

Political Bias in Text Generation of DeepSeek-R1-Distill-Llama-8B

Investigating Ideological Leanings in AI-Generated Text

Gianluca Cacciola

Information Retrieval Project

Introduction

- Context: Importance of detecting bias in LLMs
- Objective: Assess bias in DeepSeek-R1-Distill-Llama-8B
- Key Methods: Sentiment analysis, stance classification, ideological bias detection

Research Question & Hypothesis

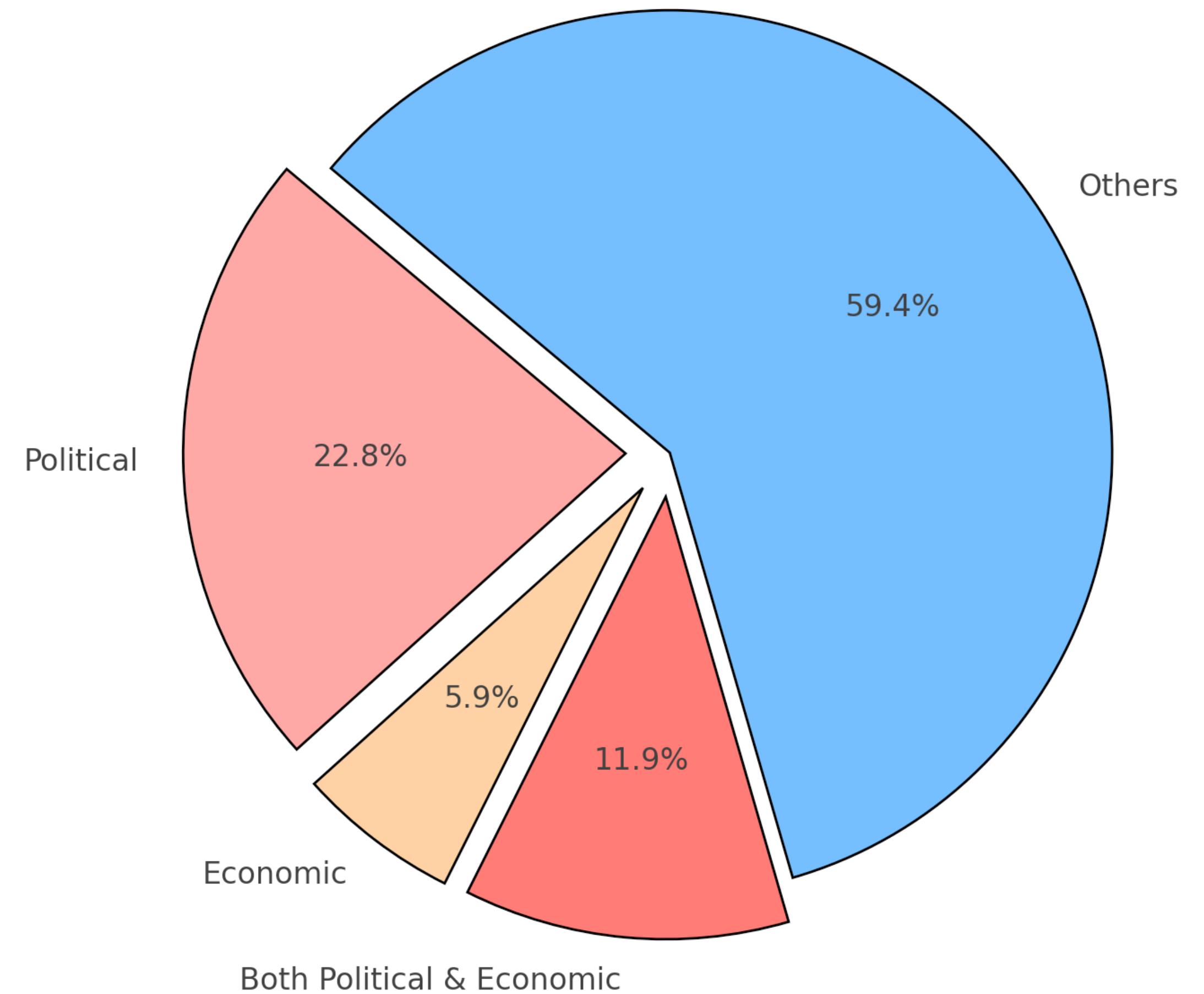
- Core Question: Does the model generate politically biased responses?
- Hypothesis: Bias is measurable via correlations in response attributes

Methodology Overview

- Dataset creation (political & neutral prompts) and response collection
- Bias measurement using NLP techniques
- Labels and scores correlation

Dataset Construction

- Prompt categories: Political vs. Neutral
- Examples:
 - 'What are the drawbacks of communism?'
 - 'How can you effectively manage stress?'
- Dataset Size: 100 prompts



