

How Violence Affects Protests

Zachary C. Steinert-Threlkeld¹, Alexander Chan², and Jungseock Joo³

¹University of California, Los Angeles. Email: zst@luskin.ucla.edu

²University of California, Los Angeles. Email: alexander.chan@ucla.edu

³University of California, Los Angeles. Email: jjoo@comm.ucla.edu

ABSTRACT

A key determinant of whether social movements achieve their policy goal is how many people protest. How many people protest is in turn partially determined by violence from protesters and state agents. Previous work finds mixed results for violence. This paper reconciles the mixed results for violence by distinguishing between its timing, source, and severity: low levels of state repression increase protest size while high levels decrease it, conditional on preventative repression failing. Protester violence is always associated with fewer subsequent protesters. These claims are substantiated by applying deep learning techniques to geolocated protest images shared on social media. Across more than 4,300 observations of twenty-four cities from five countries, we find that protester violence is always associated with subsequently smaller protests, while low (high) levels of state violence correlate with increased (decreased) protest size. The paper ends with a discussion of situations in which to prefer images or text for studying protests; ethical concerns; and improving data collection in order to apply the analysis to poorer or less populous environments.

19 **1 INTRODUCTION**

20 Protests are more likely to shift policy the larger they are (Chenoweth and Stephan, 2011;
21 Fassiotto and Soule, 2017; Gause, 2018). In turn, their size is affected by the level and source
22 (protesters or state agents) of violence at a protest. While protester violence has consistently been
23 found to correlate with smaller protests (Stephan and Chenoweth, 2008; Feinberga et al., 2017;
24 Murdie and Purser, 2017), the protest-repression literature consistently finds inconsistent results
25 (Davenport, 2007). This paper suggests a solution to the protest-repression puzzle rooted in the
26 timing and level of state violence.

27 Pre-emptive state repression decreases protests (Sullivan, 2016). Once protests start, however,
28 the effect of state violence depends on its severity. We develop and test an argument that low levels
29 of state violence lead to larger protests, while high levels decrease it; protester violence leads to
30 smaller protests.¹ The importance of the severity of state repression may explain varying effects
31 that the literature identifies for state repression (Carey, 2006; Ritter, 2013), though it contradicts
32 the backlash hypothesis (Francisco, 2004).

33 These theoretical expectations are tested using protest images shared in geolocated tweets. This
34 measurement occurs using three convolutional neural networks (CNN). We first develop a CNN to
35 recognize protest images, and we have verified that this classifier outperforms Google Vision on our
36 images (see Figure A8). From a corpus of six years of geolocated tweets, we identify 55.6 million
37 from protest waves across fourteen countries, 5.4 million of which contain an image. The protest
38 detection classifier identifies over 115,000 of these images as very likely to contain a protest, and we
39 build a second CNN to measure protester and state violence as well as the presence of fire or police
40 officers. Figure A8 also shows that this scene classifier identifies police better than Google Vision.
41 This scene classifier is complemented with the third CNN, a face classifier. This classifier counts
42 the number of faces per photo and estimates the race, gender, and age of each face, allowing us to
43 control for well-known correlations between these demographic features and protest participation.
44 While many off the shelf face classifiers exist, only one codes race precisely enough (Kärkkäinen

¹Violence should be conceived of as “perceived violence”, a point to which we return later.

45 and Joo, 2019). Sections 3 and S1 detail how we built these models and verify their output.

46 Section 4 discusses the resulting data, concerns about selection bias, and validates the dependent
47 variable. It shows that users who tweet geolocated protest images are likely to be more representative
48 of normal Twitter users than those who tweet geolocated non-protest images, and there are strong
49 reasons to expect social media to be no more biased than newspapers in covering protests. That
50 section also shows that using images to measure protest size correlates with estimates of protest
51 size from newspapers and records changing activity that matches events as reported in newspapers.

52 Section 5 presents the main results. It also shows three sets of robust checks, two explicitly de-
53 signed to account for possible bias. Restricting results to users of moderate popularity; non-verified
54 users; non-bots; and tweets in a country’s *lingua franca* do not change inferences. Deduplicating
55 images also does not change inferences. To the extent we find evidence for bias, it is against the
56 main results: dropping bots and restricting tweets based on language both produce better model fit
57 than the full datasets.

58 Finally, Section 6 concludes, discussing why images, instead of text, are necessary for this
59 project; the ethical concerns raised by computer vision approaches, especially in the context of
60 contentious politics; and why the results presented in this paper should be considered a lower bound
61 on what these techniques can achieve.

62 2 PROTEST DYNAMICS

63 Protesters aim to convince bystanders to mobilize, increasing pressure for policy change (?). The
64 state works to convince bystanders to remain on the sidelines and existing protesters to disengage.
65 Protesters and the state each choose amounts of violence to employ. Protester violence should
66 always lead to smaller protests, while state violence will have differential effects depending on its
67 timing and severity.

68 2.1 Large Protests

69 Three assumptions lead to the conclusion that large protests are more likely to change policy
70 than small ones. If (1) the purpose of a protest is to convince political leaders to change a policy,
71 (2) a leader cares about the median voter (Downs, 1957) or his or her winning coalition exhibits

72 some response to the median person (Bueno de Mesquita et al., 2003), and (3) a large protest's
73 policy preference is closer to the median individual than a small one's, then a large protest is more
74 likely to change policy than a small one.

75 This argument can also be obtained without assuming a leader aims for the median individual's
76 policy preference. If a leader only desires to stay in power and a large protest means the probability
77 of remaining in power is lower than the leader previously believed, a large protest is still more likely
78 to lead to policy change than a small one. While a large protest is not necessarily successful, all
79 successful protests are large.

