# Sentiment analysis of Social media discourse:

## a case study of Cecilia Sala on BlueSky and YouTube

M. D'Amato, F. Guardione, V. Mascellaro, D. Montepara, N.Vuolo

# ✴ Key context

**> Focus:** the case of Cecilia Sala, a 29-year-old Italian journalist detained in Tehran on December 19, 2024, for alleged violations of Iranian laws

**> Arrest context:** occurred shortly after Italy detained an Iranian engineer, sparking diplomatic efforts led by PM Giorgia Meloni to secure Sala's release

# ✴ Themes in analysis

**>** Sala's work highlights international negotiations and the risks faced by journalists in crisis zones

**>** Social media discourse reveals:
- Gender-based biases and stereotypes
- Professional critiques faced by women in public roles

# Exploring social media platforms

## Research goals

1. Analyze social media framing of international negotiations
2. Examine gendered language and sentiment towards Sala as a woman journalist
3. Study perceptions of journalists in high-risk environments

## Platforms Compared

### BlueSky

1. **Decentralized platform** fostering **text-based,** participatory dialogue
2. **Early-stage community** reduced algorithmic influence
3. Focus: sentiment patterns and cultural framing of discourse

### YouTube

1. **Established platform** emphasizing **visual and performative content**
2. **Amplifies narratives** via video engagement and viewer interactions
3. Focus: sentiment shaped by visuals and gendered professional critiques

# Data collection and preprocessing

## ✸ Data sources

**YouTube: 1,113 comments** from 10 videos related to Cecilia Sala's case, selected based on keywords such as **"Cecilia Sala," "Meloni," "Abedini," "journalism," and "woman."**

**BlueSky: 1,480 posts** retrieved via BlueSky's API, covering a timeframe from 15/12/2024 to 15/01/2025, using the keyword **"Cecilia Sala."**

## ✸ Preprocessing steps

- Applied **text cleaning** with regex to remove punctuation, links, emojis, and anomalies
- Used **spaCy NLP** for lemmatization to standardize text
- Calculated relative **term frequency** to highlight recurring linguistic patterns and their association with sentiment and emotions

# Analytical methods

**Sentiment analysis**

Model: **tabularisai/multilingual-sentiment-analysis**
- Categories: Very Positive, Positive, Neutral, Negative, Very Negative
- Goal: measure overall public opinion polarity

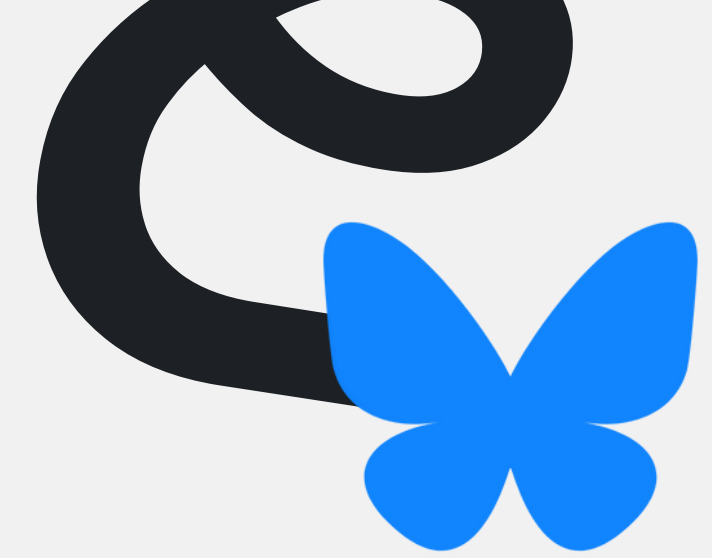**Emotion analysis**

Model: **MilaNLProc/feel-it-italian-emotion**
- Categories: Joy, Sadness, Anger, Fear
- Goal: capture emotional engagement with themes of gender, journalism, and international events