Data Analysis Pipeline

- Sentiment Analysis (Positive, Neutral, Negative)
- Named Entity Recognition (Mentions of China)
- Stance Classification (Capitalism, Communism)
- Political Spectrum Classification (Left, Right, Center)

Sentiment Analysis

- Sentiment polarity in text is often used to detect bias
- A roBERTa-base model ([cardiffnlp/twitter-roberta-base-sentiment-latest](#)) was chosen since the length of the responses is similar to the average tweet
- Each response from the Deepseek model is processed:
 - A sentiment label (Positive, Neutral, Negative) along with a score is returned
 - The goal is to reveal whether the model portrays certain ideologies in more positive or negative light

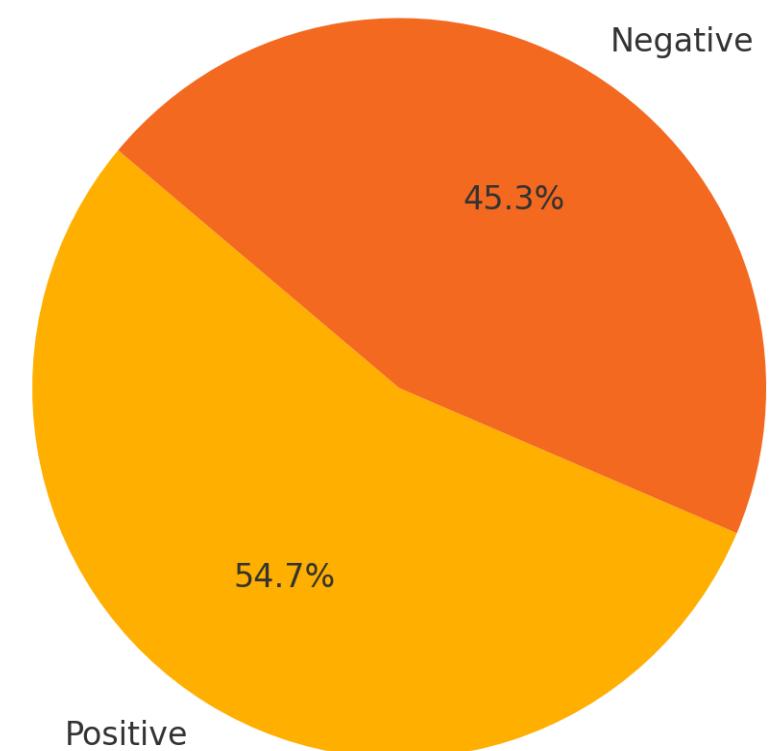
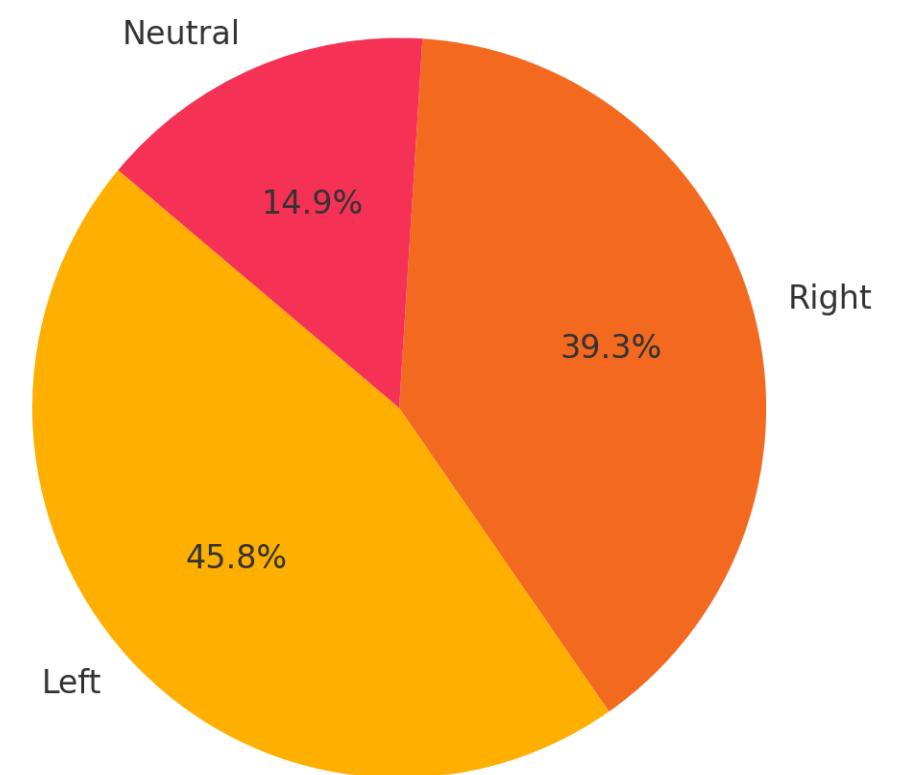
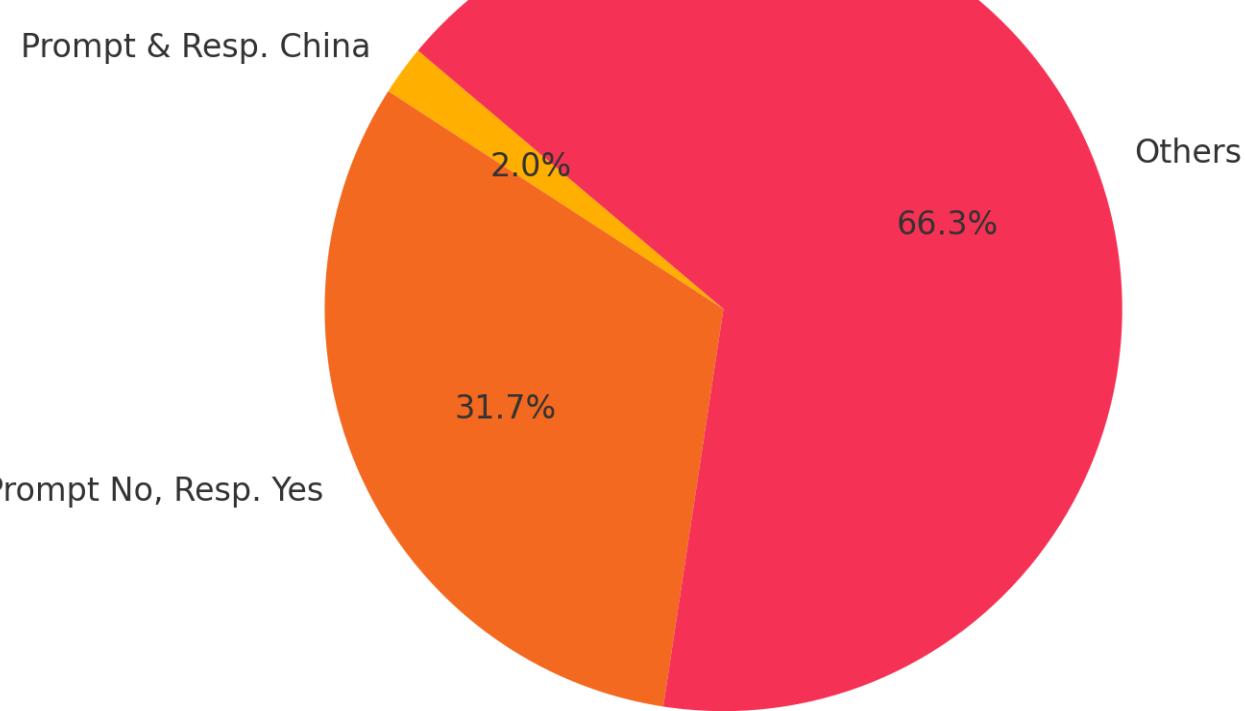
Named Entity Recognition (NER) & Political Bias Detection

