

Anthropic	<p><b>Transparency</b></p> <p>In March 2024, Anthropic shared the full system prompt alongside the release of Claude 3 as a one-off <a href="#">[Fast Company, 2024]</a>.</p> <p>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 <a href="#">[TechCrunch, 2024; X, 2024; Anthropic]</a>.</p> <p>Simon Willison reported that the publicly shared version does not include the description of various tools available to the model <a href="#">[Simon Willison, 2025]</a>.</p>
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.
OpenAI	No transparency on system prompts for Frontier Systems.
x.AI	<p><b>Transparency:</b></p> <p>Following the incident in May 2025 listed below, x.AI published their system prompts for Grok (on xAI &amp; X) on Github and promised these will be regularly updated <a href="#">[Github, 2025]</a>.</p> <p><b>Incidents:</b></p> <p>February 2024:</p> <p>A change to the system prompt, Grok briefly censored responses about Elon Musk and Donald Trump spreading disinformation. After the issue received public attention, xAI quickly reverted the changes and publicly stated that the problem was caused by an unnamed employee conducting unauthorized modifications <a href="#">[Fortune, 2024]</a>.</p> <p>May 2025:</p> <p>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 AI Chatbot told users it was ‘instructed by my creators’ to accept ‘white genocide as real and racially motivated’ <a href="#">[Guardian, 2025]</a>. x.AI quickly apologized for this incident and rolled back the changes. They reported that unauthorized modifications by an employee caused the incident <a href="#">[X, 2025]</a>.</p>
Zhipu AI	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.