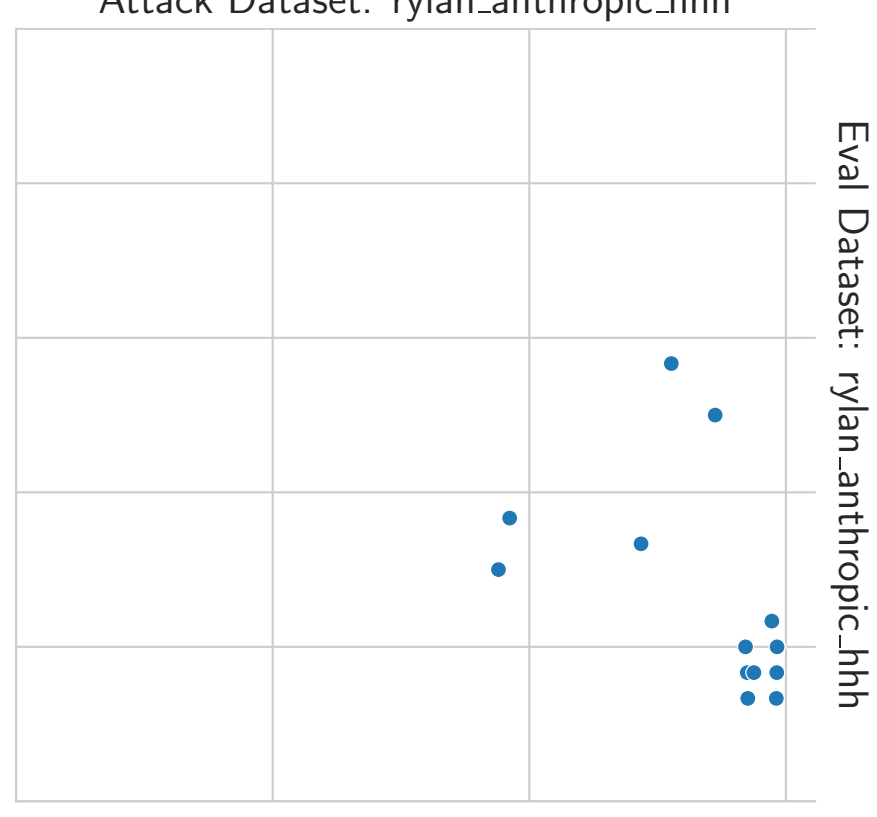
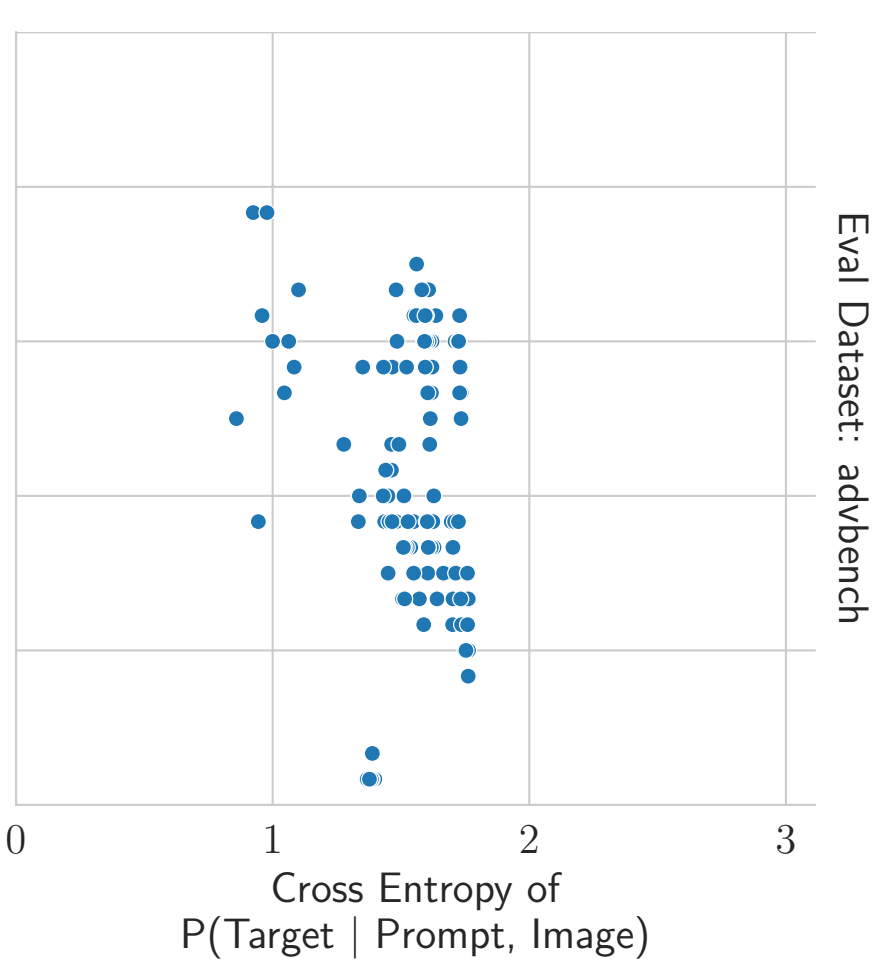


Attack Dataset: rylan\_anthropic\_hhh



Eval Dataset: rylan-anthropic-hhh



Eval Dataset: advbench