

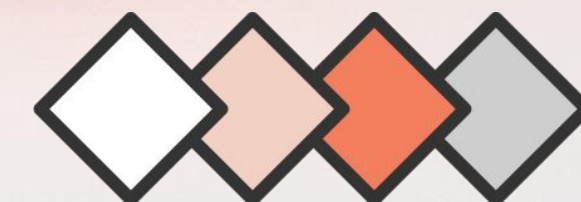
Countries in the News

FDAC 545

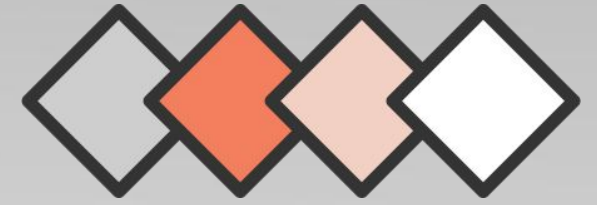


Faith Chernowski, Harshvardhan, Colin
Canonaco, Jeremiah Augustine

December 3, 2024

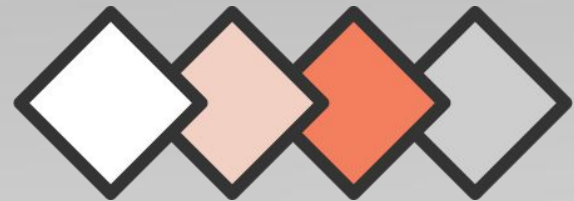


Motivation



- Analyze the portrayal of countries in international media using sentiment analysis
- Identify patterns in news coverage that can help to understand their influence on global perception and international relations
- Reveal biases and disparities in news coverage



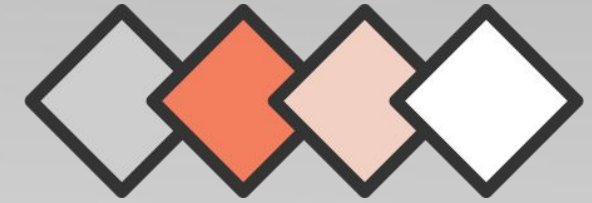


Dataset and Sampling



- GDELT 1.0: Open-access global database tracking events since 1979
 - 2013-2024 = 8 GB, 2.5% random sample -> 200 MB
 - Google BigQuery and Google Cloud Storage
- Attributes used:
 - Event date, country affiliation, stabilizing score, sentiment score, source
- Data Scope
 - Articles from 2013-2024 from publishers like The New York Times, BBC, Xinhua, etc
 - 1,868,683 articles included

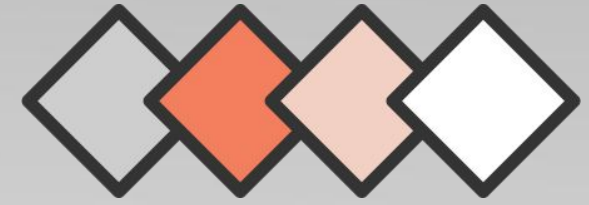
Methods



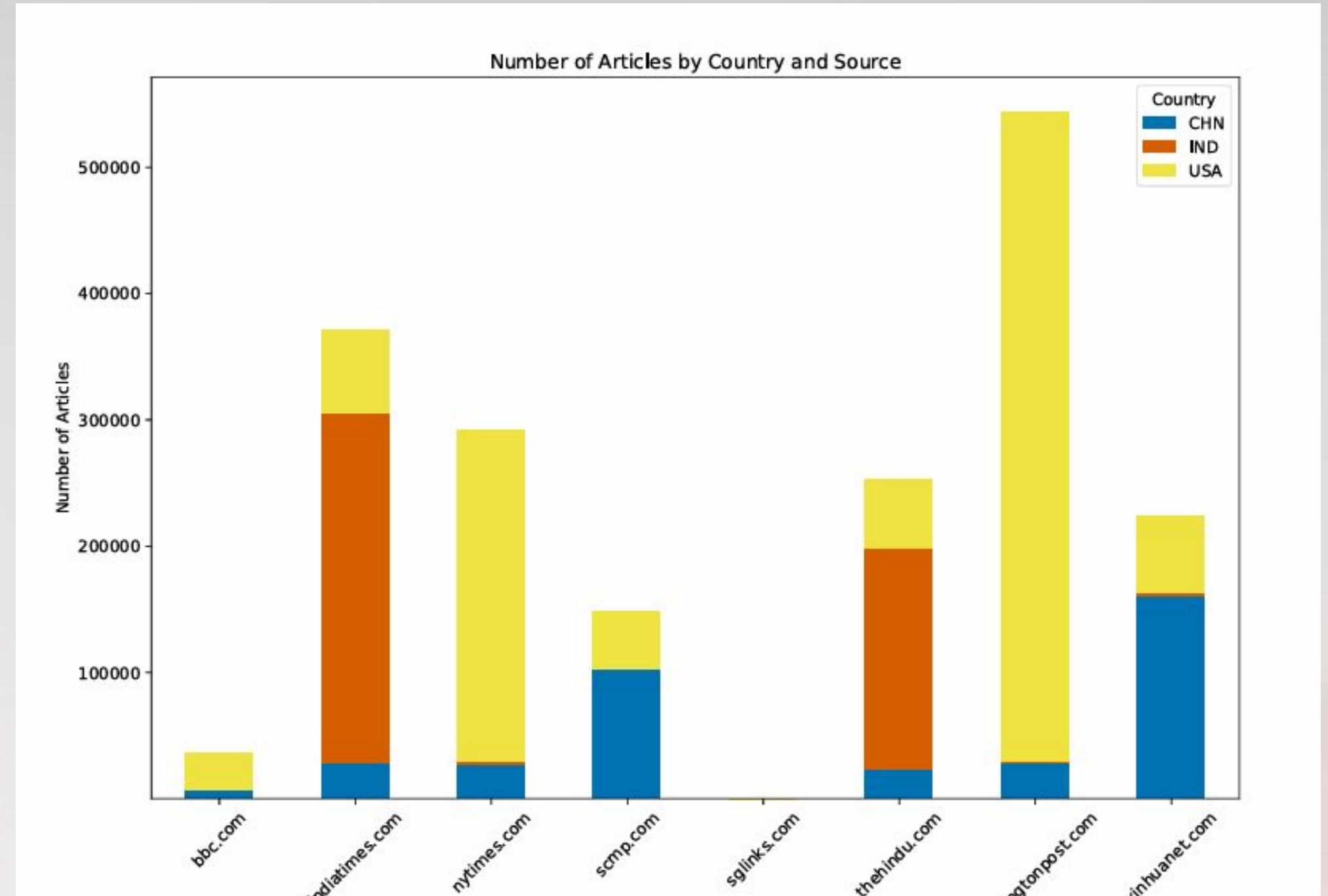
- Data Preparation:
 - SQL query using Google BigQuery
 - CSV with date, countries, tone, Goldstein, source
 - Filtered query to focus on specific sources
- Visualization
 - Color-coded matrices for individual relationships
 - Bar and line graphs for broader trends
 - Sorting axis labels to distinguish patterns

