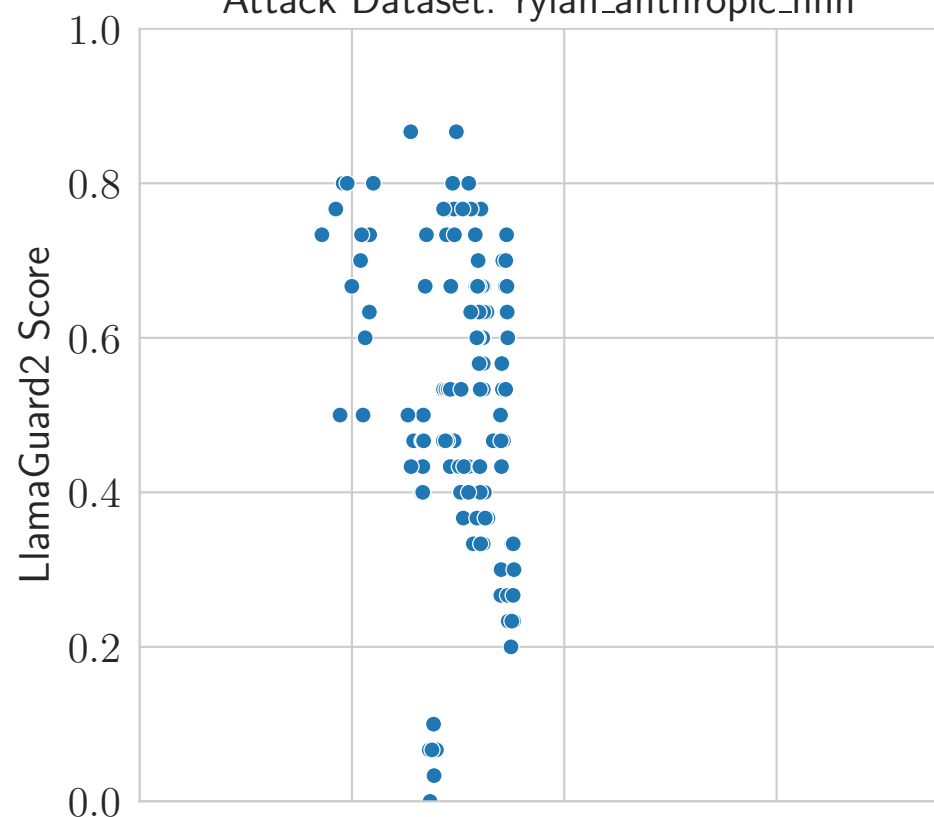
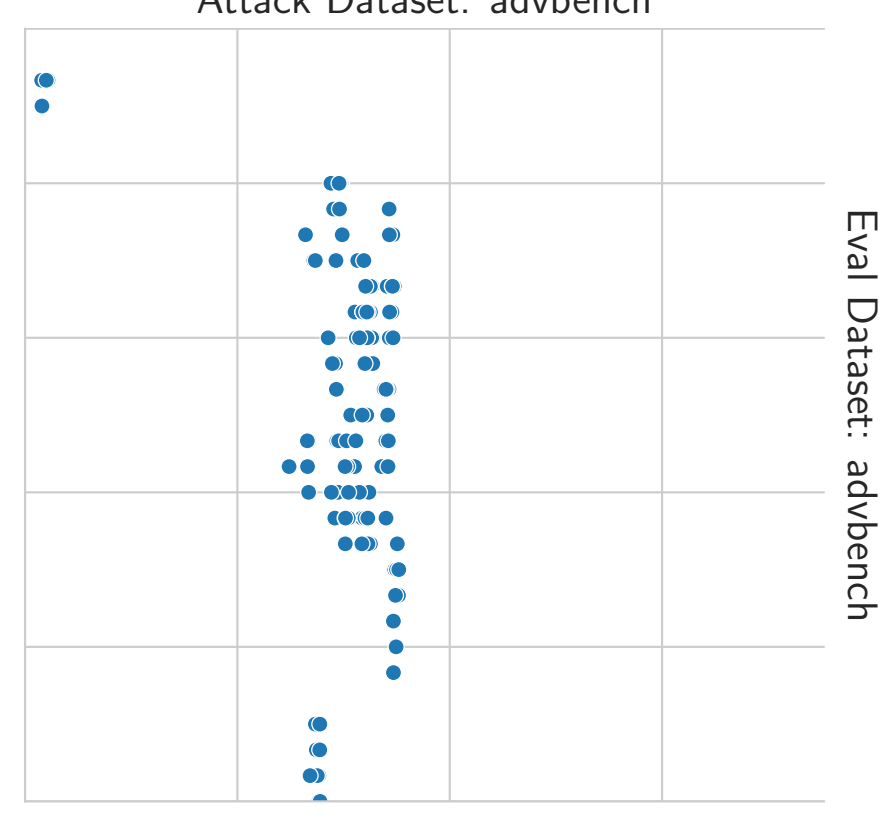


