

Llama 3.3 70B

MedHELM: medi_qa (rescale, max=5)
 MedHELM: medhallu
 MedHELM: medication_qa (rescale, max=5)
MedHELM: pubmed_qa
 HealthBench: 02.Communication quality
 Answered with Evidence: 02.Red (reverse scored)
 HealthBench: 02.Instruction following
 HealthBench: 01.Emergency referrals
 HealthBench: 02.Accuracy
 CPC
 HealthBench: 01.Expertise-tailored communication
 HealthBench: 01.Health data tasks
 HealthBench: 01.Response depth
 HealthBench: 01.Responding under uncertainty
 HealthBench: 00.Overall
 HealthBench: 02.Context awareness
 HealthBench: 01.Global health
 HealthBench: 01.Context seeking
 Answered with Evidence: 01.Green
 HealthBench: 02.Completeness

