

# Attacked Model(s)

'prism-siglip-controlled+7b'

## Evaluated Model

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

advbench

model\_to\_eval

'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

