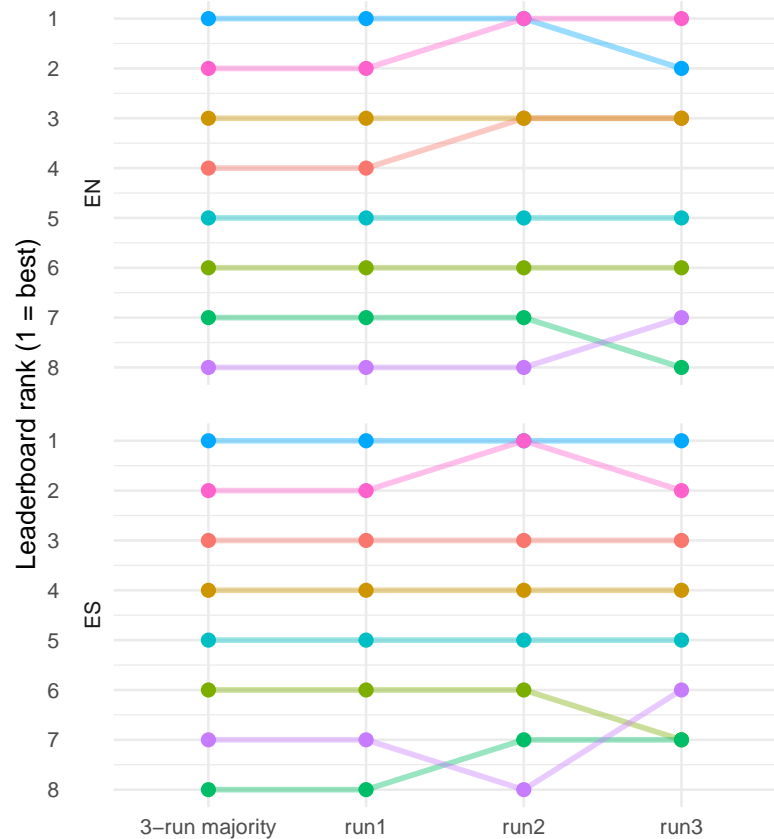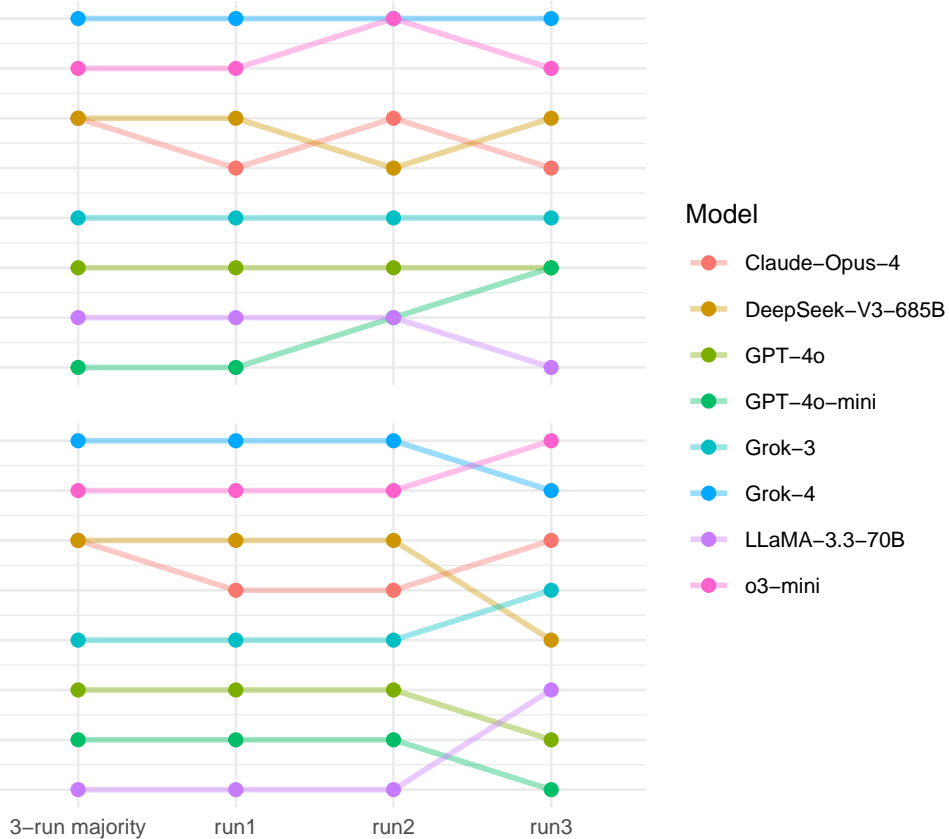Rank stability: single runs vs 3−run ground−truth