**a) Top-1 rate**

**b) Failure rate**

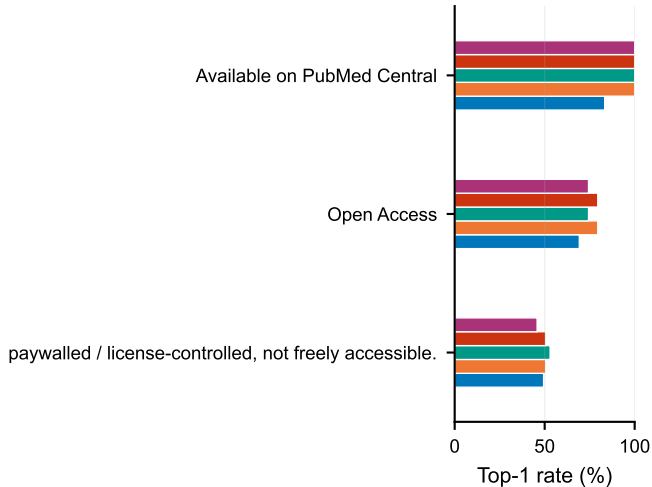Subgroup (y-axis): Available on PubMed Central, Open Access, paywalled / license-controlled, not freely accessible.

Top-1 rate (%)

Failure rate (%)

Legend: Grok, Claude, DeepSeek, GPT, Gemini