**Tools and insights**

Pretrained models from Hugging Face tailored for Italian text classification
- Combined sentiment and emotion analysis for a comprehensive understanding of public discourse:
  1. **Sentiment analysis: evaluated evaluative stances**
  2. **Emotion analysis: explored psychological and affective reactions**
  3. Enabled nuanced insights into audience perceptions across YouTube and BlueSky platforms
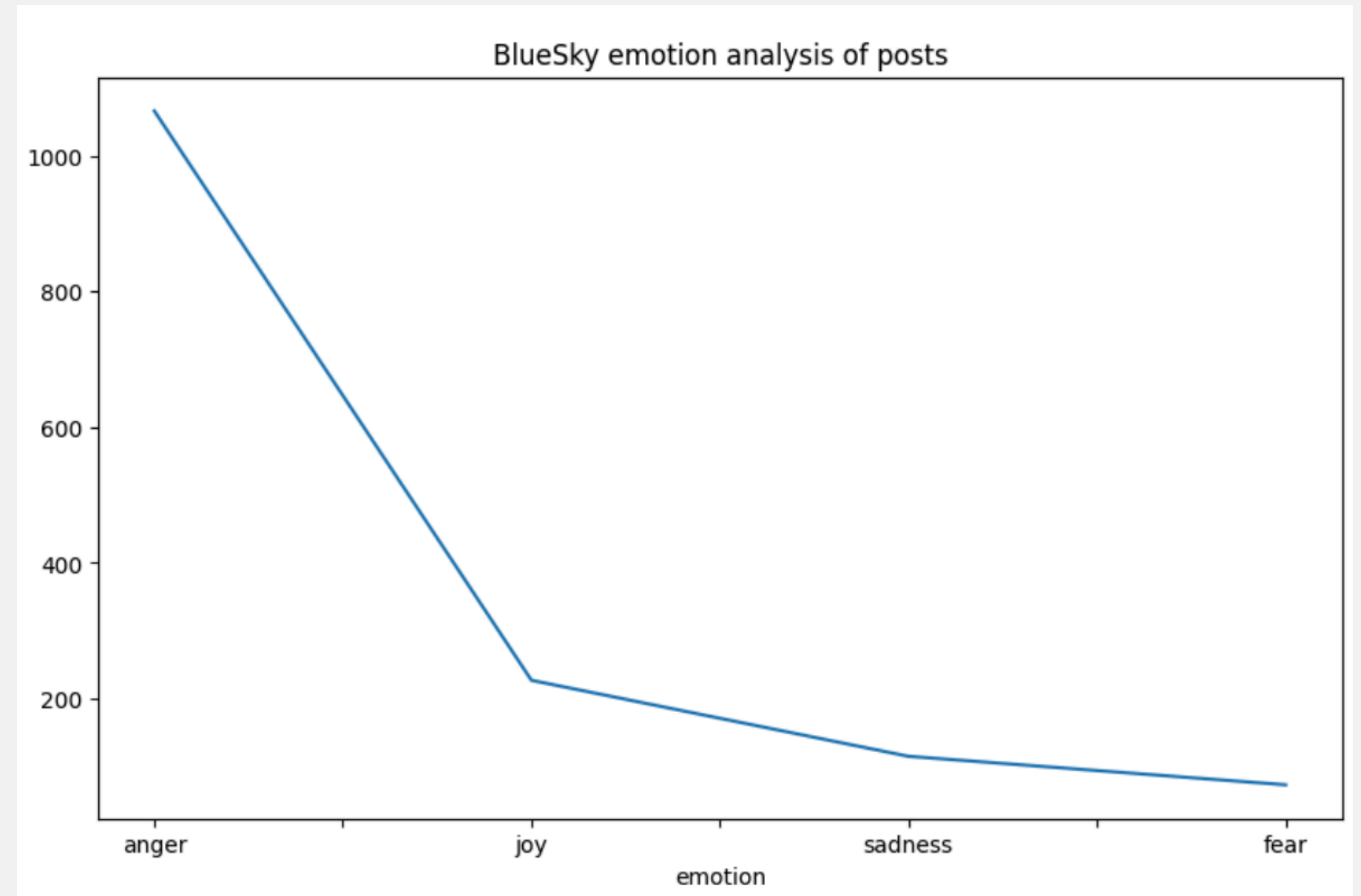
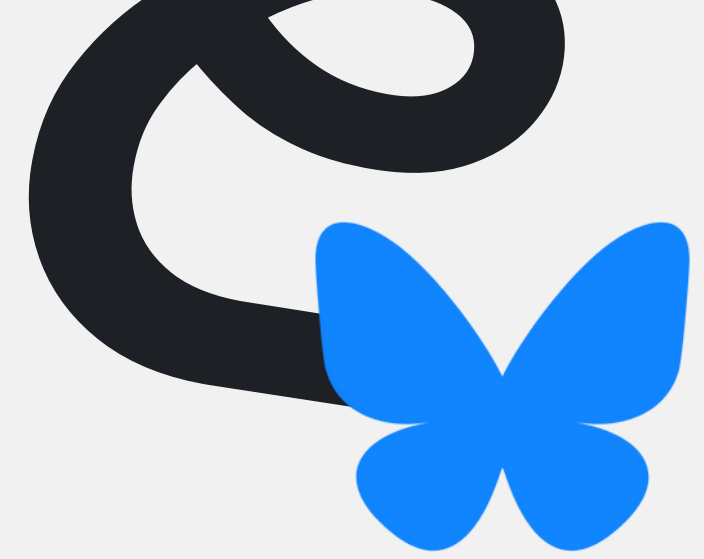# *B*lueSky analysis and results

Emotion analysis results

- **Anger dominates**: over 1,000 instances recorded, reflecting frustration or criticism
- **Joy**: second most frequent, indicating limited positive engagement
- **Sadness and Fear**: highlight concerns, grief, and apprehension, adding complexity to emotional responses