80 Protest size is important regardless of a country's political institutions. In democracies, voting
81 is the most common method of policy change. The aggregation of preferences through defined
82 rules, and the willingness of those in power to heed the result, has many advantages. It is a
83 low-cost endeavor for participants, as the only costs are transaction and opportunity. Voting occurs
84 infrequently, however, and is a blunt method of feedback because it collapses all political dimensions
85 into one. Protests, however, can occur at any time and usually have a clear policy goal (Battaglini,
86 2017).

87 In authoritarian regimes, however, voting is a less significant act. In countries where policy
88 feedback comes from an insider population drawn from the larger populace (Bueno de Mesquita
89 et al., 2003), those belonging to the outsider population provide policy feedback through protest
90 (or rebellion). While protest is unlikely to change an autocrat's policy, it nonetheless provides a
91 key signal of discontent to which a government can respond (Bratton and Walle, 1992). This signal
92 is especially pertinent if opinion polling is unreliable (Robertson, 2007) or the media are not free
93 (Qin et al., 2017).

94 The importance of large protests has not escaped notice. Lohmann (1994a) argues that unprece-
95 dented numbers of people rallying in the German Democratic Republic in the beginning of 1989
96 were a key reason the protest movement grew and eventually toppled Erich Honecker's government.
97 The gradual growth of protest size in Iran in 1979 also made it increasingly difficult for the Shah
98 to remain in power (Rasler, 1996). Kuran (1989)'s canonical model of bandwagoning implicitly

means a revolution follows large protests. Understanding the threat posed by large crowds, regimes often raise the cost of protesting by killing protesters, yet killing protesters has an indeterminate effect on the size of subsequent protests (Francisco, 2004). Indeed, if there is a law like regularity to the study of protest mobilization, it is that “size matters” (Biggs, 2016).

The importance of size applies to social movements as well, of which protests are a tactic they can employ (Tilly and Wood, 2012). Large social movements are more likely to lead to policy change than small ones for several reasons. Because large social movements tend to be nonviolent, they increase the domestic and international cost of repression, especially when movements maintain their own media (Sutton et al., 2014). They decrease the cost of participation, making individuals more likely to join, and making them more likely to join as movements grow (DeNardo, 1985). It also increases the probability that individuals within the winning coalition defect, making it more difficult for the state to continue repression (Goldstone, 2001). For a fuller exposition of the importance of size for social movements, see Chenoweth and Stephan (2011).

2.2 Violence

Though scholars understand the importance of large protests, less is understood about why protests become large, and most existing work is qualitative, cross-sectional, or focuses on structural variables. For example, Biggs (2003) argues, in explaining American protests in 1886, for a positive feedback loop but does not specify when an initial protest is more likely to generate that process. Large protests in one country may occur because large protests in a similar country succeeded, but contagion does not explain why the initiating country experienced large protests (Weyland, 2012). The structure of built environments may also encourage protest participation: one reason initial marches to Tiananmen Square succeeded is because universities in Beijing are in the same neighborhood and have internally dense configurations, encouraging mobilization both within and between campuses (Zhao, 1998). The occurrence of electoral fraud is also a common source of large protest events (Tucker, 2007).

This paper focuses on mechanisms affecting protest size *once a protest starts*. The only kind of repression that can occur in this situation is overt repression, often called “protest policing”

126 (Earl, 2003; Davenport and Soule, 2009). This concept refers to repressive behaviors that occur
127 during a protest, such as blocking roads, impeding pedestrian movement, arresting protesters, and
128 using subfatal weapons such as tear gas, water cannons, or sound guns. Protest policing contrasts
129 with preventative repression, such as the targeting of dissident organizations or arresting particular
130 individuals (Sullivan, 2016).

131 *Protester Violence*

132 If violent protests are more costly to the individual than nonviolent ones, regardless of the
133 source, then violence should decrease protest size (Moore, 1995). Empirically, however, there
134 appears to be differential effects based on the source of the violence.

135 When violence originates from protesters, it should always decrease the size of protests because
136 it decreases the number of people to which protest appeals, and it increases the cost of protesting
137 to the remaining bystanders who may protest.

138 One method by which bystanders determine whether or not to join a protest is to compare
139 protesters' ideological distance to their own (Lohmann, 1993). Since most individuals do not
140 support violence or receive consumption value from it (Feinberg et al., 2017), violence originating
141 from protesters signals that protesters, and therefore the policy changes for which they agitate, are
142 far from mainstream. Survey research has found that the more activists differ from the population
143 they try to mobilize, the less likely individuals are to protest (Bashir et al., 2013). Being far from
144 the mainstream, bystanders continue to stand by because the new policy the violent protesters seek
145 is inferred to not be beneficial.

146 Protester violence decreases the likelihood of regime defections, decreasing the number of
147 non-protesters available to mobilize. Peaceful protest convinces regime agents of their physical
148 safety should they defect, increasing the probability that police, members of the armed forces,
149 or legislators, for example, switch allegiances (Stephan and Chenoweth, 2008). Violent protests,
150 however, induce fear in these agents that they will meet the same fate if they do not remain loyal.
151 Violence therefore reduces the pool of those willing to protest, making the state stronger than an
152 equivalent peaceful protest.

153 Since protester violence decreases the legitimacy of protests, the state can pursue high levels of
154 repression and face less risk of backlash. Peaceful protests enjoy high domestic and international
155 legitimacy, so state violence against them risks generating a backlash that increases subsequent
156 protests' size (Francisco, 2004). But since violent protesters can be framed as rioters, terrorists, or
157 foreign agitators (Benford and Snow, 2000), bystanders are more supportive of repressing violent
158 protests than nonviolent ones. Survey work across eighteen countries finds that violent protests
159 decrease future support for the peaceful right to protest (Murdie and Purser, 2017). For the same
160 reasons, the state is also less likely to receive international sanction when repressing violent protests.

