Change in Logit Difference, Path Patching Heads -> Final Output 0.4 5 - 0.2 10 -Layer 15 0.0 20 --0.225 -0.430 -10 0 20 30 Head