Key insights

- Emotional discourse reflects sensitive topics like:
  - **Gender issues**
  - **Journalism in conflict zones**
  - **International negotiations**
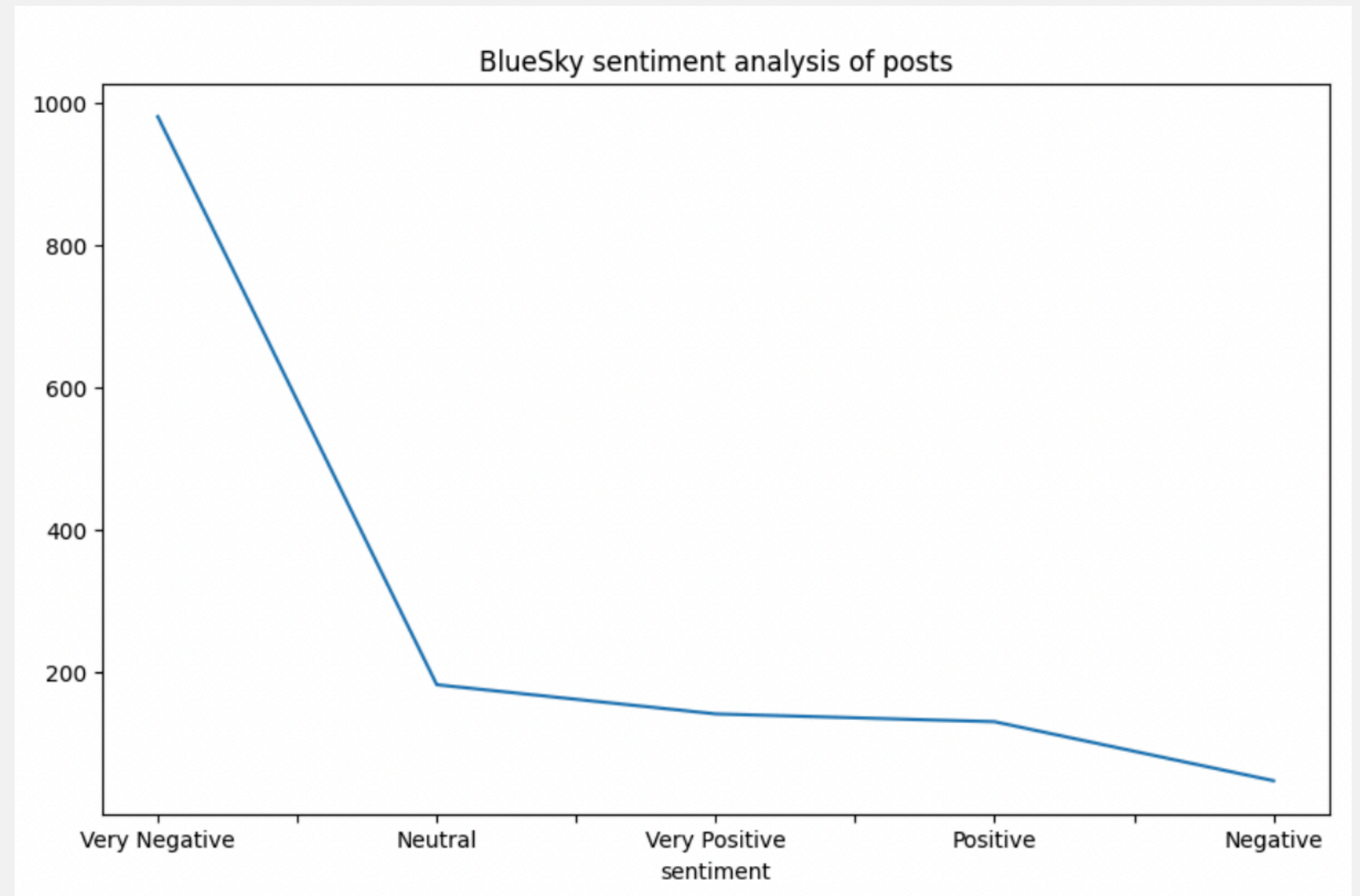- Multifaceted emotional landscape with anger as a central theme



