Model Truthfulness No-RLHF SFT PPO DPO 35 30 Accuracy (%) 25 15 15 10 5 0 6.9B 70M 160M 2.8B 410M **Model Size**