MLP Layers Effect (Noising)
llama27b - Letter-String Analogy Task ('+1' vs No Rule)