

RLHF Chosen vs Rejected Win Rates on Train Examples

