





An analysis of emotions and the prominence of positivity in #BlackLivesMatter tweets

Anjalie Field^{a,1,2}, Chan Young Park^{a,1}, Antonio Theophilo^{a,b,c,1}, Jamelle Watson-Daniels^d, and Yulia Tsvetkov^e

Edited by Molly Crockett, Yale University, New Haven, CT; received April 12, 2022; accepted July 24, 2022 by Editorial Board Member Susan T. Fiske

Emotions are a central driving force of activism; they motivate participation in movements and encourage sustained involvement. We use natural language processing techniques to analyze emotions expressed or solicited in tweets about 2020 Black Lives Matter protests. Traditional off-the-shelf emotion analysis tools often fail to generalize to new datasets and are unable to adapt to how social movements can raise new ideas and perspectives in short time spans. Instead, we use a few-shot domain adaptation approach for measuring emotions perceived in this specific domain: tweets about protests in May 2020 following the death of George Floyd. While our analysis identifies high levels of expressed anger and disgust across overall posts, it additionally reveals the prominence of positive emotions (encompassing, e.g., pride, hope, and optimism), which are more prevalent in tweets with explicit pro-BlackLivesMatter hashtags and correlated with on the ground protests. The prevalence of positivity contradicts stereotypical portrayals of protesters as primarily perpetuating anger and outrage. Our work offers data, analyses, and methods to support investigations of online activism and the role of emotions in social movements.

emotion analysis | BlackLivesMatter | Twitter | natural language processing

The term #BlackLivesMatter originated in posts made by activists Alicia Garza and Patrisse Cullors in 2013 following George Zimmerman's acquittal over the killing of Trayvon Martin, an unarmed Black teenager (1).* The term has since become popularized as referring to movements against police brutality and the extrajudicial killing of Black people. These movements have continually grown and evolved, garnering widespread attention following the deaths of Michael Brown in Ferguson and Eric Garner in New York in 2014 (2, 3) and more recently, George Floyd in Minneapolis (2020). The death of George Floyd, in addition to the deaths of Ahmaud Arbery and Breonna Taylor, led to widespread protests against police violence and racism.

Social media has been an integral part of these movements. In addition to #BlackLives-Matter, millions of tweets were posted with hashtags like #Ferguson, #JusticeForGeorge-Floyd, and #ICantBreathe. While forms of "digital protest" and "hashtag activism" can occur organically, they are often a tool used by community activists who may plan hashtag campaigns, promote in-person activism, and intentionally bypass traditional media (1, 2, 4-6). Thus, social media not only provides an avenue for analyzing modern social movements, but understanding social media messaging is also essential for providing insight into these events.

In this work, we analyze a dataset of tweets related to Black Lives Matter protests from 24 May to 28 June 2020 using a domain adaptation model for measuring emotions perceived in tweets about specific events. In the past few decades, social psychologists have recognized the important role emotions play in activism; "moral shocks" can facilitate people joining a movement, while hope and pride are necessary to sustain involvement (7-10). Understanding the dynamics between emotions (such as what balance between anger and optimism produces a "hopeful anticipation of impact" that motivates continued action) can provide both insight into past movements and guidance for future efforts

Furthermore, projected emotions have been used to falsely characterize Black people, leading to tangible harms. For example, the "angry Black woman" stereotype can result in negative physical, social, and economic impacts, such as facilitating workplace discrimination (12, 13). In the context of social movements, negative stereotypes of Black protesters

Significance

In a corpus of 34 million tweets about Black Lives Matter from June 2020, although negative emotions, like anger and disgust, occur commonly, positive emotions, like hope and optimism, are more prevalent in tweets with pro-BlackLivesMatter hashtags and significantly correlated with the presence of on the ground protests. These results contrast "angry Black" stereotypes and portrayals of protesters as perpetuating anger and outrage. This work demonstrates how natural language processing techniques can shed insight into social movements and counter harmful stereotypes, and also, it offers methodology for extracting social meaning from text data.

Author affiliations: ^aLanguage Technologies Institute, Carnegie Mellon University, Pittsburgh, PA 15289; ^bInstitute of Computing, University of Campinas, Campinas 13083-852, Brazil; ^cElectronics and Mechatronics Infrastructure Division, Center for Information Technology Renato Archer, Campinas 13069-901, Brazil; dHarvard John A. Paulson School of Engineering and Applied Sciences, Harvard University, Cambridge, MA 02138; and *Paul G. Allen School of Computer Science & Engineering, University of Washington, Seattle, WA 98195

Author contributions: A.F., C.P., A.T., and Y.T. designed research; A.F., C.P., and A.T. performed research; A.F., C.P., A.T., and J.W.-D. analyzed data; and A.F., C.P., A.T., J.W.-D., and Y.T. wrote the paper.

The authors declare no competing interest.

This article is a PNAS Direct Submission. M.C. is a guest editor invited by the Editorial Board.

Copyright © 2022 the Author(s). Published by PNAS. This open access article is distributed under Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 (CC BY-NC-ND).

¹A.F., C.P., and A.T. contributed equally to this work.

²To whom correspondence may be addressed. Email: anjalief@cs.cmu.edu.

This article contains supporting information online at https://www.pnas.org/lookup/suppl/doi:10.1073/pnas. 2205767119/-/DCSupplemental.

Published August 23, 2022.

^{*}We generally use "Black Lives Matter" to refer to the broad movement against police brutality rather than the official organization (https://blacklivesmatter.com/about/).

[†]Information is available at https://acleddata.com/2020/09/03/demonstrations-political-violence-in-america-new-data-forsummer-2020/.

as violent angry "thugs" have long been used to derail civil rights activism (14).‡ Analyzing emotions in tweets about protests can provide evidence refuting these types of negative portrayals.

However, measuring emotions is nontrivial, and computational models that overestimate expressions of emotions, like "anger," can reinforce negative stereotypes. Previous examinations of emotions and affect expressed in tweets about the Black Lives Matter movement have relied on lexicon (Linguistic Inquiry and Word Count, LIWC) scores (3), and analyses of other protest events have similarly relied on lexicon-based approaches (15, 16). While recent research has led to the development of more powerful deep learning-based models and annotated datasets, these models nevertheless are prone to overfitting to shallow lexical cues and often perform poorly in new domains (17-19). Thus, in this work, we leverage recent natural language processing (NLP) techniques, including in-domain pretraining and few-shot learning, to improve emotion analysis model performance across domains in an easily adaptable framework. We evaluate our model using two annotated datasets of emotion classification in two different domains, Reddit and Twitter, and for six different emotion categories: anger, disgust, positivity, surprise, fear, and sadness (20).

