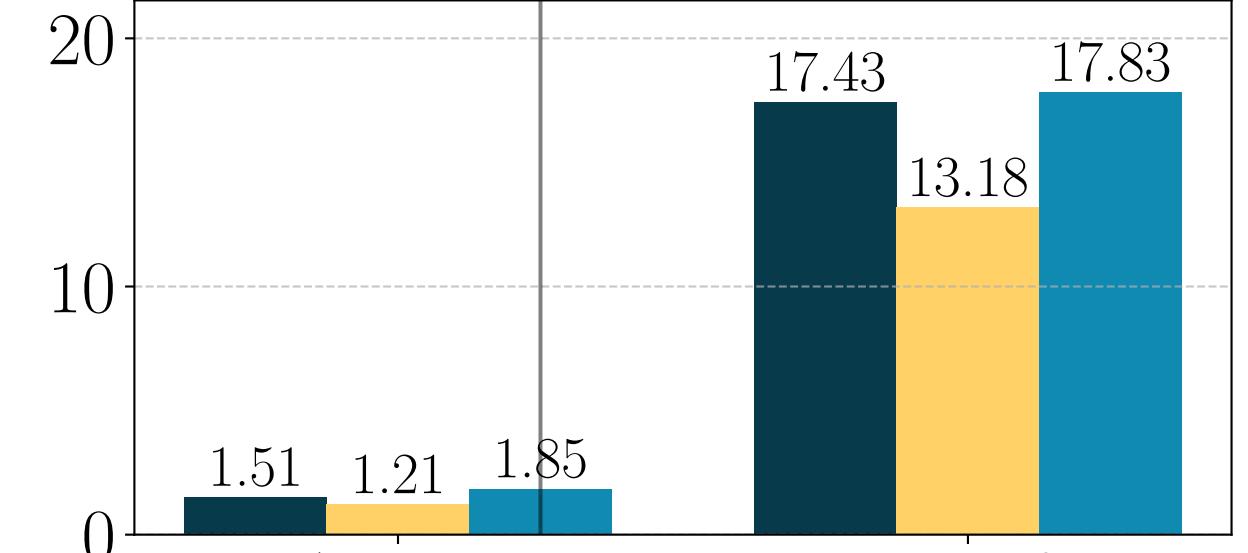


DeepsSeek-R1

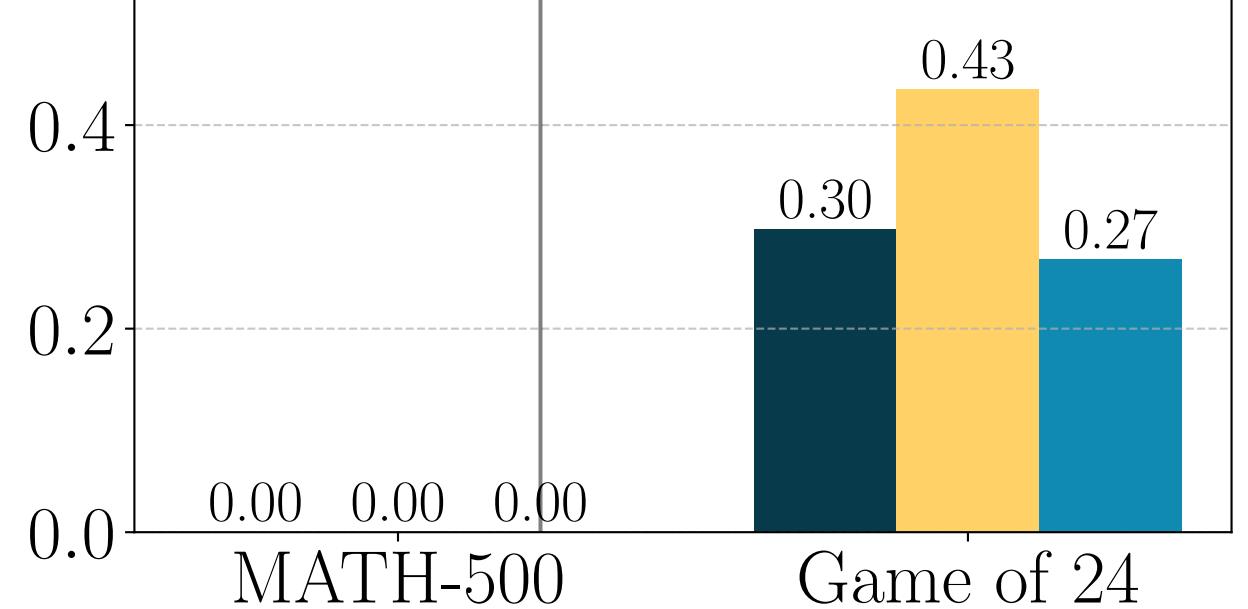
Grok 3 Mini Beta

QwQ-32B

