

Model performance across benchmarks with different freezing strategies

Accuracy

Model

- OLMo-2-1124-7B-Instruct
- SFT
- SFT (frozen MLPs)
- SFT (unfrozen Q and K)

