

proportion of choosing original word

wikipedia

news

nonsense

randomseq

Qwen 1.5B Qwen 1.5B IT OLMo 7B OLMo 7B-seed1 OLMo 7B-seed2 OLMo 7B-seed3 Qwen 7B Qwen 7B IT Llama3.1 8B Llama3.1 8B IT OLMo 13B OLMo 13B-seed1 OLMo 13B-seed2 OLMo 13B-seed3 Llama3.1 70B Llama3.1 70B IT Llama3.3 70B IT Qwen 72B Qwen 72B IT Mistral 123B IT Llama3.1 405B

Qwen 1.5B Qwen 1.5B IT OLMo 7B OLMo 7B-seed1 OLMo 7B-seed2 OLMo 7B-seed3 Qwen 7B Qwen 7B IT Llama3.1 8B Llama3.1 8B IT OLMo 13B OLMo 13B-seed1 OLMo 13B-seed2 OLMo 13B-seed3 Llama3.1 70B Llama3.1 70B IT Llama3.3 70B IT Qwen 72B Qwen 72B IT Mistral 123B IT Llama3.1 405B

Direct (log prob) Meta (prompting) Meta mean

base instruct