

# Topicality Boosts Popularity: A Comparative Analysis of NYT Articles and Reddit Memes

Kate Barnes<sup>1</sup>, Péter Juhász<sup>2</sup>, Marcell Nagy<sup>1</sup>, Roland Molontay<sup>1,3\*</sup>

<sup>1</sup>Department of Stochastics, Budapest University of Technology and Economics, Műegyetem rkp. 3, Budapest, H-1111, Hungary.

<sup>2</sup>Department of Mathematics, Aarhus University, Aarhus, Denmark.

<sup>3</sup>HUN-REN-BME Stochastics Research, Budapest, H-1111, Hungary.

\*Corresponding author(s). E-mail(s): [molontay@math.bme.hu](mailto:molontay@math.bme.hu);

Contributing authors: [kbarnes@edu.bme.hu](mailto:kbarnes@edu.bme.hu); [peter.juhasz@math.au.dk](mailto:peter.juhasz@math.au.dk);  
[marcessz@math.bme.hu](mailto:marcessz@math.bme.hu);

## Acknowledgments

We thank József Pintér for collecting Reddit data; Donát Kóller and Levente Murgás for their research insights. The authors declare they have no competing interests. The research was supported by Fulbright Hungary and the European Union project RRF-2.3.1–21-2022–00004 within the AI National Laboratory.

## Abstract

This study sheds light on interconnected topic dynamics across traditional news sources and social media platforms, emphasizing the influential role of topicality in shaping content popularity in social media. Using the Latent Dirichlet Allocation and BERTopic models, we define sets of 120 New York Times (NYT) topics to compare with 899,766 image-with-text memes from Reddit, showing that social media content aligns with many of the same topical patterns observed in news outlets. Topicality is formalized based on the temporal distributions of the topics over the past 5 years. Using these topicality features, the investigation reveals significant correlations between the rising popularity of NYT topics and increased average upvotes on Reddit, particularly evident in “innovator” memes posted during the early stages of a topic’s prevalence in the NYT. Furthermore, topicality features show significant predictive power over other content-based control features in a CatBoost classifier prediction of viral Reddit memes.

**Keywords:** memes, popularity prediction, machine learning, New York Times, Latent Dirichlet Allocation, BERTopic

# 1 Introduction

Popular content from social media sites such as Reddit provides a window into the opinions and interests of internet users. As people are spending more time on the internet, the question of what they are engaging with has come to the forefront. However, the forces underlying the emergence of popular internet content are still little understood. Particularly, the influence of memes' topics on their popularity has not yet been analyzed. Do social media sites discuss the same issues as verified news sources? Are topical subjects that are currently prevailing in the news more likely to go viral on social media too? Is it worth "riding the wave" of topicality, or does interest in a topical subject decline rapidly on social media sites?

The widespread prevalence of image-with-text memes online shows how well-suited the internet is to multi-modal reasoning and communication. Although video, sound, and image data are more complicated to analyze than pure text data, making it less studied in the literature, our understanding of online communication would be incomplete without these analyses (Halevy et al., 2022). To address this shortage of multi-modal analyses, here we study image-with-text memes. These are usually images with overlaid text that are copied and shared extensively across the internet. Due to higher information density, images are more succinct than pure text content. They grab our attention and are able to rapidly communicate complicated messages (Nelson et al., 1976). Like many forms of online media, the topic of a meme is not only communicated through the text caption but also through the underlying image.

In the digital age, many are concerned about losing the depth of conversation in the short-form discourse of the internet. However, social media content may simply achieve nuance in different ways than longer-form written communication. Multi-modality is one way of conveying nuance online. The level of nuance conveyed in online communication is further enhanced by referencing relatable and popular topics in society. The mechanism of reference is especially important to memes, which are elsewhere defined exactly as a reference to reality followed by a punchline (Spitzberg, 2014). By identifying meme topics, we connect the individual pieces of media to the broader discourse they are a part of. Using big data, we analyze overall trends in the discourse across 5 years of data.

Memes often reference cultural and political themes (Brodie, 2009; Du et al., 2020). They can be used to consolidate group identities (Dynel and Messerli, 2020) and influence opinions (McClure, 2016). The work of Zannettou et al. (2018) labeled memes from four social media sites with annotations from the Know Your Memes (KYM) website, indicating whether a meme template has political content. They found that memes on many social media sites contain political content, and on Reddit political memes are more popular than non-political memes. Today, over half of U.S. adults get news from social media (Jacob Liedke and Luxuan Wang, 2023). Furthermore, political memes have the potential to influence social actions and perceptions. For example, one study found that images shared online to promote the Black Lives Matter movement led to greater political action (Leach and Allen, 2017). More recently, TikTok has been used to increase public engagement in the Israeli-Palestinian conflict (Yarchi and Boxman-Shabtai, 2023). Importantly, online information is not always trustworthy and memes containing quotes and factual information have been shown to change

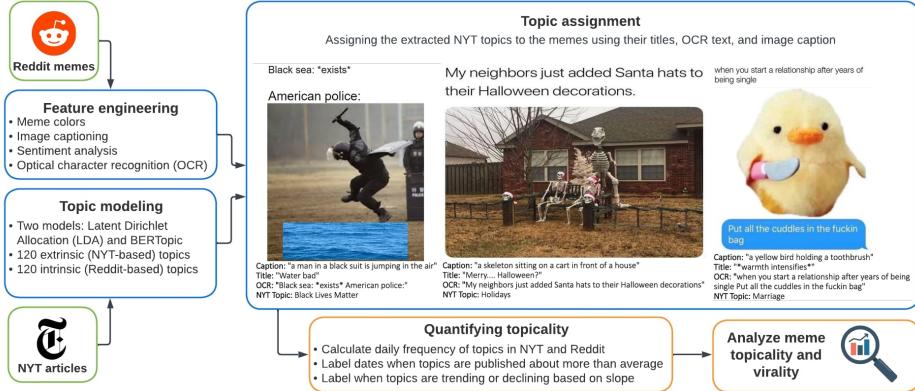
as they replicate across the internet (Simmons et al., 2011). Often memes do not make unbiased references but are imbued with opinions. In memes, the intertwining of informational and normative influence is abundantly clear (Cialdini and Goldstein, 2004).

Although many studies highlight the fact that online content references current societal events, topicality has yet to be studied more broadly. We introduce formal definitions of topicality in Section 3. Informally, when a particular topic is widely prevailing, it is said to be topical.

Topic analysis can be used to connect internet memes to the broader societal themes that they reference. Memes are not made in a tabula rasa environment but can be grouped by their topics, which in turn makes them relevant. Liking and re-posting are implicit endorsements of the memes’ topics and opinions (Spitzberg, 2014), implying that topics may influence meme diffusion. Indeed, previous research shows that individuals are more likely to spread information about topics of interest to them. Twitter users are more likely to adopt hashtags that align with their own topical interests (Grabowicz et al., 2016) and retweet hashtags that they have already tweeted about in the past (Weng et al., 2012). Typically, the diffusion of memes through the internet is modeled as a simple contagion process, similar to the spread of disease (Wang and Wood, 2011), however, topical content clearly exhibits more complex diffusion mechanisms. Poux-Médard *et al.* studied the interactions between topics in the diffusion of information on Reddit news pages, showing that the topics of previous posts have minor explanatory power for predicting the topics of subsequent posts (Poux-Médard et al., 2022). These analyses suggest that topics are interrelated and exhibit popularity trends.

