

**Upstream bias**  
Avg. negative  
sentiment

**Fine-tuning  
dataset bias**  
Prevalence of  
toxic mentions