BlueSky emotion analysis of posts
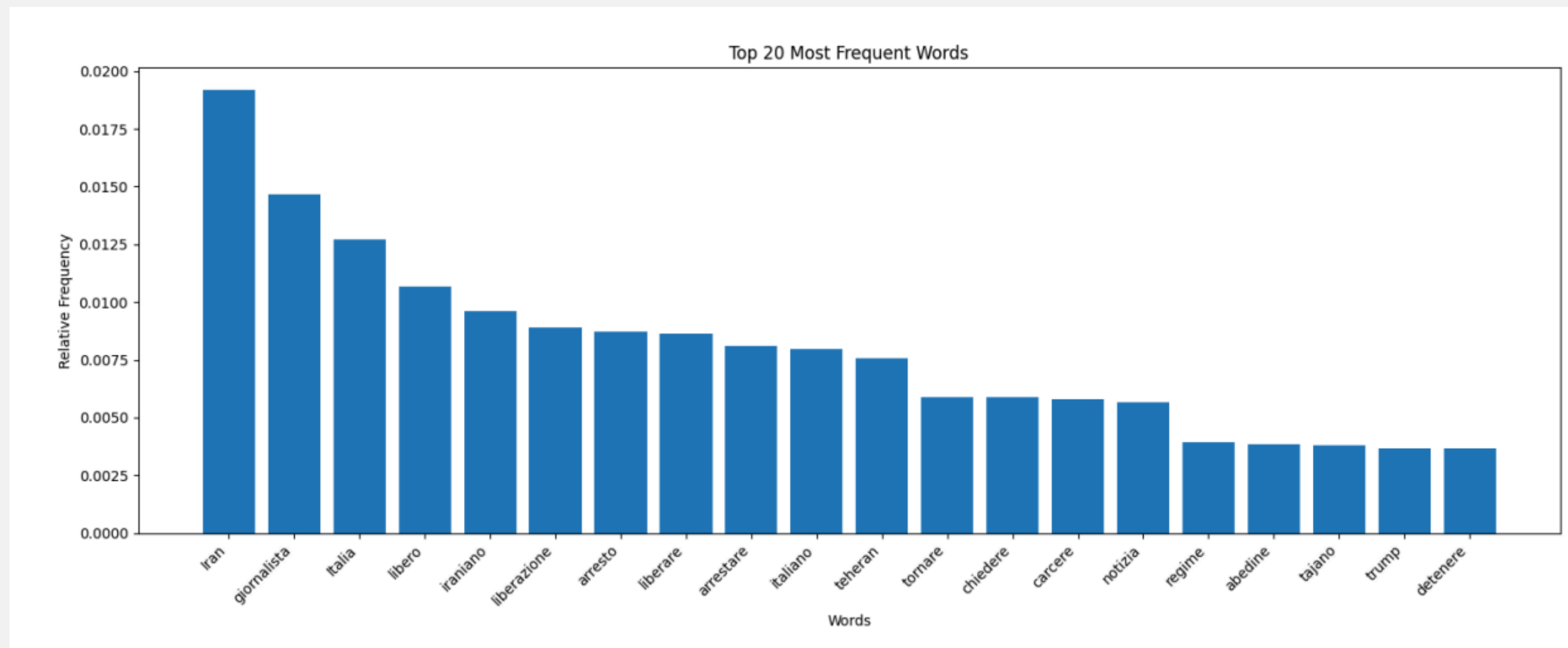
# *B*lueSky analysis and results

Sentiment analysis results

- **Very negative sentiment is predominant** (~950 posts), showing strong public reaction to Cecilia Sala's arrest
- **Contextual analysis is crucial:** Negative sentiment reflects genuine **responses to a high-profile event, not systemic bias**



