

**Patterns in the Most Influential U.S. News YouTube Channels: Becoming the Negative
Generation**

Yicong (John) Deng, Jingwen (Wendy) Xiao, Maurice Onyonyi

Georgetown University

Data Science I: Foundations PPOL 5203

Professor Dr. Tiago Ventura

16 Dec. 2025

Abstract

This report analyzes how major U.S. news YouTube channels generate influence and how audiences respond emotionally, using engagement metrics and large scale text analysis for October 2025. The project aimed to include the top 10 channels by subscribers, but the final dataset includes nine because TMZ could not be scraped. Using the YouTube Data API v3, we collected video metadata and large batches of comments for up to 100 videos per channel. We measure influence with a weighted engagement score and log scaling, then apply BERTopic for title themes, RoBERTa for title and comment sentiment, and TF IDF plus bigrams to identify dominant discussion frames. Results show clear differences in headline themes across channels, but consistently negative comment sentiment and a discourse dominated by political conflict terms. Overall, the findings suggest that engagement on YouTube news is often driven by disagreement, and we conclude with limitations and directions for future work.

1. Introduction

YouTube has become a central venue for news consumption and political conversation, but it differs from traditional media because influence is visible through platform signals and public audience text. Every upload produces measurable behaviors—views, likes, and comments—that reflect different levels of attention and participation. For news organizations, this creates both an opportunity and a challenge: the same mechanisms that expand reach can also reward emotionally charged topics and conflict-heavy discussion, especially in comment sections.

This project asks how the most-subscribed U.S. news YouTube channels generate influence and what audience sentiment reveals about that influence. Rather than treating views as the only outcome, we treat influence as engagement-weighted reach, because commenting and liking require more effort than clicking. We also study the language of titles and comments to connect publisher framing with viewer reaction. To make channel comparisons fair, we standardize the sample to a single period (October 2025) and limit the dataset to up to 100 videos per channel. This design reduces distortions from unequal posting frequency and helps the analysis reflect differences in content strategy and audience response within the same news cycle. The original plan was to include the top 10 channels by subscribers, but TMZ could not be scraped, leaving nine channels in the final sample.

Methodologically, the project combines three layers of analysis. First, we construct an influence metric from video metadata by weighting views, likes, and comments (1–3–5) and applying log scaling to reduce the impact of outliers. Second, we model title content using BERTopic to summarize dominant themes by channel and transformer-based sentiment classification to capture the tone of headlines. Third, we analyze audience discourse by running sentiment classification on hundreds of thousands of comments and extracting salient terms and bigrams to identify recurring actors, institutions, and policy frames.

Together, these steps provide an interpretable picture of YouTube news influence: what channels emphasize, how viewers react, and whether influence is associated with neutral headline strategies, polarized discussion, or both. Understanding these dynamics can inform platform moderation, newsroom packaging decisions, and media literacy by showing when engagement reflects approval versus frustration, sarcasm, or partisan conflict.

2. Data and Methods

2.1 Data source.

We used the **YouTube Data API (v3)** to collect (a) channel/video metadata and (b) comment text at scale. The metadata included video titles, publish dates, view counts, like counts, comment counts, and durations. The nine channels used throughout the analysis are ABC News, CNN, Fox News, NBC News, CBS News, The Young Turks, USA Today, The Wall Street Journal, and The New York Times.

2.2 Sampling strategy.

The dataset was filtered to **October 2025** and “limited up to 100 videos” per channel to ensure apples-to-apples comparisons. This matters because engagement metrics are highly sensitive to both time and volume: Channels that post more frequently can accumulate more total engagement while channels that publish fewer but higher-budget videos might dominate per-video engagement.

2.3 Comment collection.

Upon sampling our data, we retrieved large batches of comments for sampled videos and stored raw text for downstream processing and analysis. For every comment, the video ID, comment ID, author, published date, and like-count are draThe dataset contains roughly 300,000 total comments.

