

Test-time Scaling Baseline

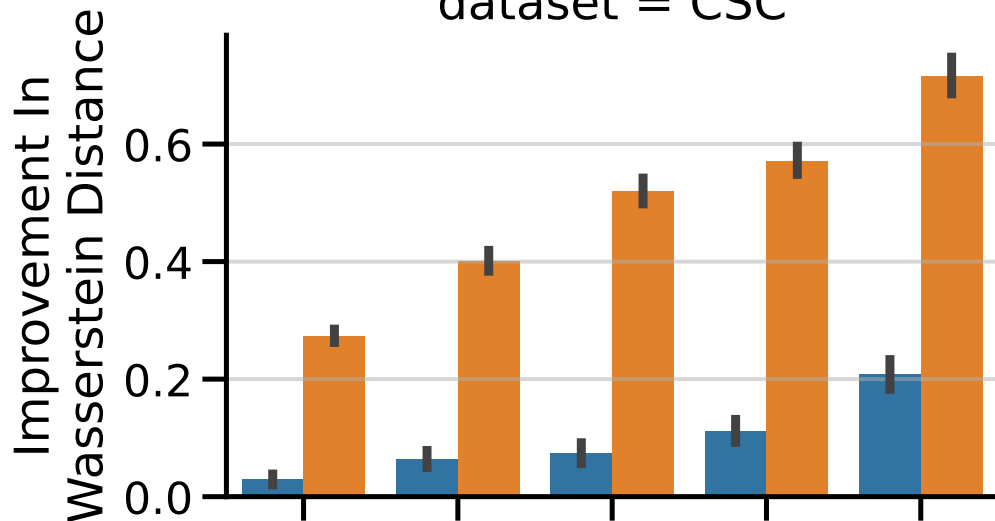


Model Averaging

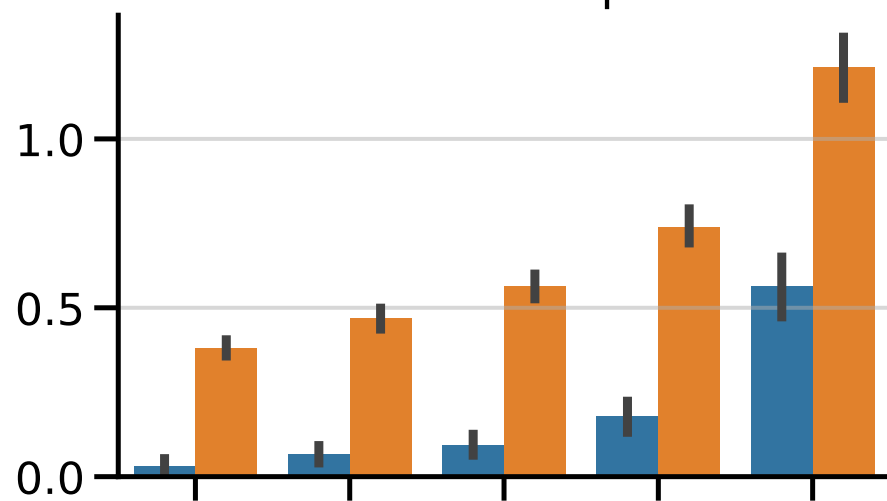


BoN Oracle

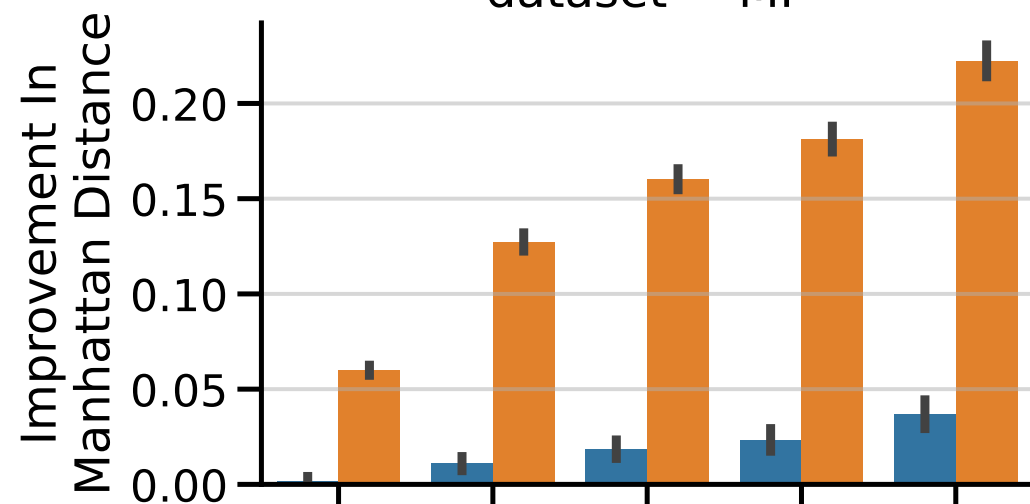
dataset = CSC



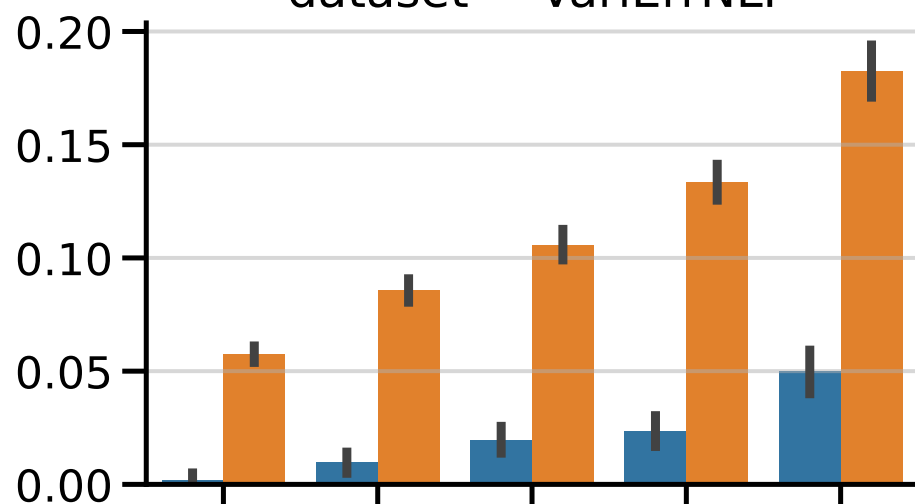
dataset = Paraphrase



dataset = MP



dataset = VariErrNLI



Prediction Diversity

Prediction Diversity