

Attacked Models

Cross Entropy of P(Target | Prompt, Image)

'prism-siglip+7b'

10^3

Gradient Step

10^4

'prism-clip+7b'

10^3

Gradient Step

10^4

prism-reproduction-llava-v15+7b

Evaluated Model

eval_model_str

prism-reproduction-llava-v15+7b