BlueSky sentiment analysis of posts

# *B*lueSky analysis and results

Bar Chart analysis
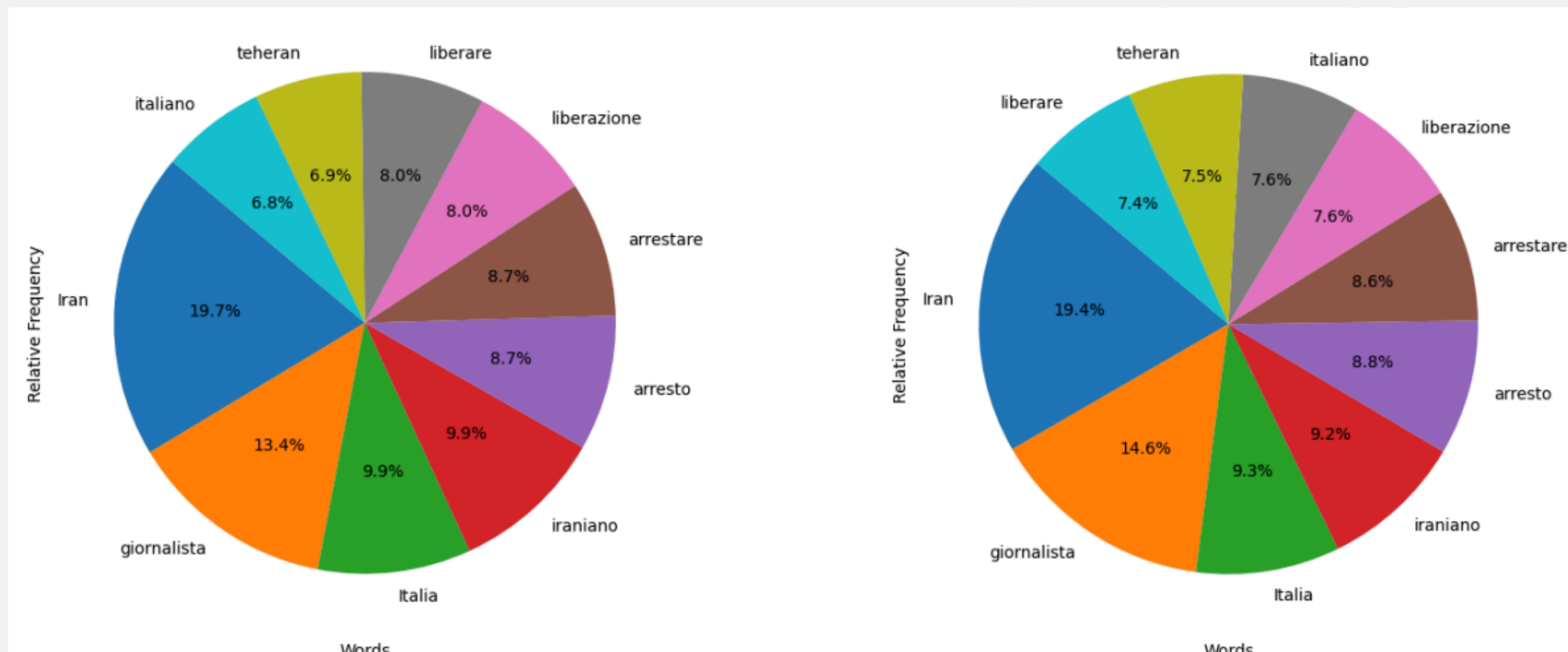


Top 20 Most Frequent Words

**Top 20 words reflect thematic focus:**

- Geopolitical context: **"Iran," "Teheran," "regime."**
- Journalistic role: **"giornalista", "notizia"**
- Detention and freedom: **"arresto", "liberare"**

# *B*lueSky analysis and results

Pie Chart analysis



**Anger and very negative aentiment words:**
- **"Iran" (19.4–19.7%)** and **"giornalista" (13.4–14.6%)** dominate, reflecting focus on journalism and the Iranian context
- Legal terms like **"arresto" (arrest)** and **"liberazione" (freedom)** underscore discussions on detention and human rights

## Conclusions:

1. Key themes: media freedom, geopolitical tensions, and calls for justice

2. Data provides a quantitative foundation to understand public priorities and framing in sensitive international and gendered issues.
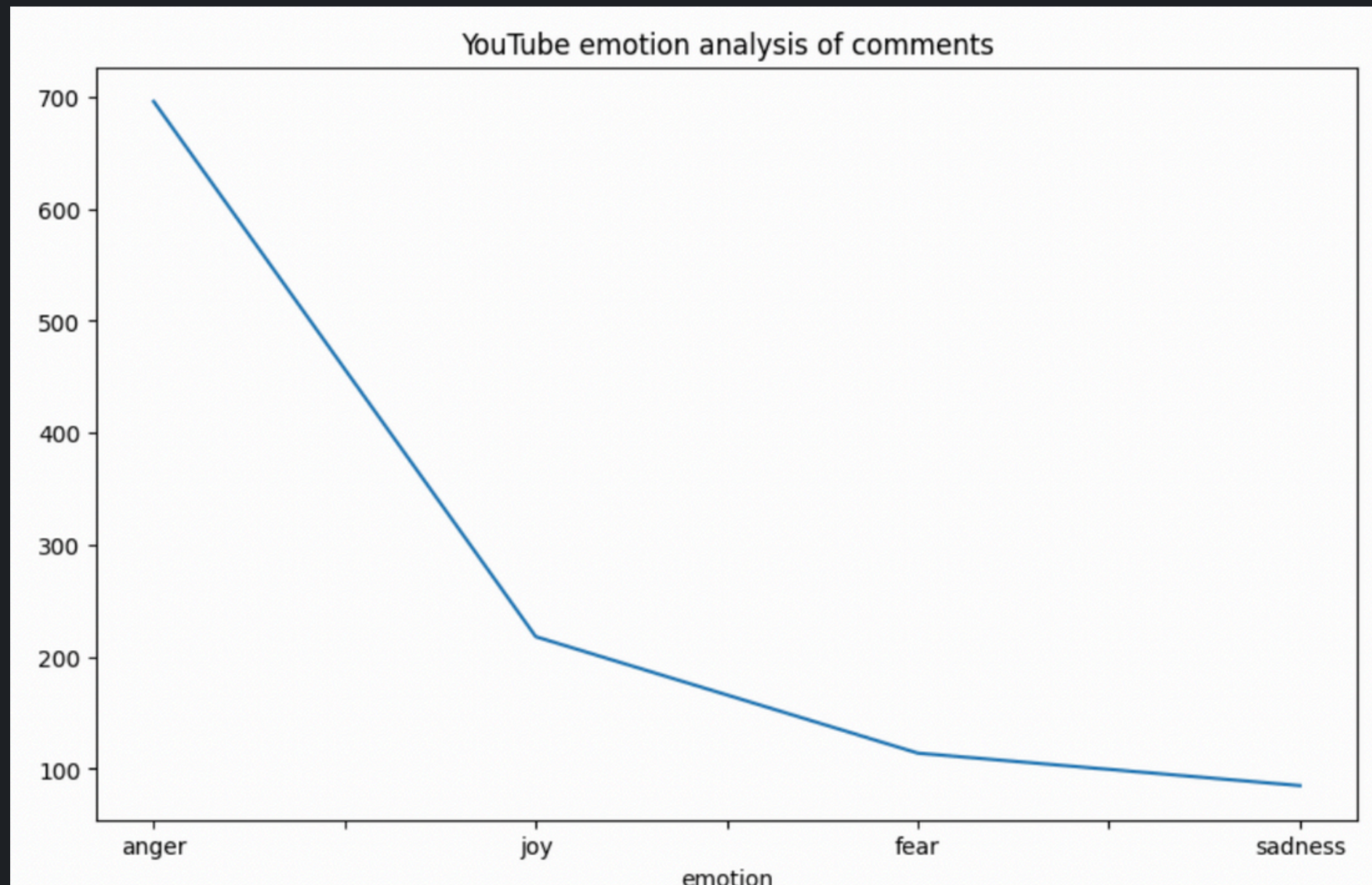
## ✳ Focus area

