Generalization Gap (Testset vs. Best Validation Score) on LiveBench Math for GPT