2.4 Constructing the influence metric.

Because “influence” is not directly observable from one number, we defined influence as a function of engagement signals that represent different levels of effort:

- **Views:** lowest-effort exposure signal
- **Likes:** lightweight endorsement
- **Comments:** higher-effort participation

Weighted influence is specified as: **Views (1), Likes (3), Comments (5)**.

The implementation then applies log scaling to stabilize extreme values. In the title analysis code, the scaled influence variable is constructed as:

$$\text{Influence} = (3 * \text{Likes} * 5 * \text{Comments}) * \text{Views}$$

$$\text{Influence_scaled} = \log_{1p} (\text{Views} + 3 * \text{Likes} + 5 * \text{Comments})$$

While the exact multiplicative form is one design choice (others could be additive), the guiding principle is consistent: **Comments receive the largest weight**, reflecting the idea that commenting indicates deeper engagement than passively viewing. The log transform (“expressed in log scale”) is explicitly used to “overcome the effects of extremely large numbers.” We used this influence metric in multiple ways: ranking the most influential videos, examining the distribution of influence by channel, and aggregating influence across videos to compare total influence by channel.

2.5 Duration conversion

Video durations were recorded in ISO 8601 format (e.g., PT9M15S). The channel-analysis notebook converts duration strings into total seconds using regular expressions and sums metrics by channel for descriptive comparisons.

2.6 Title themes using BERTopic.

To summarize channel-level patterns in what publishers talk about, we applied BERTopic to video titles (per channel), producing a set of topics and identifying a dominant topic per channel by count. The topic analysis also flagged cases where the dominant topic was “-1,” interpreted as an outlier/no clear theme.

2.7 Sentiment analysis for titles and comments.

Sentiment analysis was implemented using transformer-based models and batching. In the text-analysis notebook, sentiment is computed in batches of 64 using tokenization with truncation and maximum length constraints, then mapped to labels NEGATIVE/NEUTRAL/POSITIVE.

2.8 Topic extraction in comments: TF-IDF and bigrams.

To identify what people talk about beyond sentiment, we extracted salient terms from comments and recurring phrases from combined title+comment text. The notebook computes bigrams using a CountVectorizer with `ngram_range = (2,2)` and reports the top phrases by frequency.

3. Analysis and Results

3.1 Baseline engagement and activity patterns

A first step is to understand scale and activity—how much content each channel published (within the capped sample) and how much reaction it generated. The project summarizes core totals (views, likes, comments, and duration) by channel from the October video dataset. A key engagement signal is **total comments** on sampled videos. In the October dataset, the channels with the highest total comments were **Fox News (83,095)** and **The Young Turks (81,270)**, followed by **CNN (63,630)** and **The New York Times (36,580)**. These totals matter because comments represent active participation rather than passive exposure.