The majority of previous research using topic analysis models such as Latent Dirichlet Allocation (LDA) (Blei et al., 2003) and BERTopic (Grootendorst, 2022) on social media data focuses on the topical interests of internet users. To the best of our knowledge, this is the first study to address topicality comprehensively. Weng and Menczer show that if a particular hashtag on Twitter is adopted by people with diverse topical interests, the hashtag is more likely to go viral, due to exposure to a diverse, and presumably broadly connected, community (Weng and Menczer, 2015). In other cases, homogeneous communities with strong central figures promote content virality when the topic of the shared content matches the interest of the community members, implying the importance of shared beliefs and attitudes (Cialdini and Goldstein, 2004). Preferential attachment can be described via topical analysis as well. For example, Instagram users tend to connect with other users with the same interests as them (Ferrara et al., 2014). All of these topic analyses only used one document corpus to define and analyze topics, whereas our analysis uses a benchmark dataset (NYT) to identify the topics of Reddit memes, rather than simply extracting topics from the memes themselves. This comparative aspect leads to better-identified topics and can provide insight for other cross-corpus topic analyses.

In our previous work, we showed that Reddit memes with COVID-19-related content were more likely to go viral in March 2020 at the beginning of the pandemic than memes without COVID-19-related content (Barnes et al., 2021). Other events likely exhibit similar properties. We assume that, in general, topical posts about subjects



**Fig. 1:** Research workflow including examples of OCR text extraction, BLIP image captioning, and meme topic assignments.

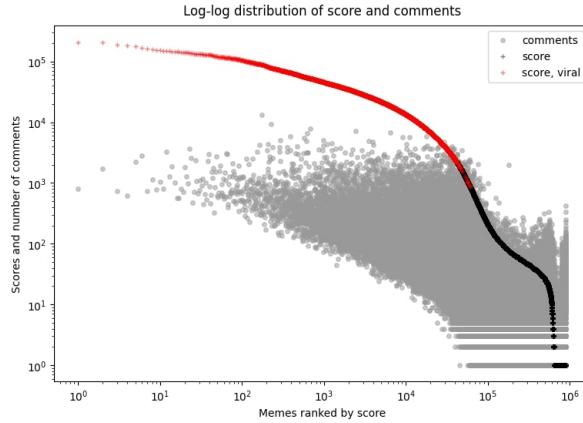
that are currently widely discussed in society are more likely to go viral online. In the present work, we investigate this hypothesis using the New York Times (NYT) as a reference document for topicality. Specific contributions are listed below.

- Identification of Reddit- and NYT-based topics.
- Comparing the temporal distributions of 120 topics on Reddit and NYT.
- Showing that Reddit posts about topics that are gaining prevalence in the NYT also receive more upvotes on Reddit, including an innovator’s advantage for being among the first to post about a topic.
- Demonstrating topic-based features have significant predictive power in meme virality prediction.

In order to establish these findings, we performed a topical analysis of NYT and Reddit data sets, engineered topic-based features, and used these to predict viral memes with a CatBoost classifier (Prokhorenkova et al., 2018). The research workflow is summarized in Figure 1. Section 2 discusses the data sets and engineering of control features describing the Reddit memes. Topicality features, describing topics extrinsic and intrinsic to the Reddit data, are discussed in Section 3. Section 4 discusses the training of a Catboost classification model to predict viral memes using these topic-based features. We analyze the incremental predictive power of topicality features over other content-based control features, as well as the importance of intrinsic vs. extrinsic topicality features in the same section.

## 2 Data Description and Preparation

Data was collected from two sources: image-with-text memes from a popular social media site called Reddit, discussed in Section 2.1, and archived article metadata from



**Fig. 2:** Distribution of upvotes in Reddit data and corresponding number of comments each meme received.

the New York Times (NYT), see Section 2.2. Sections 2.3 and 2.4 discuss the engineering of content-based control features such as the color content of meme images and the sentiment of meme texts. Both data sets can be found on Github ([HSDS, 2024](#)).

## 2.1 Reddit

Reddit, nicknamed "the front page of the internet", is the source of a lot of viral internet data, making it likely that the viral content analyzed here also circulates elsewhere on the internet ([Sanderson and Rigby, 2013](#)). Data from r/Memes, the largest community dedicated to sharing memes on Reddit, were collected using the Pushift API ([Baumgartner et al., 2020](#); [Boe, 2016](#); [Podolak, 2021](#)). We limited the collection to a maximum of 1000 randomly sampled posts per day. There were 37 missing dates from which we were unable to collect data due to issues with the API. After removing items with broken image links and gifs, the Reddit data set contained 899,766 memes from between January 1, 2018 and November 14, 2022. The r/Memes subreddit is the 12th largest community on Reddit, with more than 26 million subscribers. Based on the above characteristics, we believe that the collected data is representative and our conclusions can be applied to memes appearing elsewhere.

In addition to the meme images, the API provides several metadata features including the title and caption posted with the meme, publication date, an over-18 content indicator variable, the number of comments, and the score. On Reddit, the score of a post is calculated as the number of upvotes minus the number of downvotes it received. These attributes provide opportunities to engineer various other features.

As visualized in the memes' score distribution in Figure 2, few memes receive a lot of attention, while the vast majority go unnoticed. In general, memes with high scores also receive more comments, reflected by the Spearman rank-order correlation coefficient of  $r = 0.65$  ( $p < 0.001$ ). In many places on the internet including Reddit, popular content whether measured by likes, followers, or views, is sorted to the top of

| Feature            | Type        | Description                                 |
|--------------------|-------------|---|
| Score              | Int         | Number of upvotes minus downvotes           |
| Comments           | Int         | Number of comments                          |
| Viral              | Binary      | 1 if score in top 5% of memes               |
| Date               | String      | Date on which meme was posted               |
| day-of-week        | Categorical | Day of week when meme was posted            |
| all-text           | String      | Title, BLIP caption and OCR text            |
| Title length       | Int         | Number of characters in title               |
| OCR length         | Int         | Number of characters in OCR text            |
| Over 18            | Binary      | Reddit metadata, content warning            |
| Emoji              | Binary      | 1 if memes' all-text feature contains emoji |
| Sentiment          | Float       | Text valence score                          |
| Height             | Float       | Thumbnail height                            |
| Width              | Float       | Thumbnail width                             |
| HSV                | Float       | Average HSV image components                |
| RGB                | Float       | Average RGB image components                |
| 10 colors          | FLOATS      | Percentage of color in meme image           |
| Face               | Binary      | 1 if the meme image contains a face         |
| Topic num.         | Categorical | Topic assigned to meme                      |
| Topic prob.        | FLOATS      | Probabilities associated with topic         |
| Topic entropy      | FLOATS      | Shannon's entropy of memes' topics          |
| Monthly Topicality | Float       | Monthly average topicality                  |
| Daily Topicality   | Float       | Daily average topicality                    |
| Slope Topicality   | Float       | Slope of topicality distribution            |

**Table 1:** Features table. Shaded rows indicate features used for the Catboost identification of viral memes including topic-related features (red) and control features (gray).

users' news feeds (Figueiredo et al., 2014). Reddit's default content sorting algorithm, used to curate the "hot" tab, prioritizes content based on a mixture of how recently it was posted and the logarithm of the score (Amir Salihefendic, 2015). Content ranking methods such as this create a feedback loop between what is popular and what is visible online. While unpopular memes are visible briefly when they are first posted, popular memes are widely viewed. One study manipulated the score of Reddit posts, finding that by providing 10 upvotes soon after the content was posted they could greatly increase the posts' chance of going viral (Carman et al., 2018). This can explain the heavy-tailed distribution of popularity measures on many sites, including in our Reddit data.

It is a common choice to model the heavy-tailed popularity distributions on social media as binary variables (Weng and Menczer, 2015; Ling et al., 2021). Here, we define virality as a percentile. If a given meme is in the top 5 percent of memes posted within a  $\pm 7$ -day period around it, it is labeled viral (1), and otherwise, it is labeled non-viral (0). By measuring virality locally, in a two-week period, we avoid replicating seasonal popularity trends on Reddit. Additionally, as the Reddit userbase is constantly growing (Tiago Biachi, 2023), this assures we do not overstate the influence of memes from recent years. The viral/non-viral target variable, also visualized in Figure 2, is used in the binary classification task discussed in Section 4.

Table 1 summarizes all features describing the Reddit meme data. Shaded rows indicate features used for training the predictive model, including topicality-related features (red) and control features (gray).

