Change in Logit Difference, Path Patching Heads -> Final Output 0.4 5 - 0.2 10 Layer 0.0 15 -0.220 -0.425 5 10 0 15 Head