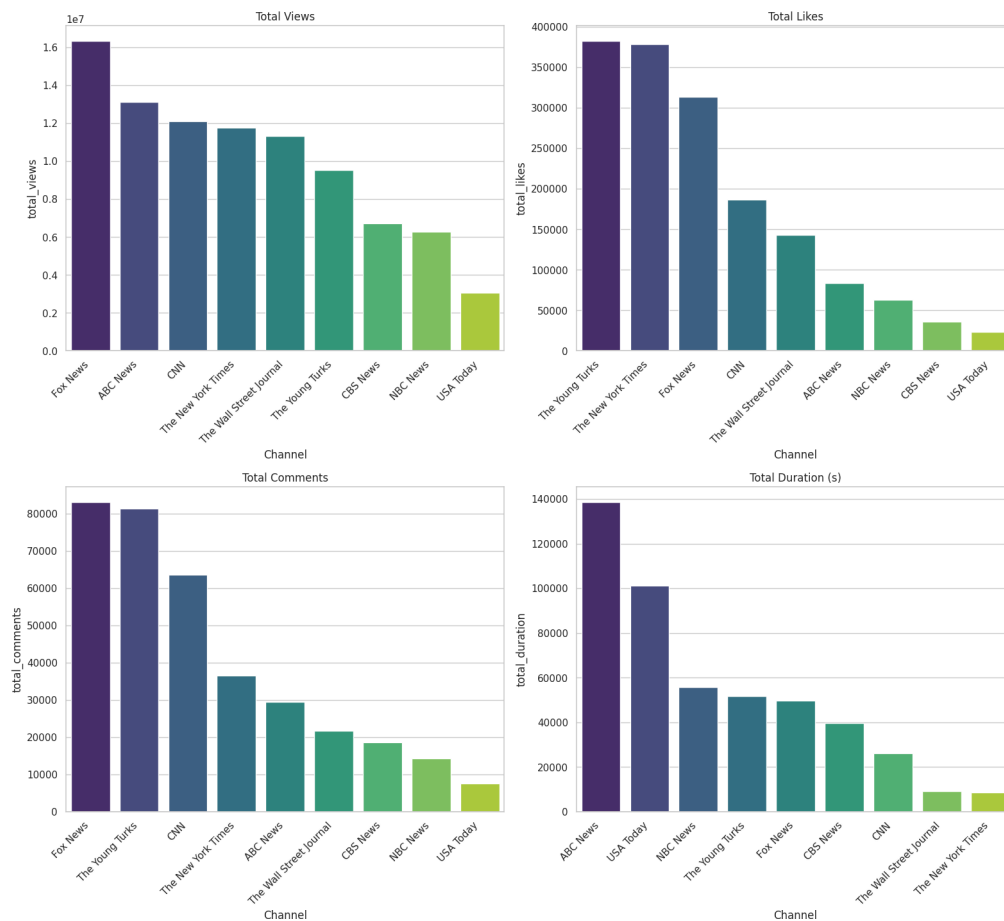


Figure 1: Summary of total scores (views, likes, comments, and duration)

3.2 Influence distribution and what “influence” highlights

Using the scaled influence score (log-transformed engagement index), we can rank individual videos and compare channels.

At the video level, we discovered that the top-ranked items include several high-performing videos from The Wall Street Journal and The New York Times, along with a high-performing The Young Turks video. For example, the ranked list includes: “I Tried the First Humanoid

Home Robot...” (WSJ) and “How the Louvre Jewelry Heist Unfolded” (NYT), each with Influence_scaled values in the mid-30s on the log scale. This illustrates a key property of the influence measure: it elevates videos that combine high viewership with substantial likes and comments, even if the channel is not the most-commented overall.

At the channel level, the project aggregates scaled influence across videos and visualizes it as centered horizontal bars. This view shifts attention from “who got one viral upload” to “who consistently generated engagement across the month.”