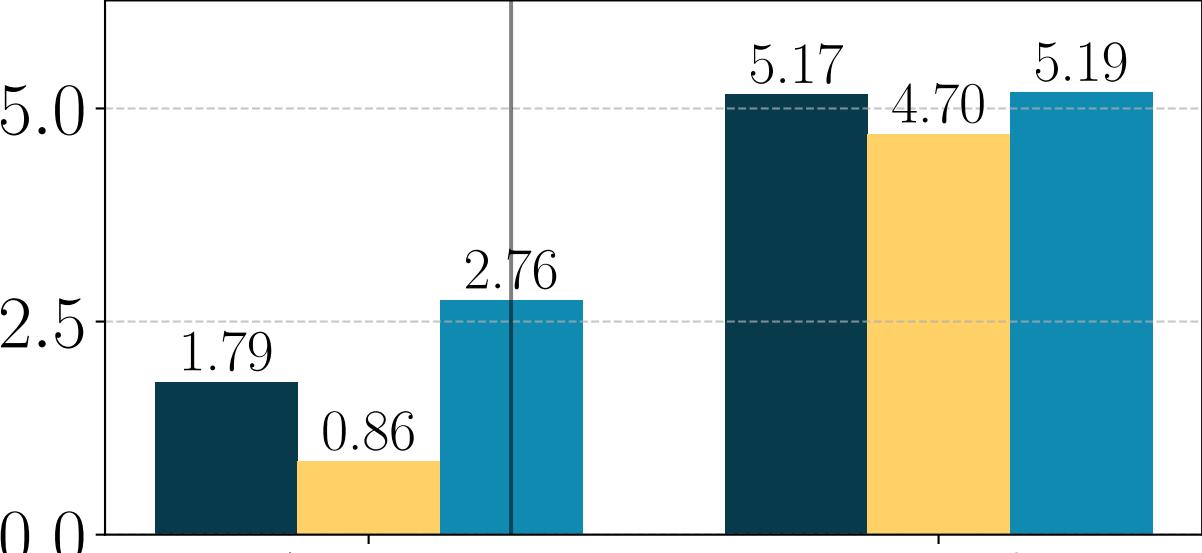
Average Solution Count



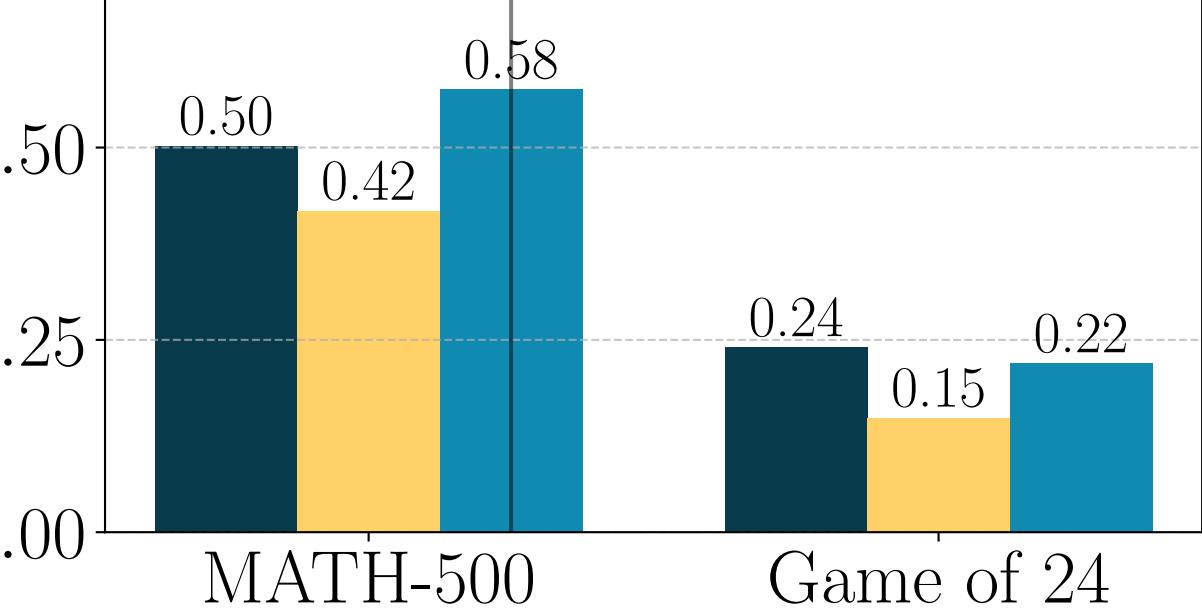
Average Overthinking Rate