## 2.2 New York Times

Data from the same timeframe as the Reddit memes, January 1 2018 to November 14, 2022, was collected from the New York Times Archives API<sup>1</sup>. The API returns all article metadata for a given month, including the abstract, snippet, lead paragraph, headline, publication date, keywords, and type of material. Most (72%) of the collected data was from the news section, but other types of material such as crosswords, obituaries, sports, and book reviews were also included. The NYT published roughly 200 articles per day on weekdays and 100 articles per day on weekends, and publishing rates showed a slight decline over the 5-year period we examined. On a few dates, the archives contained more than 500 items on a single day due to the NYT updating the archived podcast episodes, or, during the first year of the COVID-19 pandemic and during US elections, due to state-by-state statistics reports. These statistics reports and podcasts were removed to prevent the topic models from forming topics specific to this type of post and to ensure an even daily distribution of articles. No other statistics that should be handled similarly were identified. After removing these entries, our data set contained 255,783 NYT articles in total.

## 2.3 Text features

We extracted text from the Reddit meme images using Optical Character Recognition<sup>2</sup>, and generated image captions using Bootstrapping Language-Image Pre-training (BLIP) (Li et al., 2022). Examples of OCR and image caption results can be seen in Figure 1. The OCR text and image captions were combined with the title and caption posted with the meme image on Reddit into one feature containing all text associated with the meme. This all-text feature was subsequently used for topic extraction and assignment in the Reddit data. In the NYT data set, the abstract, snippet, lead paragraph, headline, and keywords were all combined into one all-text feature for topic extraction.

The all-text features for the Reddit and NYT data were cleaned identically to assure alignment when modeling the topics in both document corpora. We removed punctuation, made the text lowercase, removed stop words, and stemmed the words to their root forms using NLTK (Bird et al., 2009). Additionally, we made custom edits to the NYT text data based on observed differences between the NYT and Reddit vernaculars. For example, the NYT referred to the COVID-19 pandemic with the word “coronavirus” whereas posts on Reddit tended to use “covid”, “rona”, “corona”, and “pandemic”. This is important as the LDA algorithm identifies and assigns topics based on shared words. Therefore, if a post on Reddit referred to “covid”, it would not necessarily be associated with an NYT topic that only included the word “coronavirus”. To solve this issue, we added the words “covid”, “rona”, “corona”, and “pandemic” to the all-text feature of every NYT article that was tagged with the NYT keyword “Coronavirus”. Other examples include adding “RGB” to articles tagged with Ruth Bader Ginsberg and “BLM” to articles tagged with Black Lives Matter. Additionally, acronyms of the form “G.O.P.” and “N.F.L.” were edited to “GOP” and “NFL” in

---

<sup>1</sup><https://developer.nytimes.com/apis>

<sup>2</sup>OCR.space Free OCR API and Online OCR, <https://ocr.space/>

order to not lose this information when punctuation was removed in the text cleaning process.

In addition to text cleaning, we extracted numerical features from the Reddit text data. These features were used as control features describing the Reddit memes for classification. As seen in Table 1, we recorded the number of characters in the title posted with the meme and the OCR text extracted from the meme image. Other studies show that on Twitter, the length of the post is strongly correlated with popularity ([Tsur and Rappoport, 2015](#)). A binary variable indicating whether or not the all-text feature contained emojis was engineered using the Emoji 2.8.0 identification Python library ([Kim and Wurster, 2023](#)). The valence of the Reddit title and OCR-extracted text was analyzed with the NLTK sentiment model ([Shreyas, 2019](#)). Sentiment scores close to 1 are more positive while sentiment scores close to 0 represent more negative sentiment.

## 2.4 Image features

In addition to text features, we extracted numerical features from the meme images following a similar procedure to our previous work ([Barnes et al., 2021](#)). In total, we used 19 image-related features, including high-level features such as the BLIP image descriptions discussed in the previous section and low-level features such as the color content of the images.

The pre-trained Multi-Task Cascading Convolutional Neural Network (MTCNN) was used to detect faces in the meme images ([Zhang et al., 2016](#)). The model returned the probability that a face was present and the number of faces present in the image. For simplicity, we elected to use a binary variable indicating whether or not the meme image contained a face. Using the OpenCV image segmentation technique to mask the meme images, we calculated what percentage of the image area contained each of 10 color ([Stone, 2018](#)). The average hue, saturation, and value components of the HSV representation of the images and the average red, green, and blue components of the RGB image were also used as control features.

Previous work has shown that low-level image features such as those mentioned above can have a large effect on popularity ([Russakovsky et al., 2015; Khosla et al., 2014](#)). However, while an earlier study on the aesthetics of images shows that high-definition, bright-colored images are more appealing ([Datta et al., 2006](#)), in the case of memes it is the opposite. Popular memes generally contain dull colors ([Barnes et al., 2021](#)). For memes, these low-level features could encode the template image used to create memes, indicating that meme templates rather than the low-level features themselves, have an impact on popularity ([Coscia, 2014](#)).

## 3 Topic Modeling

Using two topic models, LDA and BERTopic, we analyzed the topics of Reddit memes and NYT articles. For comparison, we used top models from both paradigms of topic modeling: probabilistic (LDA) and embedding (BERTopic) approaches ([Egger and Yu, 2022](#)). Each topic is represented by a list of its most common words, as shown in the example topics in Tables 2 and 3. To identify “extrinsic topics” (NYT-based),

| %   | Corpus | Type         | Topic and its top 10 words  |
|-----|--------|--------------|---|
| 6.3 | NYT    | NYT-based    | <b>US Politics &amp; Government</b><br>“state”, “unit”, “govern”, “polit”, “president”,<br>“trump”, “nation”, “american”, “washington”, “parti” |
| 6.3 | NYT    | NYT-based    | <b>Coronavirus</b><br>“coronaviru”, “pandem”, “vиру”, “covid”, “corona”,<br>“rona”, “ncov”, “quarantin”, “reopen”, “mask”                       |
| 4.4 | Reddit | NYT-based    | <b>Technology</b><br>“compani”, “industri”, “compute”, “internet”,<br>“social”, “medium”, “facebook”, “online”, “technolog”                     |
| 3.7 | Reddit | NYT-based    | <b>Cuisine</b><br>“food”, “restaur”, “cook”, “recip”,<br>“cookbook”, “chef”, “farm”, “wine”, “drink”, “tabl”                                    |
| 0.8 | Reddit | Reddit-based | <b>Video Games</b><br>“game”, “play”, “video”, “drive”, “minecraft”,<br>“car”, “skeleton”, “mustach”, “control”, “player”                       |
| 0.8 | Reddit | Reddit-based | <b>Cuisine</b><br>“milk”, “thano”, “slow”, “soup”, “mushroom”,<br>“popcorn”, “spoil”, “hunger”, “drank”, “waiter”                               |

**Table 2:** Most common LDA-identified topics in NYT and Reddit corpora. Words are listed in order of importance to the topic.

we trained the topic models on the NYT data and performed inference on the Reddit data. To determine “intrinsic topics” (Reddit-based), the topic models were trained on and inference was performed on the Reddit text data. From the distributions of each topic over time we defined topicality features for the Reddit data. In total, 12 topicality features were added for each topic model. This is twice the number listed in Table 1 because the features were calculated using both the extrinsic and intrinsic topic models.

In this section, we discuss the LDA topic model (Section 3.1), the BERTopic model (Section 3.2), and the formalization of topicality into numerical features (Section 3.3). At the end of Section 3.3 we report statistics describing the relation between topicality features and the popularity of the Reddit memes. We display plots comparing topicality trends from the the past 5 years of NYT and Reddit memes data. The interested reader can see the complete sets of NYT-based topics identified by LDA and BERTopic on our GitHub page ([HSDS, 2024](#)).