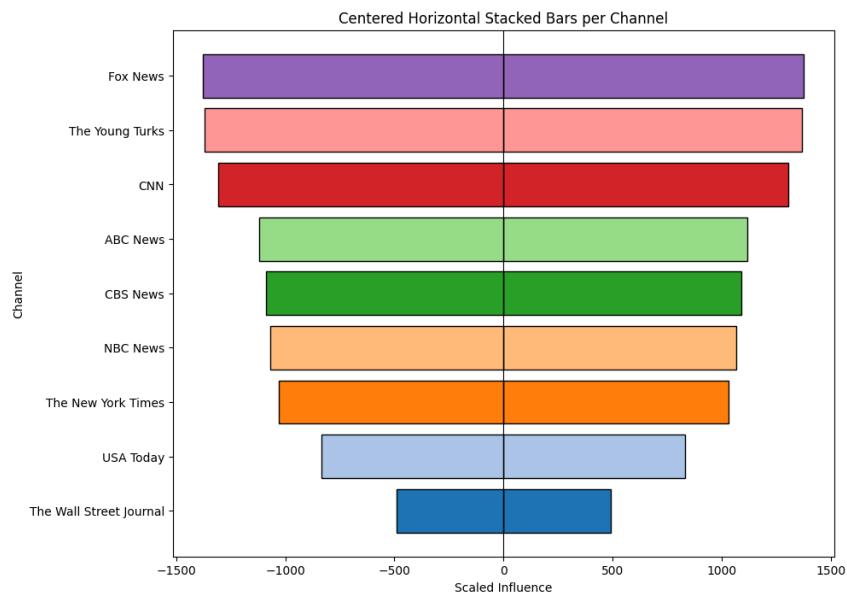


Figure 2: Aggregate Scaled Influence by Channel

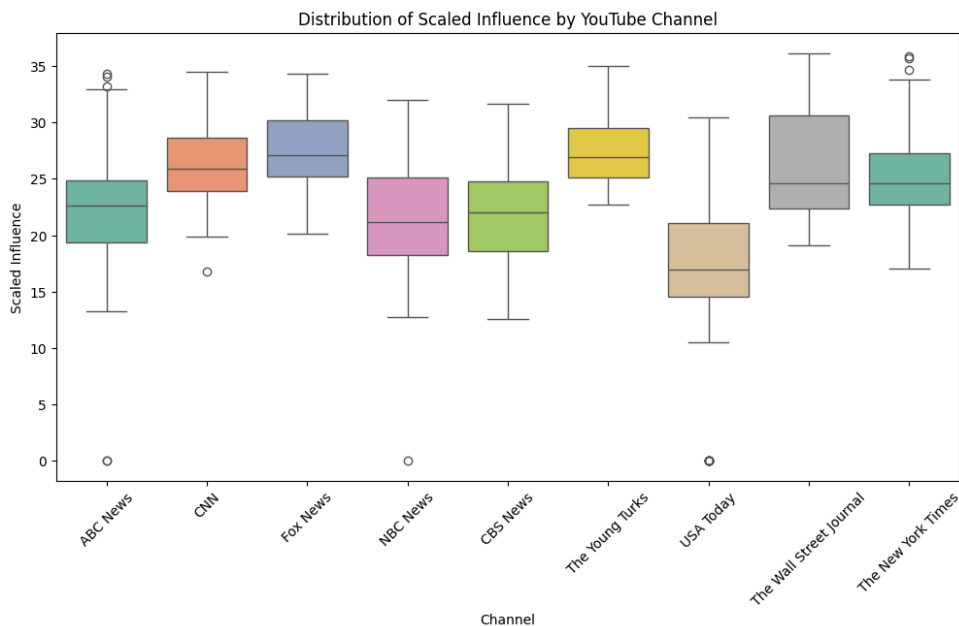


Figure 3: Influence Distribution Boxplot

3.3 Publishing themes and headline sentiment

Channel identity is partly reflected in recurring coverage themes. The project’s BERTopic-based thematic summary reports:

ABC News: General news / mixed coverage
CNN: Disaster / weather / humanitarian coverage
Fox News: Mixed / varied coverage
NBC News: Disaster / weather / humanitarian coverage
CBS News: Politics / government coverage
The Young Turks: Politics / U.S. politics
USA Today: Politics / government & international affairs
The Wall Street Journal: Politics / government / policy
The New York Times: Mixed / varied coverage

3.3.1 Headline sentiment

The title sentiment analysis asks whether publishers frame headlines as negative, neutral, or positive. A key takeaway from the project is that most channels keep headlines relatively neutral, likely to maintain broad appeal and avoid sounding extreme.

An additional lens comes from the emotion classifier run on 818 video titles: The most common predicted emotion was neutral (341 titles), followed by fear (179) and anger (108), with smaller counts for sadness, disgust, joy, and surprise. This does not contradict the “neutral sentiment” insight; instead, it suggests that even when sentiment is not strongly negative/positive, a meaningful share of headlines are emotionally keyed toward threat and conflict framing (fear/anger), which is consistent with political and crisis-oriented news cycles.