○ Emotion analysis: distribution of anger, joy, fear, sadness
○ Sentiment analysis: spectrum from Very Negative to Very Positive

## ✳ Methodology

Frequency analysis of key terms in an Italian-language corpus

# Youtube comments and results

# Youtube comments and results



YouTube emotion analysis of comments

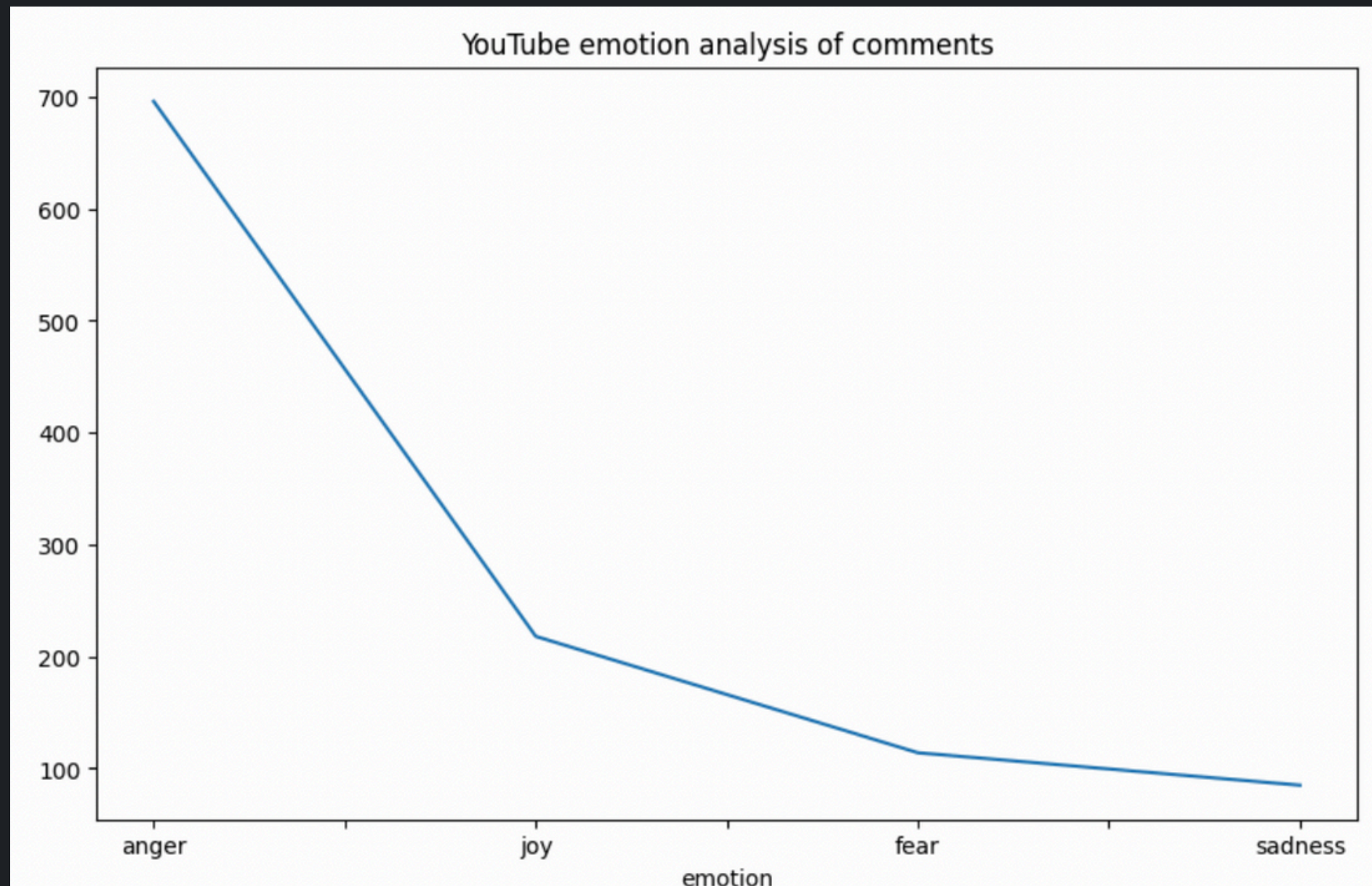Emotion analysis results

Four primary emotions analyzed:

- **Anger: predominant** with ~700 instances
- **Joy: limited presence**, fewer than 200 instances
- **Fear and sadness: minimal occurrence**, below 150 instances

Key insight

Discourse reflects strong **frustration and negativity**, with limited positive emotions

# Youtube comments and results



Sentiment analysis results

Distribution pattern:

- **Very Negative: dominates** (~600 comments).
- **Neutral: moderate presence** (~200 comments).
- **Positive & Very Positive: Minimal presence** (100–150 comments).
- Negative: Lowest (~<100 comments).

