Universality of Image Jailbreaks Attack Dataset (Train Split) AdvBench Cross Entropy Claude 3 Opus Anthropic HHH Same Data Distribution True 6×10^{-1} False Attack Failure Rate 3×10^{-1} 3×10^{-1} 10^{0} 89×10^{-1} 8×10^{-1} 7×10^{-1} 2×10^{-1} 6×10^{-1} HarmBench LlamaGuard2 6×10^{-1} Attack Failure Rate 4×10^{-1} 4×10^{-1} Attack Failure Rate 3×10^{-1} 3×10^{-1} 2×10^{-1} 3×10^{-1} 10^{3} 10^{4} 10^{3} 10^{4} Gradient Step Gradient Step