

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