

# AUROC

Test dataset

GPT-3

0.87

0.79

0.70

GPT-4

0.66

0.73

0.55

Gemini

0.92

0.80

0.85

Llama2

0.99

0.80

0.98

Mistral

0.94

0.80

0.85

PaLM2

0.99

0.80

0.97

Qwen

0.99

0.80

0.90

o1

0.56

0.72

0.50

Accuracy

Low FNR

Low FPR

Threshold