### 3.1 Latent Dirichlet Allocation

LDA is a statistical language model used to distill topics from a corpus of documents using a bag-of-words approach in which word order is not considered. Words that commonly co-occur in documents are grouped into topics. Topics are then represented as a probability distribution over all of the words in the document corpus,

| %   | Corpus | Type         | Topic and its top 10 words  |
|-----|--------|--------------|---|
| 0.9 | NYT    | NYT-based    | <b>Arts</b><br>“artist”, “work”, “art”, “sculptur”, “museum”,<br>“exhibit”, “galleri”, “artist”, “play”, “novel”            |
| 0.5 | NYT    | NYT-based    | <b>Sports</b><br>“team”, “soccer”, “leagu”, “game”, “player”,<br>“basebal”, “footbal”, “tenni”, “basketbal”, “sport”        |
| 0.7 | Reddit | NYT-based    | <b>Police Violence</b><br>“polic”, “murder”, “shoot”, “homicid”, “attempt”,<br>“murder”, “homicid”, “offic”, “shot”, “kill” |
| 0.6 | Reddit | NYT-based    | <b>Cuisine</b><br>“restaur”, “cook”, “recip”, “cookbook”, “cook”,<br>“chef”, “food”, “wine”, “dish”, “chicken”              |
| 0.6 | Reddit | Reddit-based | <b>Animals</b><br>“dog”, “cat”, “sit”, “cat sit”, “banana”, “dog sit”,<br>“bear”, “pet”, “cartoon cat”, “pictur dog”        |
| 0.6 | Reddit | Reddit-based | <b>Technology</b><br>“phone”, “cell phone”, “cell”, “share”, “vote”,<br>“reddit”, “meme”, “youtub”, “comment”, “post”       |

**Table 3:** Most common BERTopic-identified topics in NYT and Reddit corpora. Words are listed in order of importance to the topic.

with the words most highly associated with the given topic receiving the highest probabilities. In the inference step, the word distributions of the topics are compared with the word distributions of each document to assess the probabilities with which each document discusses each single topic. Documents are then represented as probability distributions over the set of topics (Blei et al., 2003).

One limitation of the LDA algorithm is the requirement to specify the number of topics in advance. We used the *keywords* metadata provided by the NYT to estimate an appropriate number of topics in the NYT data set. The NYT indexed 1,583,868 unique keywords covering 92% of the articles, the most used keywords being “US Politics and Government”, “Trump, Donald J.”, and “Coronavirus (2020-nCov)”. We filtered the keywords list based on three criteria. The keyword should be used, 1. at least 100 times total, 2. on at least 50 distinct dates, and 3. at least 5 times on at least one day. Only, 127 unique keywords fit these criteria, providing an estimated number of topics for the NYT corpus, which we rounded down to 120 for simplicity.

The NYT *all-text* feature was used to train the LDA model parameterized to find 120 topics in 1500 iterations over the corpus. We labeled articles in the NYT corpus according to their highest probability topic. Table 2 shows that the most prevalent topics identified in the NYT data are about US politics & government and Coronavirus. This aligns with the most common NYT-applied keywords, lending confidence to our model. Using the trained LDA model we then performed inference on the Reddit *all-text* feature, resulting in a topic distribution assigned to each meme in the Reddit data,

which described the probability that the meme is about each of the 120 NYT-based topics. Later, we will refer to this topic distribution as the vector,  $\tau_m$ .

We labeled the Reddit memes according to their highest probability topic, saving that topic's probability of being assigned to the meme as well. Additionally, we calculated the Shannon's entropy of the memes' topic distributions. Memes uniformly assigned to all topics receive high entropy values while memes with few dominant topics receive low entropy values. These are the "topic", "prob" and "topic entropy" features described in Table 1.

In the same manner, a set of 120 intrinsic topics (500 iterations over the Reddit corpus) was extracted by training the LDA algorithm on the Reddit *all-text* feature. In this case, we chose 120 topics to be consistent with the NYT-based topic set. Again, the highest probability intrinsic topic, its associated probability, and the Shannon's entropy of the topic distribution were saved for the Reddit memes.

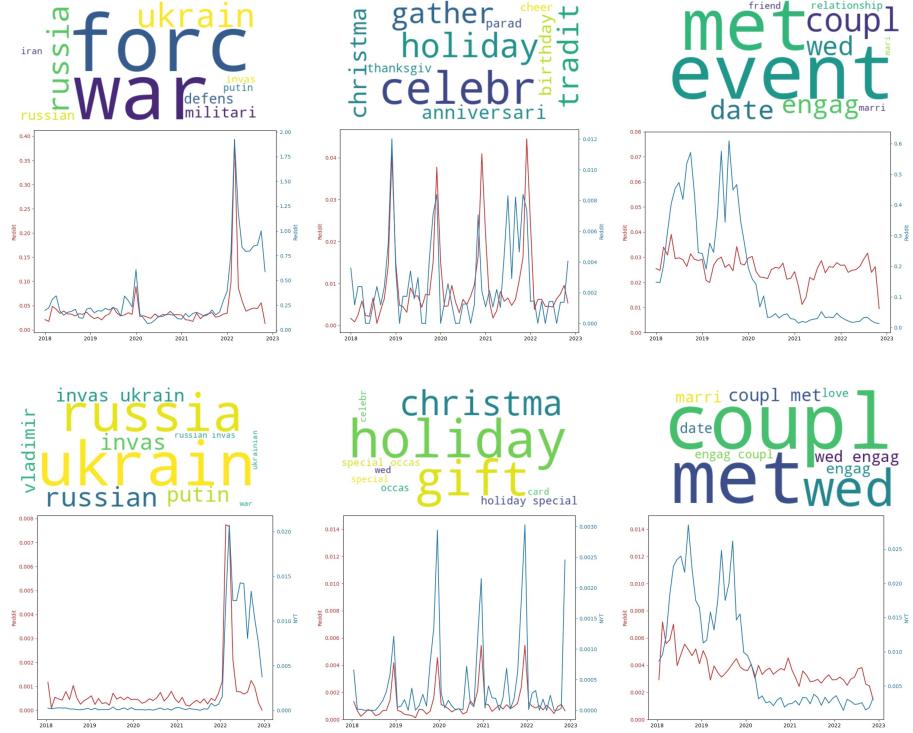
Both the intrinsic and extrinsic topic models created one miscellaneous topic in which the most important words in the topic were common words like "already", "life", "way", "time" and "person". This topic was the most prevalent in both the NYT and Reddit corpora, but it did not dominate the data sets. In both cases, the miscellaneous topic was assigned to less than 10% of the documents. Table 2 shows the top 2 most frequent topics in the Reddit and NYT data excluding this miscellaneous topic.

### 3.2 BERTopic

BERTopic is a newer topic modeling technique that uses a transformer based approach (Grootendorst, 2022). The model generates high-dimensional document embeddings, reduces the dimensionality of these embeddings with UMAP, and then clusters the vectors into topic groups using HDBSCAN. Word order is important for the transformer-based model to understand the context of texts therefore notably, word order was not changed during the text cleaning steps described in Section 2.3. To ensure comparability, we used the same text data for both the BERTopic and LDA models.

A set of 120 extrinsic topics was extracted using BERTopic from the NYT *all-text* document corpus, and inference was run on the Reddit *all-text* corpus to assign these extrinsic topics to the Reddit memes. Although it is not necessary to specify the number of topics before fitting BERTopic, the number of topics can be reduced after training. Our BERTopic model was parameterized to cluster topics containing at least 100 documents, resulting in 194 NYT topics. This was reduced to a set of 120 topics, to match the LDA model, using automatic and manual approaches. We applied BERTopic's default topic-reduction technique which merges the most similar HDBSCAN clusters. Additionally, we qualitatively assessed the similarity of topics' word distributions and manually merged those that were most similar. We trained a second BERTopic model to find 120 intrinsic topics, using only the Reddit *all-text* features, following the same procedure.

As with the LDA model, BERTopic initially generated miscellaneous topics with common words such as "like" and "also". These topics were manually combined into one miscellaneous topic which, again was the most prevalent topic assigned to both



**Fig. 3:** Monthly average topicality in the NYT (blue) and Reddit (red) as identified by the LDA (top row) and BERTopic (bottom row) topic models. Word clouds above the distributions show the top ten words associated with the topic.

the NYT and Reddit data. However, again this topic accounted for less than 10% of the data sets.

