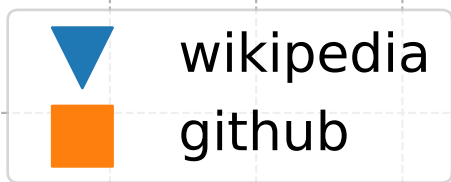


Loss Difference From Intervention



Mean Attention Weights on Extreme Tokens