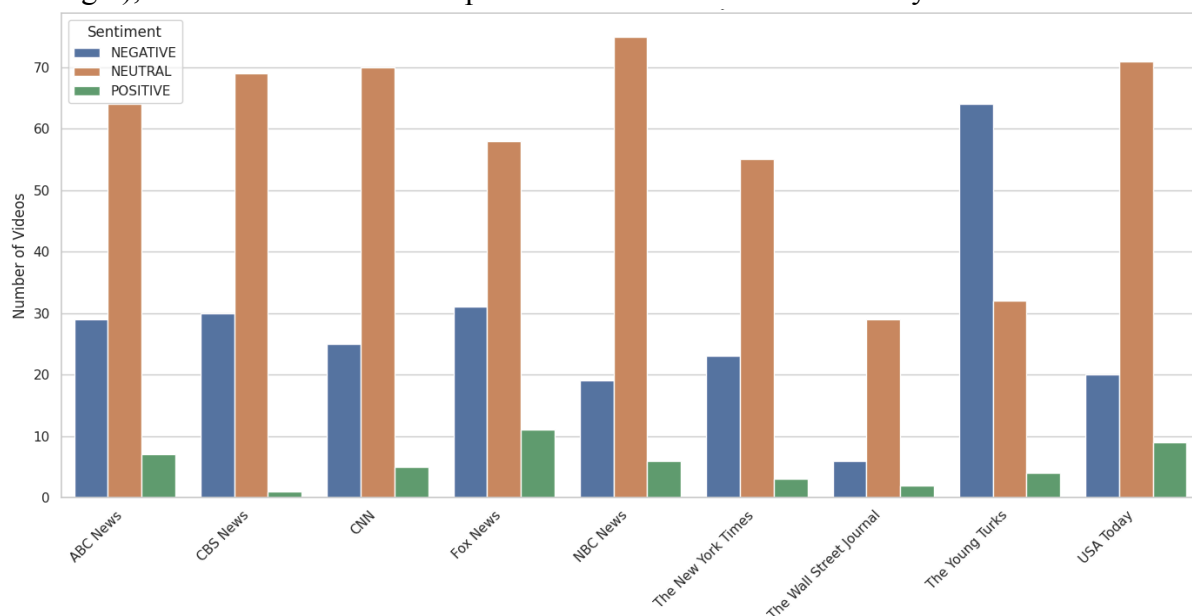


Figure 4: Title Sentiment Distribution by Channel

3.4 Audience discussion topics and comment sentiment

From our analysis, the comment corpus behaves like “political debate club + emotional debate.” We first ran a single-word TF-IDF analysis to identify the top related keywords in the comment after removing common stop-words. Most of them are strongly political words: most notably “trump” (top term), followed by words like “democrats,” “america,” “president,” and “government.” However, when titles and comment text are combined and analyzed via bigrams, the most frequent phrases include “mike johnson” (15,625), “government shutdown” (11,745), “snap benefits” (8,773), and “white house” (8,421), among others. These repeated phrases indicate that (at least in October 2025) the comment ecosystem is heavily oriented around U.S. political elites and institutional conflict.

