limits – evaluators could automate or batch queries with no additional throttling or hard caps.		Elevated but capped – evaluators had higher quotas than the public/enterprise tier but were still subject to explicit caps (e.g., requests-per-minute or daily token limits). Query rates can depend on technical feasibility in some cases.
Does your organization log and retain the model interactions of external evaluators during predeployment testing?	 Yes - Inputs and outputs are logged and retained. No - Inputs and outputs are neither logged nor retained, protecting evaluator IP. Other (please describe): 	Other (please describe): We will communicate with the evaluators to confirm whether it is permissible to retain relevant records.		Other (please describe): Zero Data Retention available upon request, if technically feasible during pre-deployment periods (for some new models or products, ZDR is not always possible during pre-deployment testing).

Internal Deployments (3 Questions)

Deployment levels:

- 1. Broad deployment: Many teams within the company have access for normal use.
- 2. Development access: Access limited to specific teams or projects that are actively testing the model or developing it further.

Question Title	Available options	Zhipu Al	xAI	OpenAl
If you specified external pre- deployment safety evaluations in the previous section, were these performed before or after broad internal deployment? (Select one)	Before - External safety tests were completed before broad internal deployment. Partial - All external evaluations on situational awareness, scheming, and cyber-offense were conducted before broad internal deployment. After - External safety tests were completed after broad internal deployment. Other (please explain briefly):	Partial - All external evaluations on situational awareness, scheming, and cyber-offense were conducted before broad internal deployment.		After - External safety tests were completed after broad internal deployment.
What level of safety testing does your company require for broad internal deployment of frontier AI models? (Select one)	 No formal risk management requirements for internal deployments Formalized risk management for internal deployments with less stringent requirements than external deployment framework for the following risks/capabilities: situational awareness, scheming, Al R&D, cyber-offense. Formalized risk management for internal deployments with the same requirements as external deployment framework for the following risks/capabilities: situational awareness, scheming, cyber-offense. Company requires the same risk management effort for internal and external deployments. Other (Please briefly describe): 	Formalized risk management for internal deployments with less stringent requirements than external deployment framework for the following risks/capabilities: situational awareness, scheming, AI R&D, cyber-offense.	No formal risk management requirements for internal deployments	As described in our public Preparedness Framework, we believe that models that have reached or are forecasted to reach Critical capability under our framework will require additional safeguards (safety and security controls) during development, regardless of whether or when they are externally deployed. We do not currently possess any models that have Critical levels of capability, and we expect to further update this Preparedness Framework before reaching such a level with any model.
Does your company require any of the following safeguards for broad internal deployments of frontier AI models? (Select all that apply)	 Inference time safety mitigations for misuse risks (including cyber & bio risks) Restricting access to helpful-only models and only granting time-bound access to staff that apply with a legitimate research need Logging all inputs and outputs from internal use and retaining them for at least 30 days Not currently logging, but introduced an *official, written* plan to start doing so after models reach a specified capability threshold Analyzing all internal model interactions for abnormal activity, including harmful use or unexpected attempts by AI systems to take real-world actions Live monitoring and automated editing/resampling of suspicious outputs None of the above Other (please describe briefly): 	Inference time safety mitigations for misuse risks (including cyber & bio risks),Restricting access to helpful-only models and only granting time-bound access to staff that apply with a legitimate research need,Logging all inputs and outputs from internal use and retaining them for at least 30 days,Analyzing all internal model interactions for abnormal activity, including harmful use or unexpected attempts by Al systems to take real-world actions,Live monitoring and automated editing/resampling of suspicious outputs	Logging all inputs and outputs from internal use and retaining them for at least 30 days	See answer to Q24, above.



Safety Practices, Frameworks, and Teams (9 Questions)