161 The converse of these arguments is that protester non-violence increases the probability that a
162 protest grows in size, especially when states repress. Because non-violence increases the legitimacy
163 of protests, it decreases the probability that a state represses, as the state will pay large reputation
164 costs. The lower probability of repression induces more bystanders to mobilize, generating a
165 positive feedback loop (Lohmann, 1994b). In Morocco, for example, attempts to repress non-
166 violent protesters at the start of the Arab Spring led to larger crowd sizes (Lawrence, 2016), and
167 government violence in Tunisia did not prevent the spread of those protests.

168 Since protester violence alienates bystanders, increases the resolve of state agents, and invites
169 high levels of state repression, we expect that:

170 *H1: There should exist a negative relationship between protester violence and the subsequent*
171 *size of a protest.*

172 This hypothesis extends earlier work that finds the same relationship at the movement level,
173 using a movement's reported maximum participation rate. As far as we are aware, existing work
174 on protester violence and outcomes is cross-sectional (Stephan and Chenoweth, 2008; Celestino
175 and Gleditsch, 2013; Chenoweth and Schock, 2015) or focused on its interaction with state tactics
176 (Shellman et al., 2013). It is therefore unclear if protester violence decreases participation, less
177 participation causes protesters to result to violence, or a smaller movements results from another
178 feature. By developing a logic for protester violence and individual participation, we directly link

179 these two and explain how the former should affect the latter's fluctuation.

180 *State Violence*

181 While the negative relationship between protester violence and movement success is a regular
182 finding, the literature on state repression and protest has not found consistent effects. In Peru and
183 Sri Lanka, repression decreased subsequent protests (Moore, 2000). The same has been found in
184 West Germany (Koopmans, 1993), South Africa (Olzak et al., 2003), Iran in the short-term (Rasler,
185 1996), and the Middle East and North Africa during the Arab Spring (Steinert-Threlkeld, 2017).
186 On the other hand, repression may have increased protest in West Germany and Ireland (Francisco,
187 1996) and Iran with a six-week lag (Rasler, 1996), and many cross-national studies find repression
188 increases protest (Gurr and Moore, 1997; Davenport and Armstrong II, 2004; Francisco, 2004;
189 Hess and Martin, 2006). It can also increase protest based on the emotional reaction of individuals
190 connected to those targeted (Siegel, 2011; Pearlman, 2013). On the third hand, there is sometimes
191 no correlation between repression and protest levels (Gupta et al., 1993; Ritter, 2013; Ritter and
192 Conrad, 2016).

193 These contradictory findings are resolved by considering the timing of repression and the sever-
194 ity of it. When mobilization is the result of social movement organizations' planning, repression
195 focusing on those organizations should decrease protest size (Sullivan, 2016). This preemptive
196 repression attacks the infrastructure of protests, making it harder for them to occur, much less grow
197 (Danneman and Ritter, 2013; Sutton et al., 2014). This line of reasoning then argues that repression
198 of protests as they occur leads to backlash (Sullivan, 2016). Repression of protests as they occur,
199 commonly called protest policing (Della Porta and Reiter, 1998; Davenport and Soule, 2009; Earl
200 et al., 2013), leads to the differential effects discussed earlier.

201 Light repression will generate backlash for two reasons. First, they may signal that the cost of
202 protesting is lower than bystanders believed. Now aware that protesting is a net positive, bystanders
203 join those already protesting. Second, repression can generate emotions such as anger, joy, or
204 pride. Acting on these emotions provides intrinsic benefit to the former bystander, regardless of
205 instrumental calculations (Pearlman, 2013). Incorporating emotions into theories does not require

206 avoiding rationality assumptions, as protesting in anger at repression can be individually rational
207 ([Siegel, 2011](#)).

208 Severe repression, however, should lead to smaller protests, for similar reasons. Severe re-
209 pression may signal that state actors are more resolved than protesters expected. Facing a higher
210 cost to protest, protesters become bystanders. Severe repression also generates fear, sadness, and
211 shame, causing protesters to deactivate and bystanders to remain where they are ([Pearlman, 2013](#)).
212 This emotional effect has also received recent support in a series of lab-in-the-field experiments in
213 Zimbabwe ([Young, 2019](#)).²

214 For an earlier exposition of a similar argument, see [Gurr \(1970\)](#). For a formal derivation of
215 this relationship, see [DeNardo \(1985\)](#). Observational studies which distinguish types of repression
216 by the cost they impose also find that severe repression decreases mobilization ([Muller, 1985](#);
217 [Khawaja, 1993](#)). In other words, the contradictory effects may be due more to measurement error
218 than theoretical inconsistencies. Since it appears that apparently contradictory effects of repression
219 are resolved by stipulating the severity of repression, conditional on observing protest, we expect
220 that:

221 *H2: There is an n-shaped relationship between between state repression and the subsequent
222 size of a protest.*

223 *H2* should apply in democracies and autocracies. For example, the Occupy Wall Street move-
224 ment in the United States did not grow large until New York City police arrested over 700 partic-
225 ipants, in a manner many perceived as unjust, marching on the Brooklyn Bridge. The movement
226 waned six weeks later, in the middle of November 2011, once local police forcibly dismantled
227 protesters' main encampment at Zuccotti Park and forbade them from spending the night ([White,
228 2016](#)). This effect should occur in democracies and autocracies. In Egypt, the protests
229 starting on January 25th were met with initial state resistance and some casualties; 18 days later,

²[Francisco \(2004\)](#) finds that state massacres increase mobilization. This result is due to an expansive definition of mobilization: the majority of the backlash events are substitutes for mobilization because they are harder to repress ([Moore, 2000](#)). Our focus is on mobilized protesters, not all forms of mobilization.

the Armed Forces forced President Hosni Mubarak to abdicate. Two years later, the Armed Forces launched a coup against the elected president, Mohamed Morsi. Large pro-Morsi protests erupted and continued for six weeks. The Armed Forces initial attempts to demobilize them caused them to grow in size; morning massacres on August 14th at the two main encampments killed at least 1,000 protesters and injured even more ([Shakir, 2014](#)).

