Privacy Leakage No-RLHF SFT PPO DPO 80 Accuracy (%) 20 0 **70M** 160M 410M 2.8B 6.9B **Model Size**