We ultimately use our model to examine emotion trends in a dataset we collected containing ~34 million tweets related to the Black Lives Matter movement. In examining estimated perceived emotions over time, in tweets with specific hashtags, and in comparison with on the ground protests, our results consistently identify the prominence of positivity (e.g., pride, optimism, excitement), which supports social theories about the importance of emotions like hope and pride and offers evidence countering "angry Black" stereotypes.

Results

Data. Our primary corpus consists of tweets about Black Lives Matter. We gathered English tweets posted between 24 May and 30 June 2020 using the Twitter search API. SI Appendix, section 1 contains the full list of terms used for data collection, which includes terms likely to be used by both supporters and critics of the Black Lives Matter protests. Our final dataset, which we refer to as #BLM2020, consists of 250 million tweets (34.7 million excluding retweets) by 18.9 million users. Fig. 1 presents the volume of tweets and users through the time span. There is high Twitter engagement in the first 10 d followed by a slow decrease

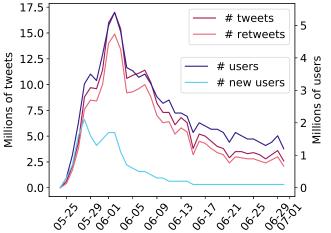


Fig. 1. Distribution of tweets, retweets, users, and new users in #BLM2020.

in the subsequent 4 wk. We ceased data collection at the end of June given the substantial decline in tweet volume by the end of the month.

In general, our work involves analysis of a sensitive social issue, and while all data was publicly available at the time of collection, Twitter users did not explicitly consent to this analysis. In order to facilitate reproducibility while preserving anonymity and privacy as much as possible, we do not make the raw data freely available, but we will make tweet identifications available for academic research purposes only upon request in accordance with Twitter terms of service.

Detecting Emotions Expressed in Tweets. In order to analyze emotions expressed in #BLM2020, we develop and evaluate models for identifying six emotion categories: anger, disgust, fear, positivity, surprise, and sadness, which are the primary core emotions according to Ekman's taxonomy (20). This approach assumes that emotions identified by annotators in tweets can be represented in discrete categories, and taking a psychological constructionist perspective of measuring emotions [e.g., focusing on the dimensions of valence and arousal (21-23)] may have different results (24). We follow prior work in considering these six Ekman emotions to be supersets of finer-grained emotions (18):

- anger: anger, annoyance, disapproval, and rage;
- disgust: disgust, loathing, and boredom;
- fear: fear, nervousness, vigilance, and apprehension;
- positivity[§]: amusement, approval, excitement, gratitude, love, optimism, relief, pride, admiration, desire, caring, acceptance, anticipation, serenity, trust, and ecstasy;
- surprise: realization, confusion, curiosity, amazement, and distraction; and
- sadness: disappointment, embarrassment, grief, remorse, and pensiveness.

Throughout this work, we treat emotions as nonexclusive (e.g., a tweet may contain both anger and sadness). We also aim to capture emotions that Twitter users choose to express or solicit on the platform, which may not reflect their actual emotional state, and we discuss this distinction in our analysis.

Traditionally, social scientists have used lexicon-based approaches to measure emotions in tweets, determining whether or not a tweet expresses anger based on whether or not it contains any words from a list of "angry" words. While lexiconbased approaches remain popular because of their ease of use, they can be brittle and fail to adapt to new domains. Word connotations change in different contexts (21), particularly in protest movements, which often aim to subvert the status quo. For example, the National Research Council Canada (NRC) Word-Emotion Association Lexicon (EmoLex), which contains words associated with eight emotions, associates "police" with fear, positive, and trust, which are contradictory to the connotations of "police" in protests against police brutality (15, 25). More recently, machine learning-based NLP models have outperformed traditional lexicon approaches at identifying affect in text (18, 19, 26, 27). Neural models are trained on annotated datasets and used to infer affect in unseen text. However, a model trained on a precollected dataset may still perform poorly on data from a different domain where connotations differ. Collecting new

[‡]Information is available at https://www.nbcnews.com/news/us-news/not-accident-falsethug-narratives-have-long-been-used-discredit-n1240509.

 $[\]S$ Ekman's taxonomy refers to this dimension as "joy" but defines it as encompassing all positive emotions (20). We use the term "positivity" instead of joy throughout this work to reflect the breath of emotions encompassed by this dimension and avoid suggesting that happiness is a dominant emotion in #BLM2020.

annotated datasets for every domain of interest is prohibitively time consuming and expensive, especially for tasks that require in-domain knowledge or involve subjective judgements.

Instead, we take a domain adaptation approach; given a set of source data annotated for perceived emotional content (for example, tweets with binary present/not present labels for emotions, like anger, surprise, and fear), our goal is to infer emotion labels for a set of target data from a different domain using explicit methods to adapt the model to this new domain. Different domains could include text about a different event or from a different social media platform. Domain adaptation allows us to reuse annotated datasets rather than collecting new annotated data for every domain of interest. We train and evaluate a base classifier for inferring emotions with two variants of domain adaption:

Base classifier (BASE). In the simplest setting, we train a prediction model over the annotated source data and infer labels on the target data without any explicit domain adaptation. We specifically use a pretrained language model (Bidirectional Encoder Representations from Transformers, BERT) fine-tuned over the source data (details are in Materials and Methods).

Task-adaptive pretraining (+TGT). NLP has recently seen large performance improvements through masked language model pretraining; models are pretrained by optimizing them to predict words that have been obfuscated from input sentences (28). The same model can then be fine-tuned for a specific task. Following prior work, we use masked language model pretraining over unannotated sentences from the target data to encourage domain adaption and then, fine-tune the model to infer emotions using the annotated source data, as in BASE (19, 29, 30).

Few-shot learning (+FSL). While collecting a large annotated dataset for every new domain can be infeasible, collecting annotations over a small number of in-domain labeled data is often practical. In this model, we fine-tune the classifier over small sets of annotated target data (300 instances), after training over the larger source dataset.

