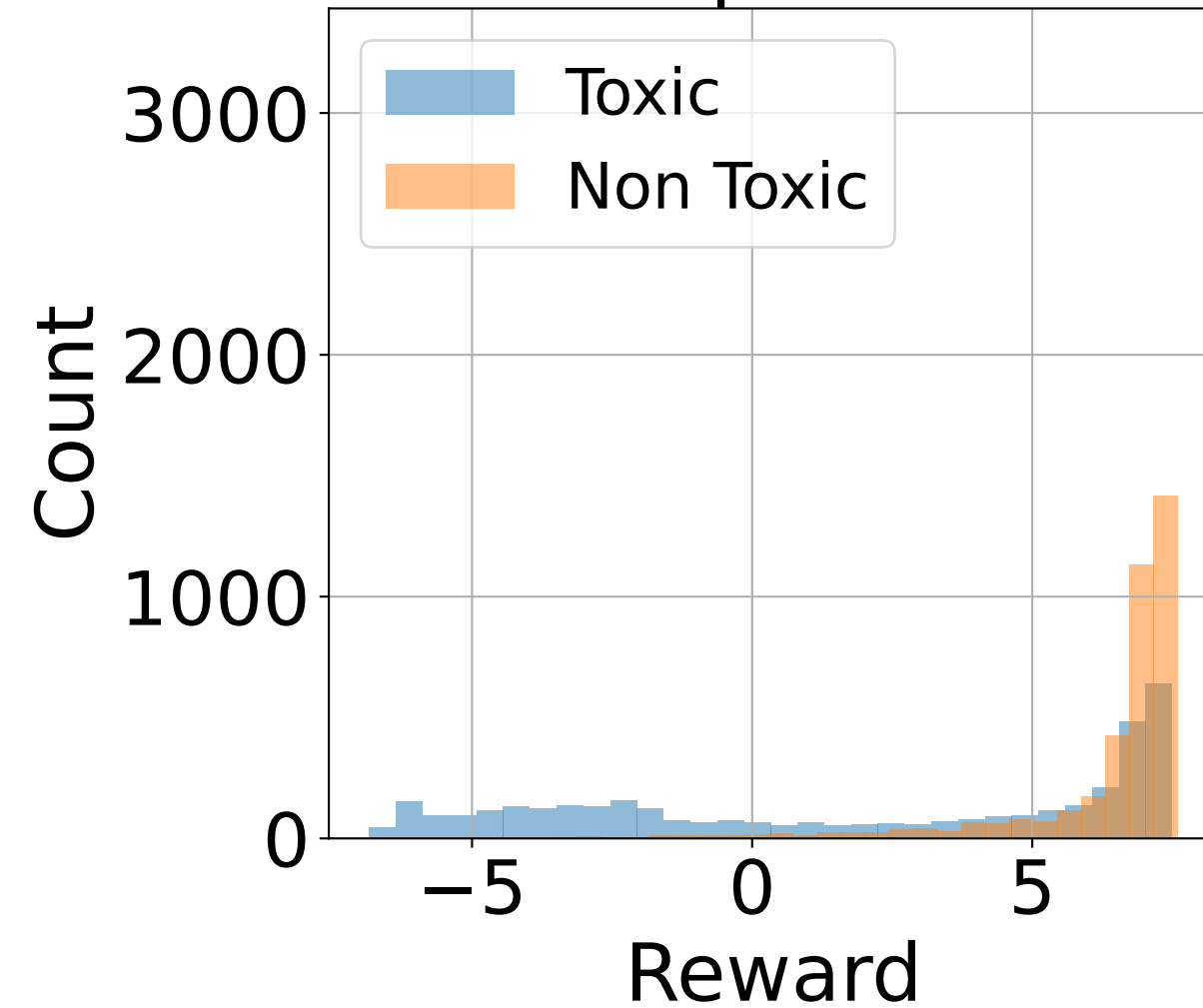
