

Anthropic	DeepSeek	Google DeepMind	Meta	OpenAI	x.AI	Zhipu AI
<p>Bug bounty on universal jailbreaks</p> <ul style="list-style-type: none"> - 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	<p>Abuse Vulnerability Reward Program:</p> <p>Accepts certain abuse-related discoveries:</p> <ul style="list-style-type: none"> - Prompt Attacks - Training Data Extraction - Manipulating Models - Adversarial Perturbation - Model Theft <p>(excludes jailbreaks) [Google]</p>	<p>Bounty programs are restricted to privacy or security issues, like extracting training data through tactics like model inversion or extraction attacks. [Meta]</p>	<p>Early access for safety testing (December 2024)</p> <p>One-off programs allowed safety researchers to apply for early access to frontier models to help surface novel risks. No payments announced. [OpenAI, 2024]</p>	None	None