Our primary training data are drawn from two sources, GoEmotions and HurricaneEmo (18, 19). GoEmotions consists of 58,000 English Reddit comments manually labeled for emotion categories or neutral (18). We randomly divide these data into train (80%), validation (10%), and test (10%) splits. HurricaneEmo consists of 15,000 English tweets about hurricanes annotated for 24 emotions according to Plutchik's scheme (19, 31), which we map to the Ekman scheme (described in SI Appendix, section 2). The original dataset provided a different train–test split for each emotion; thus, we created our own instance-level data split of train (70%), validation (10%), and test (20%). To facilitate few-shot learning and evaluation, we additionally collect emotion annotations over 700 randomly sampled tweets from #BLM2020 using the six Ekman emotions, and we use 300 as training data, 100 as development data, and 300 as test data. We provide further details in Materials and Methods and SI Appendix, section 2.

Fig. 2 shows evaluation results over the annotated #BLM2020 test data, where we use both GoEmotions and HurricaneEmo as training data and use 300 of the annotated #BLM2020 for few-shot learning. We provide additional validation metrics over larger test datasets in SI Appendix, section 3.

In addition to classification models, we provide LIWC as a baseline since it is a popular dictionary-based analysis method and has previously been used in analyzing tweets about Black Lives Matter (3, 32). We map the LIWC dimensions of "anger," "positive emotion," and "sadness" to anger, positivity, and sadness, respectively, since they are the only emotions that directly map to LIWC dimensions, and we map the floating-point scores produced by LIWC to binary labels using the best-performing threshold over the validation dataset.

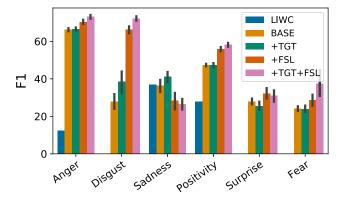


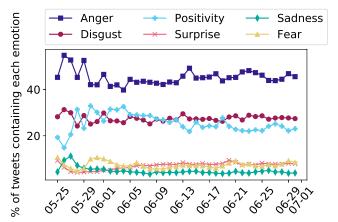
Fig. 2. F1 scores of emotion classifiers evaluated over #BLM2020. Error bars indicate the 95% CIs.

The machine learning classifiers generally outperform LIWC, few-shot learning brings a large performance improvement, and +TGT+FSL achieves the best overall performance. As +TGT+FSL outperforms other models, we use it to obtain perceived emotion labels for all tweets in #BLM2020, which we analyze in the following section. We generally focus our analysis on the emotions that our model identified with highest F1 and that had the highest interannotator agreement in our annotated data (reported in SI Appendix, section 2): anger, disgust, and positivity. Performance of the +TGT+FSL model is poor for sadness in Fig. 2; however, sadness is very sparse in the #BLM2020 test set, and +TGT+FSL outperforms other models when evaluated over a larger test set (SI Appendix, section 3). In contrast, surprise has poorer model performance and lower interannotator agreement over all test sets. Thus, we avoid extended discussion of this emotion, although we do display metrics for all emotions.

Analysis of Emotions in #BLM2020. We first use our inferred emotion labels to examine how emotions expressed in #BLM2020 change over time and with different hashtags. Because retweets are not written independently nor displayed as separate posts to Twitter users and because we do not expect model performance to be reliable over very short tweets, we exclude retweets and tweets with fewer than five tokens, leaving 34.1 million tweets for analysis. Given critique of Ekman's taxonomy (33–35) and the potential for classifier error, we provide analysis metrics using sentiment models and probabilistic aggregation in SI Appendix, sections 6 and 7. We also provide anonymized examples from our data and additional visualizations in SI Appendix, section 5 and a comparison with tweets from 2012 to 2015 in SI Appendix, section 8.

Changes in emotions over time. In Fig. 3, we plot the percentage of tweets that contain each emotion over time estimated using our model. Although positivity captures a broader range of emotions than anger, anger is the most prevalent emotion throughout, consistently occurring in >40% of tweets. Positivity and disgust are also prevalent, with positivity gradually decreasing over time, while anger and disgust gradually increase after an initial peak. A small peak in sadness occurs early on but is quickly eclipsed. A peak in fear occurs from Sunday 31 May to Monday 1 June, directly following the first weekend of protests (see Fig. 6).

Anger and positivity have a strong negative correlation over time (-0.79), while anger and disgust have a strong positive one (0.69). We also note that annotators who labeled emotions in #BLM2020 described anger and disgust as difficult to distinguish in this setting (Materials and Methods), which is consistent with the identification of "moral outrage" as involving anger and disgust (36).



Percentage of tweets that contain each emotion over time (24 May and 30 June). Emotion categories are drawn from Ekman's taxonomy (20) and inferred over a dataset of 34.1 million tweets using a neural classification model with domain adaptation components (+TGT+FSL).

Fig. 3 presents emotions over the entire dataset, which contains tweets both supportive of the Black Lives Matter movement and opposed to it. Thus, it provides no insight into how emotions are directed and does not distinguish between, for example, protesters' anger and anger at protesters. In Fig. 4, we display emotion levels only for tweets that contain a pro-BLM hashtag (defined in SI Appendix, section 1). Over a subset of ~ 600 tweets that annotators labeled for stance (SI Appendix, section 2), using these hashtags to recover tweets annotated as "pro-BLM" obtained a precision of 82.7% and a recall of 29.4%. In Fig. 4, the initial peak in sadness is even more apparent as is the high peak of anger, both of which predate the first weekend protests. Positivity rises shortly before the first weekend and continues through the second weekend before declining. A later peak in positivity occurs on 19 June 2020, which is Juneteeth, a holiday celebrating the emancipation of people who had been enslaved in the United States. On this day, #Juneteenth was the second-most common hashtag in the total dataset after #BlackLivesMatter.

Common hashtags for each emotion. In Table 1, we report hashtags that are most overrepresented in tweets that our model identifies as containing each emotion, calculated using log odds with a Dirichlet prior (37). These hashtags are highly indicative of the predicted emotions. Tweets labeled with positivity commonly contain #love and #pride; tweets labeled with sadness contain #sad and #RIP. Importantly, hashtags associated with the same emotions often reflect opposing viewpoints; tweets labeled with anger frequently contain both #MAGA (Donald Trump's campaign slogan) and #TrumpResignNow. In SI Appendix, section 5, we additionally provide associated hashtags when the data are divided into tweets with pro-BLM and anti-BLM hashtags, as well as word clouds of words associated with each inferred emotion.

