

MMLU, llama-3.1-8b to GPT-4.1 mini

Accuracy

0.85
0.80
0.75
0.70
0.65
0.60

0.00

0.25

0.50

0.75

1.00

Routing Ratio

- average-token-prob
- verbalization-1s
- verbalization-2s
- p(true)
- trained-probe
- perplexity
- jaccard-degree
- ood-probe