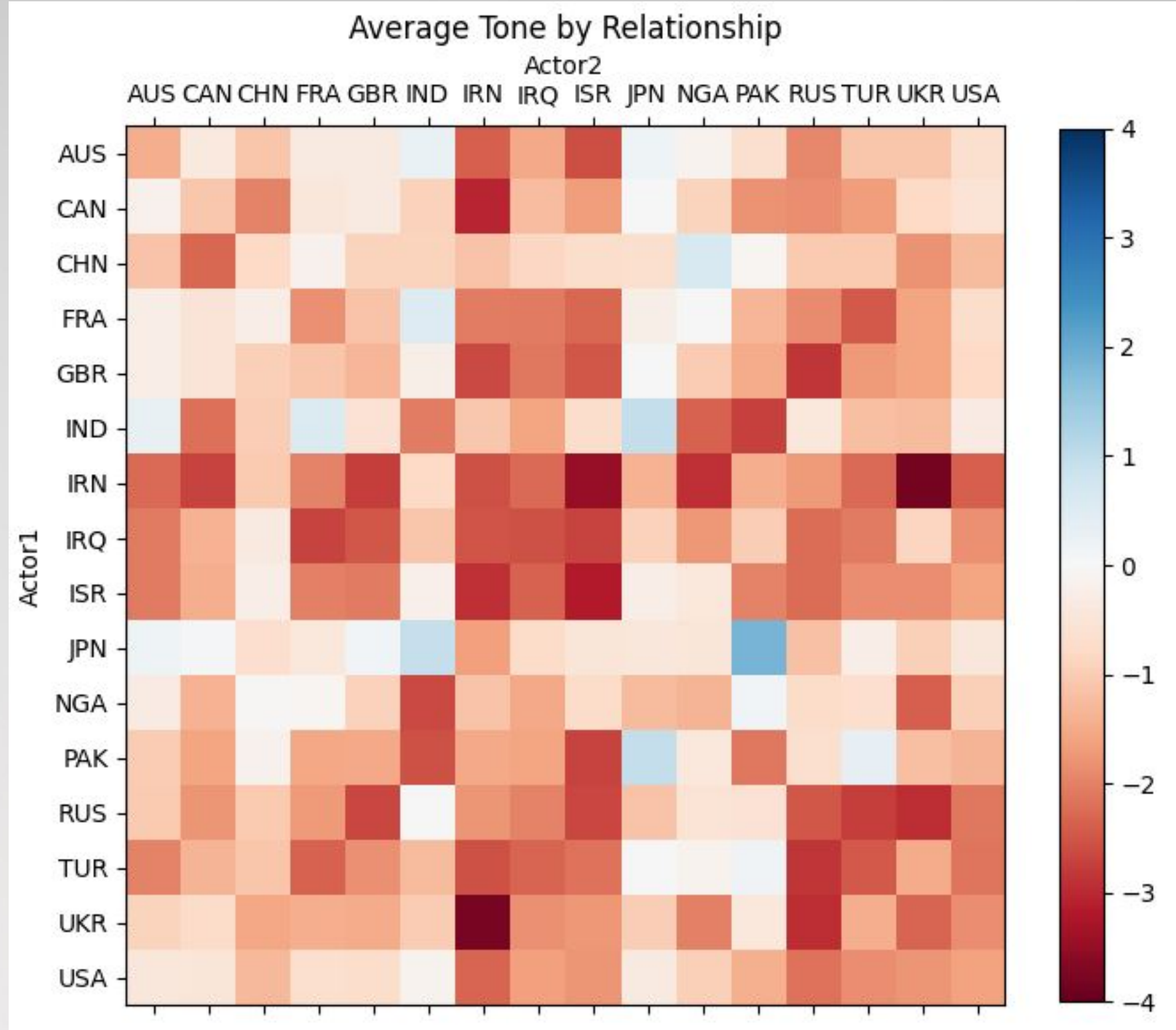
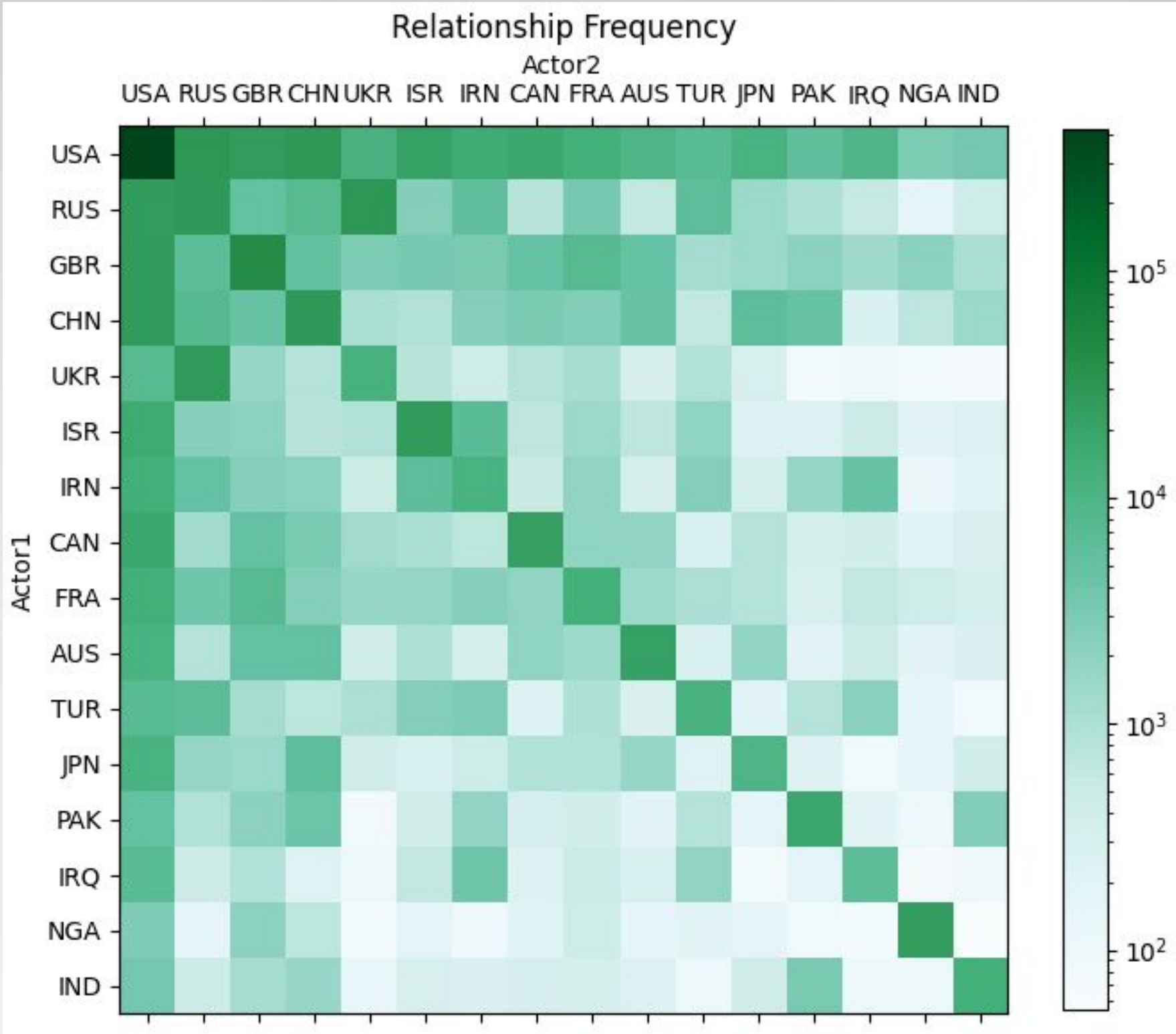


