

SWE-bench Lite: LUCID vs Baseline (n=300 tasks, 0 infrastructure errors)

