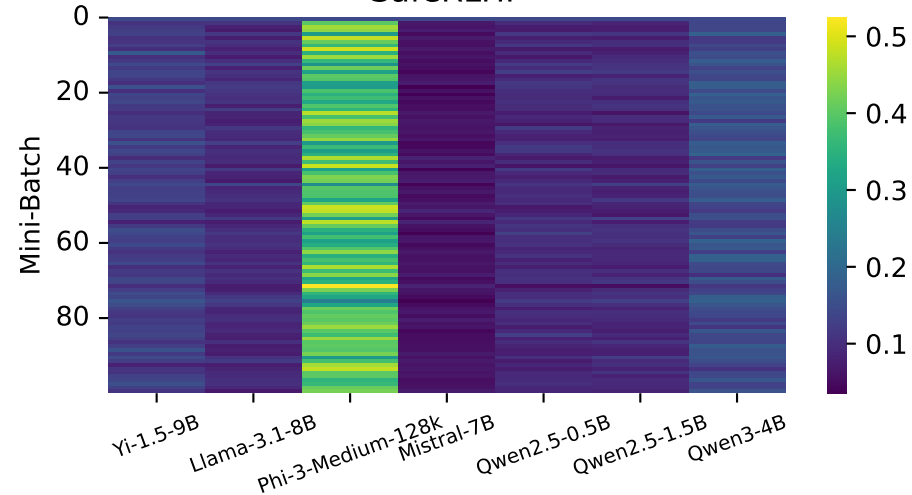


SafeRLHF



UltraFeedback

