

Histogram of Perplexity vs Average Harmfulness Evaluation  
Risk Area: Discrimination, Exclusion, Toxicity, Hateful, Offensive