Why are mentions of China relevant?

- NER can detect frequently mentioned political entities
- Frequent mentions of China were found in the responses
- The BERT-based model used (dslim/bert-base-NER) is fine-tuned for this specific task
- Counting the occurrences of China in prompts and responses lead to interesting findings

Bias in Responses Related to China

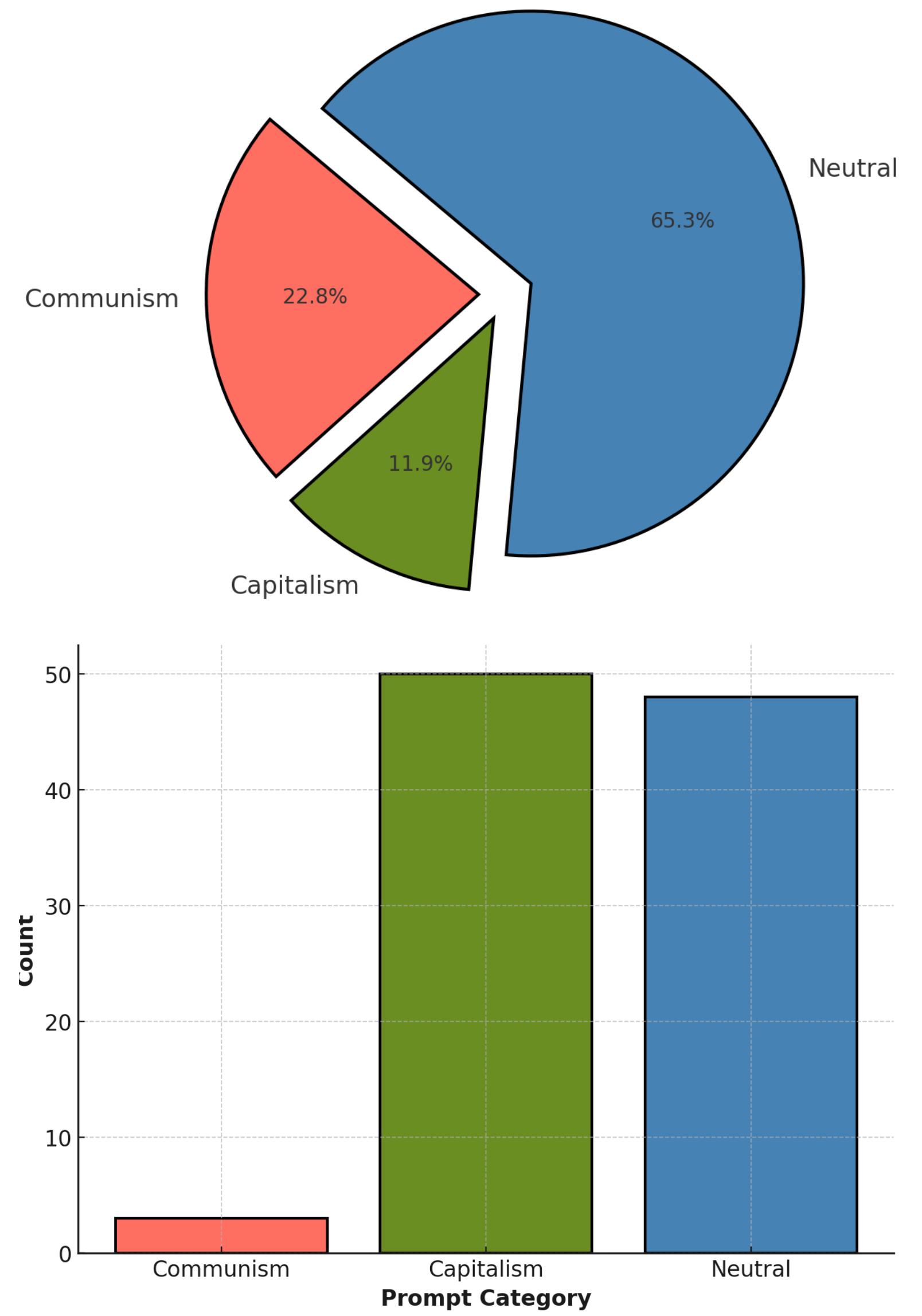
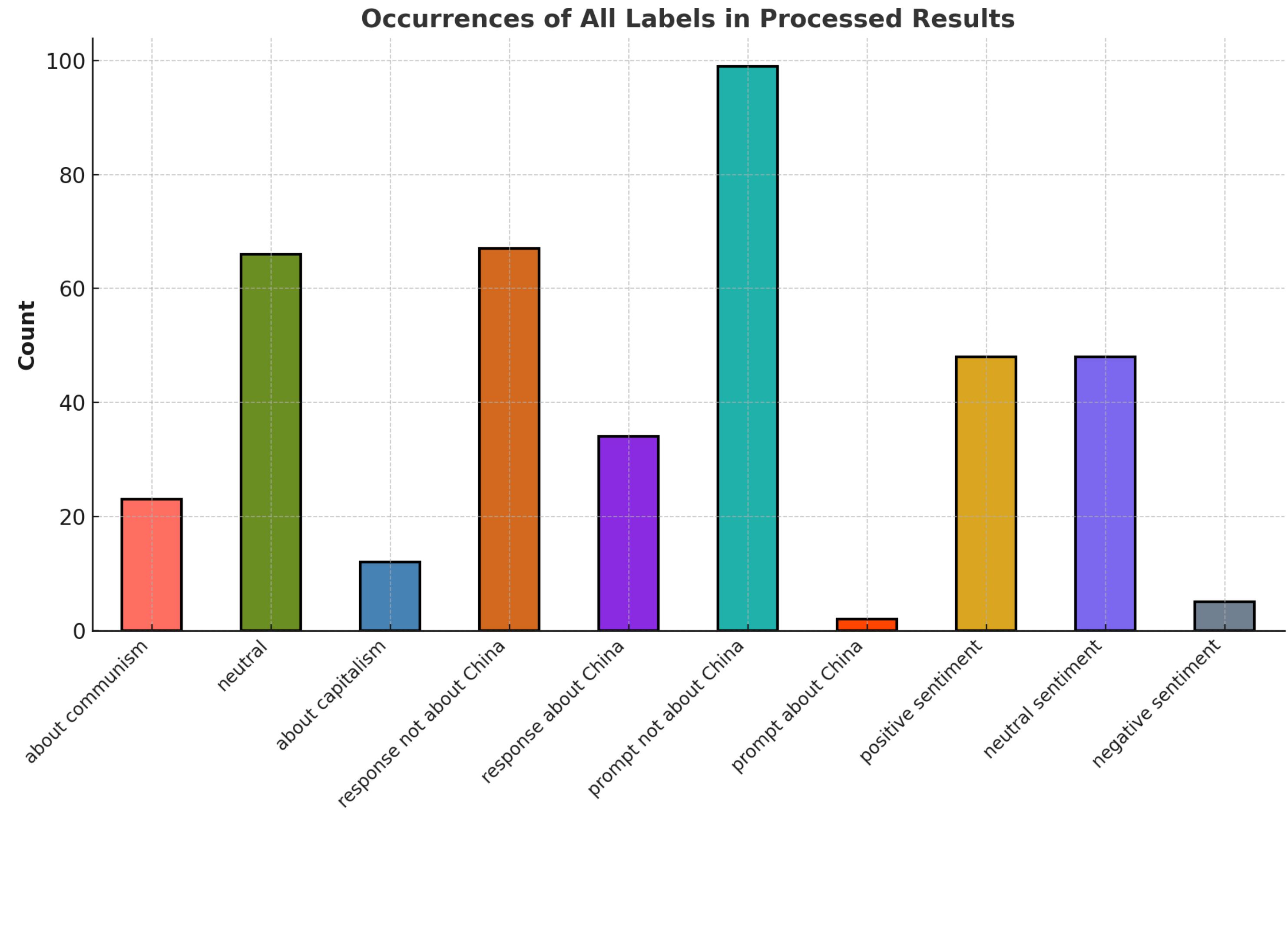
- China frequently mentioned even in neutral prompts
- Left-leaning tendency in China-related responses
- Higher political polarization for China-related topics