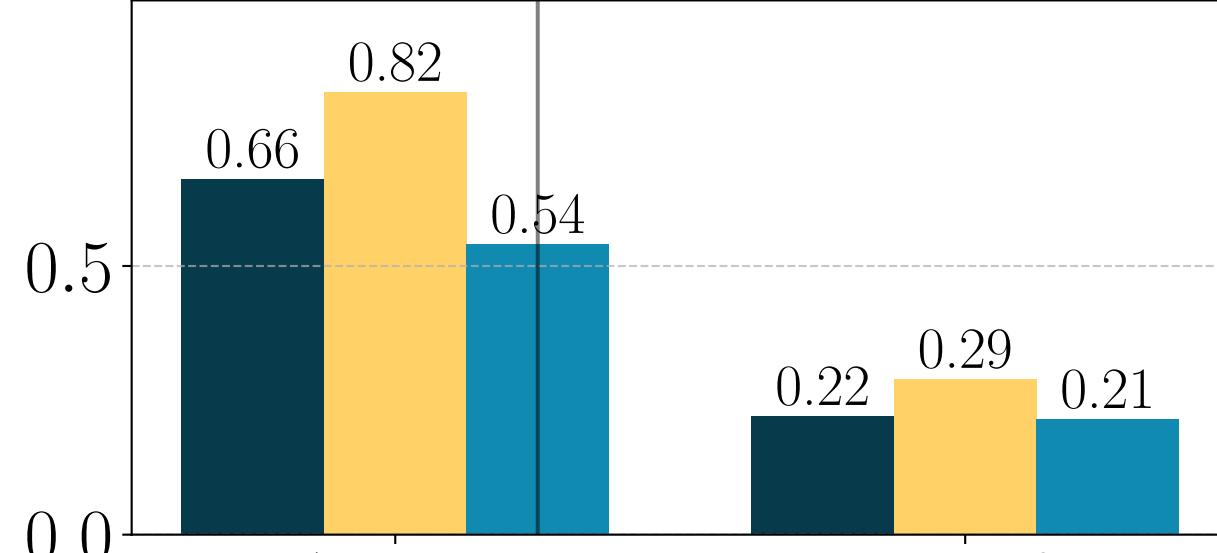
Average Jump Distance



Average Verification Rate



Average Success Rate



Forgetting Rate