Results

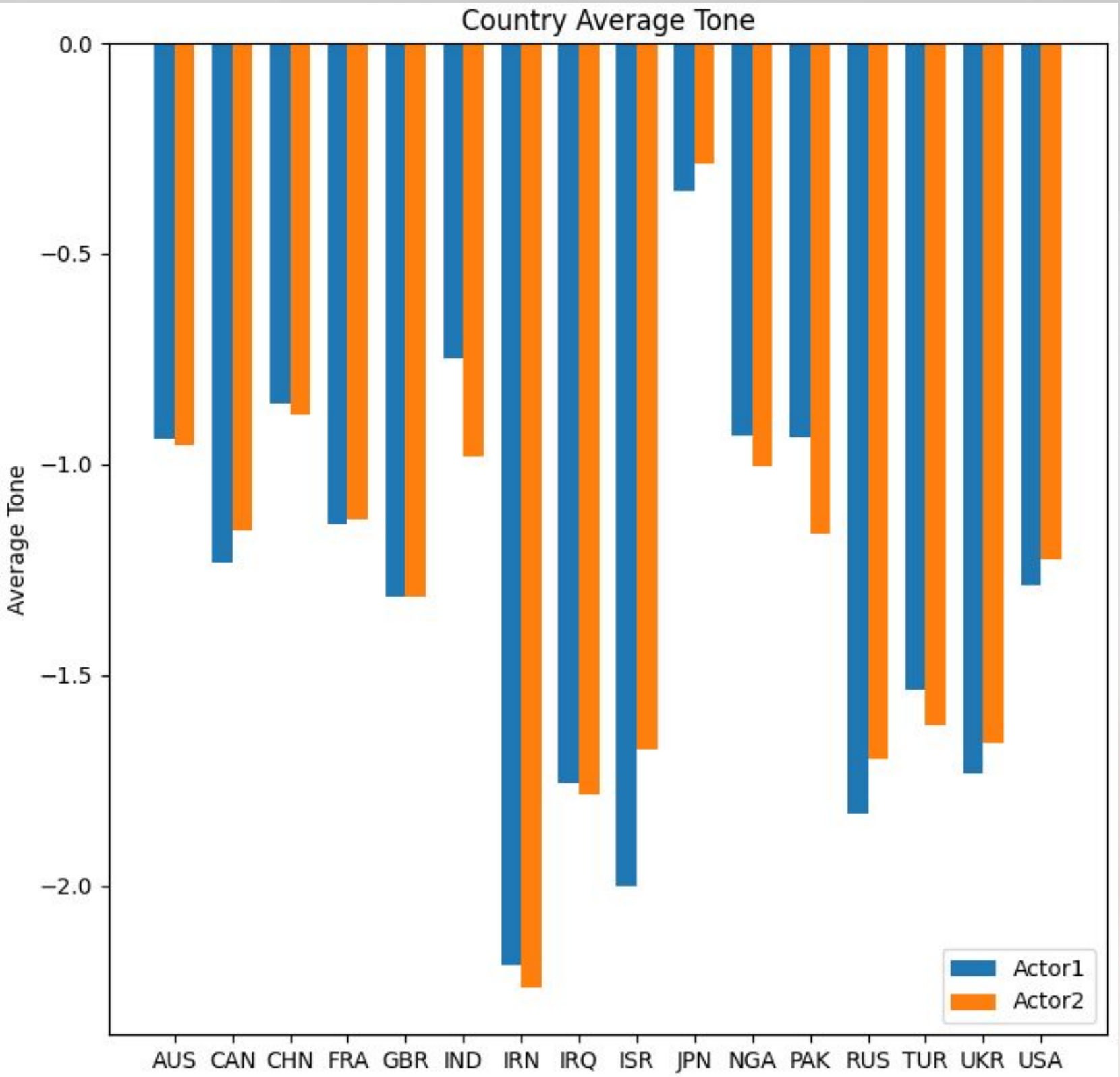
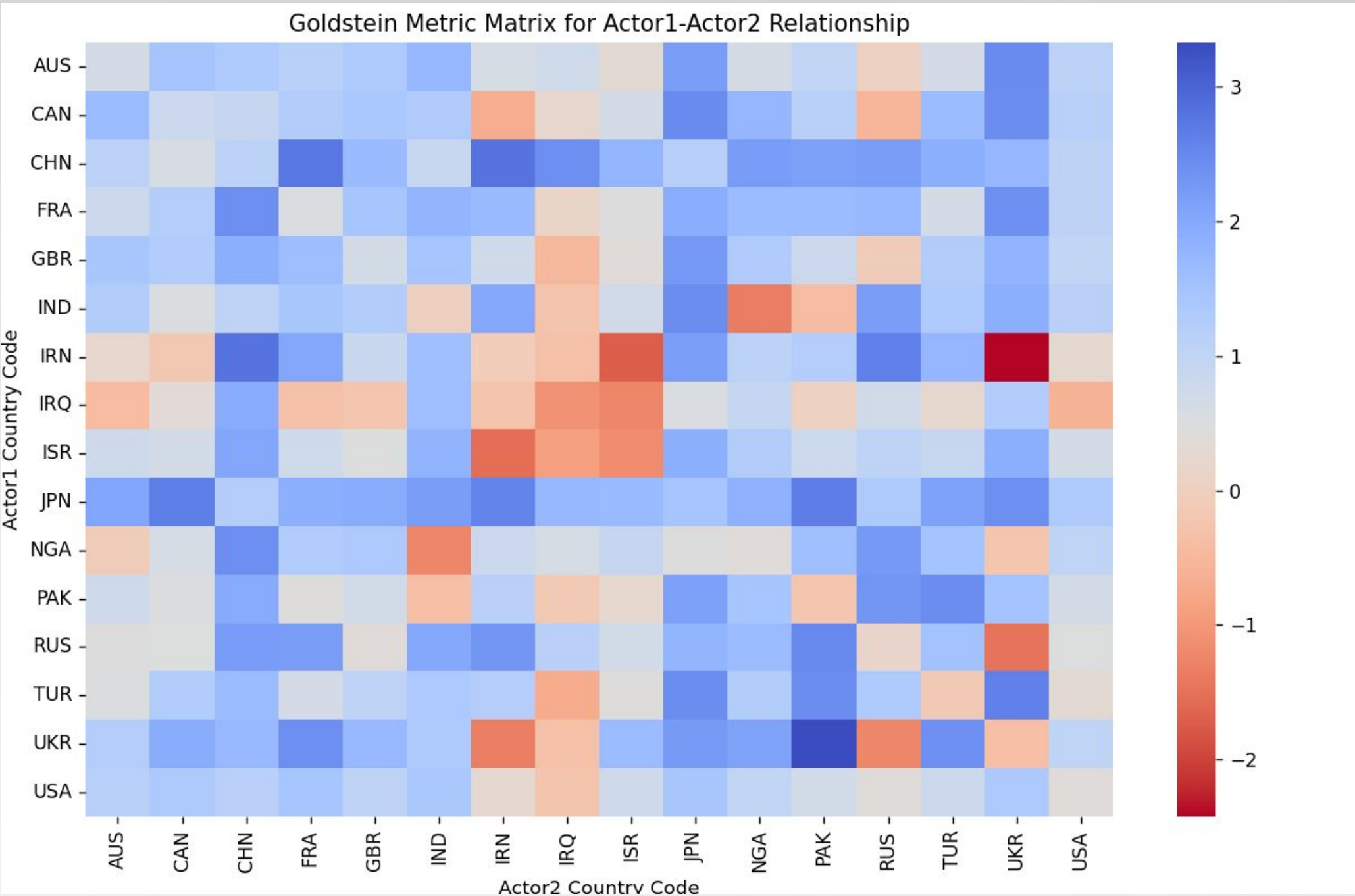
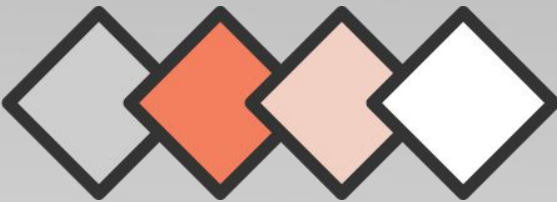


- Coverage differs greatly by news source
 - Quantity
 - Proportion
- Motivated analysis of country relationships for broader perspective



