

We observe that ROME is much more successful than FT+L at generating text that is consistent with the counterfactual; this finding is consistent with results in Table 4 that show that ROME generalizes better than FT+L. Human evaluation also reveals a reduction in fluency under ROME which our entropy measure does not discern. Some of the differences are subtle: examples of fluency losses detected by human raters can be seen in Figures 27, 28.