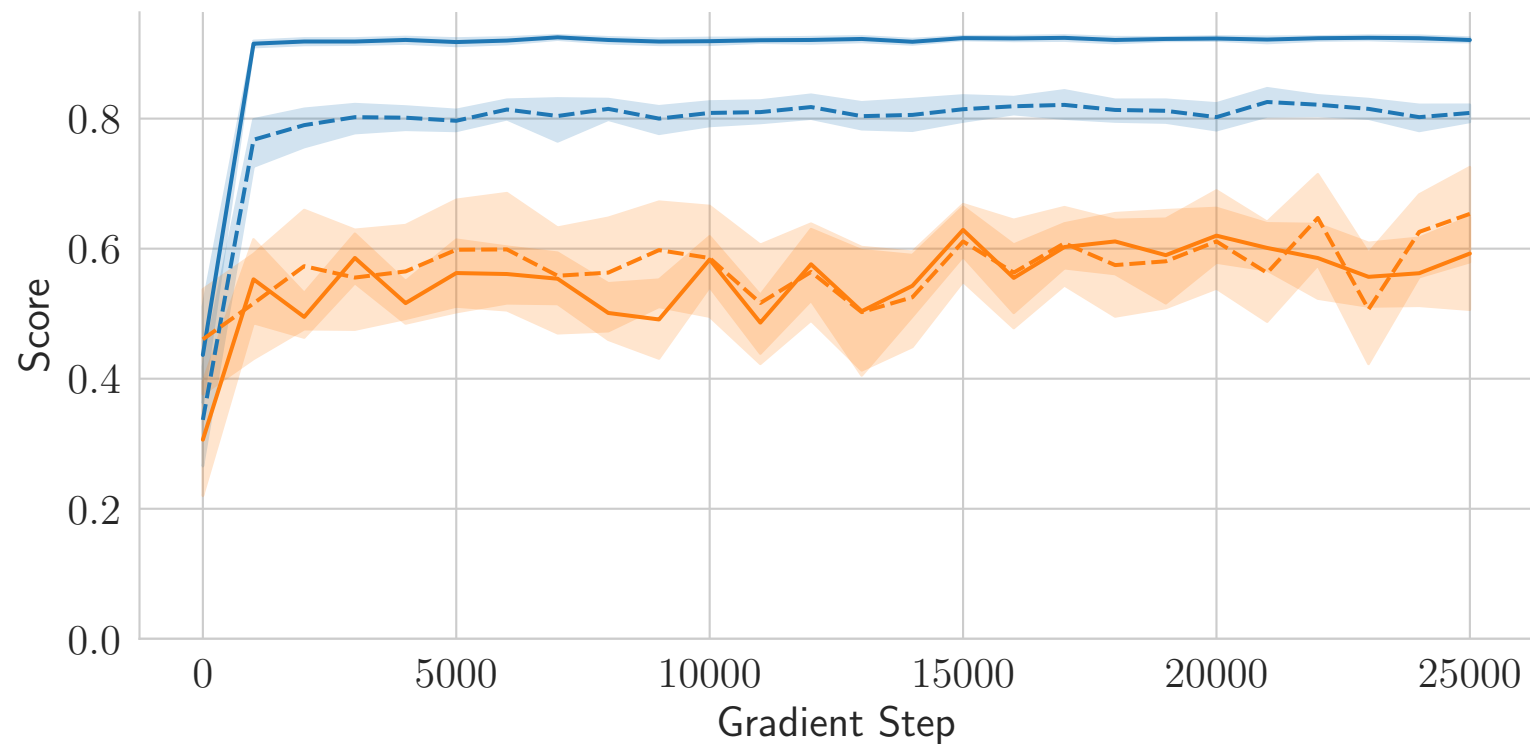


Universality of Image Jailbreaks

LlamaGuard2 Score



Attack Dataset (Train Split)

advbench

rylan_anthropic_hhh

Same Data Distribution

True

False