Emotions by keywords. Fig. 5 shows the percentage of tweets our model identifies as containing each emotion, where tweets are divided as containing pro-BLM hashtags, anti-BLM hashtags, terms related to police, and terms related to protests as enumerated in SI Appendix, section 1. In all cases, positivity, anger, and disgust occur much more than fear, sadness, and surprise. Both positivity and sadness occur more often in tweets with pro-BLM hashtags than in any of the other subsets. Notably,

anger and disgust are lower in tweets with explicitly pro-BLM hashtags than in tweets with explicitly anti-BLM hashtags, while positivity is higher. As users often use hashtags to engage in public narratives and direct content to particular streams (38), these data offer counterevidence to the narrative of BLM protesters as angry "thugs". There is more positivity and less anger and disgust in tweets with pro-BLM hashtags (i.e., that are explicitly directed toward streams about the movement) than in tweets discussing these events more generally, including tweets with reactionary #AllLivesMatter hashtags. The highest percentage of anger occurs in tweets mentioning police, which encompass both anger over police brutality and calls for reform as well as reactionary propolice posts expressing anger at protesters. The highest percentage of fear occurs in tweets mentioning protests, which capture direct references to events that occurred during protests, like aggressive police responses.

Correlations between Emotions in Tweets and on the Ground Protests. Finally, we compare tweet volume and emotions with the volume of on the ground protests during the same time period. To estimate on the ground protests, we use data collected from two sources: The Armed Conflict Location & Event Data Project (ACLED)# (39) and the Crowd Counting Consortium (CCC) (40). The ACLED contains records of political violence, demonstrations, and strategic developments across the United States. Entries are hand coded by ACLED researchers and based on media reports by 2,400 sources. The CCC contains records of political crowds reported in the United States, including marches, protests, strikes, demonstrations, riots, and other actions, and is maintained by a dedicated project manager and research assistants.

Fig. 6 shows the number of protests across the United States per day as reported by the ACLED and the CCC. Data from both sources show similar patterns, although the CCC consistently reports slightly more protest events than the ACLED. The first peak in protests occurs from 30 May 2020 to 31 May 2020, the weekend directly following George Floyd's death. The highest peak in Twitter activity occurs after this weekend, which may suggest how early protests called attention to George Floyd's death. The peak volume of protests occurs after the highest peak

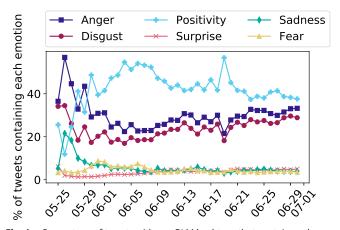


Fig. 4. Percentage of tweets with pro-BLM hashtags that contain each emotion over time (24 May and 30 June). Emotion categories are drawn from Ekman's taxonomy (20) and inferred using a neural classification model with domain adaptation components (+TGT+FSL). The dataset is restricted to 6.5 million tweets that contain pro-BLM hashtags (defined in SI Appendix, section 1).

 $[\]P$ Precision and recall for using anti-BLM hashtags to recover anti-BLM stances are 62.5% and 4.5%, respectively, and we caution that these results reflect the use of these hashtags, not necessarily stance.

Information is available at https://acleddata.com/special-projects/us-crisis-monitor/. Information is available at https://sites.google.com/view/crowdcountingconsortium/

Table 1. Most common hashtags for tweets labeled for each emotion computed using log odds with a Dirichlet prior (37)

Anger	Disgust	Positivity	Surprise	Sadness	Fear
BreonnaTaylor	Trump	BlackLivesMatter	BLM	RIPGeorgeFloyd	NYCScannerDuty
GeorgeFloyd	Racist	RaiseTheDegree	AllLivesMatter	JusticeFor GeorgeFloyd	NYCProtests
DefundThePolice	MAGA	Love	WhiteLivesMatter	GeorgeFloyd	PDX911
PoliceBrutality	TrumpResignNow	BlackOutTuesday	AskingForAFriend	ICantBreathe	COVID19
ACAB	AllLivesMatter	Juneteenth	AlmostBrokeMy HeartAtTheEnd (Thai)	BlackLivesMatter	BlackLives MatterNYC
Riots2020	BunkerBoy	PrideMonth	Confused	RIP	NYCProtest
GeorgeFloydWas Murdered	DefundThePolice	MatchAMillion	BlackOutTuesday	Sad	DCProtest
MinneapolisRiots	Trump2020	Music	Nkurunziza	JusticeForFloyd	DCProtests
JusticeForGeorge Floyd	RacistInChief	Art	달빛보다_찬란 한_준휘야-생 일축하해	RestInPower	Breaking
DerekChauvin	DemocratsAre Destroying America	Pride2020	HNGInternship	WeAreTired	GeorgeFloyd Protests
Trump	AntifaTerrorists	Juneteenth2020	Dollar (Arabic)	RIPHumanity	FoxNews
BreonnaTayor	Democrats	2MforBLM	AmUnbroken	Palestinian LivesMatter	Coronavirus
FakeNews	BLM	Pride	달빛보다_찬란 하_준휘	BlackLivesMatters	SeattleProtest
TrumpResignNow	TrumplsARacist	NYCScannerDuty	Kalu	ShootATweet	NYCScanner
DemocratsAre Destroying America	ACAB	Equality	365DNI	JusticeForGeorge	Protests
AntifaTerrorists	Antifa	Peace	BlueLivesMatter	JusticeForJeyaraj AndFenix	NYPD

Emotion categories are drawn from Ekman's taxonomy (20) and inferred using a neural classification model with domain adaptation components (+TGT+FSL). Hashtags are deduplicated after case normalization.

in Twitter activity on 6 June 2020, the second Saturday. While definite conclusions cannot be drawn from these few data points, this pattern suggests a possible symbiotic relationship between online and offline protests; the first peak of in-person protests encouraged increased engagement on Twitter, which in turn, resulted in even more protests the following weekend. After this weekend, the volume of protests steadily declines, with regular peaks on subsequent weekends.

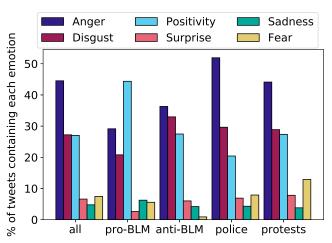


Fig. 5. Percentage of tweets that contain each emotion, where tweets are divided by keywords and hashtags. Emotion categories are drawn from Ekman's taxonomy (20) and inferred using a neural classification model with domain adaptation components (+TGT+FSL).

