

Anthropic	<p>Constitutional AI:</p> <p>Method 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.</p> <p>What's it for:</p> <p>1) Supervised learning phase: Model self-critiques and revises its outputs based on constitutional principles, creating a supervised learning dataset</p> <p>2) RLAIF phase: Model compares response pairs using constitutional principles to generate preference labels, then trains via RL on these AI-generated preferences</p> <p>Timeline & Development:</p> <p>December 2022: Original Constitutional AI paper published</p> <p>May 2023: Claude's constitution made public (58 principles)</p> <p>Constitution (May 2023):</p> <p>58 principles (1.2k words) drawn from:</p> <ul style="list-style-type: none"> - UN Declaration of Human Rights - Apple's Terms of Service - DeepMind's Sparrow principles - Non-Western perspectives - Anthropic's own research <p>Example principle: "Please choose the response that most supports and encourages freedom, equality, and a sense of brotherhood."</p> <p>Benefits:</p> <p>Readable, transparent, and explicitly formulated principles, as opposed to RLHF, which leverages implicit values.</p> <p>Limitations:</p> <p>Version uncertainty: Only the May 2023 constitution is public; the current production versions are unknown</p> <p>Anthropic uses a "variety of techniques including human feedback, Constitutional AI [..], and the training of selected character traits." Given that other approaches are incorporated in post-training, the impact of any one of them is unclear.</p> <p>Since 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 values conflict.</p> <p>Source: [Anthropic, 2025]</p>
	<p>DeepSeek</p> <p>No detailed specification available, but frontier model weights are public, so models can be modified.</p>
	<p>Google DeepMind</p> <p>No detailed specification available</p>
	<p>Meta</p> <p>No detailed specification available, but frontier model weights are public, so models can be modified.</p>
	<p>OpenAI</p> <p>OpenAI Model Spec:</p> <p>OpenAI's Model Spec is a detailed (~28k words), public, living rule-book that defines the objectives, safety rules, and default behaviours OpenAI trains its models —via human feedback and deliberative alignment—to follow.</p> <p>What's it for:</p> <p>1) Human RLHF guidance – provides a single, public rule-book that labelers follow when creating preference data.</p> <p>2) Deliberative Alignment – o-series models (o1, o3, o4-mini) are explicitly taught to read and reason over the Spec before answering.</p> <p>3) Automated evaluation – OpenAI ships a challenge-prompt suite to measure adherence.</p> <p>Timeline & Versions:</p> <p>1st May 2024</p> <p>2nd Feb 2025</p> <p>3rd Apr 2025</p> <p>Framework:</p> <p>Three principal types:</p> <p>1) Objectives – broad goals such as "assist the developer & end user" and "benefit humanity."</p> <p>2) Rules – hard, platform-level constraints (e.g., comply with law, prohibit or restrict certain content, protect privacy, uphold fairness).</p> <p>3) Defaults – stylistic and behavioural norms that developers/users may override.</p> <p>Sections: Stay in bounds · Seek the truth together · Do the best work · Be approachable · Use appropriate style.</p> <p>Includes specific guidance on specific policy areas such as potential, medical, or harmful content.</p> <p>Risk taxonomy: Misaligned goals · Execution errors · Harmful instructions.</p> <p>Chain of command:</p> <p>Platform (OpenAI) → Developer → User → Guideline → Untrusted text.</p> <p>Within any level, explicit > implicit, later > earlier.</p> <p>(OpenAI's Usage Policy overrides the Spec if the two conflict.)</p> <p>Ongoing Development:</p> <p>Released under CC0 license (public domain)</p> <p>Changelog and version history maintained on GitHub</p> <p>OpenAI commits to regular updates as the spec evolves</p> <p>Key Benefits</p> <p>Greater transparency of intended model behavior.</p> <p>Finer-grained steerability via the chain of command</p> <p>Reduced reliance on implicit human values; models can show interpretable reasoning steps grounded in the Spec.</p> <p>Transparency & Limitations</p> <p>Production models don't fully reflect the spec yet.</p> <p>OpenAI states: "While the public version of the Model Spec may not include every detail, it is fully consistent with our intended model behavior."</p> <p>Source: [OpenAI, 2025]</p>
x.AI	<p>No detailed specification available</p>
Zhipu AI	<p>No detailed specification available, but frontier model weights are public, so models can be modified.</p>