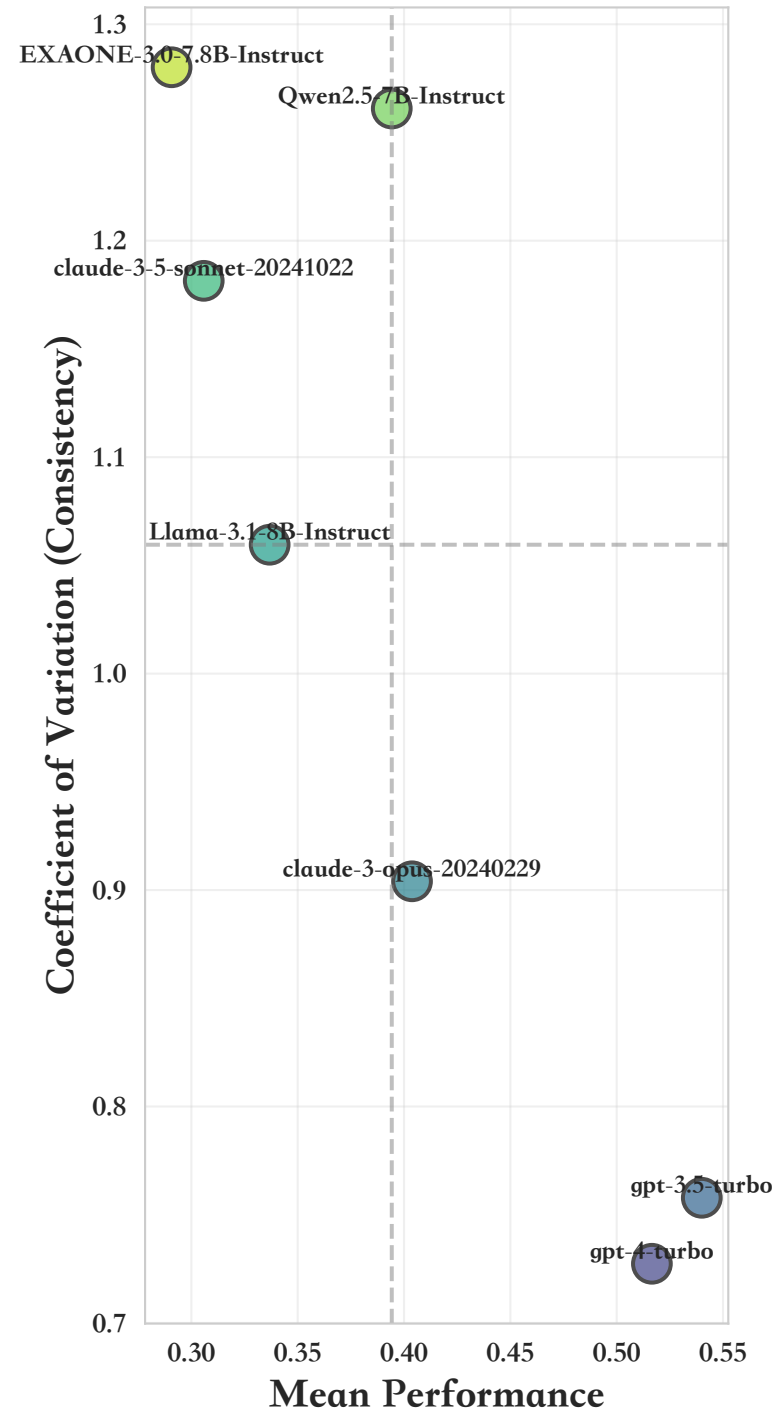


Model Performance vs Consistency
(Lower CV = More Consistent)



Model Performance Range Across Tasks
(Diamond = Mean)