A disadvantage of BERTopic is the number of outliers produced. These are documents that are not assigned to any topic. To address this, we used the BERTopic *calculate probabilities* parameter to calculate the probability each document belonged to every topic. Then, we automatically assigned every document to its highest probability topic, resulting in no outlier documents. Furthermore, this parameter was used to match the LDA results, in which every document is represented by a probability distribution over the topics,  $\tau_m$ . As in the case of the LDA model, we supplemented the meme data with the highest probability topic, its probability, and the Shannon's entropy of the topic distribution.

The BERTopic and LDA topic models obtained similar results. As shown in Figure 3, there was a considerable amount of overlap between the topics found by the LDA and BERTopic models, and their 5-year temporal distributions. Table 3 shows the top 2 most common topics identified by BERTopic in the NYT and Reddit data sets. In the NYT, the most prevalent topics were those that the NYT publishes about regularly, such as arts events, real estate, and sports. COVID-19 and the invasion of



**Fig. 4:** Stream charts showing 5-year trends for 20 NYT-based topics identified by BERTopic.

Ukraine were the only two current events that appeared in the top 10 NYT topics. Notably, topics about police brutality and sexual assault were among the top 10 most prevalent NYT-based topics in the memes data as identified by both models, but these topics were not as prevalent in the NYT data.

Although NYT and Reddit follow many of the same topicality trends, the proportions with which they publish about these topics differs. This can be seen in the stream charts in Figure 4. The NYT publishes much more about the arts, economics and real estate than Reddit, whereas Reddit posts are more likely to discuss animals and technology than the NYT. Notably, current events are discussed at length in the NYT, but on Reddit interest in these topics declines rapidly after the date of the event. The impact of Covid-19 on NYT publishing is also evident in Figure 4 as live events such as theater performances, weddings and fashion week were cancelled due

to the pandemic. Moreover, the effect of Reddit outages, when the website was down for 5-10 hours in 2019 and 2021, can also be observed.

### 3.3 Topicality Features

In addition to saving the topic, probability and topic entropy for each Reddit meme, we developed variables to assess whether trending, “topical”, topics were more popular than non-trending topics. These variables were calculated for both the LDA and BERTopic models and for both the NYT-based, extrinsic topics and the Reddit-based, intrinsic topics.

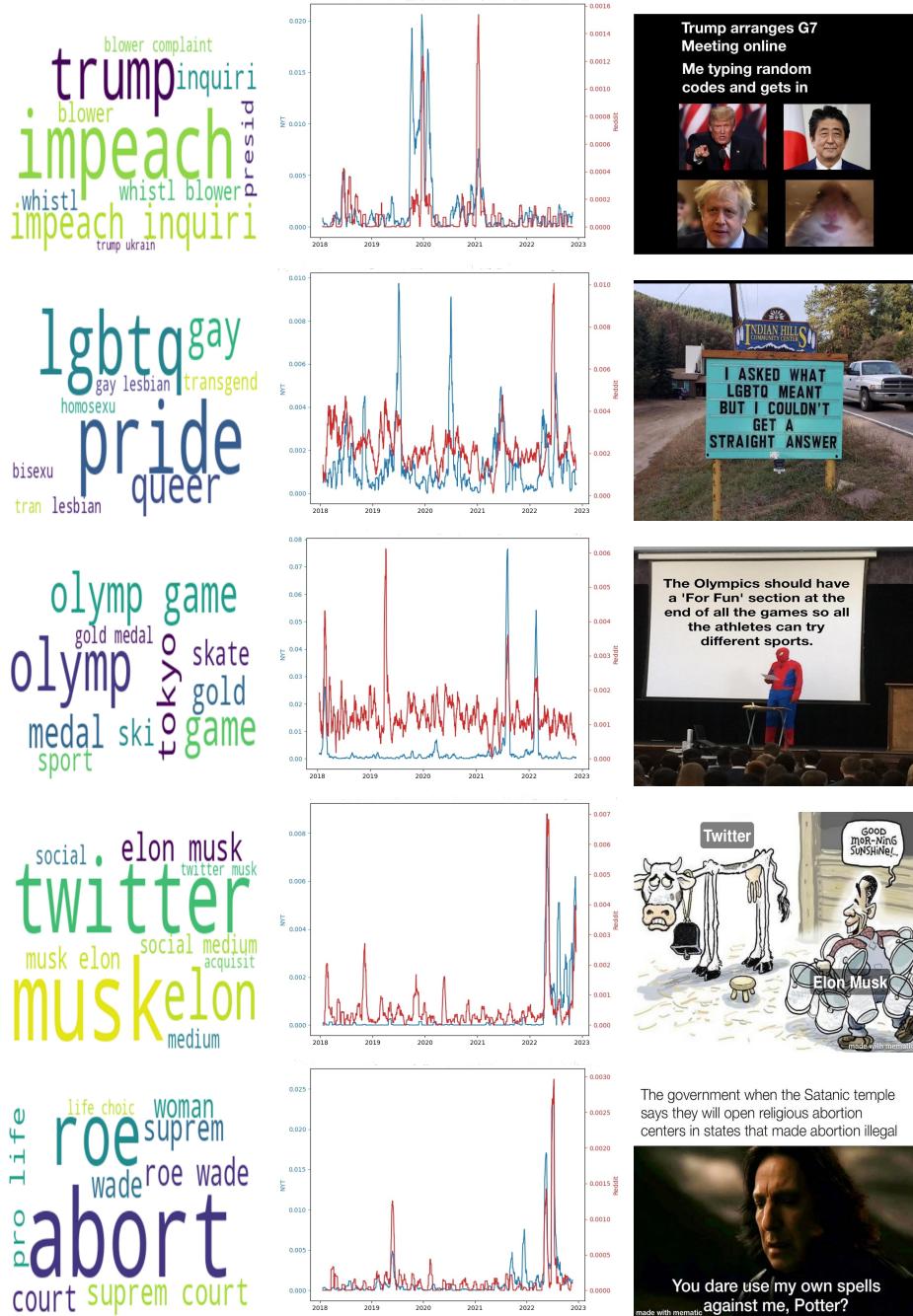
#### 3.3.1 Defining Trending Topics

First, we define trending topics in terms of how frequently they were published about on given dates. This topicality variable,  $\tau_t(d)$ , was calculated as the normalized sum of the probability with which documents were assigned to a given topic  $t$  on a given date  $d$ . Formally,

$$\tau_t(d) = \frac{\sum_{a \in A_d} p_a(t)}{|A_d|}, \quad (1)$$

where  $A_d$  is the set of articles published on date  $d$  and  $p_a(t)$  is the probability that document  $a$  is about topic  $t$ . Note that  $\tau_t(d)$  can be calculated for either daily or monthly granularity. Figures 3 and 4 show monthly topic trends, meaning  $d$  in equation (1) was the month and year in which the document was posted. Figures 5 and 6, on the other hand, show daily topic trends, meaning  $d$  was the date on which the document was posted. Furthermore,  $\tau_t(d)$  was calculated for both NYT articles and Reddit memes. We use “document” to refer to *either* an article or a meme. The blue lines in the above-mentioned figures show NYT distributions in which  $a$  in equation (1) represents NYT articles, whereas for the red lines,  $a$  represents Reddit memes. The plots serve to compare topical trends in the NYT and Reddit and thus are only based on the sets of extrinsic topics.