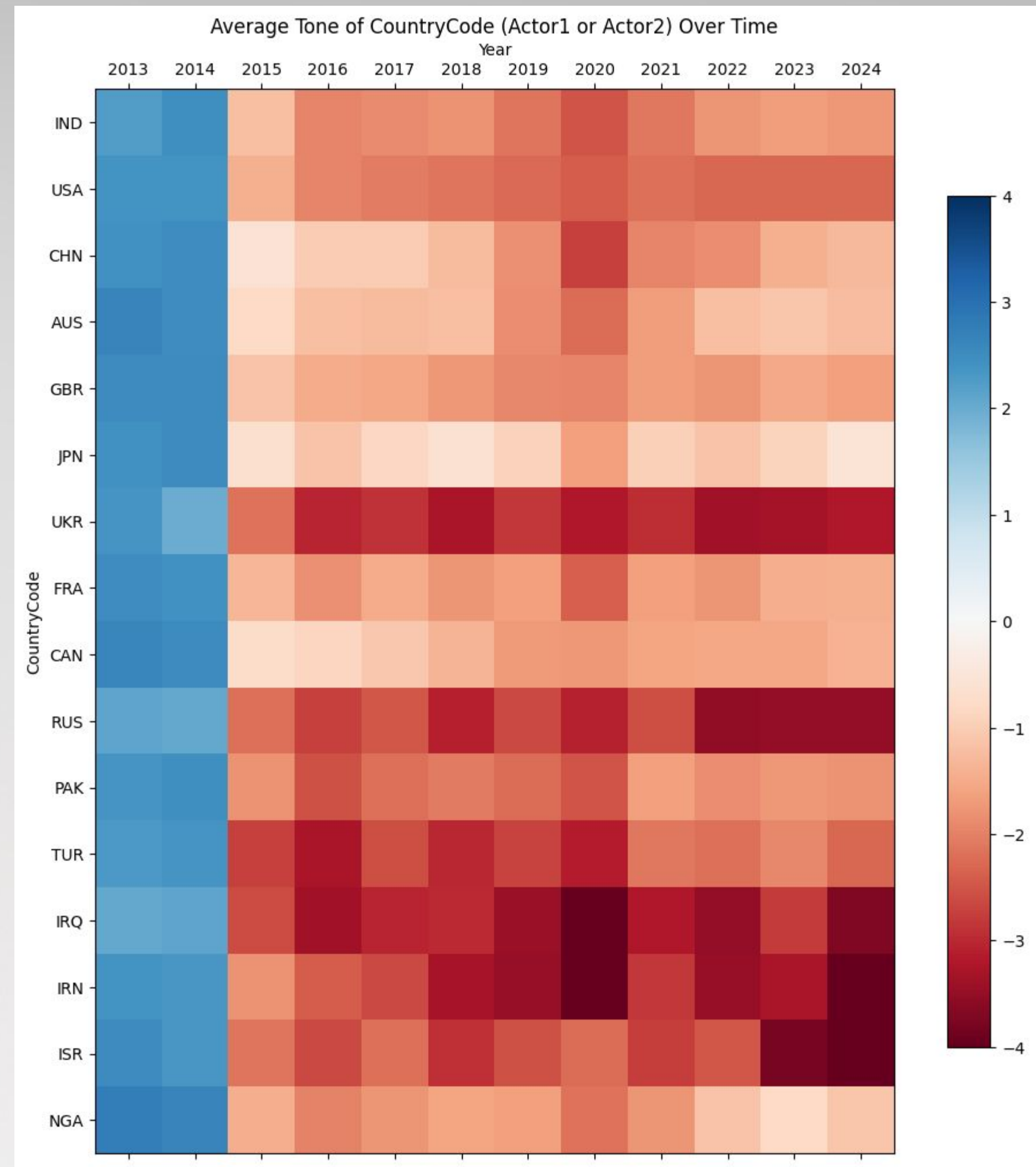


Results

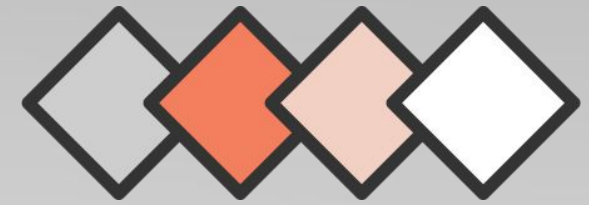


The chart displays the average tone of 16 country codes over time. The y-axis represents the average tone, ranging from -4 to 3, and the x-axis represents the year, ranging from 2014 to 2024. The legend lists the countries: IND, USA, CHN, AUS, GBR, JPN, UKR, FRA, CAN, RUS, PAK, TUR, IRQ, IRN, ISR, and NGA. The chart shows a general downward trend for most countries, with a sharp drop around 2015. The average tone for most countries starts between 2.0 and 2.8 in 2013 and drops to between -2.5 and -3.5 by 2016. After 2016, the average tone for most countries fluctuates between -1.0 and -3.5, with some countries showing a slight upward trend towards the end of the period.

CountryCode	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024
IND	2.1	2.5	-1.2	-2.5	-2.2	-1.8	-1.8	-2.5	-1.8	-1.8	-1.7	-1.8
USA	2.4	2.5	-2.8	-3.3	-2.8	-2.2	-2.2	-2.5	-2.2	-2.2	-2.2	-2.3
CHN	2.6	2.5	-1.5	-1.0	-1.0	-1.2	-1.2	-2.8	-1.8	-1.8	-1.4	-1.3
AUS	2.4	2.4	-1.8	-1.2	-1.2	-1.2	-1.2	-2.2	-1.8	-1.2	-1.1	-1.2
GBR	2.5	2.5	-1.2	-1.5	-1.5	-1.8	-1.8	-1.8	-1.8	-1.8	-1.5	-1.6
JPN	2.8	2.6	-1.4	-2.0	-1.8	-1.6	-1.6	-1.6	-1.0	-1.1	-0.8	-1.1
UKR	2.3	2.4	-2.2	-3.0	-2.8	-3.3	-3.3	-3.2	-3.0	-3.3	-3.2	-3.2
FRA	2.5	2.5	-1.8	-1.8	-1.5	-1.6	-1.6	-2.2	-1.8	-1.8	-1.5	-1.4
CAN	2.6	2.5	-1.6	-0.8	-0.8	-1.4	-1.4	-1.7	-1.5	-1.5	-1.5	-1.4
RUS	2.1	2.1	-2.2	-2.7	-2.7	-3.1	-3.1	-3.1	-2.5	-2.5	-2.5	-2.5
PAK	2.4	2.4	-1.2	-2.0	-1.8	-1.8	-1.8	-2.5	-1.8	-1.8	-1.7	-1.8
TUR	2.4	2.4	-2.8	-3.3	-2.8	-2.2	-2.2	-2.5	-2.2	-2.2	-2.2	-2.3
IRQ	2.6	2.5	-1.5	-3.4	-3.4	-3.0	-3.0	-4.1	-3.2	-3.5	-2.8	-3.7
IRN	2.4	2.4	-1.8	-2.4	-2.6	-3.3	-3.5	-4.2	-2.8	-3.5	-3.3	-4.4
ISR	2.5	2.5	-2.2	-2.6	-2.2	-2.9	-2.9	-2.5	-2.5	-2.5	-3.8	-4.4
NGA	2.7	2.6	-1.4	-2.0	-1.8	-1.6	-1.6	-1.6	-1.0	-1.1	-0.8	-1.1



Potential Biases



Analysis on Potential Biases:

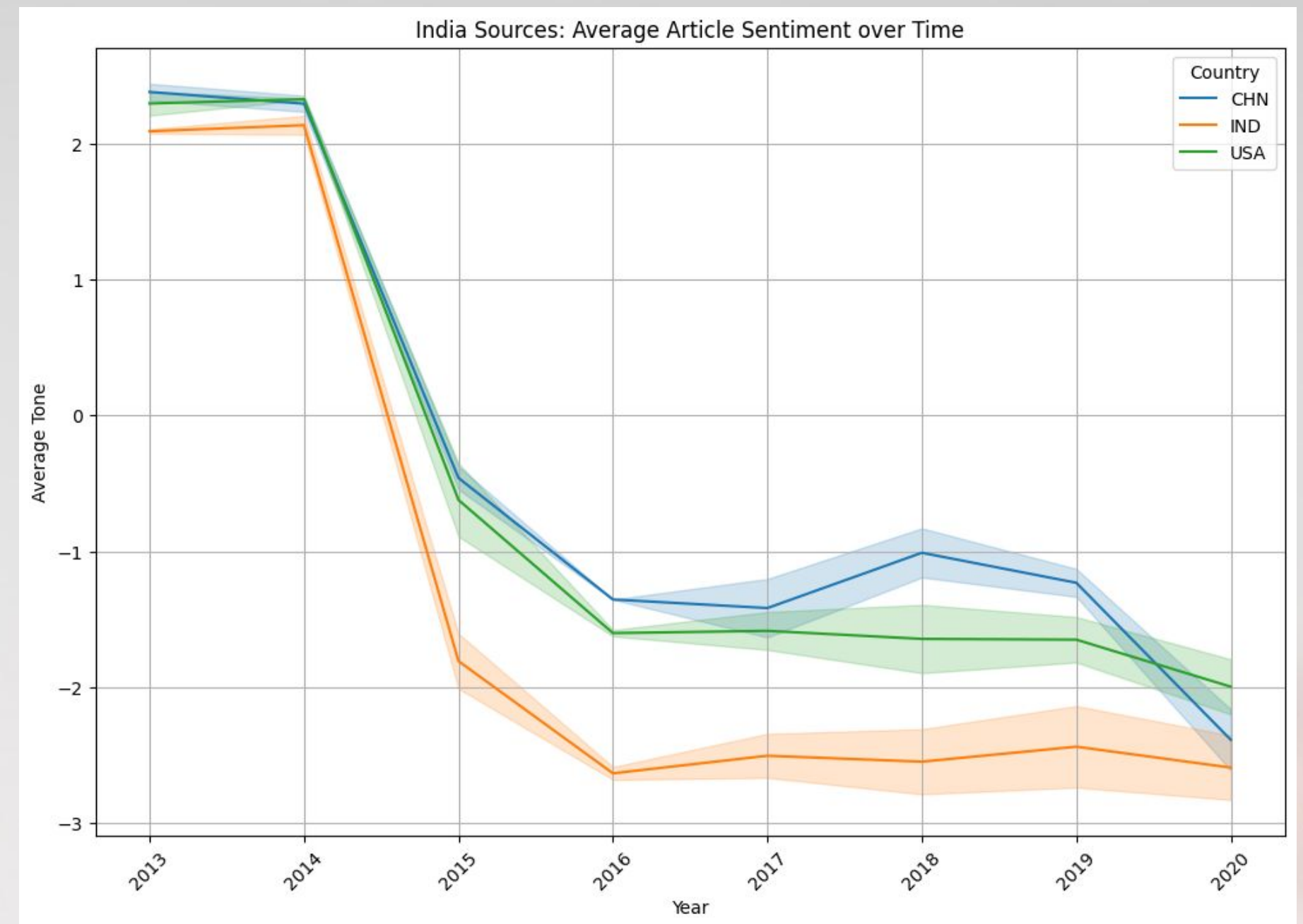
- Must know sources Country of Origin
- This allows us to see how other countries news bias

Chosen Sources:

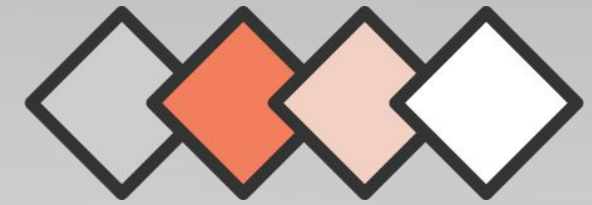
- NYtimes & WashingtonPost: USA
- SCMP & XinhuaNet: China
- IndiaTimes & TheHindu: India