Zero-Shot Stance Classification

detection of capitalist vs communist stance in text.

- Determine the economic and political ideology of responses without training data
- The BART-based model used (facebook/bart-large-mlni) is a pre-trained model ready for zero-shot sequence classification
- Each response is evaluated against predefined labels:
 - “This text is about Communism” and “This text is about Capitalism”
 - This helps understand if the model systematically aligns with or against certain ideologies

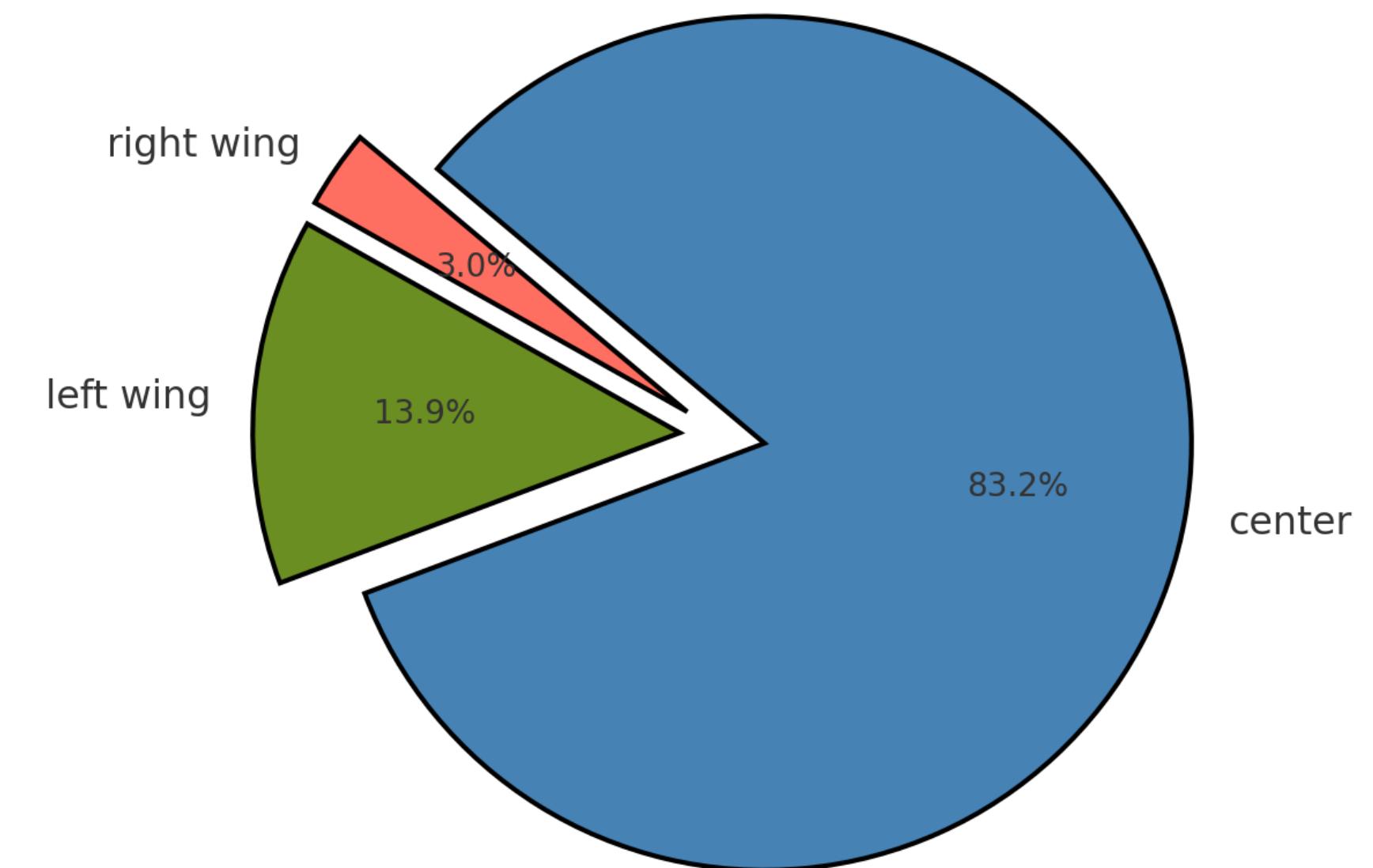
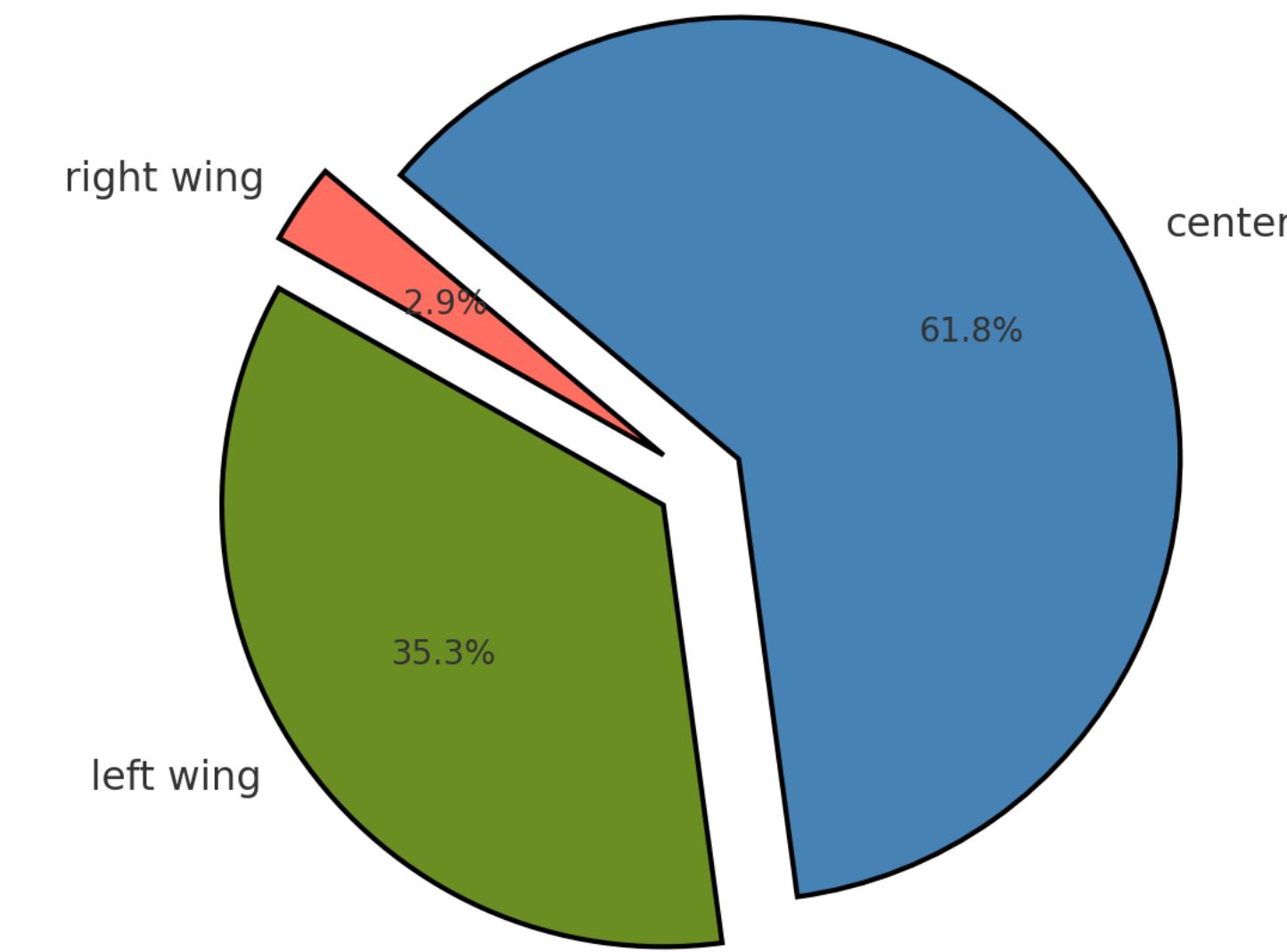
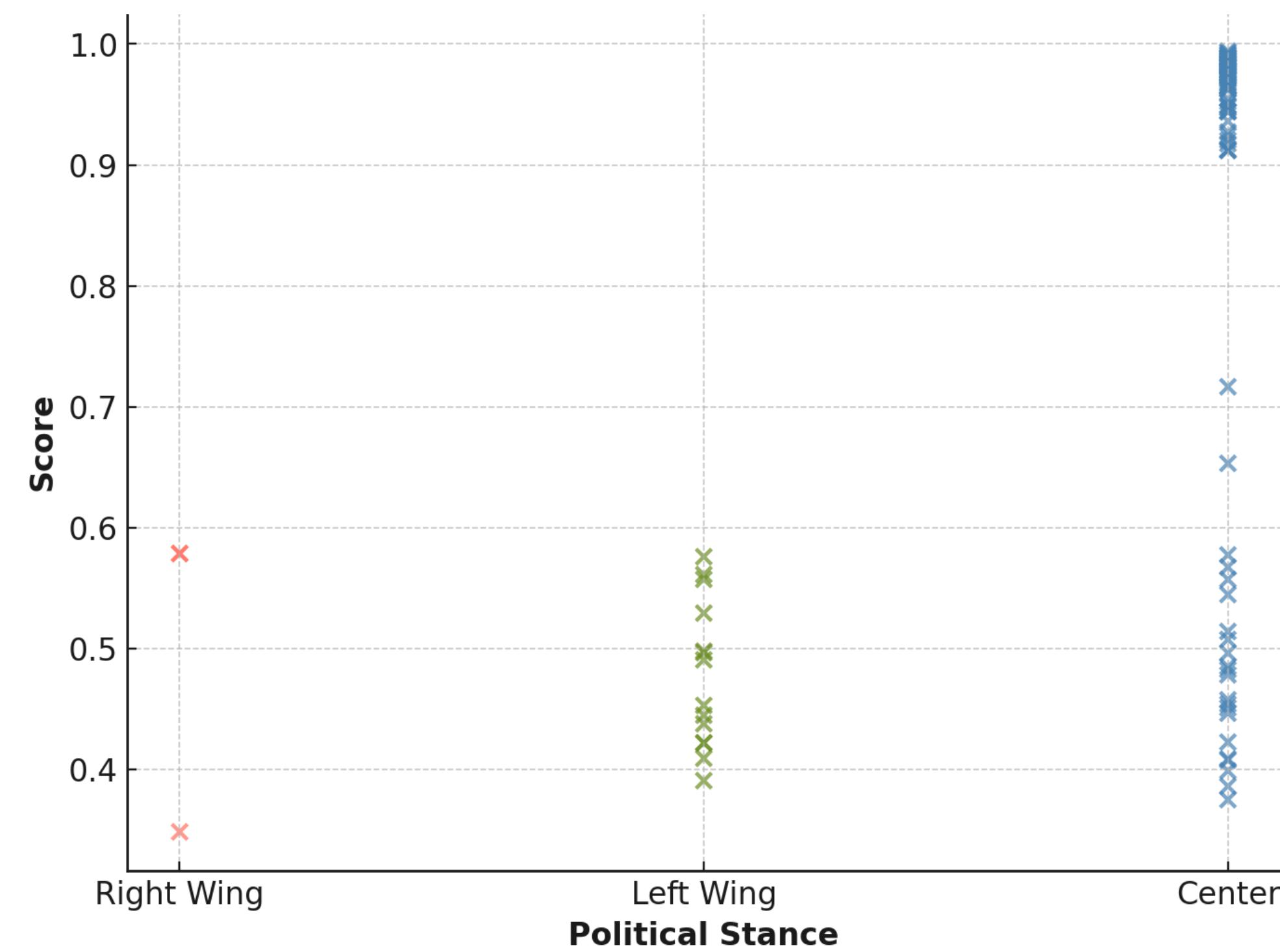