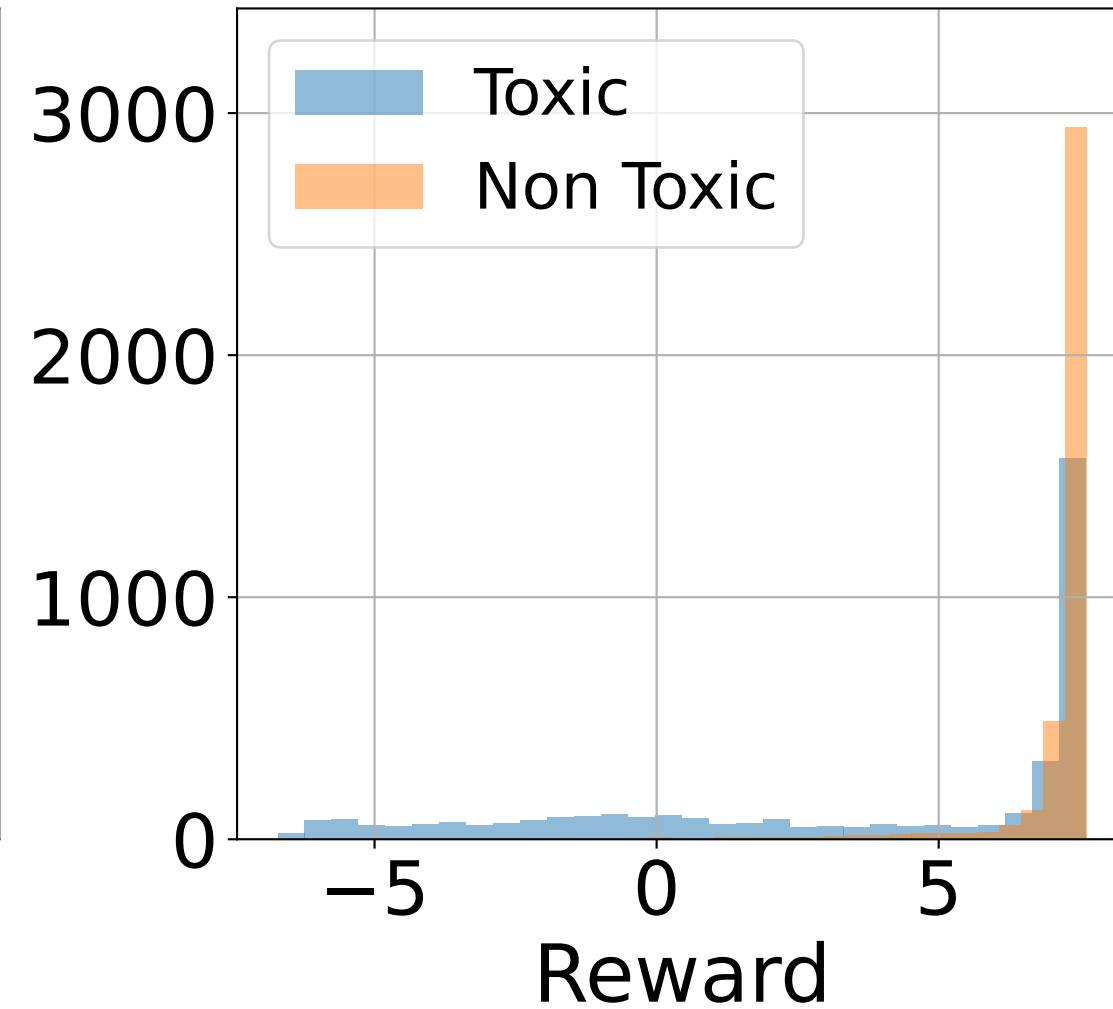