Sentiment scale:

**five-point system** (Very Negative to Very Positive) for detailed tonal analysis
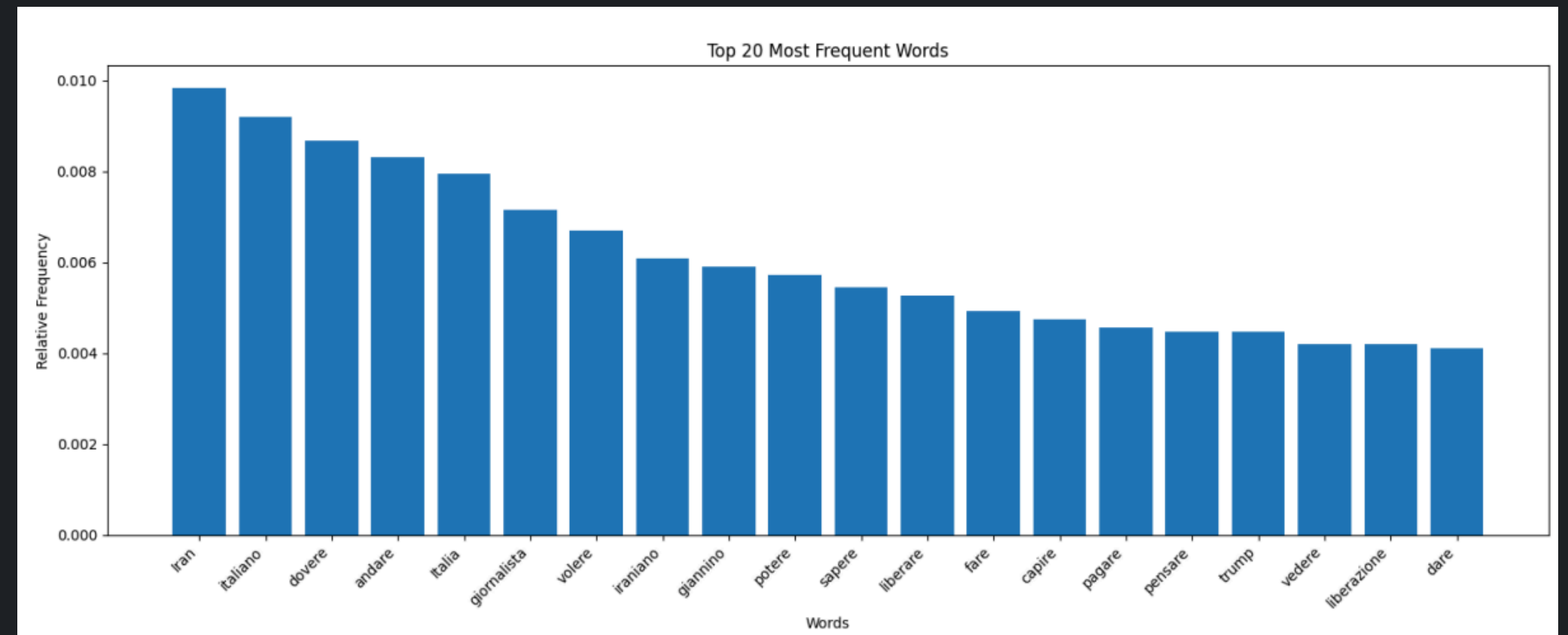
# Youtube comments and results

Bar Chart analysis

**Most frequent terms**

- **"Iran"** (≈ 0.0095): aligns with Sala's focus on Middle Eastern affairs
- **"Giornalista"** and **"Italia"**: reflect professional identity and national context

**Action-oriented verbs:**

- "Liberare", "dovere", "andare", "volere", "pensare"
- Indicates themes of advocacy, movement, and press freedom



