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