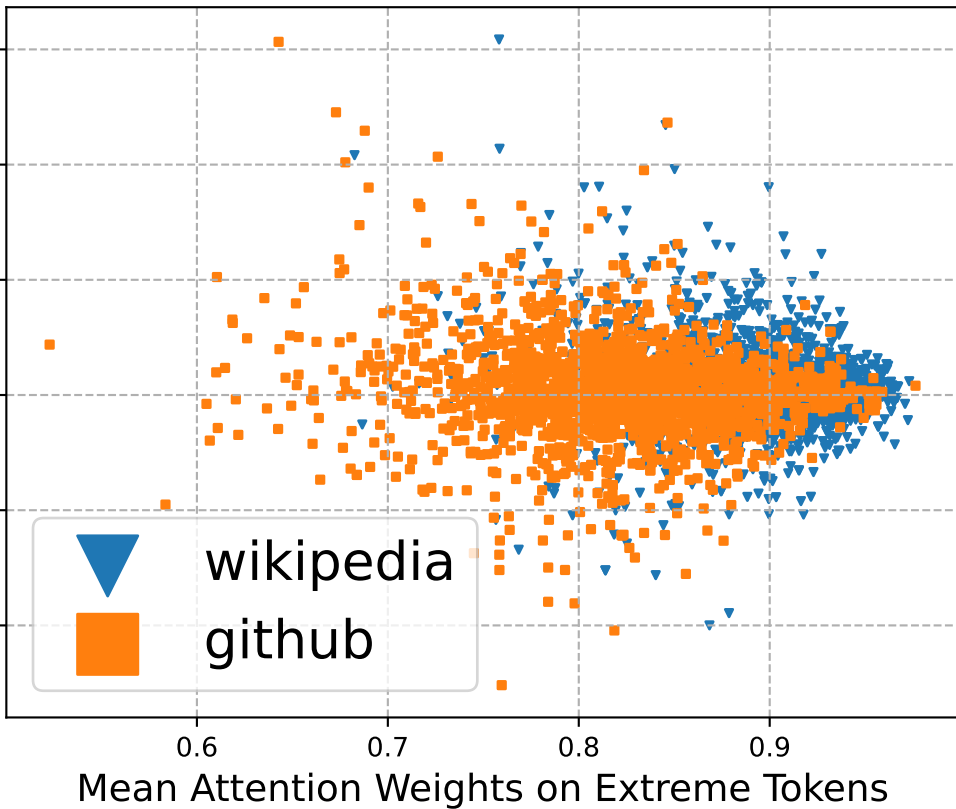


Loss Difference From Intervention



wikipedia



github

Mean Attention Weights on Extreme Tokens