Topics about current events typically exhibited a rapid spike followed by a slower decline. Oscillating topics appeared for sports events, holidays and other regular celebrations such as Pride. In many topics, like the invasion of Ukraine and Christmas holidays shown in Figure 3, the distribution of Reddit memes posted about the given topics aligns very well with the distribution of NYT articles. In other topics, such as the marriage and romance topic in Figure 3, the two sources differ. We can also see that while the topic models identify the documents associated with each topic well, they are not perfect. For example, some NYT articles about Russian interference in the US 2020 presidential elections are associated with the topic of the invasion of Ukraine by the LDA topic model, as seen in the spike in 2020 in the LDA Ukraine invasion topic in Figure 3. This is likely due to words about Russia being prevalent in both topics. Despite this disadvantage, the identified topics describe the majority of the memes well for both algorithms used.



**Fig. 5:** BERTopic-identified topical alignments between the NYT and Reddit.



**Fig. 6:** LDA-identified topical alignments between the NYT and Reddit.

### 3.3.2 Defining Topical Memes

Equation (1) and the corresponding charts describing trending topics can be used to identify which memes are topical and which are not. Here, we consider Reddit memes topical if they are about events which are currently on the rise in NYT publishing. Visually, memes published about a given topic when there is a spike in the corresponding NYT distribution (*e.g.* on the date of the current event) are considered topical.

Applying equation (1) for each topic on the NYT document corpus results in topicality vectors with a length of 120 for each date  $d$ ,  $\underline{\tau}_d$ , describing how frequently each topic was being discussed on a given month or day in the 5 years that our data span. This vector can be formalized as,

$$\underline{\tau}_d = \sum_{t \in T} \tau_t(d) \cdot \underline{i}_t, \quad (2)$$

where  $T$  is the set of topics, and  $\underline{i}_t$  is the unit vector on  $|T|$  coordinates with the coordinate corresponding to topic  $t$  being 1.

Using the topicality vector  $\underline{\tau}_d$  and the topic distribution  $\underline{\tau}_m$  assigned to a meme  $m$  in the model inference step, we define the topicality of a Reddit meme  $m$  as,

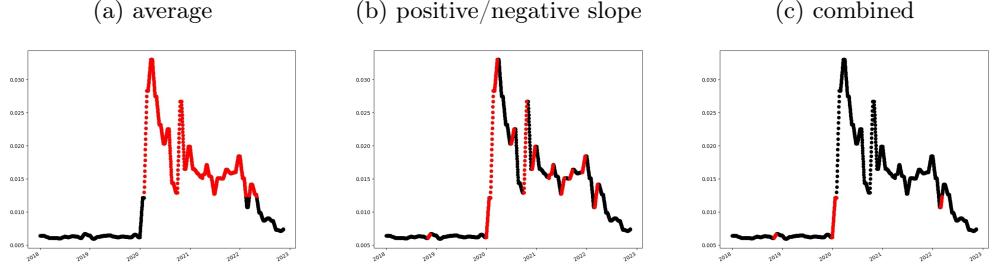
$$\text{top}_m = \underline{\tau}_{d_m} \cdot \underline{\tau}_m \quad (3)$$

where  $\underline{\tau}_{d_m}$  is the topicality distribution of the reference data (NYT for extrinsic topics, Reddit for intrinsic topics) on the given date  $d_m$  that the meme  $m$  was posted, and  $\underline{\tau}_m$  is the topic distribution assigned to the particular meme instance.

Note that again, this can be calculated for either a daily or monthly granularity, meaning  $d_m$  is either the day on which the meme  $m$  was posted or the month in which the meme was posted. Thus, equation (3) was used to define the "daily topicality" and "monthly topicality" features listed in Table 1. We suspect the monthly topicality variable to have reduced noise. For example, the monthly average topicality for the holidays topic, seen in Figure 3, is higher for all memes posted in December due to the effect of Christmas, but the daily average topicality variable has higher values only on the few days around Christmas and the other holidays.

In addition to the daily and monthly topicality features, reflecting how frequently a given topic was being discussed in the NYT or Reddit, we designed a feature describing the slope of the daily topicality distributions (*e.g.* Figure 6). This was calculated as the average slope of the reference data distribution for the  $\pm 2$  days around when the meme was posted. Given that topics about current events tend to spike rapidly, we used a small time frame. This is the "slope topicality" feature in Table 1. We suspect that a topic will be more popular on Reddit when the topicality distribution has a positive slope, indicating increasing topicality trends.

Two binary topicality features were engineered based on the slope and daily topicality features. These binary variables were used for statistical tests showing that topical memes are more popular than non-topical memes. The first binary topicality variable, visualized in Figure 7a, indicates whether the meme's topic is being posted



**Fig. 7:** Reddit extrinsic COVID-19 topicality distribution colored according to three binary topicality features (red=topical, black=non-topical).

| Model | Type          | Topical Mean Score | Non-topical Mean Score | p-value |
|-------|---------------|--------------------|------------------------|---------|
| LDA   | average       | 577                | 590                    | 0.059   |
| LDA   | pos/neg slope | 632                | 541                    | <0.001  |
| LDA   | combined      | 644                | 571                    | <0.001  |
| BERT  | average       | 564                | 594                    | 0.901   |
| BERT  | pos/neg slope | 694                | 672                    | <0.01   |
| BERT  | combined      | 609                | 575                    | <0.001  |

**Table 4:** Pooled t-tests results comparing memes' scores, groups determined by three binary topicality features.

about more (labeled 1) or less (labeled 0) than average in NYT on the date the meme was posted. The second binary variable, visualized in Figure 7b, was calculated based on the slope topicality. This variable takes the values 1 when the slope of the reference topicality distribution is positive, and 0 when the slope of the distribution is negative.

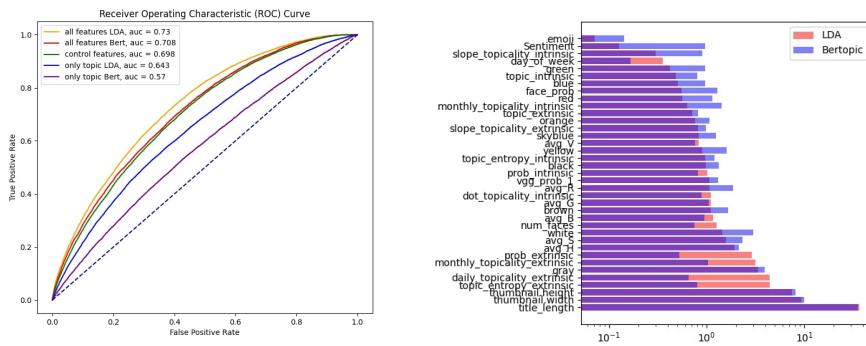
The cross-tabulation of these two binary topicality variables makes a categorical topicality feature with 4 categories. This categorical feature indicates both whether the slope of the reference distribution was positive or negative on the date when the meme was posted and whether the topic was being published about more or less than average on the date when the meme was posted. Figure 7c shows category 1 when the slope of the reference distribution is positive but the topic is not yet being posted about more than average.

Table 4 reports t-tests comparing the number of upvotes received by topical and non-topical memes by these binary variables. On average, memes about topics that show increasing, positive slope, topicality trends receive a higher score than memes about topics with decreasing topicality trends. There is an additional innovators advantage for memes posted right when a topic is beginning to gain popularity in the NYT (figure 7c). Differences between the other groups, for which T-test results were not reported, were not significant. Taken together, these results indicate that memes posted about a topic which has recently hit the news receive more upvotes on average. Topicality has a positive effect on Reddit popularity, however this effect doesn't

endure long after the date of the event, as indicated by the non-significant results of the figure 7a variable.

## 4 Identifying Viral Memes

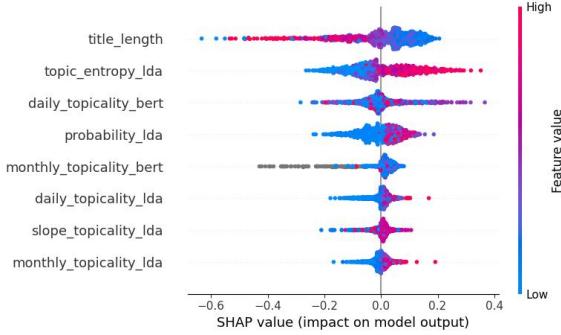
In order to gain a better understanding of how topicality influences the virality of a meme, we trained a CatBoost classifier. CatBoost, an ensemble learning algorithm for gradient boosting on decision trees, has also been shown to outperform other algorithms on a variety of problems (Ibrahim et al., 2020). The classifier is able to parse categorical features, meaning the 120-category topic features did not have to be one-hot-encoded before training. We used the image, text, and topicality features summarized in the shaded rows of Table 1 for the binary classification task of identifying viral memes in our Reddit data. Figure 8 shows the algorithm’s performance in identifying viral memes.



