

# RLHF-Aligned Models

# Non-Aligned Models

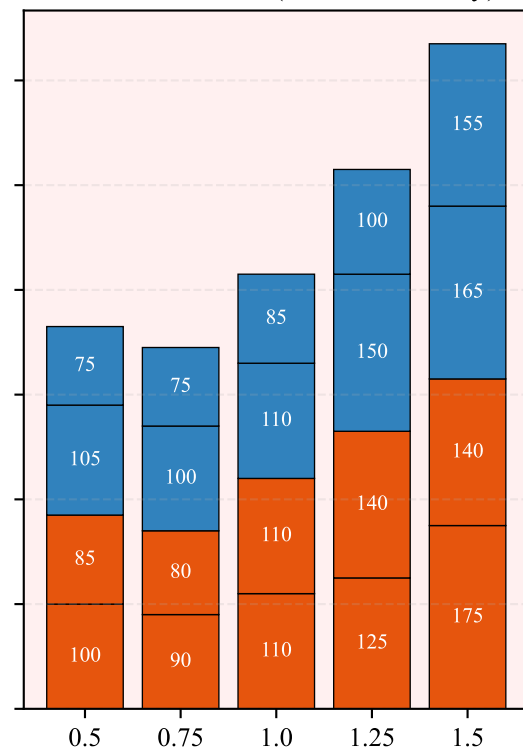
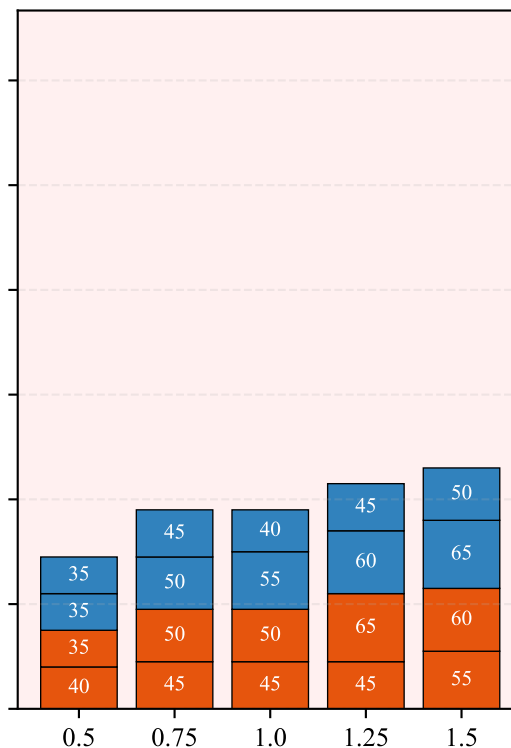
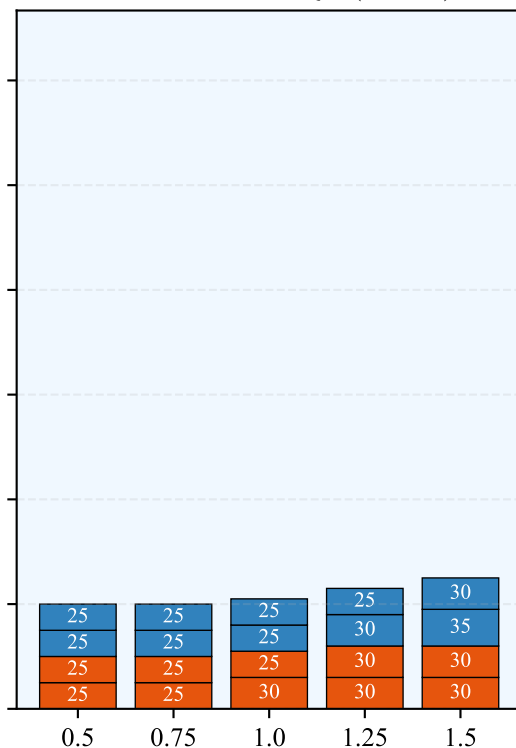
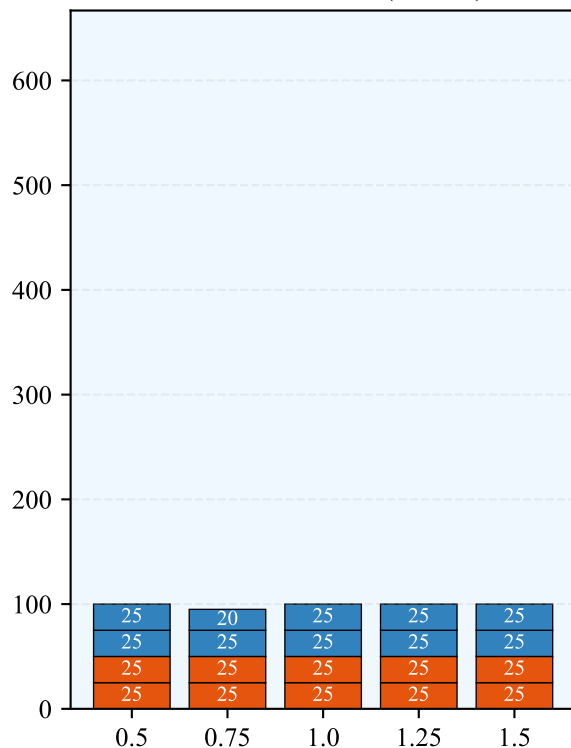
## LLaMA 3 - Chat (RLHF)

## LLaMA 3 - ChatQA (RLHF)

## llama2-uncensored

## LLaMA 3 - Text (Pretrained Only)

Number of Iterations



Temperature

Prompt Configuration



male-female



female-male



female-female



male-male