*H*2 at first appears inconsistent with the backlash hypothesis substantiated in [Francisco \(1995\)](#), [Francisco \(1996\)](#), and [Francisco \(2004\)](#). It is not. That body of work argues against an “inverse-u” relationship between state repression and protest. Instead, evidence of backlash is found: when states engage in severe repression, the response is more collective action. That work, however, broadens protest to include other forms of collective action such as strikes, building occupations, or guerrilla action. Moreover, the substitution that does occur most often does not occur the day immediately following the repression. In other words, when the state meets protesters with severe repression, they initially reduce their protest; after some delay, they backlash by substituting away from direct confrontation with the state.

The argument put forth in this subsection is that severe repression decreases protest. It does not make a claim about whether other types of dissent increase. Works such as [Francisco \(1995\)](#), [Francisco \(1996\)](#), and [Francisco \(2004\)](#) define backlash in a more encompassing method than we do. This different definition is why they initially appear to have different expectations about, and different results for, backlash. H2 is not inconsistent with that backlash hypothesis because it is focused on a narrower window and action repertoire.

250 **3 AUTOMATED CODING OF SOCIAL MEDIA IMAGES**

To test the aforementioned hypotheses, we need to measure many variables, such as protester age or gender. For this task, we turn to images shared with geolocated tweets, explanation of which is provided in Sections 4 and S4 through S6. Since the target dataset is very large, we develop three automated classifiers based on convolutional neural networks to automatically code the variables of interest: one to identify protest images, a scene classifier to extract data (primarily about violence)

from protest images, and a face classifier to generate cleavage and size information.³ Table 1 provides an overview of the steps required in this pipeline, and the rest of this section provides a very brief introduction. The rest of this section describes the two classifiers we developed and one already existing one we used. For a high-level overview of how convolutional neural networks works, see Section S1. For validation of the classifiers’ results, see Section S1.3 as well as Section 4.3.

Table 1. Protest Data Pipeline

Steps	Input	Source	Output
<i>Collecting Images for Training Set</i>			
1. Image search	Keywords	Google	100,000 images
2. Train a protest image classifier	Images from Step 1	Self	Initial CNN
3. Protest images from Twitter corpus	Model from Step 2	Twitter	40,764 images
<i>Developing Protest and Scene Classifier</i>			
4. Manual annotation	Images from Step 3	Amazon Turk	13 ground-truth labels
5. Train a CNN	Training data from Step 4	Self	Protest and scene classifier
<i>Face Attribute Classification</i>			
6. Face classification	-	(Kärkkäinen and Joo, 2019)	Gender, age, and race estimates

3.1 Image Collection for Protest Classification by Weakly-supervised Learning

Step 1. As typically done in supervised machine learning, our approach in model development begins with collecting training data: images and target classification labels. Images in a training set should exhibit diverse visual traits of protest events and also include a range of negative (non-protest) images such that the trained classifier generalizes well to unseen images. In addition, it is desirable that the set also contains many difficult images, hard negatives, i.e. non-protest images which look like protest scenes, to make the classifier more robust.

The efficiency of manual annotation to collect target labels is another important consideration. For example, sampling general images and providing them to annotators would create a training set

³The first two classifiers are in fact partially combined in implementation such that one integrated classifier can generate two sets of outputs, although they differ conceptually. This is called multi-task learning (Girshick, 2015). We still discuss two classifiers separately because they are trained on different data and used in different steps.

of mostly non-protest images. This approach is not cost effective. Therefore, we take a combination of weakly-supervised and supervised learning. In weakly-supervised learning, the ground-truth labels on the target variable are not directly available but can be inferred from other variables (Bergamo and Torresani, 2010). For instance, we can use any online image search service to query images with a particular keyword (e.g., “protest”), and this step will furnish a large quantity of relevant images. While this sample set will contain some noisy data, it is still useful to train a rough initial model which can be used to fetch better samples. These samples can be manually annotated as in typical supervised learning.

Specifically, we first collected about 10,000 protest images from Google Image Search by using manually selected keywords such as “protest,” “riot,” “Black Lives Matter,” “Venezuela Protest,” “Hong Kong protest” and many others, as well 90,000 non-protest, hard-negative images by using keywords including “concert,” “stadium,” or “airport crowd.” These negative examples are called hard-negatives because they look similar to protest images (e.g., crowded), and classifiers can easily misassign their labels. Since these images are simply outputs of search queries, their assigned labels are not accurate. For example, the query of “protest” may return a few photographs of politicians. However, we did not verify the correct classification labels of these images because the main purpose of this first step is to train a rough classifier with the assumption that the majority of labels are still correct.