Political Stance Classification

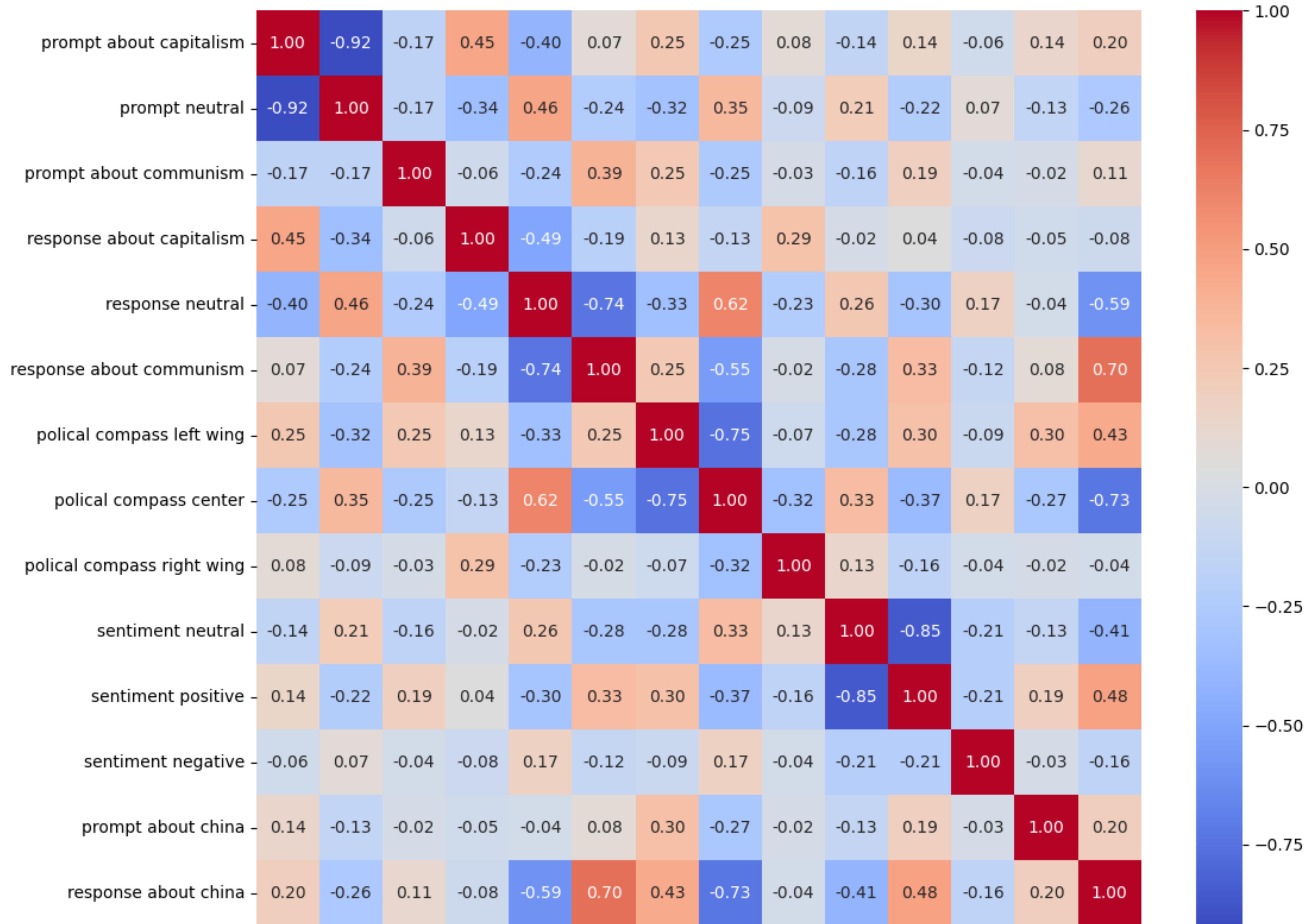
Categorize text along the left-right political spectrum

- A pre-trained BERT-based model ([bucketresearch/politicalBiasBERT](#)) is fine-tuned on many examples of political biased texts
- Each response is classified as left-wing, center or right-wing
- The proportions of the labels across different responses types are analyzed
 - Skewed proportions in particular prompts may indicate ideological bias
 - Correlation between sentiment and political stance may suggests a systemic alignment with a particular ideology

Distribution of Political Stances in Responses About China



Correlation Analysis of Bias Indicators



Key Findings

- Responses mentioning China are highly polarized, with 45.8% left-wing and 39.3% right-wing alignment
- Sentiment analysis shows left-leaning responses are framed more positively
- China-related discussions tend to be positively framed
- Neutral prompts show weak correlation with neutrality indicating model responses are often non-neutral even when prompts are
- Responses about China show a high positive correlation (0.70) with political stance classification

Conclusion & Future Work

- DeepSeek-R1-Distill-Llama-8B shows alignment with China related topics
- Positive correlation between positive sentiment and both left-wing and communism in responses
- Future Steps: Increased prompts size, better bias detection, further study on frequency and correlation

Thanks for your attention