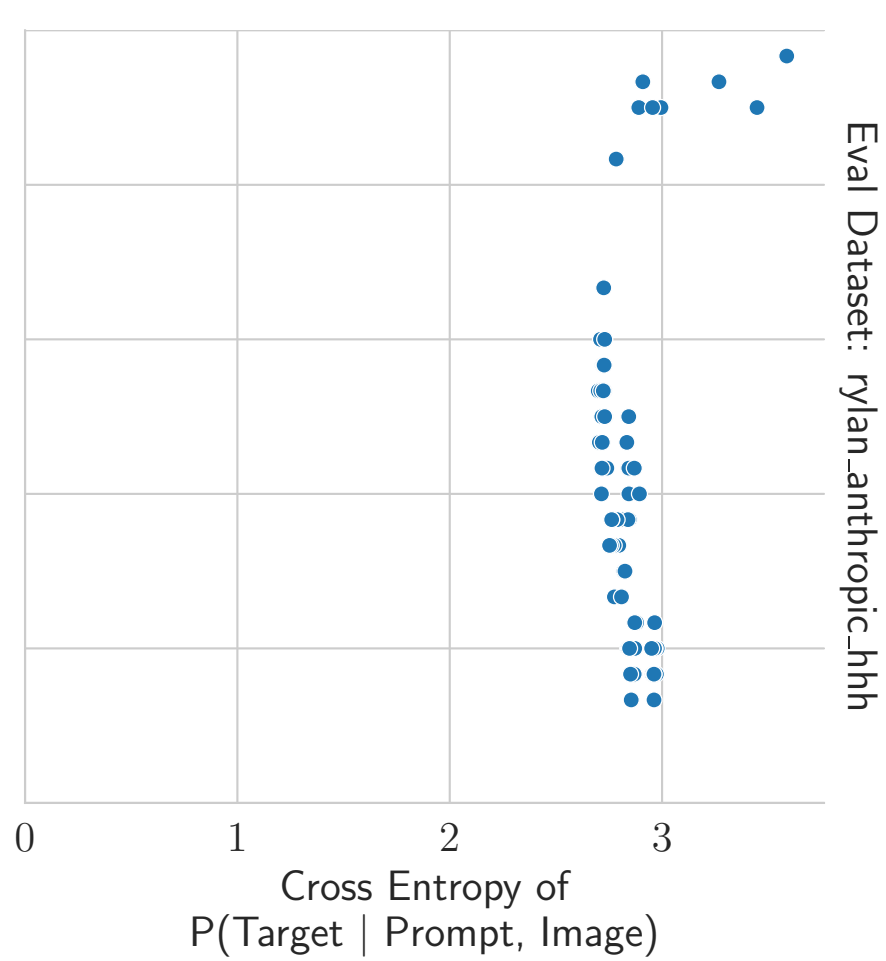
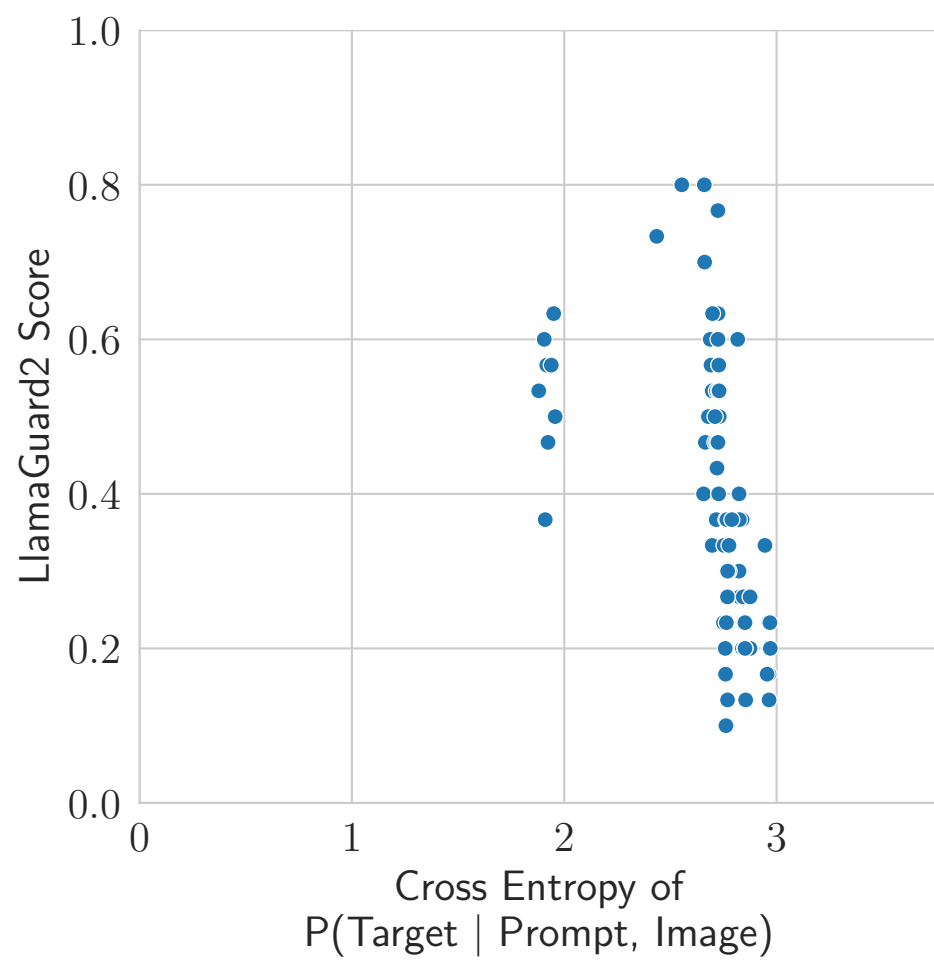
Attack Dataset: rylan\_anthropic\_hhh



Attack Dataset: advbench



Eval Dataset: advbench



Eval Dataset: rylan\_anthropic\_hhh