To the right is a good baseline of what we found

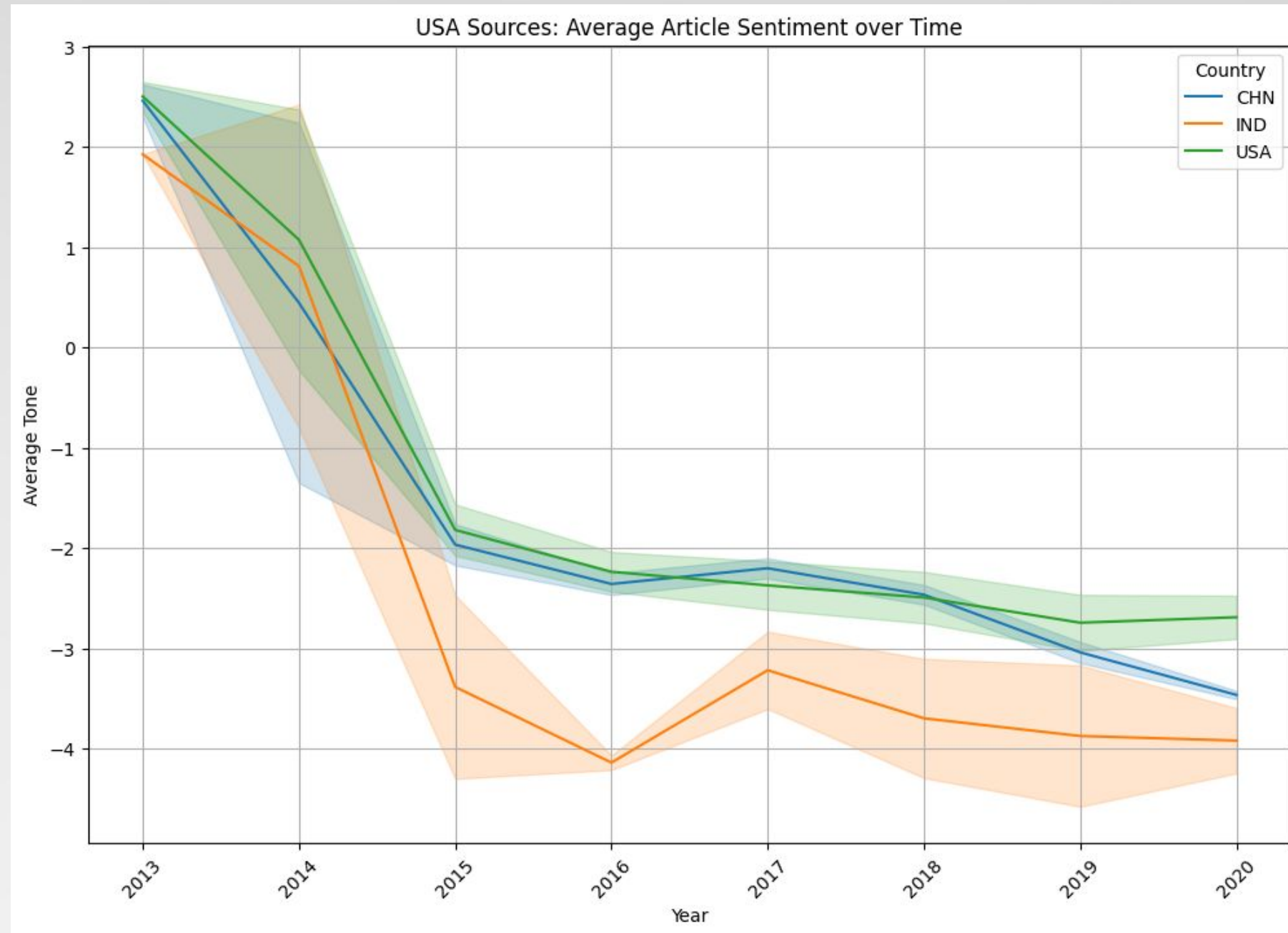
This graph was seen as fairly unbiased



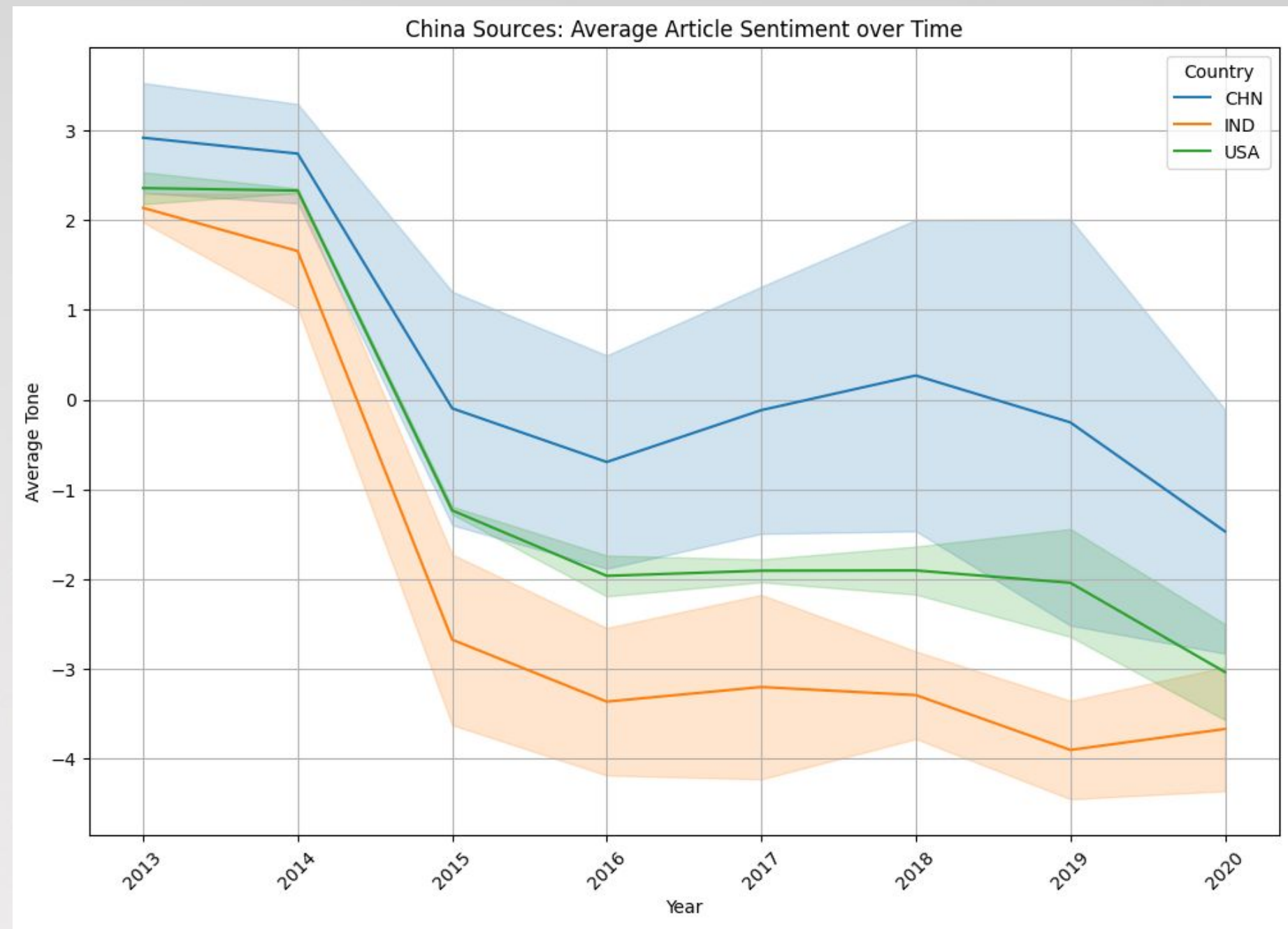
Potential Biases: Results



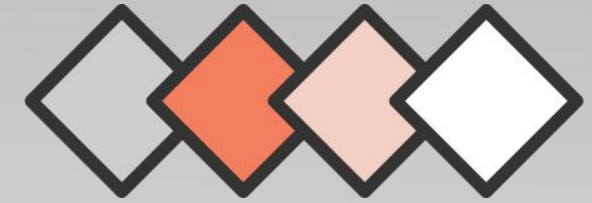
Fairly Unbiased



Clear Bias (China Sources -> China Avg Tone)



Word Analysis

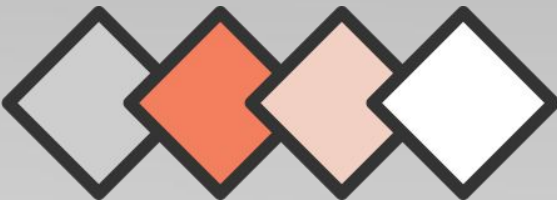


Analyze Article Titles extracted from URLs to Identify words with the strongest correlation to the average tone of the articles.

