Attention Heads Effect (Denoising)
gptj6b - Letter-String Analogy Task ('+1' vs No Rule)