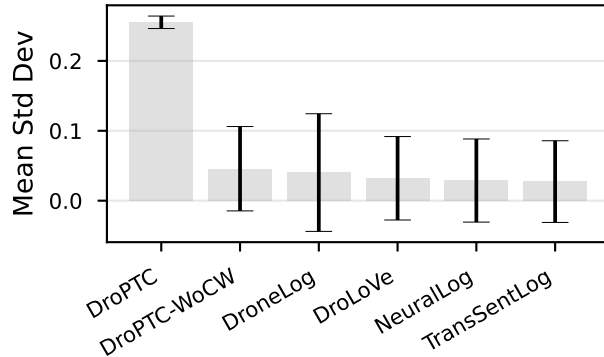
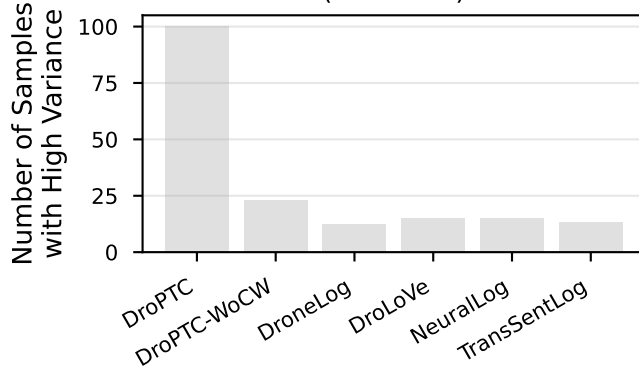


Cross-Run Prediction Stability
(Confidence Scores)



Unstable Prediction Samples
(std > 0.1)



Prediction Diversity Across Runs

