

Attacked Model(s)

Cross Entropy of
 $P(\text{Target} \mid \text{Prompt, Image})$

'prism-dinosiglip+13b'

advbench

Evaluated Model

'prism-clip+7b'

'prism-clip-controlled+7b'

'prism-reproduction-llava-v15+7b'

'prism-siglip+7b'

'prism-siglip-controlled+7b'

Attack Dataset

advbench

----- rylan_anthropic_hhh

10^3

10^4

Gradient Step