While protests broke out across the United States, they were more widespread and lasted longer in certain areas than in others, which allows us to compare emotions expressed on Twitter and on the ground protests by comparing tweets by users in different locations. We identified location for users in our dataset based on the user-populated location string in their profiles (details are in SI Appendix, section 4). This value was nonempty for 62.36%, of users in our dataset, and we were able to map 20.66% of users to a US state and 12.3% of users to US cities listed in the ACLED data.** Our results in this analysis are limited to users who specified locations in their Twitter profiles, and we cannot conclude how well they generalize to users who did not, although prior work has suggested that geolocated tweets provide accurate measures of protest events, even though geolocation data are typically sparse

In Tables 2 and 3, we show the Pearson correlations between the number of protests in each city or state and the percentage of tweets containing each emotion as measured by our model for tweets posted by users in those locations. Because we can expect larger and more populous states to have more protests, we normalize the number of reported protests in each state by the number of counties in the state (a US administrative/political/geographic subdivision of a state with some level of governmental authority),

^{**}Given the sparsity of city-level data, we only compute results for cities for which we were able to identify at least 500 users. We do not observe substantial differences in results if we change the cutoff to 100 or 1,000 users.

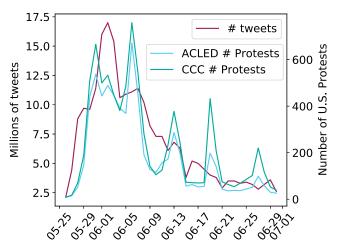


Fig. 6. Volume of US protests and collected tweets. Protest data are drawn from the ACLED (39) and the CCC. Twitter data are collected in this work.

obtaining county counts from US Census data reported by Safe-Graph.†† We believe that counties are a reasonable normalization term because they reflect factors that influence protests, which typically take place in a single geographic area and are often targeted toward local government. At a city level, where we do not expect as substantial geographic barriers and given the importance of size in social movements (44, 45), we weight protest events by ACLED and CCC size estimates in order to compute total protest volume (SI Appendix, section 4 has details and discussion).

At both the city and state levels, positivity is positively correlated with more protest events, and anger, disgust, and sadness are negatively correlated. Fear is additionally positively correlated at a city level. Importantly, our results demonstrate geographic correlations, not temporal ones. We cannot distinguish if expressions of positivity precede protests and are thus predictive of on the ground activism or if they are reactionary (posted during or after a protest). Given the potential misuses of technology for predicting protests [actors have sought to discourage collective action, destabilize movements, or promote polarization (46, 47)] as well as the nonlinear temporal trends in our data (Fig. 6) (protest volume declines over time with spikes on weekends), we do not compute temporal trends or make any attempt to predict protest or tweet volume.

Discussion

Political and social psychology research has identified anger as a politically motivating emotion using survey data, laboratory experiments, and theoretical analyses in protest movements and political involvement generally (48-50) as well as specifically for Black people (51-53). Our results do not directly contradict this research, in that we do find anger and disgust as the most commonly expressed emotions in our dataset, and also, we see initial peaks in these emotions (Fig. 3). However, we find that these emotions are negatively correlated with in-person protests (Tables 2 and 3), whereas positivity is positively correlated. This difference in result from prior work could result from differences between actual emotional state and what users choose to post on Twitter. Negative stigma around angry Black people could disincentive people from posting expressions of anger or disgust on Twitter (54). Additionally, as we focus on geographic rather than temporal relations, our model captures emotions expressed

before, during, and after protests, and feelings of camaraderie and pride resulting from protests could outweigh other emotions expressed on Twitter.

Relatedly, our results also show positive correlations between fear and protests at a city level, which seemingly contradict prior identification of fear, anxiety, and sadness as dispiriting emotions that deter political engagement (8, 9, 55). However, both sadness and fear are uncommon in our data, which is consistent with these theories, as posting on Twitter is itself an act of engagement and people feeling sadness or fear may choose not to post at all. An examination of the relatively small percentage of tweets that our model does identify as reflecting fear suggests that they often focus on events specifically related to protests, including community monitoring of police activity, like severe crowd control tactics during protests (references to "ScannerDuty" in Table 1 and a higher prevalence of fear in tweets referring to protests in Fig. 5). These results are consistent with the discussion of fear in the analysis in ref. 50 of Arab Uprisings, which notes that protesters express fear and suggests that identifying conditions under which people press on despite fear is more relevant than identifying conditions under which fear disappears. Unlike fear, sadness is negatively correlated with protest activity, which is consistent with prior identification of this emotion as dispiriting (48, 50, 55).

Overall, our results consistently identify the role of positive emotions in Black Lives Matter social media posts. In addition to the correlations with on the ground protests, tweets with pro-BLM hashtags contain more positivity than other tweets in our dataset, such as ones with anti-BLM hashtags. These results support social psychology theories suggesting that positive emotions are an important component of social movements (8, 9). While outrage and anger can encourage people to become involved, participants must also have optimism and hope for change, or they will not have the motivation to act (8, 11). Similarly, joy and camaraderie (e.g., feeling affective bonds as a member of a group) encourage sustained involvement (8, 9). Our findings additionally also offer evidence countering the narrative of protesters as perpetuating anger. However, our analysis is limited to specifically tweets from June 2020, and we cannot conclude how they may generalize to other data sources or time periods, especially given that our work highlights the importance of context in emotion analyses.

Prior work on the Black Lives Matter movement has also examined emotions. One study uses LIWC lexicons to measure several dimensions, including positive/negative affect, anger, anxiety, sadness, and swear, in a dataset of tweets about Black Lives Matter protests in 2014 and 2015 (3). The authors find that

Table 2. Pearson correlations between the percentage of tweets with each emotion and the number of protests in each state

Emotion	CCC	P value	ACLED	P value
Anger	-0.38	0.0072	-0.42	0.0020
Disgust	-0.19	0.1869	-0.30	0.0356
Joy	0.48	0.0004	0.48	0.0003
Surprise	-0.31	0.0267	-0.18	0.2009
Fear	-0.02	0.8867	0.11	0.4517
Sadness	-0.21	0.1435	-0.40	0.0040

Tweets are associated with US states based on locations listed by users in their profiles, where 20.66% of users were aligned to states. Emotion categories are drawn from Ekman's taxonomy (20) and inferred using a neural classification model with domain adaptation components (+TGT+FSL). Protest data are drawn from two initiatives: the ACLED and

^{††}Information is available at https://www.safegraph.com/.

