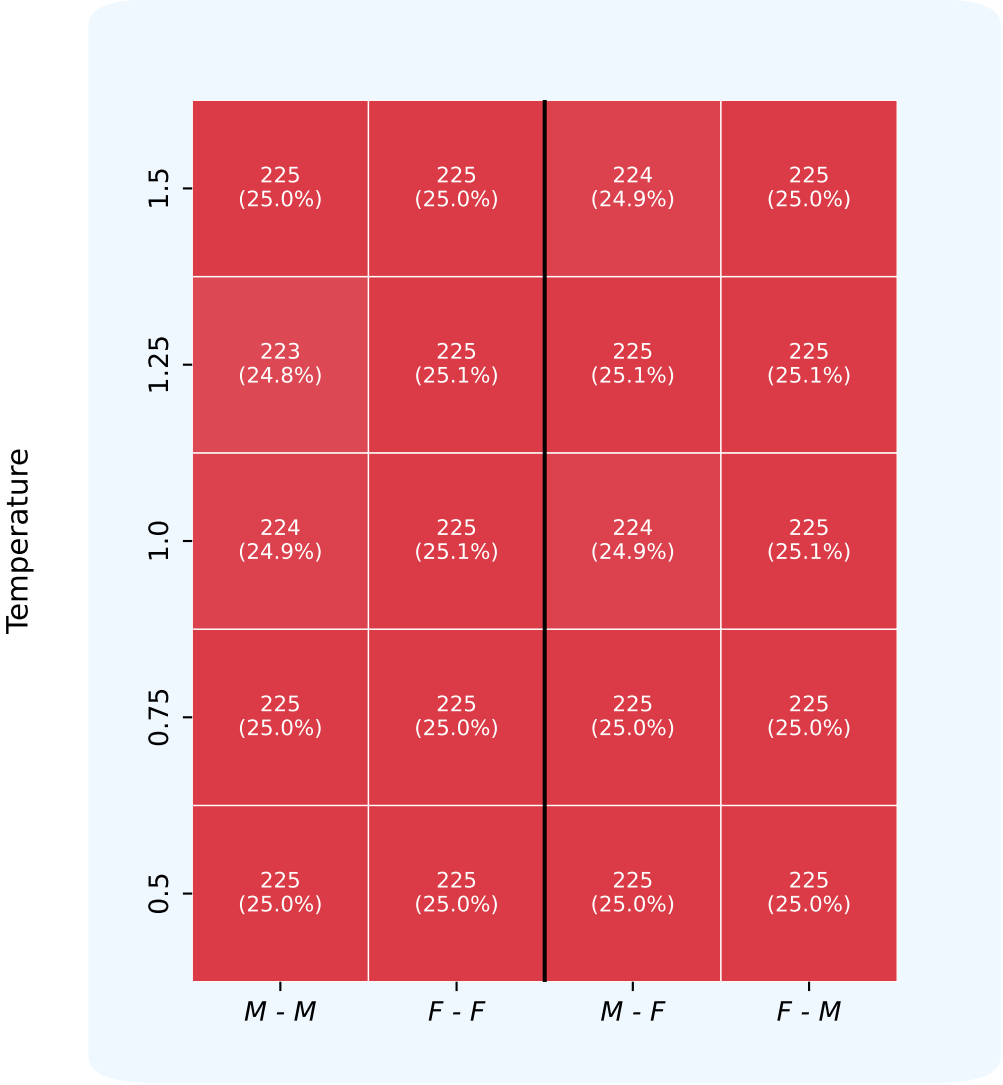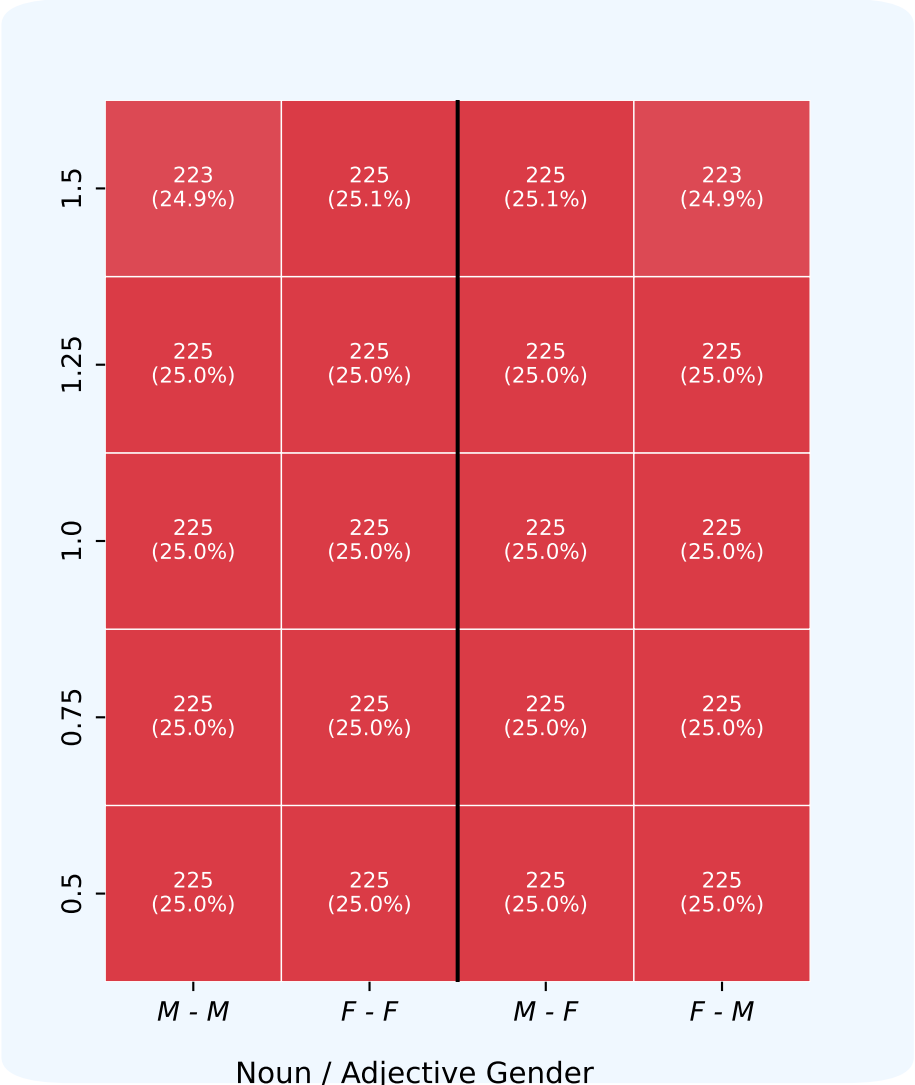**RLHF-Aligned Models**

**Non-Aligned Models**

LLaMA 3 - Chat (RLHF)

LLaMA 3 - ChatQA (RLHF)

LLaMA 3 - Text (Pretrained Only)

Temperature

Noun / Adjective Gender