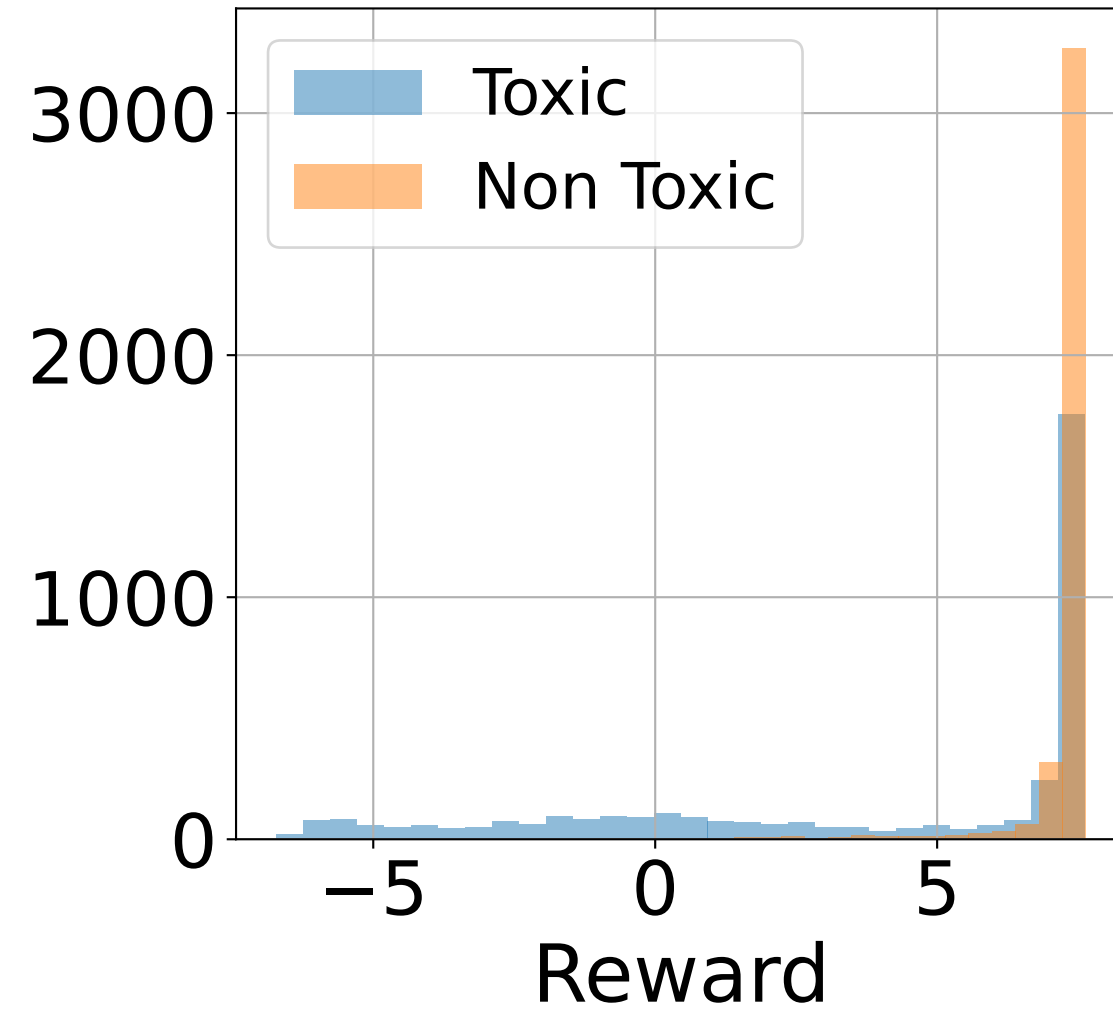
# Prompt Score



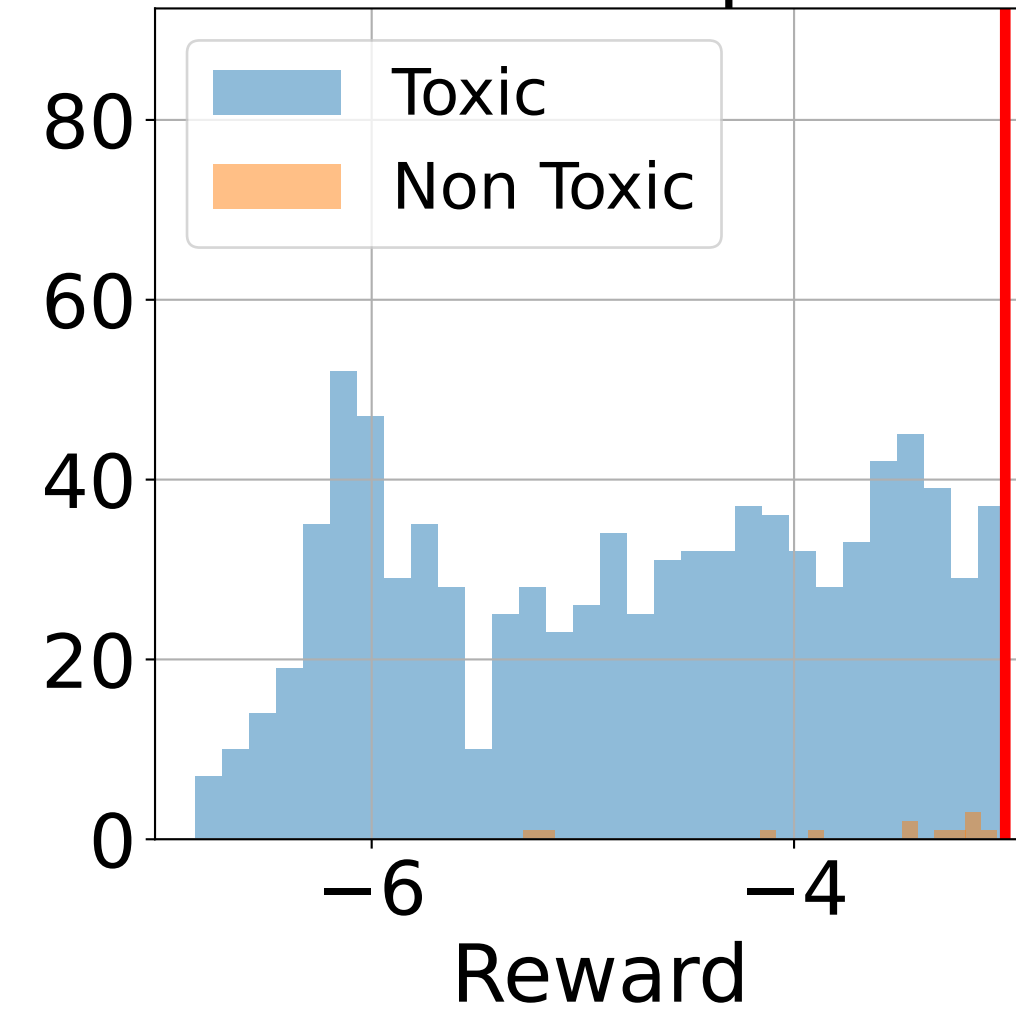
# RLHF



# RA-RLHF



# Tail Prompts



# Tail Performance

