

Generalization Gap (Testset vs. Best Validation Score) on HotpotQA for GPT-4

