Attention Heads Effect (Noising)
gpt2 - Letter-String Analogy Task ('+1' vs No Rule)