- **Extract Article Titles**
- **Tokenize Titles:** Break down article titles into individual words
- **Pair Words with Their Article's Tones**
- **Calculate Correlation:** Group data and calculate the correlation between all of the word occurrences and the average tone of articles.



Word Analysis



Word	Correlation
pm	0.050511
modi	0.048074
to	0.047394
cooperation	0.046950
summit	0.046828
prince	0.046653
celebrates	0.045348
visit	0.044791
ties	0.039746
c90000	0.037597
meets	0.036833
president	0.036323
with	0.035501
china	0.035466
king	0.035315
celebrate	0.035041
charles	0.034916
harry	0.034842
11ef	0.034663
content	0.033035

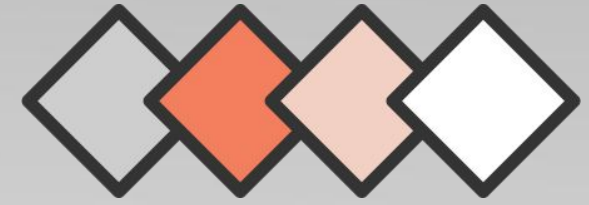
Positive Correlation

Word	Correlation
israel	-0.112326
attack	-0.094831
killed	-0.086772
man	-0.085358
iran	-0.078577
israeli	-0.075639
gaza	-0.074471
murder	-0.073554
russian	-0.071784
police	-0.070379
hezbollah	-0.068671
hamas	-0.068544
war	-0.067689
shooting	-0.067227
attacks	-0.065330
kill	-0.063166
accused	-0.061525
charged	-0.060675
arrested	-0.060009
death	-0.058084

Negative Correlation

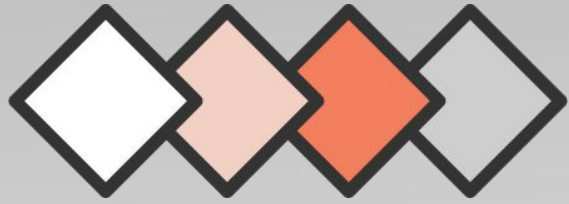


Primary Issues

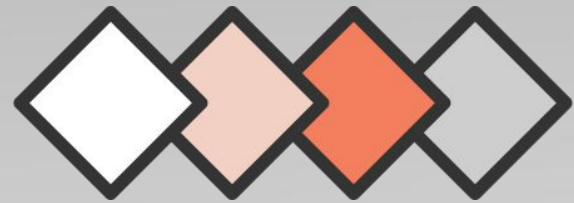


- Size of data set was overwhelming so storing it was a pain and we had to pay for storage at one point
- We had to get a new dataset halfway through to use more recent articles
- Topic modeling didn't work for the 1st GDELT data set because there was no text or content
 - Also unable to scrape all the websites for text despite numerous attempts (they didn't like it)
 - Used an LDA algorithm for the topic modeling

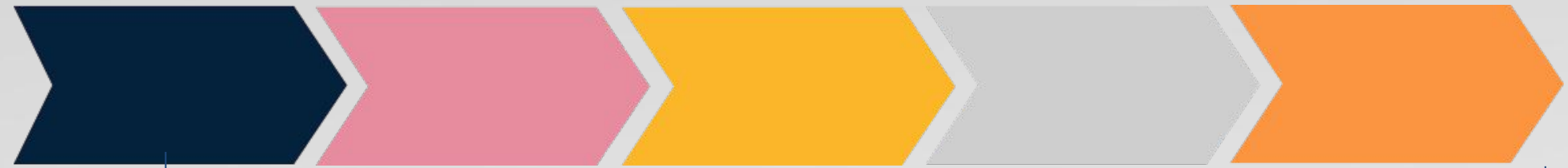




Team Member	Responsibilities
Colin Canonaco	Exploring, investigating, and cleaning original GDELT data set and performing sentiment analysis. Use data/results to identify biases and potential conclusions.
Harshvardhan	Helping with data set cleaning and original analysis of news sources, then overseeing analyses. Use data/results to identify biases and potential conclusions.
Jeremiah Augustine	Worked with sentiment analysis and creating visuals to represent results and conclusions. Use data/results to identify biases and potential conclusions.
Faith Chernowski	Worked with other analysis methods such as topic modelling and clustering to represent scraped data. Use data/results to identify biases and potential conclusions.



Timeline



Sprint 1 - Sep 15 - Sep 30

Project proposal and
original GDELT data
scraping

Sprint 3: Oct 15 - Oct 30

Secondary data set
retrieval and original
sentiment analysis
averages

Sprint 5: Nov 15 - Nov 30

Presentation and final report
creation, finishing results
discussion/analysis

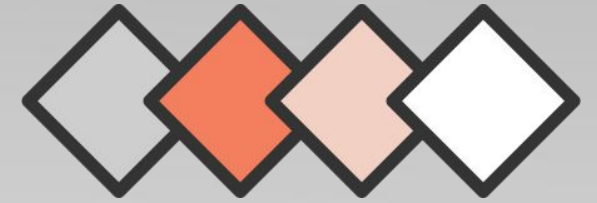
Sprint 2 - Oct 1 - Oct 15

Topic modeling attempts
and analysis by news
source numbers

Sprint 4: Oct 30 - Nov 15

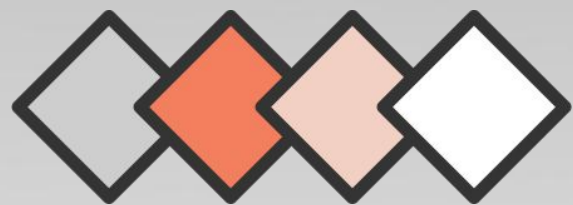
Tone visualization by
country and over time
(sentiment analysis)

Future Work



- Expand Dataset Scope: Include additional timeframes, more countries, and diverse publishers to capture broader media perspectives.
- Deepen Analysis: Investigate sub-national portrayals, cross-domain topics like environmental or human rights issues, and cultural influences on media narratives.
- Advance Techniques: Apply fine-grained sentiment analysis, network analysis for global media dynamics, and event chronology to track narrative shifts.
- Real-world Applications: Develop policy recommendations, media literacy tools, and interactive visualizations for public and academic use.
- Comparative Studies: Analyze platform-specific portrayals, cultural framing differences, and contrasts between traditional and modern media.





THANK
YOU

