Benchmark: MMLU Logical Fallacies
Performance Metric: Target Greedily Decoded
Correlation Metric: Kendall