**Fig. 8:** ROC and features importances results identifying Reddit meme virality with CatBoost model

The binary target variable was highly unbalanced because it was defined as a percentile: only the top 5 percent of Reddit memes were considered viral. We used the CatBoost weighting parameter to improve performance on the imbalanced class prediction. The CatBoost classifier was trained using data from 763,388 memes and tested using 84,821 memes. The model was fine-tuned with 5-fold cross-validation, and sklearn’s GridSearchCV for parameter selection. Notably, the title length feature is the most important for classifying viral memes. Memes posted with short titles perform better than memes accompanied by a lot of text. This mirrors earlier work showing brevity is the most important feature for viral tweets (Tsur and Rappoport, 2015). The reason could be that messages that are easier to convey are understood better. Given the importance of the text length variable, we used CatBoost’s per-float-feature-quantization parameter to increase the number of decision boundaries allowed for this feature, taking advantage of its predictive power.

The classifier was trained separately for the LDA-based and BERTopic-based topicality feature sets. In order to estimate the importance of the topic-based features for predicting meme virality, we trained the CatBoost model with different subsets of features. The first models were trained using only topicality features, highlighted red in Table 1. Alone, the topicality features (AUCs: LDA=0.64, BERTTopic=0.57) are already able to predict meme virality better than a random guess (AUC=0.5).



**Fig. 9:** SHAP values for extrinsic topicality features indicating that topicality positively affects the CatBoost probability memes are predicted viral.

It could be the case that the explanatory power of the topicality features is already contained by information in the control features set. To test this, we trained the model with only control features, described in Sections 2.3 and 2.4. Then, we assessed the incremental predictive power of topicality features over control features. Models trained with the entire feature set, all shaded rows in Table 1, were better able to identify viral memes than the model trained with only the control features set. We tested the statistical significance of these AUC differences using a procedure from previous research (Hanley and McNeil, 1982), finding the engineered topicality features have significant predictive power over the control features for both BERTopic ( $AUC = 0.71, p < 0.01$ ) and LDA ( $AUC = 0.73, p < 0.001$ ). Results were more significant for the LDA topic model, and topicality features ranked higher in the LDA features importances than for BERTopic, as seen in Figure 8.

Figure 9 shows the SHapley Additive exPlanations (SHAP) values for given extrinsic topicality features (Lundberg and Lee, 2017). The text length feature, which is negatively correlated with meme virality, is provided for reference – memes with fewer words have a higher probability of being viral. The topicality features on the other hand correspond positively with the CatBoost’s estimation of meme virality. Memes with higher daily and monthly topicality values are generally labeled 1, viral, with a higher probability by the CatBoost model than memes with lower topicality values. Interestingly, higher values of topic entropy and topic probability also have a positive effect on the CatBoost outcome. This suggests that memes with a strong dominant topic, and high probability, but otherwise mixed topic distributions, and high entropy

have higher probabilities of going viral. The topicality features positively impact the CatBoost's estimation of meme virality.

## 5 Conclusion

While NYT and Reddit are qualitatively very different sources of information, they do show significant alignment. NYT-based topics are good descriptors and predictors of popular Reddit memes. Here we showed:

- Many topics on Reddit and the NYT show similar topicality distributions.
- Memes that are about topics with increasing topicality trends in the NYT receive more upvotes on average.
- “Innovator” memes posted when a topic is just beginning to be prevalent in the news receive the most upvotes.

These findings can be interpreted from the perspective of agenda-setting theory, which states that media sources shape the interests of society by selecting which topics to publish about. Our finding that news topicality has a positive effect on social media virality also indicates that it may be harder to get attention for and spread non-topicality information. This study does not unveil a causal link between topics in NYT and on Reddit, however, the great overlap between the sources indicates that topically they exert similar influence on their audiences. Future research could delve into differences between the opinions purported by news and social media about the same topics.

Articles from many disciplines explore the diffusion of internet memes through online sites. Our results showing that topicality plays a role in meme popularity, suggest that topics and topicality should be taken into account to model the more complex diffusion process of online media.

Most topic analyses have analyzed the relationship between user topical interests and content. Our analyses of topics in general and the use of a reference document to assess the topicality of the topics is a unique contribution. Defining topicality based on only the NYT can also be seen as a limitation of this work. It would be an interesting challenge to develop a more complete description of what subjects are topical on the internet by combining multiple reference documents. While NYT does report on global events and on issues other than news, using NYT as a source document skewed our set of topics to be more USA- and politics-centered. The USA-skewed results were appropriate for our purposes to match the large proportion of Reddit users based out of the USA, 48% of all users ([Tiago Biachi, 2023](#)). However, this means our results are less relevant to other parts of the globe.

Finally, the results presented here contribute to the growing body of research modeling multi-modal data. With the rise of internet communication, multi-modality is a paradigm shift in the forms discourse takes in general. Predicting the popularity of multimedia content is more difficult than text-only data. For example, impressive results have been achieved predicting viral tweets ([Weng et al., 2012; Tsur and Rapoport, 2015](#)), but we have not yet seen such success predicting the popularity of

image-with-text memes. One article achieved an AUC of 0.86 predicting image-with-text memes, but they worked with a small data set labeled by hand by humans, and kept only the highest and lowest scoring memes, thereby eliminating all moderately-popular memes and reducing the difficulty of the problem (Ling et al., 2021). Recent improvements in generative AI will soon make modeling multi-modal content more tractable. For now, predicting multi-modal content remains a challenge, and our AUC of 0.73 stands out in the literature.

With the increasing availability of internet data, social media has become a popular source of data for academics. Reddit is an especially popular data source, and the ethical guidelines around using this data are not yet firmly established. In light of this, we heed recommendations from studies surveying social media users about research data use. Namely, we collected data from a large, public community, assured no personal identifying information was associated with the data, and we plan to share our results back to Reddit upon completion (Proferes et al., 2021).

## References

- Halevy, A., Canton-Ferrer, C., Ma, H., Ozertem, U., Pantel, P., Saeidi, M., Silvestri, F., Stoyanov, V.: Preserving integrity in online social networks. *Communications of the ACM* **65**(2), 92–98 (2022)
- Nelson, D.L., Reed, V.S., Walling, J.R.: Pictorial superiority effect. *Journal of experimental psychology: Human learning and memory* **2**(5), 523 (1976)
- Spitzberg, B.H.: Toward a model of meme diffusion (m3d). *Communication Theory* **24**(3), 311–339 (2014)
- Brodie, R.: Virus of the Mind: The New Science of the Meme. Hay House, Inc, ??? (2009)
- Du, Y., Masood, M.A., Joseph, K.: Understanding visual memes: An empirical analysis of text superimposed on memes shared on Twitter. In: Proceedings of the International AAAI Conference on Web and Social Media, vol. 14, pp. 153–164 (2020)
- Dynel, M., Messerli, T.C.: On a cross-cultural memescape: Switzerland through nation memes from within and from the outside. *Contrastive Pragmatics* **1**(2), 210–241 (2020)
- McClure, B.: Discovering the discourse of internet political memes. (2016)
- Zannettou, S., Caulfield, T., Blackburn, J., De Cristofaro, E., Sirivianos, M., Stringhini, G., Suarez-Tangil, G.: On the origins of memes by means of fringe web communities. In: Proceedings of the Internet Measurement Conference 2018, pp. 188–202 (2018)
- Jacob Liedke and Luxuan Wang: Social Media and News Fact Sheet.

