

AUROC

Model family (test)

Claude

0.75 0.74 0.77 0.78 0.77 0.74 0.73 0.80

GPT-3

0.84 0.83 0.85 0.85 0.85 0.84 0.81 0.80

GPT-4

0.80 0.85 0.80 0.78 0.80 0.79 0.84 0.82

Llama2

0.70 0.71 0.74 0.78 0.78 0.74 0.73 0.72

Mistral

0.63 0.64 0.66 0.68 0.67 0.65 0.63 0.64

o1

0.43 0.61 0.38 0.37 0.42 0.42 0.61 0.87

GPT-3

GPT-4

Gemini

Llama2

Mistral

PaLM2

Qwen

o1

Model family (training)