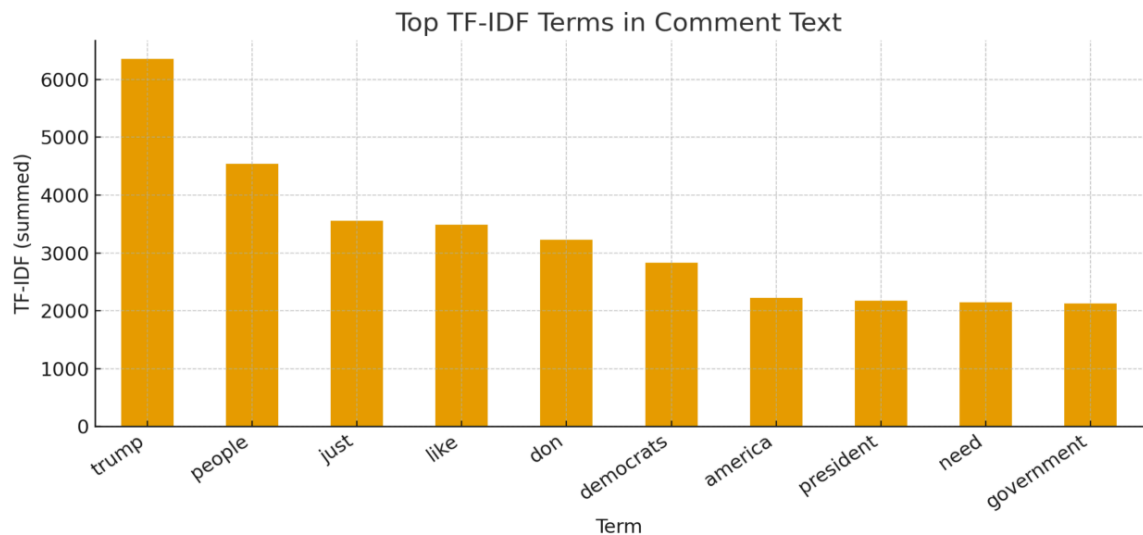


Figure 5: Top Related Single Words in Comments

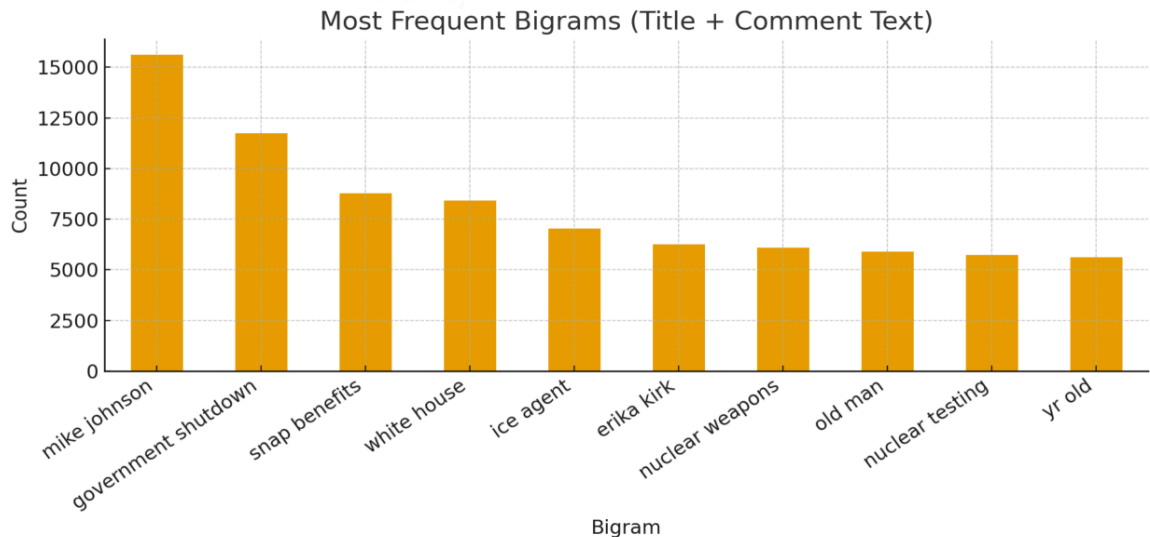


Figure 5: Top Related Bi-gram Words in Title and Comments

3.4.1 Comment sentiment: negative across every channel

Diving into the sentiments of 300,000 comments, our analysis demonstrates that almost two-thirds of the comments are negative (184,635), while 74,910 are neutral and only roughly 10% are positive (34,303). This is an unequivocal pattern. Negative sentiment is the dominant mode of audience participation in the comment sections.

