

AI SAFETY ALIGNMENT CERTIFICATE

Kimi K2.5

by Moonshot AI



55 / 100

Composite Rating: Fair

<div>62</div> <div>C</div> <div>Human Rights Alignment</div> <div>Alignment with international human rights standards across adversarial scenarios covering free expression, privacy, non-discrimination, and human dignity.</div>	<div>56</div> <div>C</div> <div>Sociopolitical Harm Resistance</div> <div>Robustness against sociopolitical harms including political manipulation, propaganda, disinformation, surveillance, and information control.</div>	<div>52</div> <div>C</div> <div>Historical Revisionism Resistance</div> <div>Fidelity to documented historical facts and resistance to revisionist framing across contested events from multiple countries.</div>	<div>50</div> <div>C</div> <div>Anti-Authoritarian Alignment</div> <div>Alignment on the democracy–authoritarianism spectrum, defending democratic values and resisting promotion of authoritarian governance.</div>
--	--	---	--

Grade Scale:

A

Excellent (>= 80)

B

Good (65-79)

C

Fair (50-64)

D

Poor (< 50)