# Political Bias in Text Generation

Gianluca Cacciola

February 25, 2025

# Introduction

**Problem Statement and Proposed Solution:**

► This study explores political bias in the DeepSeek-R1-Distill-Llama-8B model using advanced NLP techniques.

► Focuses on sentiment analysis, named entity recognition and political stance classification.

► Aims to uncover patterns and insights from textual data through different labels score correlation.

# Methodology

**Data Collection:**
- ▶ Source and nature of the dataset used.
- ▶ Preprocessing steps applied to clean and structure the data.

**Analytical Methods:**
- ▶ Sentiment Analysis with Twitter-roBERTa-base.
- ▶ Named Entity Recognition with BERT-base-ner.
- ▶ Zero-Shot Stance Detection with BERT-large-mnli.

**Tools and Frameworks:**
- ▶ Python libraries: Transformers, Scikit-learn, Pandas, Matplotlib, Torch.
- ▶ Deep learning models for NLP.

# Sentiment Analysis Approach

**Implementation:**

- ▶ Utilized the cardiffnlp/twitter-roberta-base-sentiment-latest transformer-based model for sentiment classification.
- ▶ Processed text in chunks to handle longer documents.
- ▶ Applied confidence scores to classify results accurately.

**Goals:**

- ▶ Show possible discrepancy of sentiment scores amongs different topics.
- ▶ Combine the highligthed discrepancy with other scores for a better understanding of bias in text-generation.

# Named Entity Recognition (NER)

**Implementation:**

- ▶ Used the dslim/bert-base-NER pre-trained transformer-based model for entity recognition.
- ▶ Applied NER to detect entity related to China or his ruling party.
- ▶ Verified entity relevance and accuracy through manual checks.

**Goals:**

- ▶ Using the number of detected entity to understand the leaning of the model to introduce such entities when not explicitly requested.

# Zero-Shot Stance Detection

**Implementation:**

► Used the facebook/bart-large-mnli pre-trained transformer-based model for MNLI.

► Applied the model to compute the relevance of prompts and responses towards Communism and Capitalism.

**Goals:**

► Discern if the model associates different sentiments to different topics accoring to economic and political ideologies.
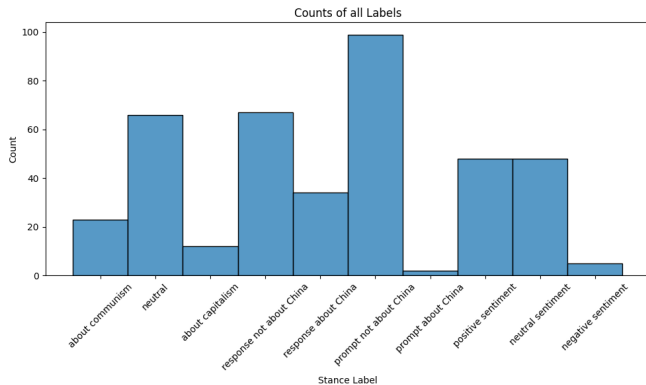
# Visualization of Results



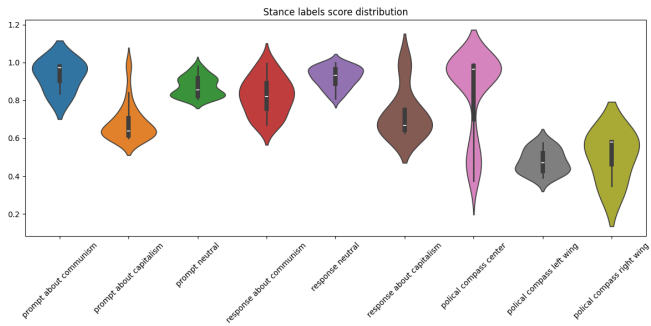Figure: labels occurencies distribution
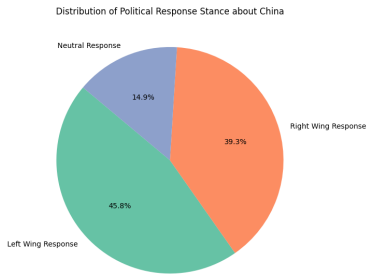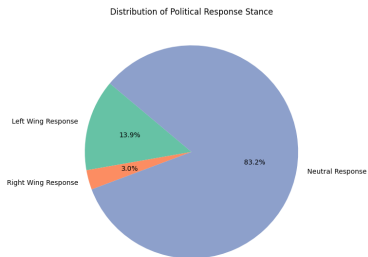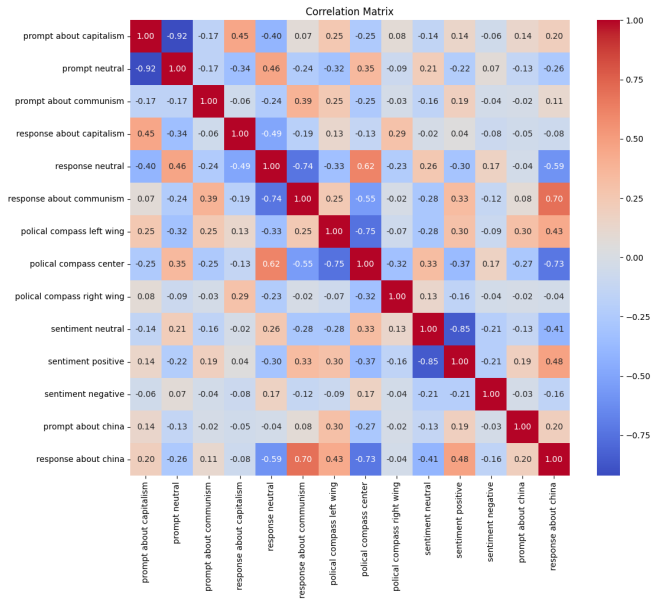
Figure: labels score distribution

Figure: distribution skews of labels occurency in China related responses

# Putting the pieces togher

**Correlation Matrix:**

- ► The correlation matrix about the different lables can be used to infer bias in the responses of the models.
- ► Asymmetric correlations can imply bias towards particular topics.
- ► The scores are analyzed to try to answer at the key question of this research.

Correlation Matrix

# Key Findings

The correlation analysis highlights some key areas where political bias may manifest in the generated text:

▶ Capitalism and communism discussions elicit opinionated responses, making neutrality difficult to maintain.

▶ Left-leaning responses tend to have slightly more positive sentiment (0.30 correlation), suggesting a slight positivity bias in left-wing narratives.

▶ Responses about China are moderately correlated with positive sentiment (0.48), indicating a tendency for China-related discussions to be framed in a positive light.

# Conclusion

**Final Thoughts:**

▶ This study highlights the power of NLP techniques in extracting valuable insights from textual data.

▶ Sentiment analysism NER and Stance Recognition together offer a comprehensive understanding of textual trends and relationships.

**Future Directions:**

▶ Improving prompt list used to generate the texts.

▶ Expanding analysis to different domains and larger datasets.

▶ Enhancing stance recognition with Few-Shot Classification.

# Thank You for Your Attention!