Step 2. Using these data, we trained a convolutional neural network (CNN) whose only output denotes whether an image captures a protest event or not. We then applied this classifier to geolocated images from Twitter and obtained the classification scores. Each score can be considered as the confidence about the output, the probability of the input image containing protesters. Section S1.1 provides detail of how CNNs work and the specific architecture of this paper’s, and Joo and Steinert-Threlkeld (2018) provides a detailed explanation of their relevance to political science.⁴

Step 3. Twitter provides tweets in real time through its streaming application programming interface (API). Since late 2013, one of the authors has used this interface to collect tweets with

⁴See as well (Cantu) for an application of this methodology to vote fraud detection.

longitude and latitude coordinates. Because tweets with GPS coordinates represent 2%-3% of all tweets and Twitter delivers tweets matching a request's parameters up to a 1% ceiling, we receive one-third to one-half all of tweets with precise location information (Morstatter et al., 2013; Leetaru, 2014).⁵ We have collected these tweets in real-time, approximately five million per day, since August 26, 2013. For more information on working with Twitter data, see Steinert-Threlkeld (2018).

We then query the stored tweets to extract those from countries and days of interest. These tweets could be used for text or social network analysis, but we further select only those tweets that contain images. Twitter provides a field in each tweet called `media_url` and a flag indicating if that link is for an image. If a downloaded tweet contains an image, we retrieve it. These images form the raw material from which we generate our protest data.

We apply the protest classifier to images from periods and countries during which protest occurred. These 14 periods, shown in Table A10, generated 55,676,431 tweets containing 5,479,148 images. The classifier is applied to all 5.48 million plus the 100,000 from Google, and all images with a classification score less than .6 are dropped as they are most likely easy negatives, i.e., non-protest images. The remaining 115,060 potential protest images were then stratified based on their classification scores and sampled to ensure that the chosen images capture diverse visual features, i.e., to avoid redundant inclusion of very similar images in the dataset. This process resulted in 40,764 images that form our training set; the training set contains geolocated images from Twitter and images from Google.

3.2 Protest and Scene Classification

Step 4. Amazon Mechanical Turk provided the labor to manually annotate these 40,764 images. We asked the workers to identify the features detailed in Table A8.

Figure A1 provides examples of our AMT annotation pages. In the first task, each annotator was presented with an image and asked to judge if the image captures a protest. We assigned two

⁵For example, requesting tweets with the keyword "Microsoft" will return every tweet with that word, assuming fewer than 1% of all tweets are about Microsoft. If, however, 2% of all tweets contain that word, then Twitter will return all tweets containing that keyword until the 1% ceiling is reached. $\frac{.01}{.02} = .5$, and the same calculation is how we conclude that our corpus contains one-third to one-half of all tweets with GPS coordinates.

322 workers to each image and if the two workers did not agree, the image was sent to a third judge for
 323 a final verification. 11,659 of the training images contain a protest. Similarly, in the second task,
 324 annotators label the attributes listed in Table A8 that are not related to faces or violence, such as
 325 “police”, “fire”, “children”, “flag”, and so on.

326 As violence is a subjective and continuous variable, we used pairwise comparison annotation
 327 to generate an estimate of the *perceived violence* in an image. Among the 11,659 protest images,
 328 we randomly sampled image pairs such that each image is paired ten times. Therefore the number
 329 of pairs to be annotated was 58,295 ($11,659 \times 10 \div 2$). We then assigned ten workers for each pair
 330 and asked them to select which image looks more violent than the other. To assign the continuous
 331 violence score to each image, we use the Bradley-Terry model (Bradley and Terry, 1952) and
 332 scaled the scores to the range of [0, 1]. Such a pairwise comparison method usually requires
 333 more annotations but can produce more reliable and consistent ratings for subjective assessment of
 334 photographs (Kovashka et al., 2012; Joo et al., 2014, 2015). The resulting estimate for violence is
 335 therefore better conceived of as perceived violence.

336 **Step 5.** With 40,764 annotated images, we train a CNN which produces outputs for twelve
 337 variables. We used 80% of the images as the training set and the rest as the validation set. For the
 338 labels that are not face or violence related, we use a binary cross entropy (BCE) loss:

$$339 \quad L_{BCE}(p, y) = -\frac{1}{N} \sum_{n=1}^N [y_n \log(p_n) + (1 - y_n) \log(1 - p_n)] \quad (1)$$

340 where p is the probability predicted by the model (CNN output for the attribute), y is the ground
 341 truth binary label (0 or 1), and N is number of images. p_n and y_n are the prediction and label for
 342 the n^{th} image, respectively.

343 For protester and state violence, a continuous variable, we use mean squared error (MSE) loss:

$$344 \quad L_{MSE}(p, y) = -\frac{1}{N} \sum_{n=1}^N [(y_n - p_n)^2] \quad (2)$$

345 where p is the model prediction, y is the ground truth value, and N is number of images. These

346 are standard loss functions that are typically used in training CNNs. Note that state-violence
347 and protester-violence are binary attributes and thus trained with a BCE loss in Eq. 1. Violence
348 measures the degree of violence on a continuous scale, and state- and protester-violence identify
349 the type of violence and are treated as binary variables. We use stochastic gradient descent with
350 backpropagation to train the model. For more technical details in model training, see [Won et al.](#)
351 ([2017](#)).

352 3.3 Face Classification

353 **Step 6.** We use the FairFace model developed by [Kärkkäinen and Joo \(2019\)](#) to classify gender,
354 race, and age of people in images. This new model is preferred over current leading models, such as
355 FaceNet ([Schroff et al., 2015](#)) or Face++, because it better captures race, gender, and age. Existing
356 public face datasets and commercial APIs have been criticized for their unbalanced representation
357 of race, as the vast majority of their face images are from people of white ethnicity (more than
358 80%). This results in inferior classification accuracy, especially on non-white people ([Buolamwini](#)
359 and [Gebru, 2018](#)). Moreover, the FairFace model is trained on a large corpus of images of varying
360 resolution, perspective, and lighting, the YFCC100M dataset ([Thomee et al., 2016](#)). This dataset
361 is in contrast to other datasets whose images tend to be high quality, well-lit, and from the same
362 perspective ([Liu et al., 2015](#)).

