

Figure 1: All Models Performance Across All Benchmarks
(n=163 tasks: 100 Phase1 + 24 FalseReject + 39 Phase2)

