

Effect of Patching Attention Layers on Logit Difference  
(Clean: "Paris" vs Corrupted: "Berlin")