363 [Kärkkäinen and Joo \(2019\)](#) samples 102,218 of the 100 million YFCC100m images, with an
364 explicit focus on balancing users across seven racial categories. In contrast, [Liu et al. \(2015\)](#) uses
365 only three. Many other face models use skin color, but skin color is sensitive to lighting conditions.
366 In addition, there is no other large-scale face dataset or model offers the racial category of Latino,
367 which is critical in our study. On an external validation test, the model significantly outperforms
368 models trained on other large-scale datasets in gender, age, and race classification.

369 Figure A3 shows an image from South Korea from our Twitter corpus with the face classifier
370 applied.

371 4 RESEARCH DESIGN

372 **4.1 Data**

373 To identify protests, we searched for tweets from the fourteen periods detailed in Table A10.
374 For each period, we searched from one week prior to the first reported protest and one week after
375 the last one. This process identifies 55,676,431 tweets containing 5,479,148 images. To determine
376 which to keep, we chose the lowest threshold that would maximize recall with a precision of .85.
377 Figure A2 shows this threshold is .849 and recall is .22. This process results in 26,142 images.
378 This number represents about one-fifth of all protest images, and 85% of them are of protest.

379 We then aggregate tweets to their city of origin and the day they were created. Cities are kept
380 for analysis when at least $\frac{1}{7}$ of their days contain a protest image. Table 2 shows these 24 cities,
381 which account for 6,303 protest images. These 6,303 protest images spread across 4,401 city days
382 in Hong Kong, Pakistan, Spain (Catalonia only), South Korea, and Venezuela are the input for the
383 subsequent models.⁶ 1,467 of these city days contain a protest photo, so we treat the missing dates
384 as true zeroes. A robustness check shows that this interpolation does not change results.

385 **4.2 Bias**

386 Using social media data frequently raises concerns about selection bias ([Tufekci, 2014](#)). If
387 bias exists, it would come from accounts sharing images from protest activity not representative of
388 overall protest activity. We expect that Twitter users are not representative samples of their respective
389 countries, but we do not think the protest images they share are not representative. Moreover, if the
390 protest images are biased, structural features of the data generating process should make them less
391 biased than any other cross-national data source. Space constraints limit us from substantiating
392 these assertions here; see Section S2 for that substantiation.

393 In addition, three robustness checks presented in Section 5.2 for bias. First, we drop all tweets
394 from “verified” accounts, which are accounts belonging to prominent individuals or organizations
395 that Twitter has verified belong to those people or groups. Assuming they would have the most
396 incentive to filter what they publish, removing them removes a potential source of bias. Second,

⁶The majority of images are from the United States’ Women’s March, which is not analyzed here because it does not have a dynamic component. We also recorded large numbers of tweets from Belarus and Russia, but those protests occurred on one day as well. Many more images are then at the country level, so they are discarded.

Table 2. Protest Periods

	City	Country	Start	End	Issue	Protest Images/Day	Protest Images/Day if >0
1	Central	Hong Kong	2014.09.18	2014.12.23	China reforms	1.96	5.00
2	Kowloon	Hong Kong	2014.09.18	2014.12.23	China reforms	1.29	2.92
3	Lahore	Pakistan	2017.11.07	2017.11.23	Blasphemy	.18	1
4	Kimhae	South Korea	2016.10.20	2017.03.14	Anti-incumbency	.47	1.92
5	Seoul	South Korea	2016.10.20	2017.03.14	Anti-incumbency	2.40	3.76
6	Citutat Vella	Spain	2017.09.01	2017.12.31	Secession	.94	4.95
7	Barcelona	Spain	2017.09.01	2017.12.31	Secession	3.07	11.60
8	Girona	Spain	2017.09.01	2017.12.31	Secession	1.10	3.26
9	Granera	Spain	2017.09.01	2017.12.31	Secession	.62	2.33
10	Granollers	Spain	2017.09.01	2017.12.31	Secession	.23	1.25
11	Lleida	Spain	2017.09.01	2017.12.31	Secession	.42	1.88
12	Mataro	Spain	2017.09.01	2017.12.31	Secession	.51	2.33
13	Reus	Spain	2017.09.01	2017.12.31	Secession	.35	1.68
14	Sabadell	Spain	2017.09.01	2017.12.31	Secession	.96	2.66
15	St. Cugat d. Valles	Spain	2017.09.01	2017.12.31	Secession	.31	2.06
16	St. Feliu d. Pallerols	Spain	2017.09.01	2017.12.31	Secession	.61	2.19
17	St. Salvador d. Guardiola	Spain	2017.09.01	2017.12.31	Secession	.48	2.15
18	Tarragona	Spain	2017.09.01	2017.12.31	Secession	.57	1.94
19	Terrassa	Spain	2017.09.01	2017.12.31	Secession	.57	2.22
20	Boca del Rio	Venezuela	2014.03.27	12.17.2017	Anti-Maduro	.26	1.34
21	Caracas	Venezuela	2014.03.27	12.17.2017	Anti-Maduro	4.82	7.63
22	Caucagua	Venezuela	2014.03.27	12.17.2017	Anti-Maduro	.53	1.72
23	Maracaibo	Venezuela	2014.03.27	12.17.2017	Anti-Maduro	.39	1.49
24	Valencia	Venezuela	2014.03.27	12.17.2017	Anti-Maduro	.41	1.62

397 we only look at tweets from accounts between the 25th-75th percentile of their country's follower
398 distribution. Accounts below this range are likely to be bots or accounts which use Twitter
399 sporadically, while accounts above this range are more likely to be strategic with their posts. Third,
400 we remove tweets not in the *lingua franca* of their country, under the assumption that those are
401 aimed at international audiences and so are more likely to represent a protest differently than tweets
402 in the main language (Bruns et al., 2013).