Top 20 Most Frequent Words
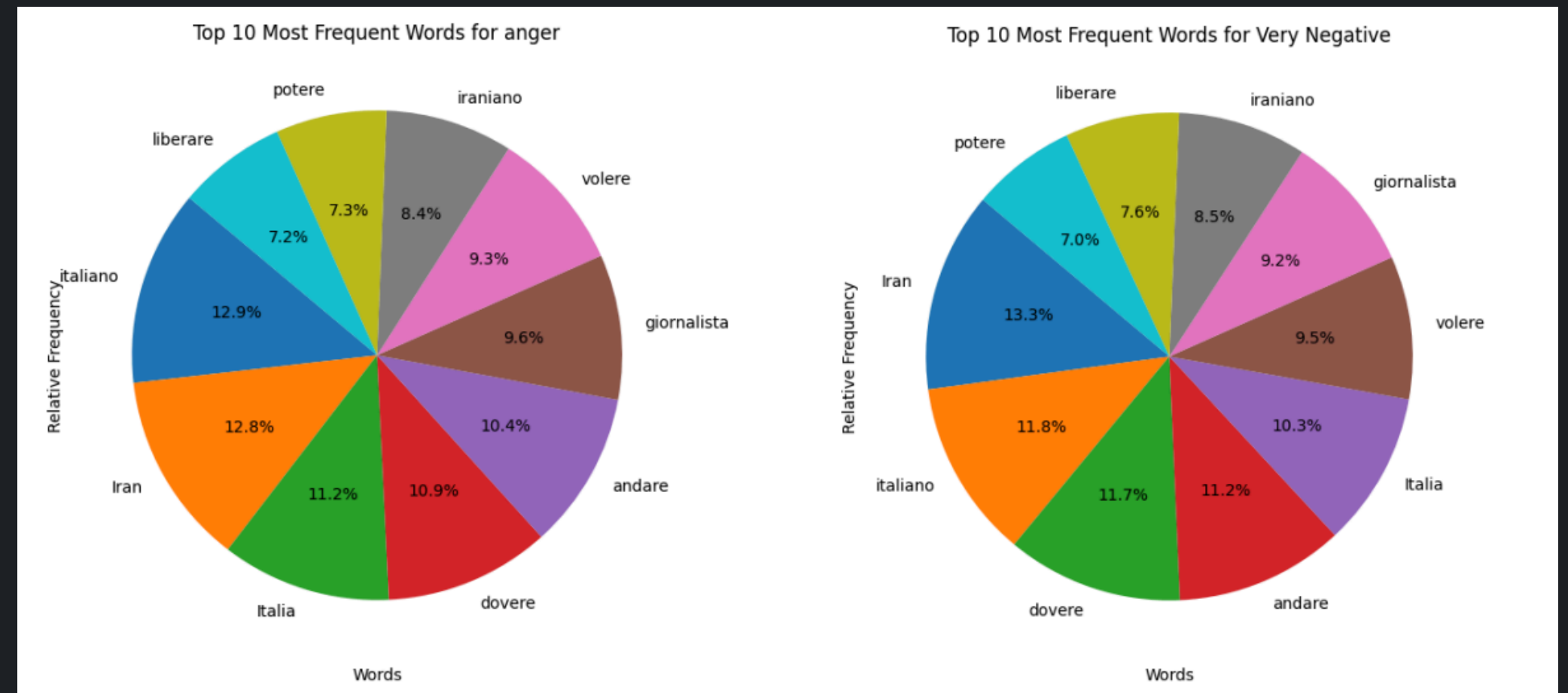
# Youtube comments and results

Pie Chart analysis

**Themes identified:**
- **Press freedom, journalist safety, and geopolitical tensions**
- Strong audience engagement with topics on **professional risks** and international reporting

**Impact of platform dynamics:**
Highlights how **gender, journalism, and audience sentiment intersect in online discussions**

# *D*iscussion

## ✸ Key findings

- **Sentiment and emotion analysis:** highlights **public attitudes** toward journalism, gender, and international affairs via BlueSky and YouTube comments

- **Platform-specific insights:**
1. <span style="color:red">**YouTube:**</span> highly **polarized discourse**, influenced by **platform algorithms**
2. <span style="color:blue">**BlueSky:**</span> limited by **restricted query functionality**, leading to **incomplete datasets**

- **API constraints:**
1. <span style="color:blue">**BlueSky**</span> lacks metadata filtering
2. <span style="color:red">**YouTube**</span> does not support date-based filtering, hindering longitudinal analysis

- **Bias in NLP models:**
1. Pretrained models misclassify complex expressions (e.g., sarcasm, irony)
2. Lemmatization errors (e.g., "Sala" → "salare") distort sentiment results

# Discussion

✸ **Proposed solutions**

- **Enhanced Data Collection:**
1. Cross-platform integration (Mastodon, Reddit, news comments)
2. Periodic scraping and historical archiving to ensure data representativeness

- **Improved NLP models:**
1. **Domain-specific training datasets** for journalism and gender discourse
2. **Customized Named Entity Recognition (NER)** to handle proper nouns

- **Advanced analytical techniques:**
1. **Topic modeling** for thematic insights
2. **Longitudinal studies** for evolving trends in sentiment and emotions

- **Ethical considerations: transparent methodologies** and **AI frameworks** to ensure responsible use of public data

**Conclusion:**
**Big Data analytics and AI-driven NLP offer robust tools to understand digital discourse, informing academic research and policy on online communication dynamics.**

# Conclusions

Key findings:

- Sentiment analysis: **predominantly negative and very negative comments**, consistent with research on **gender-based digital harassment** toward female journalists.
- Emotion analysis: **dominance of fear and anger**, reflecting concerns about **press freedom, international affairs, and political discourse**
- Lexical insights: strong focus on **geopolitical topics** (e.g., "Iran") and **journalistic identity** (e.g., "giornalista"), emphasizing the relevance of the case study

Methodological insights:

- Strengths:
  - Big Data and AI-driven NLP enable large-scale, systematic textual analysis
  - Structured sentiment, emotion classification, and lexical frequency analysis provide actionable insights
- Limitations:
  - **Algorithmic biases** in NLP models
  - **Data access restrictions** and need for continuous **model refinement**

**Big Data analytics and AI-driven NLP offer robust tools to understand digital discourse, informing academic research and policy on online communication dynamics.**

Thank you for your attention