More importantly, when visualized, the negative skew holds at the channel level. As shown in figure 7, Fox News' comments include 39,503 negative vs 16,955 neutral and 9,579 positive, while The Young Turks shows a whopping 55,582 negative vs 14,642 neutral and 5,876 positive. The same pattern is visible across all nine channels in the per-channel sentiment table.

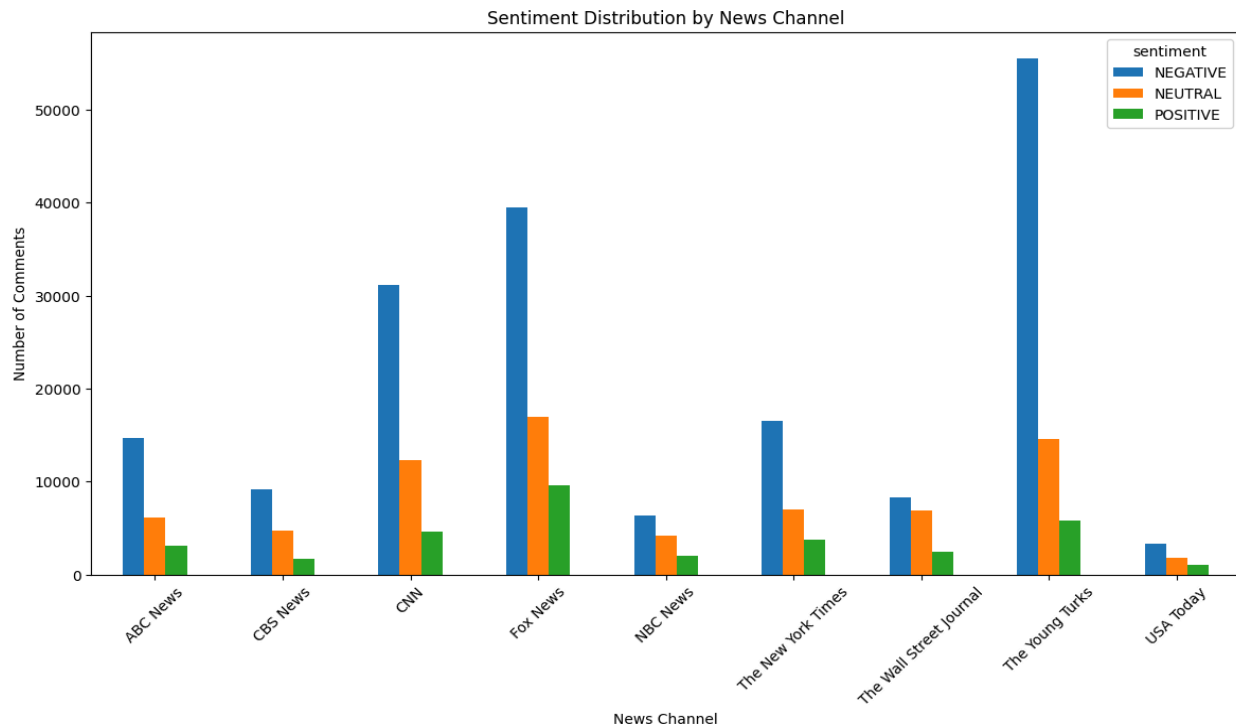


Figure 7: Comment Sentiment by Channel

4. Discussion

4.1 What the results suggest about publisher strategy vs. audience behavior

Two patterns stand out when we compare headline sentiment and comment sentiment. First, publishers appear to moderate their headline sentiment. This indicates that most channels keep their headlines relatively neutral, so that it appeals to broader audiences without sounding extreme. The emotion results reinforce this: “neutral” is the most common emotional label for titles, even though fear and anger are also frequent. One plausible interpretation is that publishers use professional norms and brand constraints to avoid overtly positive/negative phrasing, while still selecting topics (e.g., political conflict, institutional breakdown) that carry emotional charge.