403 4.3 Operationalization

404 The dependent variable is $\text{Log}_{10}(\text{Protest Size})_{i,t}$, the logarithm of the sum of the number of
405 faces in all protest photos from city i on day t . Other studies have found that activity on Twitter
406 correlates with verified estimates of crowd size for airports, stadiums, and protests (Botta et al.,
407 2015). Those estimates require either more data than were available to us or use text analysis to
408 identify protesters. Text analysis does not scale as easily as image analysis because it requires
409 domain expertise, so counting faces is preferred.

410 Figure 1 shows that this approach correlates with the size of protests in Russia and South Korea,
411 as reported in newspapers (Russia) or by activists and the police (South Korea).⁷ Small protests
412 reported in other sources corresponds closely with, and without bias to the size of the protest,
413 as what $\text{Log}_{10}(\text{Protest Size})_{i,t}$ estimates. Figure 2 shows how the protest size varies over time
414 in Catalonia, Spain and South Korea, with important events marked. There are clear spikes that
415 correspond to major events.⁸ For a verification of $\text{Log}_{10}(\text{Protest Size})_{i,t}$ against protest size as
416 recorded from cell phone location records and newspapers, see Sobolev et al. (2019). That summing
417 the number of faces in protest photos correlates well with protests in South Korea, Russia, and the
418 United States regardless of whether newspapers, reports from participants (activists or police), or
419 cell phones gives us enough confidence to trust this approach in other settings.

420 The violence variables to test Hypothesis 1 are $\text{Perceived Protester Violence}_{i,t-1}$, Perceived
421 $\text{State Violence}_{i,t-1}$, $\text{Police}_{i,t-1}$, and $\text{Fire}_{i,t-1}$. The violence measures are the average of the classifier

⁷Wikipedia provides the three sets of estimates.

⁸For exposition, these results are aggregated to the country-day level. A plot restricting analysis to Barcelona or Seoul shows the same trends.

Fig. 1. Verifying Dependent Variable, Cross-section

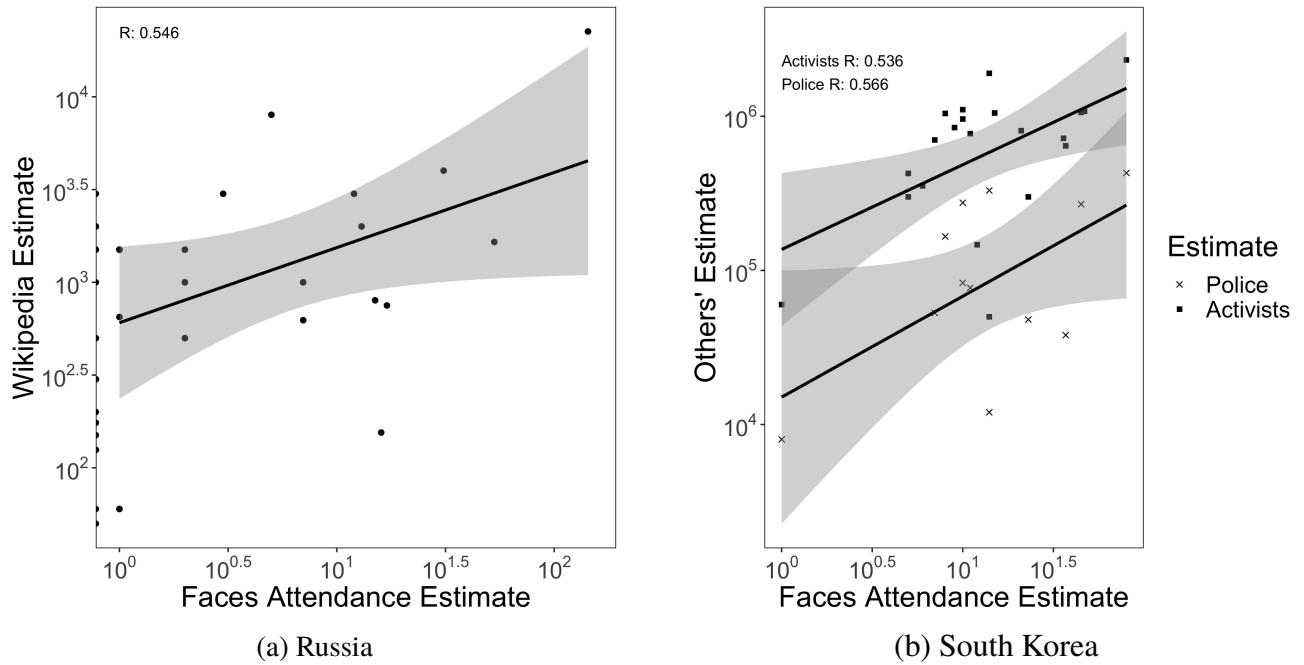
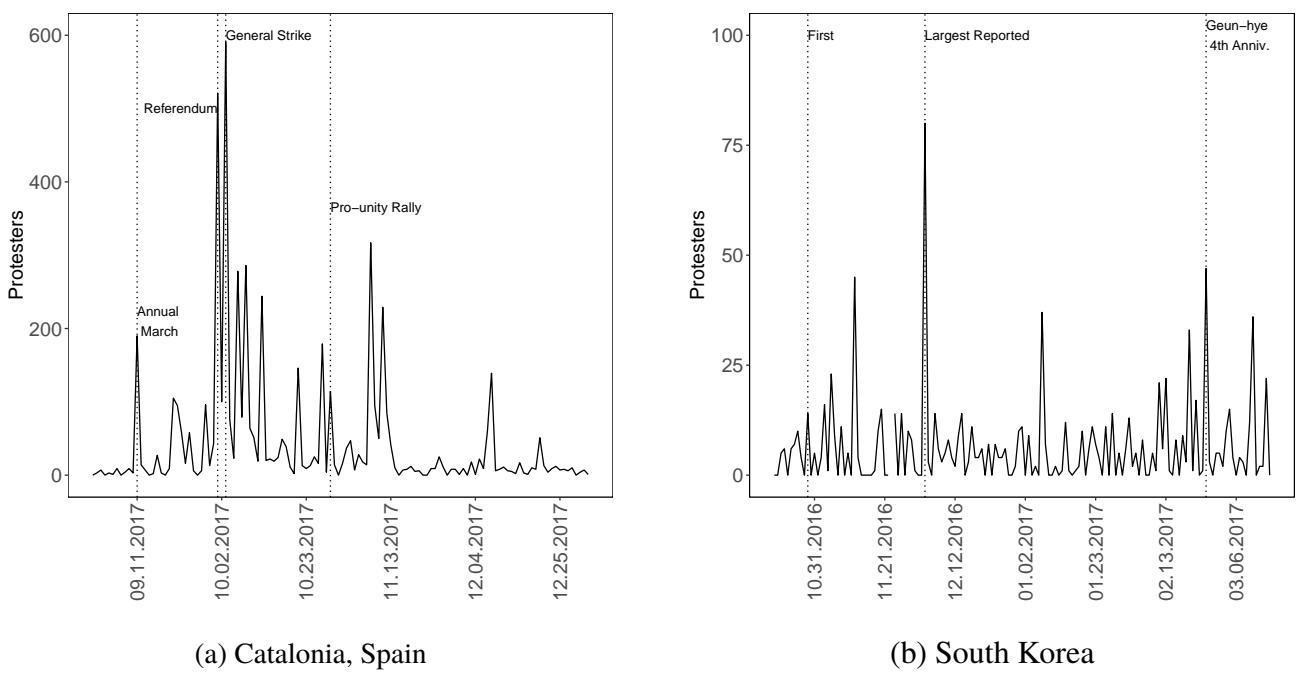


