Universality of Image Jailbreaks Attack Dataset (Train Split) AdvBench Cross Entropy Claude 3 Opus Anthropic HHH Same Data Distribution of Harmful Responses  $6 \times 10^{-1}$   $4 \times 10^{-1}$   $3 \times 10^{-1}$   $2 \times 10^{-1}$ True False  $10^{0}$  $89 \times 10^{-1}$  $8 \times 10^{-1}$  $7\times10^{-1}$  $6\times 10^{-1}$ HarmBench LlamaGuard2 solves  $6 \times 10^{-1}$  We show that  $4 \times 10^{-1}$   $3 \times 10^{-1}$   $3 \times 10^{-1}$   $3 \times 10^{-1}$ % of Harmful Responses  $4\times 10^{-1}$  % - $3 \times 10^{-1}$  $10^{3}$  $10^{4}$  $10^{3}$  $10^{4}$ Gradient Step Gradient Step