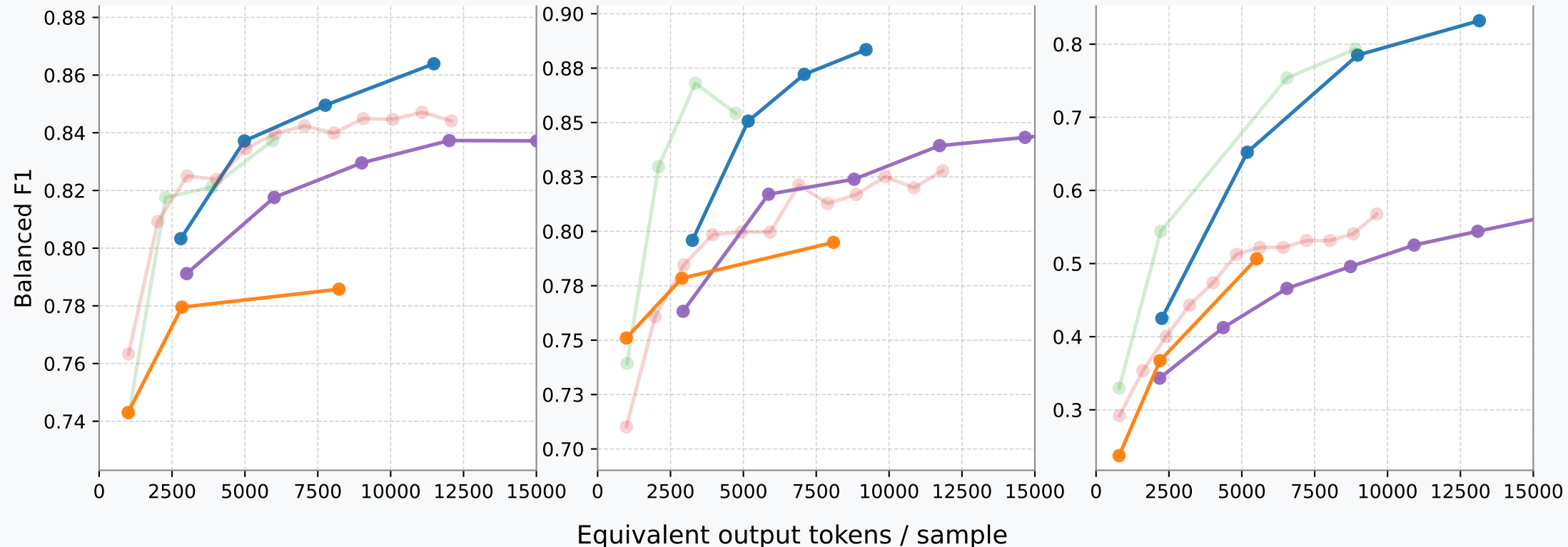


# Performance and Efficiency of Pessimistic Verification Methods (GPT-5-mini)

## IMO-GradingBench

## Hard2Verify

## QiuZhen-Bench



progressive simple pessimistic progressive (low reasoning\_effort) simple pes. (low reasoning\_effort) standard long CoT