MLP Layers Effect (Denoising)
Letter-String Analogy Task ('+1' vs 'Swap' Rule) 0.8 Average Effect on Logit Difference 0.6 0.4 0.2 0.0 -0.2-0.410 11 Layer