Accessed Dec 1, 2023. <https://www.pewresearch.org/journalism/fact-sheet/social-media-and-news-fact-sheet/> (2023)

Leach, C.W., Allen, A.M.: The social psychology of the black lives matter meme and movement. *Current Directions in Psychological Science* **26**(6), 543–547 (2017)

Yarchi, M., Boxman-Shabtai, L.: The image war moves to TikTok evidence from the may 2021 round of the Israeli-Palestinian conflict. *Digital Journalism*, 1–21 (2023)

Simmons, M., Adamic, L., Adar, E.: Memes online: Extracted, subtracted, injected, and recollected. In: *Proceedings of the International AAAI Conference on Web and Social Media*, vol. 5, pp. 353–360 (2011)

Cialdini, R.B., Goldstein, N.J.: Social influence: Compliance and conformity. *Annu. Rev. Psychol.* **55**, 591–621 (2004)

Grabowicz, P., Ganguly, N., Gummadi, K.: Distinguishing between topical and non-topical information diffusion mechanisms in social media. In: *Proceedings of the International AAAI Conference on Web and Social Media*, vol. 10, pp. 151–160 (2016)

Weng, L., Flammini, A., Vespignani, A., Menczer, F.: Competition among memes in a world with limited attention. *Scientific Reports* **2**, 335 (2012) <https://doi.org/10.1038/srep00335>

Wang, L., Wood, B.C.: An epidemiological approach to model the viral propagation of memes. *Applied Mathematical Modelling* **35**(11), 5442–5447 (2011)

Poux-Médard, G., Velcin, J., Loudcher, S.: Properties of Reddit news topical interactions, 16–28 (2022). Springer

Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. *Journal of Machine Learning Research* **3**(Jan), 993–1022 (2003)

Grootendorst, M.: BERTopic: Neural topic modeling with a class-based tf-idf procedure. arXiv preprint arXiv:2203.05794 (2022)

Weng, L., Menczer, F.: Topicality and impact in social media: diverse messages, focused messengers. *PloS one* **10**(2), 0118410 (2015)

Ferrara, E., Interdonato, R., Tagarelli, A.: Online popularity and topical interests through the lens of instagram. In: *Proceedings of the 25th ACM Conference on Hypertext and Social Media*, pp. 24–34 (2014)

Barnes, K., Riesenmy, T., Trinh, M.D., Lleshi, E., Balogh, N., Molontay, R.: Dank or not? Analyzing and predicting the popularity of memes on Reddit. *Applied Network Science* **6**(1), 1–24 (2021)

Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A.V., Gulin, A.: Catboost: unbiased boosting with categorical features. *Advances in neural information processing systems* **31** (2018)

HSDS: Topical alignments between the NYT and Reddit. Accessed Feb 1, 2024. <https://github.com/hsdslab/topicality-online> (2024)

Sanderson, B., Rigby, M.: We've Reddit, have you?: What librarians can learn from a site full of memes. *College & Research Libraries News* **74**(10), 518–521 (2013)

Baumgartner, J., Zannettou, S., Keegan, B., Squire, M., Blackburn, J.: The pushshift Reddit dataset. In: *Proceedings of the International AAAI Conference on Web and Social Media*, vol. 14, pp. 830–839 (2020)

Boe, B.: PRAW: The Python Reddit API Wrapper. <https://github.com/praw-dev/praw>. Accessed: 2022-12-15 (2016)

Podolak, M.: PMAW: Pushshift Multithread API Wrapper. <https://github.com/mattpodolak/pmaw>. Accessed: 2022-12-15 (2021)

Figueiredo, F., Almeida, J.M., Gonçalves, M.A., Benevenuto, F.: On the dynamics of social media popularity: A youtube case study. *ACM Transactions on Internet Technology (TOIT)* **14**(4), 1–23 (2014)

Amir Salihefendic: How Reddit Ranking Algorithms Work. Accessed June 1, 2023. <https://www.yale.edu/about-yale/yale-facts> (2015)

Carman, M., Koerber, M., Li, J., Choo, K.-K.R., Ashman, H.: Manipulating visibility of political and apolitical threads on Reddit via score boosting. In: *2018 17th IEEE International Conference On Trust, Security And Privacy In Computing And Communications/12th IEEE International Conference On Big Data Science And Engineering (TrustCom/BigDataSE)*, pp. 184–190 (2018). IEEE

Ling, C., AbuHilal, I., Blackburn, J., De Cristofaro, E., Zannettou, S., Stringhini, G.: Dissecting the meme magic: Understanding indicators of virality in image memes. *Proceedings of the ACM on human-computer interaction* **5**(CSCW1), 1–24 (2021)

Tiago Biachi: Reddit - Statistics & Facts. Accessed August 1, 2023. <https://www.statista.com/topics/5672/reddit/#topicOverview> (2023)

Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In: *International Conference on Machine Learning*, pp. 12888–12900 (2022). PMLR

Bird, S., Klein, E., Loper, E.: Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit. " O'Reilly Media, Inc.", ??? (2009)

Tsur, O., Rappoport, A.: Don't let me be# misunderstood: Linguistically motivated

algorithm for predicting the popularity of textual memes. In: Proceedings of the International AAAI Conference on Web and Social Media, vol. 9, pp. 426–435 (2015)

Kim, T., Wurster, K.: Emoji 2.8.0 for Python. Accessed May 1, 2022. <https://pypi.org/project/emoji/> (2023)

Shreyas, P.: Sentiment analysis for text with Deep Learning. Medium (2019). <https://towardsdatascience.com/sentiment-analysis-for-text-with-deep-learning-2f0a0c6472b5> Accessed 2020-10-01

Zhang, K., Zhang, Z., Li, Z., Qiao, Y.: Joint face detection and alignment using multitask cascaded convolutional networks. IEEE signal processing letters **23**(10), 1499–1503 (2016)

Stone, R.: Image Segmentation Using Color Spaces in OpenCv+Python (2018). <https://realpython.com/python-opencv-color-spaces/> Accessed 2020-10-01

Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., *et al.*: Imagenet large scale visual recognition challenge. International Journal of Computer Vision **115**, 211–252 (2015)

Khosla, A., Das Sarma, A., Hamid, R.: What makes an image popular? In: Proceedings of the 23rd International Conference on World Wide Web, pp. 867–876 (2014)

Datta, R., Joshi, D., Li, J., Wang, J.Z.: Studying aesthetics in photographic images using a computational approach. In: Computer Vision–ECCV 2006: 9th European Conference on Computer Vision, Graz, Austria, May 7–13, 2006, Proceedings, Part III 9, pp. 288–301 (2006). Springer

Coscia, M.: Average is boring: How similarity kills a meme’s success. Scientific reports **4**(1), 6477 (2014)

Egger, R., Yu, J.: A topic modeling comparison between lda, nmf, top2vec, and bertopic to demystify Twitter posts. Frontiers in sociology **7**, 886498 (2022)

Ibrahim, A.A., Ridwan, R.L., Muhammed, M.M., Abdulaziz, R.O., Saheed, G.A.: Comparison of the CatBoost classifier with other machine learning methods. International Journal of Advanced Computer Science and Applications **11**(11) (2020)

Hanley, J.A., McNeil, B.J.: The meaning and use of the area under a receiver operating characteristic (roc) curve. Radiology **143**(1), 29–36 (1982)

Lundberg, S.M., Lee, S.-I.: A unified approach to interpreting model predictions. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) Advances in Neural Information Processing Systems 30, pp.

4765–4774. Curran Associates, Inc., ??? (2017). <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>

Tsur, O., Rappoport, A.: Don't let me be# misunderstood: Linguistically motivated algorithm for predicting the popularity of textual memes. In: Ninth International AAAI Conference on Web and Social Media (2015)

Proferes, N., Jones, N., Zimmer, M.: Studying Reddit: A systematic overview of disciplines, approaches, methods, and ethics. *Social Media and Society* **7** (2021) <https://doi.org/10.1177/20563051211019004>