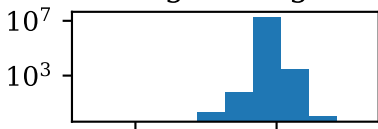
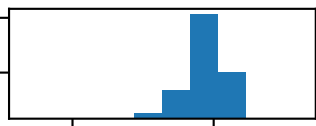


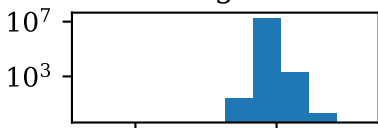
conditional
strengthening243



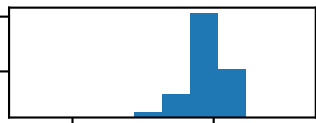
conditional
strengthening243
baseline



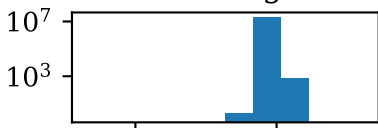
proportional
change243



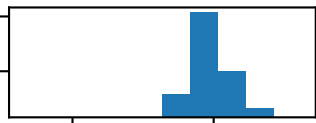
proportional
change243
baseline



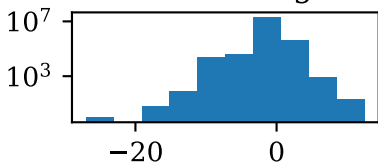
conditional
weakening243



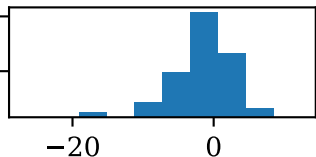
conditional
weakening243
baseline



weakening



weakening
baseline



number of model predictions

$\text{loss}(\text{clean}) - \text{loss}(\text{ablated})$