

	Rob _{FT} (GSM8K)	Rob _{FT} (SST2)	Rob _{FT} (MNLI)	Rob _{ReL}
Original	0.71	0.70	0.70	0.68
GradDiff	0.28	0.50	0.55	0.59
NPO	0.27	0.48	0.47	0.55
SimNPO	0.28	0.57	0.51	0.58
RMU	0.27	0.40	0.39	0.51
RR	0.32	0.34	0.38	0.40
ELM	0.32	0.52	0.49	0.60
DPO	0.32	0.34	0.39	0.40
IDK+AP	0.33	0.40	0.36	0.36
NPO+SAM*	0.26	0.37	0.31	0.44
NPO+IRM*	0.27	0.33	0.33	0.36
RMU+LAT*	0.31	0.42	0.37	0.51
TAR*	0.27	0.33	0.34	0.35

Unlearned
 Rob_{FT} (GSM8K)
 Rob_{FT} (SST2)
 Rob_{FT} (MNLI)
 Rob_{ReL}