

Performance on simple inferences vs. GSM-8K score

Temperature: 1, Condition: few-shot