Question Title	Available options	Zhipu Al	xAI	OpenAl
When you released your latest flagship model, did you release the same model version that the final round of safety (framework) evaluations were conducted on? (Select one)	 Yes – we released the same model version. No – we further modified the model but explicitly mentioned and described all further changes in the model documentation. No – further modifications are not described explicitly in the model documentation. 	Yes – we released the same model version.	Yes – we released the same model version.	Yes – we released the same model version. Yes. We ran our evaluations on an earlier checkpoint and then confirmed our automated evaluation results on the final checkpoint.
If your company has one or more teams focused primarily on technical AI safety research, please provide more information about the team(s) below. By technical AI safety teams, we are referring to teams researching topics such as scalable oversight, dangerous capability evaluations, mechanistic interpretability, AI control, alignment evaluations, risk-modeling, etc. Please use separate paragraphs for listing multiple teams.	1) Team name (& website URL if available) 2) Mission and scope - Briefly describe the team's focus. Please distinguish between: - immediate product safety (e.g., RLHF, jailbreak prevention, safety classifiers), and - forward-looking/fundamental research (e.g., model organisms of misalignment, mechanistic interpretability) 3) Technical FTEs - Approximate number of full-time equivalent technical staff (researchers and research engineers). Please count each individual only once, based on their primary team.	This matter is considered company confidential, and we prefer not to answer.	Team name: Al Safety Engineer Mission and scope: Forward-looking / fundamental research + model improvements such as jailbreak prevention and safety classifiers FTEs: Three Team name: Product Safety Mission and scope: Immediate product safety such as jailbreak prevention FTEs: One	We have multiple teams across safety research focused on safety, alignment, evaluations, trustworthiness and governance.
Does your organization have a formal, written policy that requires notifying external authorities when safety testing determines a model exceeds your organization's "unacceptable-risk" threshold (i.e., a risk-level that bars deployment under your own safety framework), even if the model will not be released? (Select option that best describes your policy)	 1) No policy – there is no written requirement to notify any external body. 2) Regulator-only notification – the policy mandates prompt disclosure to a competent regulatory, or supervisory authority. 3) Regulator + public transparency – as in option 2 **and** the policy provides for a public statement or summary once doing so will not exacerbate security risks. Other (please briefly describe): 	2) Regulator-only notification – the policy mandates prompt disclosure to a competent regulatory, or supervisory authority.	No policy – there is no written requirement to notify any external body.	1) No policy – there is no written requirement to notify any external body.



Question Title	Available options	Zhipu Al	xAI	OpenAl
For companies that signed the "Frontier AI Safety Commitments" at the AI Seoul Summit in 2024, and those that strive to implement equivalent safety frameworks: Which of the levels below best describes the status of your Safety Framework? Please indicate the *highest* option below that accurately describes your current state.	 No official Safety Framework published (yet). Published & Implementation in progress Published & substantially implemented Most discrete policies, processes, or technical safeguards described in the policy are fully implemented and operational. Please briefly assert which elements have not been implemented as described yet and the expected timeline for implementation: Published & fully implemented - All discrete policies, processes, or technical safeguards described in the policy are fully implemented and operational. 	Published & Implementation in progress	Published & Implementation in progress	Published & Implementation in progress
Do you have a plan for ensuring that the AGI you're trying to build will remain controllable, safe and beneficial?	No, but we're working on it Yes, internally. (Please briefly explain why you have not published it)	Yes, internally. (Please briefly explain why you have not published it) Currently, Zhipu's models have not yet reached the level of AGI, so we prefer not to release the related plans.	No, but we're working on it	Yes, internally. (Please briefly explain why you have not published it) For more on our approach to ensuring that AGI remains controllable and safe, see https://openai.com/safety/how-we-think-about-safety-alignment/



Question Title	Available options	Zhipu Al	xAl	OpenAl
Which of the following elements of an AI emergency response capability has your organization implemented? (Select all that apply)	 Maintained and tested technical capability to rapidly roll back a deployed model to a previous version globally (within 12h). Successfully tested rapid full model rollback including internal deployments within the last 12 months. Maintained and tested technical capability to rapidly tighten model safeguards and restrict specific capabilities (e.g. web-browsing) globally. Successfully tested rapid throttling or capability-restriction including internal deployments within the last 12 months. Conducted at least one full live emergency response drill/simulation in the past 12 months. Created a formal, documented emergency response plan for Al safety incidents with threshold for triggering emergency response, a named incident commander and a 24×7 duty roster. Established a risk-domain-specific (e.g. bio, cyber) 24-hour communication protocol and points of contact with relevant government agencies. None of the above Other: Please use this text-field to share URLs to relevant documentation or to clarify specific responses 	Maintained and tested technical capability to rapidly roll back a deployed model to a previous version globally (within 12h). Successfully tested rapid full model rollback including internal deployments within the last 12 months.,Maintained and tested technical capability to rapidly tighten model safeguards and restrict specific capabilities (e.g. web-browsing) globally. Successfully tested rapid throttling or capability-restriction including internal deployments within the last 12 months.,Created a formal, documented emergency response plan for Al safety incidents with threshold for triggering emergency response, a named incident commander and a 24 × 7 duty roster,Established a risk-domain-specific (e.g. bio, cyber) 24-hour communication protocol and points of contact with relevant government agencies.	Maintained and tested technical capability to rapidly roll back a deployed model to a previous version globally (within 12h). Successfully tested rapid full model rollback including internal deployments within the last 12 months., Maintained and tested technical capability to rapidly tighten model safeguards and restrict specific capabilities (e.g. web-browsing) globally. Successfully tested rapid throttling or capability-restriction including internal deployments within the last 12 months.	way. Our response capabilities include:



Question Title	Available options	Zhipu Al	xAI	OpenAl
Does your company agree with the following principles for promoting legible and faithful reasoning in advanced AI systems to ensure AI remains safe and controllable? (Select all statements you support) Leading AI companies should:	Ensure Human-Legible Reasoning - Al models should reason in ways that are accessible and understandable to humans. Developers should avoid opaque reasoning methods. Avoid Optimization That Encourages Obfuscation - Developers should exercise caution when applying optimization pressures to model reasoning, especially when removing 'undesired reasoning', to prevent fostering deceptive behavior. Disclose Optimization Pressures on Reasoning - Companies should transparently report the optimization pressures and training methods applied to model reasoning, particularly when removing 'undesired reasoning.' None of the above	**Ensure Human-Legible Reasoning** - Al models should reason in ways that are accessible and understandable to humans. Developers should avoid opaque reasoning methods,**Avoid Optimization That Encourages Obfuscation** - Developers should exercise caution when applying optimization pressures to model reasoning, especially when removing 'undesired reasoning', to prevent fostering deceptive behavior.	**Avoid Optimization That Encourages Obfuscation** - Developers should exercise caution when applying optimization pressures to model reasoning, especially when removing 'undesired reasoning', to prevent fostering deceptive behavior.,**Disclose Optimization Pressures on Reasoning** - Companies should transparently report the optimization pressures and training methods applied to model reasoning, particularly when removing 'undesired reasoning.'	**Avoid Optimization That Encourages Obfuscation** - Developers should exercise caution when applying optimization pressures to model reasoning, especially when removing 'undesired reasoning', to prevent fostering deceptive behavior. We've publicly urged against optimizing on chains of thought: https://openai.com/index/chain-of-thought- monitoring/
Task-Specific Fine-Tuning (TSFT) involves training a model to excel at potentially dangerous tasks (e.g., designing biological agents, cyber attacks). Before releasing your current frontier model, which statement best describes your TSFT safety testing? (Select one)	 None – no TSFT safety testing performed (skips follow-up). Partial – TSFT performed on ≤ 2 high-risk domains (choose below). Comprehensive – TSFT performed on ≥ 3 high-risk domains (choose below). 	Comprehensive – TSFT performed on ≥ 3 high-risk domains (choose below).	None – no TSFT safety testing performed (skips follow-up).	None – no TSFT safety testing performed (skips follow-up). None. We evaluated helpful-only models, which we believe is appropriate for the threat model of misuse for models made available via our platform and whose weights we do not release, as is codified in our Preparedness Framework.
If you selected 'Partial' or 'Comprehensive' on the previous question, Please tick the risk-domains tested with TSFT.	 Biological Persuasion Chemical Deceptive alignment / Autonomy Cyber-offense Other (please specify): 	Other (please specify): Biological, Persuasion, Chemical, Cyber-offense, Political		

Question Title	Available options	Zhipu Al	xAI	OpenAl
Question Title If you wish to provide clarifications to particular answers, you can use this textbox to do so. Please reference specific questions using their associated number. You may also share additional information about your company's policies.	Available options	Zhipu Al	xAI	Below, we include some additional information about our security work that we believe may be useful context for evaluators considering our overall posture and approach. • For additional technical detail on our security measures for Al see: Securing Research Infrastructure for Advanced Al. • Third party collaboration on security: OpenAl maintains a bug bounty program through BugCrowd (https://bugcrowd.com/openai), and welcomes responsible disclosures from third parties via our coordinated vulnerability disclosure policy (https://openai.com/policies/coordinated-vulnerability-disclosure-policy/). In addition, OpenAl runs a Cybersecurity Grant Program to support research and development focused on protecting Al systems and infrastructure. This program encourages and
				funds initiatives that help identify and address vulnerabilities, ensuring the safe deployment of AI technologies.

FLI AI Safety Index

Independent experts evaluate safety practices of leading AI companies across critical domains.

17th July 2025