Table 3. Pearson correlations between the percentage of tweets with each emotion and the number of protests in each city

Emotion	CCC	P value	ACLED	P value
Anger	-0.22	0.0001	-0.28	0.0000
Disgust	-0.21	0.0001	-0.27	0.0000
Joy	0.23	0.0000	0.26	0.0000
Surprise	-0.04	0.4534	-0.09	0.1009
Fear	0.16	0.0040	0.18	0.0012
Sadness	-0.27	0.0000	-0.18	0.0009

Tweets are associated with US cities based on locations listed by users in their profiles, where 12.3% of users were aligned to cities. Emotion categories are drawn from Ekman's taxonomy (20) and inferred using a neural classification model with domain adaptation components (+TGT+FSL). Protest data are drawn from two initiatives: the ACLED and the CCC.

anger tends to decrease over time, while friends and social tend to increase, supporting the theory that anger and outrage may cause initial participation but that joy and camaraderie facilitate sustained involvement. They also find that high negativity and sadness but low anger and anxiety on Twitter are predictive of an increased volume of future protests. Beyond language and emotion in Black Lives Matter tweets, other work has examined the motivations and identities of individuals involved, including the prominence of female activists (1), the demographics of Twitter users (56), the roles that activists take (5), communication networks and widely shared content (4), estimations of violence using images (57), and the broader implications of social media activism (6).

While our analysis focuses on tweets about Black Lives Matter, our methodology can be used in other settings, requiring only a small annotated set of in-domain data for fine-tuning and evaluation. These analyses and methodologies can enhance understanding of social movements, providing information to social scientists and activists.

Materials and Methods

Model Setup. Our primary classifier for identifying emotions uses pretrained BERT as the base network. In emotion classification, BERT has consistently outperformed other models, such as convolutional neural networks and RoBERTa (A Robustly Optimized BERT Pretraining Approach) (19, 28, 58). We append a two-layer feed-forward neural network on top of BERT, which takes the mean pooled representation of all input tokens. We train one classifier per emotion, which makes each task a binary classification task. We also experimented with multiclass classification, but we found little difference in performance and ultimately, use single-class models to ensure that any

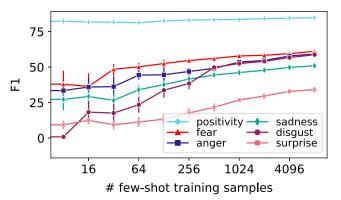


Fig. 7. F1 scores of emotion classifiers on HurricaneEmo test data using GoEmotions and varying numbers of few-shot training samples from HurricaneEmo as training data. Results are averaged across 10 random seeds, and the error bars indicate the 95% CIs.

identified correlations between emotions are not model artifacts. For BERT's hyperparameters, we used the BERT base model from the transformers library (59).

The models are optimized using the AdamW optimizer with cross-entropy loss. Ref. 60 reports that source performance on validation set is often uncorrelated with target validation performance and suggests using the target validation set for model selection even in the zero-shot setting. Following this suggestion, we used the validation split of HurricaneEmo to choose the final model in both zero-shot and few-shot learning settings. More details about the data preprocessing, model, and hyperparameters can be found in SI Appendix, section 2.

Few-Shot Training Data Size. In order to finalize the +FSL model, we use GoEmotions with small subsets of HurricaneEmo as training data and HurricaneEmo as test data to experiment with different in-domain dataset sizes. Fig. 7 reports results. Unsurprisingly, performance improves as the size of the indomain training data increases. However, the rate of improvement is not standard for all emotions. Prediction of positivity changes little with increasing dataset sizes, while prediction of disgust shows the greatest changes. The steepest rate of improvement occurs between 0 and 256 data points, after which we see diminishing returns for most emotions. Based on these results, we fix the indomain data size for the +FSL models to 300 and annotated 700 instances from #BLM2020 to facilitate few-shot learning and evaluation.

#BLM2020 Annotations. We collected an initial set of annotations over 400 tweets from #BLM2020; this was conducted by five volunteers who were living in the United States throughout the time period in our dataset. In the annotation instructions, annotators were provided with all subemotions used in GoEmotions for each high-level emotion (listed in Detecting Emotions Expressed in Tweets) and asked to select all the emotions that occurred in the tweet, either expressed by the author or solicited in the reader. For each tweet, we collected two independent judgments. If the two annotators disagreed on any label, a third independent annotation was collected. In order to ensure annotation quality in our test set, we revised the annotation scheme based on feedback from the initial annotations and collected annotations over an additional 300 tweets from six annotators, where each tweet was annotated by all six annotators. We report additional details, including instructions provided to annotators and agreement over each emotion in SI Appendix, section 2. Over the test set, Krippendorff's α is \geq 0.49 for all emotions, except surprise (0.26) and sadness (0.35), and interrater correlation is >0.5 for all emotions. Agreement for surprise and sadness is likely lower than other emotions due to the rareness of these emotions in our data. Given the subjective nature of emotions and that we avoid directing annotators on how to annotate particular types of tweets to avoid unduly influencing results, some disagreement is expected over these data. Agreement over our dataset is higher than the agreement reported for GoEmotions (18).

Data, Materials, and Software Availability. Code has been deposited in GitHub (https://github.com/chan0park/BLM-emotions) (61). Following feedback from reviewers of our initial submission as well as in accordance with Twitter terms of service, the tweet identifications for tweets used in this study will be made available for the purposes of academic research only upon request. To preserve anonymity as much as possible, they will not be posted publicly.

ACKNOWLEDGMENTS. We acknowledge feedback on and support of this work from Amanda Chen, Jamie Cho, Judeth Choi, Amanda Coston, Katherine Keith, David Orr, and Elizabeth Salesky. A.F. acknowledges support from a Google PhD Fellowship and NSF Graduate Research Fellowship Program Grant DGE1745016. A.T. acknowledges support from Google Latin America Research Awards and Sao Paulo Research Foundation (FAPESP) Grant 2019/21030-1. C.P. acknowledges support from Korea Foundation for Advanced Studies Overseas PhD Scholarship Program. J.W.-D. acknowledges support from a Ford Foundation Pre-Doctoral Fellowship and NSF Graduate Research Fellowship Program Grant DGE1745303. Y.T. acknowledges support from NSF Faculty Early Career Development Program (CAREER) Grant IIS2142739 and an Alfred P. Sloan Foundation Fellowship. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the NSF.