Second, audiences respond with negativity regardless of publisher brand. The comment sentiment totals show a large negative majority overall and within each channel. This suggests that the comment section functions less as a space for “feedback” and more as a space for identity-signaling, contestation, and emotional release, especially for political news. The TF-IDF and bigram outputs support with that story: repeated phrases focus on national political elites and conflict topics such as “government shutdown” and “snap benefits.” Together, these findings point to a structural platform dynamic: professional framing is comparatively moderated, but participatory engagement is conflict-driven.

4.2 The Young Turks’ Surprise

Among all nine channels, The Young Turks (TYT) stands out for combining the shortest average video duration with one of the highest total comment counts and the most extreme share of negative comments. Despite producing shorter clips which are often under five minutes on average, TYT’s videos attract remarkably intense viewer participation. The sentiment breakdown shows roughly 55,582 negative comments, compared to only 14,642 neutral and 5,876 positive, giving it the highest negativity ratio in the dataset. This pattern suggests that TYT’s concise, opinion-driven video format efficiently triggers engagement through controversy and affective polarization. Shorter runtimes allow for rapid reaction cycles: Audiences can watch, react, and comment multiple times within the same news window. This can be attributed to the rise of short videos (Instagram Reels, Tiktoks, YouTube Shorts), which shortened peoples’ attention span nowadays. Furthermore, the disproportionate negativity does not necessarily signal audience hostility toward the channel; rather, it reflects a community dynamic in which strong political identity and criticism coexist. In effect, TYT exemplifies a “high-intensity, short-format” engagement model, where brevity and bold commentary amplify emotional response and discussion volume, highlighting how, on YouTube, influence can be generated not by length or production scale but by resonance and conflict intensity.

4.3 Limitations

We recognize several limitations in our research. First, channel coverage is incomplete because TMZ was not included due to scraping restrictions, so our findings generalize only to the nine channels we were able to collect reliably. Second, our sampling window is narrow because the dataset is limited to October 2025, so the results may reflect month specific dynamics, including pre-election news intensity, and may not represent other periods. Future work should repeat the analysis across additional months and around major events such as elections and crises to test whether these patterns hold over time. Third, API and comment availability constraints may affect representativeness because not all videos expose the same volume of accessible comments, and retrieval can be limited by platform settings such as disabled comments and by API rate limits. Finally, there are modeling limitations because transformer based sentiment classifiers can misinterpret sarcasm, slang, and domain specific political language, meaning that negative labels may sometimes capture disagreement, irony, or rhetorical style rather than straightforward hostility.

References

CardiffNLP. “cardiffnlp/twitter-roberta-base-sentiment-latest.” Hugging Face, Hugging Face, <https://huggingface.co/cardiffnlp/twitter-roberta-base-sentiment-latest>

Google. “API Reference: YouTube Data API.” *Google for Developers*, 6 Oct. 2025, <https://developers.google.com/youtube/v3/docs>

Grootendorst, Maarten. “BERTopic: Neural Topic Modeling with a Class-based TF-IDF Procedure.” *arXiv*, 11 Mar. 2022, <https://arxiv.org/abs/2203.05794>

Liu, Yinhan, et al. “RoBERTa: A Robustly Optimized BERT Pretraining Approach.” *arXiv*, 26 July 2019, <https://arxiv.org/abs/1907.11692>

Pedregosa, F., et al. “Scikit-learn: Machine Learning in Python.” *Journal of Machine Learning Research*, vol. 12, 2011, pp. 2825–2830. <https://www.jmlr.org/papers/volume12/pedregosa11a/pedregosa11a.pdf>

Reimers, Nils, and Iryna Gurevych. “Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks.” *ACL Anthology*, Association for Computational Linguistics, 2019, <https://aclanthology.org/D19-1410/>

Wolf, Thomas, et al. “Transformers: State-of-the-Art Natural Language Processing.” *ACL Anthology*, Association for Computational Linguistics, 2020, <https://aclanthology.org/2020.emnlp-demos.6/>