

Attacked Model(s)

Cross Entropy of
P(Target | Prompt, Image)

10⁰
10⁻¹

10³

10⁴

Gradient Step

'prism-siglip+13b'

advbench

Evaluated Model

- 'prism-clip+13b'
- 'prism-clip+7b'
- 'prism-clip-controlled+7b'
- 'prism-dinosiglip+7b'
- 'prism-dinosiglip-controlled+7b'
- 'prism-reproduction-llava-v15+13b'
- 'prism-reproduction-llava-v15+7b'
- 'prism-siglip+7b'
- 'prism-siglip-controlled+13b'
- 'prism-siglip-controlled+7b'
- Attack Dataset
- advbench
- rylan_anthropic_hhh