- 1. A. V. Richardson, Dismantling respectability: The rise of new womanist communication models in the era of Black Lives Matter. J. Commun. 69, 193-213 (2019).
- Y. Bonilla, J. Rosa, #Ferguson: Digital protest, hashtag ethnography, and the racial politics of social 2. media in the United States. Am. Ethnol. 42, 4-17 (2015).
- M. Choudhury, S. Jhaver, B. Sugar, I. Weber, "Social media participation in an activist movement for 3. racial equality" in Proceedings of the International AAAI Conference on Web and Social Media (AAAI Press, Palo Alto, CA, 2016), pp. 92-101.
- D. Freelon, C. D. McIlwain, M. Clark, Beyond the Hashtags: #Ferguson,#BlackLivesMatter, and the Online Struggle for Offline Justice (Center for Media & Social Impact, American University, 2016).
- J. O. Choi, J. Herbsleb, J. Hammer, J. Forlizzi, "Identity-based roles in rhizomatic social justice movements on Twitter" in Proceedings of the International AAAI Conference on Web and Social Media (AAAI Press, Palo Alto, CA, 2020), vol. 14, pp. 488-498.
- S. J. Jackson, M. Bailey, B. Foucault Welles, HashtagActivism: Networks of Race and Gender Justice (The MIT Press, 2020)
- J. Jasper, The Art of Moral Protest: Culture, Biography, and Creativity in Social Movements (University of Chicago Press, 1997).
- J. Goodwin, J. M. Jasper, F. Polletta, Emotional Dimensions of Social Movements (John Wiley & Sons, Ltd., 2007), pp. 413-432.
- J. M. Jasper, Emotions and social movements: Twenty years of theory and research. Annu. Rev. Sociol. **37**, 285-303 (2011).
- 10. M. L. Y. Ma, Affective framing and dramaturgical actions in social movements. J. Commun. Inq. 41, 5-21 (2017)
- A. M. Allen, C. W. Leach, The psychology of Martin Luther King Jr.'s "creative maladjustment" at societal injustice and oppression. J. Soc. Issues 74, 317–336 (2018).
- P. Collins, Black Feminist Thought: Knowledge, Consciousness, and the Politics of Empowerment
- (Taylor & Francis, 1990). J. Ć. Walley-Jean, Debunking the myth of the "angry black woman": An exploration of anger in young
- African American women. Black Women Gend. Fam. 3, 68-86 (2009).
- J. Leopold, M. P. Bell, News media and the racialization of protest: An analysis of Black Lives Matter articles. Equal. Divers. Incl. 36, 720-735 (2017).
- S. M. Mohammad, P. D. Turney, Crowdsourcing a word-emotion association lexicon. Comput. Intell. 29, 436-465 (2013).
- Z. Steinert-Threlkeld, J. Joo, Protest event data from geolocated social media content. APSA Preprints [Preprint] (2020). https://preprints.apsanet.org/engage/apsa/article-details/5f594a6b1d75ae 001b0fab90 (Accessed 7 January 2022).
- 17. S. Mohammad, F. Bravo-Marquez, M. Salameh, S. Kiritchenko, "Semeval-2018 task 1: Affect in tweets" in Proceedings of the 12th International Workshop on Semantic Evaluation, M. Apidianaki et al., Eds. (Association for Computational Linguistics, New Orleans, LA, 2018), pp. 1-17.
- D. Demszky et al., "GoEmotions: A dataset of fine-grained emotions" in Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, D. Jurafsky, J. Chai, N. Schluter, J. Tetreault, Eds. (Association for Computational Linguistics, 2020), pp. 4040–4054.
- 19. S. Desai, C. Caragea, J. J. Li, "Detecting perceived emotions in hurricane disasters" in *Proceedings of* the 58th Annual Meeting of the Association for Computational Linguistics (Association for Computational Linguistics, D. Jurafsky, J. Chai, N. Schluter, J. Tetreault, Eds. 2020), pp. 5290–5305.
- P. Ekman, Are there basic emotions? Psychol. Rev. 99, 550–553 (1992).
- A. Field, G. Bhat, Y. Tsvetkov, "Contextual affective analysis: A case study of people portrayals in online #MeToo stories" in Proceedings of the International AAAI Conference on Web and Social Media (AAAI Press, Palo Alto, CA, 2019), vol. 13, pp. 158-169.
- A. Field, Y. Tsvetkov, "Entity-centric contextual affective analysis" in Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, A. Korhonen, D. Traum, L. Màrquez, Eds. (Association for Computational Linguistics, Florence, Italy, 2019), pp. 2550-2560.
- C. Y. Park, X. Yan, A. Field, Y. Tsvetkov, "Multilingual contextual affective analysis of LGBT people portrayals in Wikipedia" in Proceedings of the International AAAI Conference on Web and Social Media (AAAI Press, Palo Alto, CA, 2021), vol. 15, pp. 479-490.
- 24. L. F. Barrett, J. A. Russell, The Psychological Construction of Emotion (Guilford Publications, 2014).
- 25. S. Mohammad, P. Turney, "Emotions evoked by common words and phrases: Using Mechanical Turk to create an emotion lexicon" in Proceedings of the NAACL HLT 2010 Workshop on Computational Approaches to Analysis and Generation of Emotion in Text, D. Inkpen, C. Strapparava, Eds. (Association for Computational Linguistics, Los Angeles, CA, 2010), pp. 26-34.
- M. S. Rasooli, N. Farra, A. Radeva, T. Yu, K. McKeown, Cross-lingual sentiment transfer with limited resources. Mach. Transl. 32, 143–165 (2018).
- 27. C. Potts, Z. Wu, A. Geiger, D. Kiela, "DynaSent: A dynamic benchmark for sentiment analysis" in Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), C. Zong, F. Xia, W. Li, R. Navigli, Eds. (Association for Computational Linguistics, 2021), pp. 2388–2404.
- J. Devlin, M. W. Chang, K. Lee, K. Toutanova, Bert: Pre-Training of Deep Bidirectional Transformers for Language Understanding (NAACL-HLT, 2019).
- J. Howard, S. Ruder, "Universal language model fine-tuning for text classification" in Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), I. Gurevych, Y. Miyao, Eds. (Association for Computational Linguistics, Melbourne, VIC, Australia,
- 30. S. Gururangan et al., "Don't stop pretraining: Adapt language models to domains and tasks" in Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, D. Jurafsky, J. Chai, N. Schluter, J. Tetreault, Eds. (Association for Computational Linguistics, 2020), pp. 8342-8360.

