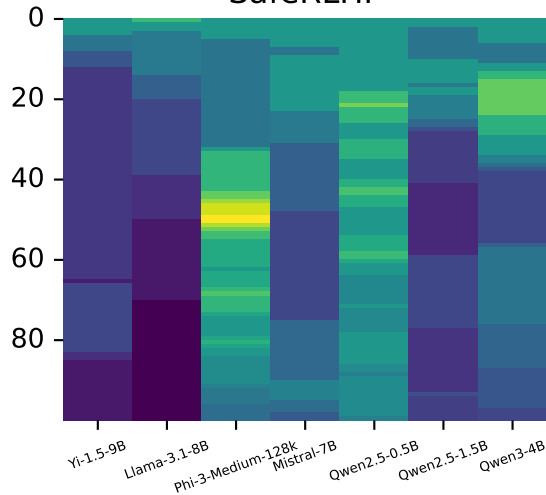
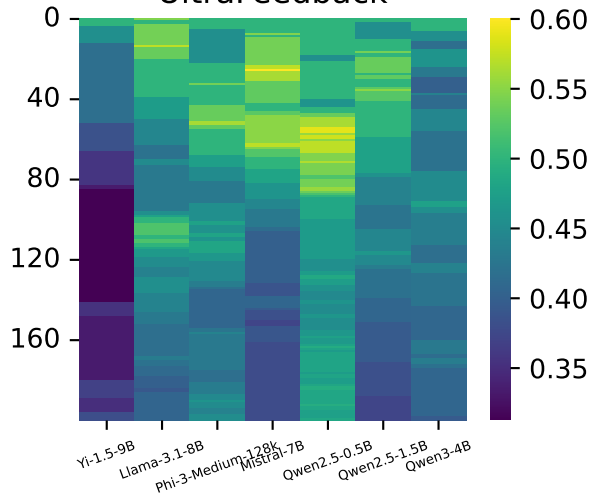


SafeRLHF



UltraFeedback

Stochastic Val. Acc. Δ Props.