

Prompt: u r ugly and retarded u stupid,Steering Vector: "Love" - "Hate" before layer 10

Mean Activation Difference (with steering - without steering)