- 31. R. Plutchik, The nature of emotions: Human emotions have deep evolutionary roots, a fact that may explain their complexity and provide tools for clinical practice. Am. Sci. 89, 344-350 (2001).
- Y. R. Tausczik, J. W. Pennebaker, The psychological meaning of words: LIWC and computerized text analysis methods. *J. Lang. Soc. Psychol.* **29**, 24–54 (2010).
- C. Chen et al., Distinct facial expressions represent pain and pleasure across cultures. Proc. Natl. Acad. Sci. U.S.A. 115, E10013-E10021 (2018).
- L. F. Barrett, R. Adolphs, S. Marsella, A. M. Martinez, S. D. Pollak, Emotional expressions reconsidered: Challenges to inferring emotion from human facial movements. Psychol. Sci. Public Interest 20, 1-68
- K. Crawford et al., Al Now 2019 Report (Al Now Institute, New York, NY, 2019).
- J. M. Salerno, L. C. Peter-Hagene, The interactive effect of anger and disgust on moral outrage and judgments. Psychol. Sci. 24, 2069-2078 (2013).
- B. L. Monroe, M. P. Colaresi, K. M. Quinn, Fightin' words: Lexical feature selection and evaluation for identifying the content of political conflict. Polit. Anal. 16, 372-403 (2017).
- J. Huang, K. M. Thornton, E. N. Efthimiadis, "Conversational tagging in twitter" in *Proceedings of the* 21st ACM Conference on Hypertext and Hypermedia, HT '10, M. Chignell, E. Toms, Eds. (Association for Computing Machinery, New York, NY, 2010), pp. 173-178.
- C. Raleigh, A. Linke, H. Hegre, J. Karlsen, Introducing ACLED: An armed conflict location and event dataset: Special data feature. J. Peace Res. 47, 651-660 (2010).
- Crowd Counting Consortium, A public interest and scholarly project to document crowds and contention in the United States (2022).
- https://sites.google.com/view/crowdcountingconsortium/home. Accessed 24 January 2022. A. Sobolev, M. K. Chen, J. Joo, Z. C. Steinert-Threlkeld, News and geolocated social media accurately
- measure protest size variation. *Am. Polit. Sci. Rev.* **114**, 1343–1351 (2020).

 B. Hecht, L. Hong, B. Suh, E. H. Chi, "Tweets from JUSTIN BIEber's heart: The dynamics of the location field in user profiles" in Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, CHI '11, D. Tan, G. Fitzpatrick, C. Gutwin, B. Begole, W. A. Kellogg, Eds. (Association for Computing Machinery, New York, NY, 2011), pp. 237–246.
- 43. B. Alex, C. Llewellyn, C. Grover, J. Oberlander, R. Tobin, "Homing in on Twitter users: Evaluating an enhanced geoparser for user profile locations" in Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16), N. Calzolari et al., Eds. (European Language Resources Association (ELRA), Portorož, Slovenia, 2016), pp. 3936-3944.
- E. Chenoweth, M. J. Stephan, Why Civil Resistance Works: The Strategic Logic of Nonviolent Conflict (Columbia University Press, 2011).
- M. Biggs, Size matters: Quantifying protest by counting participants. Sociol. Methods Res. 47, 351-383 (2018)
- G. King, J. Pan, M. E. Roberts, How the Chinese government fabricates social media posts for strategic distraction, not engaged argument. Am. Polit. Sci. Rev. 111, 484-501 (2017).
- A. Arif, L. G. Stewart, K. Starbird, Acting the part: Examining information operations within #BlackLivesMatter discourse. Proc. ACM Hum. Comput. Interact. 2, 20 (2018).
- N. A. Valentino, K. Gregorowicz, E. W. Groenendyk, Efficacy, emotions and the habit of participation. Polit. Behav. 31, 307-330 (2009).
- N. A. Valentino, T. Brader, E. W. Groenendyk, K. Gregorowicz, V. L. Hutchings, Election night's alright for fighting: The role of emotions in political participation. J. Polit. 73, 156–170 (2011).
- W. Pearlman, Emotions and the microfoundations of the Arab uprisings. Perspect. Polit. 11, 387-409 50. (2013).
- A. J. Banks, I. K. White, B. D. McKenzie, Black politics: How anger influences the political actions blacks pursue to reduce racial inequality. Polit. Behav. 41, 917-943 (2019).
- C. D. Burge, Introduction to dialogues: Black affective experiences in politics. Polit. Groups Identities 8 390-395 (2020).
- J. S. Scott, J. Collins, Riled up about running for office: Examining the impact of emotions on political ambition. Polit. Groups Identities 8, 407-422 (2020).
- 54. D. L. Phoenix, The Anger Gap: How Race Shapes Emotion in Politics (Cambridge University Press,
- 55. L. Huddy, S. Feldman, C. Weber, The political consequences of perceived threat and felt insecurity. Ann. Am. Acad. Polit. Soc. Sci 614, 131-153 (2007).
- A. Olteanu, I. Weber, D. Gatica-Perez, "Characterizing the demographics behind the #BlackLivesMatter movement" in Proceedings of AAAI Spring Symposia on Observational Studies through Social Media and Other Human-Generated Content (AAAI Press, Palo Alto, CA, 2015), pp. 4144-4154.
- D. Won, Z. C. Steinert-Threlkeld, J. Joo, "Protest activity detection and perceived violence estimation from social media images" in Proceedings of the 25th ACM International Conference on Multimedia, MM '17, Q. Liu et al., Eds. (Association for Computing Machinery, New York, NY, 2017), pp. 786-794.
- Y. Liu et al., RoBERTa: A robustly optimized BERT pretraining approach. arXiv [Preprint] (2019). https://arxiv.org/abs/1907.11692 (Accessed 1 March 2022).
- T. Wolf et al., "Transformers: State-of-the-art natural language processing" in Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, Q. Liu, D. Schlangen, Eds. (Association for Computational Linguistics, 2020), pp. 38-45.
- P. Keung, Y. Lu, J. Salazar, V. Bhardwaj, "Don't use English dev: On the zero-shot cross-lingual evaluation of contextual embeddings" in Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), B. Webber, T. Cohn, Y. He, Y. Liu, Eds. (Association for Computational Linguistics, 2020), pp. 549-554.
- A. Field, C. Y. Park, A. Theophilo, J. Watson-Daniels, Y. Tsvetkov, Code and data for "An Analysis of Emotions and the Prominence of Positivity in #BlackLivesMatter Tweets." GitHub. https://github. com/chan0park/BLM-emotions. Deposited 7 July 2022.