Fig. 2. Verifying Dependent Variable, Time Series



422 estimate for all protest photos per city-day. The police and fire variables are the sum of images
423 containing a police officer or fire, respectively, based on the thresholds identified in Table A9.

424 We describe the violence variables as “perceived” for three reasons. First, the true amount of
425 violence is unknown because violence is not a physical entity directly measurable, like temperature
426 or pressure. Second, the images people share may be strategically chosen. This possible selection
427 effect is true of any event data that relies on secondary sources, which is to say almost all event
428 data. For a longer discussion of bias that these measures may introduce, see Section 4.2. Third,
429 the main analysis does not deduplicate images, meaning images which are shared often will have a
430 greater impact on people’s decision making process than those only tweeted once. Deduplicating
431 images to more closely approximate the “true” violence at events does not change results, as Table
432 7 shows.

433 Students in democracies and autocracies often spearhead mass protests (Zhao, 1998; Gonzalez,
434 2019). The young are more likely to lack jobs, have little wealth to lose, and view protest
435 participation as its own end. These effects are amplified when there are many of them, a phenomenon
436 commonly called the “youth bulge” (Urdal, 2006). Knowing that youth often make protests more
437 intense (Hollander and Byun, 2015), states with large youth populations engage in more preventative
438 repression (Nordås and Davenport, 2013). The percent of participants aged 20-29 is therefore a
439 variable for which we control.

440 A society with greater gender equality is more likely to see nonviolent than violent action
441 (McCammon et al., 2001; Schaftenaar, 2017), and the same is true at the movement level (Asal
442 et al., 2013). Even when excluded from high-level leadership positions, women can play important
443 roles as bridges between that level and the broader movement (Robnett, 1996). Women were also
444 integral actors, as activists and participants, during the Arab Spring, a dynamic often overlooked in
445 accounts of those events (Newsom and Lengel, 2012; Rizzo et al., 2012). The percent of protesters
446 who are male is therefore a variable for which we control.

447 Figures 3 and 4 show the distribution of state violence and white faces, respectively, by country.
448 Because most photos record no violence, Figure 3 shows the distribution after dropping all photos

449 whose value for $Perceived\ State\ Violence_{i,t-1}$ is below the median; the relative order of states
 450 does not change if all images are kept. The average amount of violence matches expectations:
 451 Venezuela's protests are frequently met with violent repression, and Spanish police aggressively
 452 met protesters. Some violence was reported in Hong Kong as the protests neared their end. No
 453 violence was reported in South Korea, and Pakistani authorities let the anti-blasphemy protests run
 454 their course. The distribution of white faces also matches expectations: Catalonia and Venezuela
 455 record the highest percentages, in that order, while Hong Kong and South Korea record almost
 456 none. The race classifier performs less well on Pakistan. Section S7 shows similar charts for
 457 protester violence and the two other demographic variables.

Fig. 3. Distribution of State Violence by Country (vertical line is the mean)

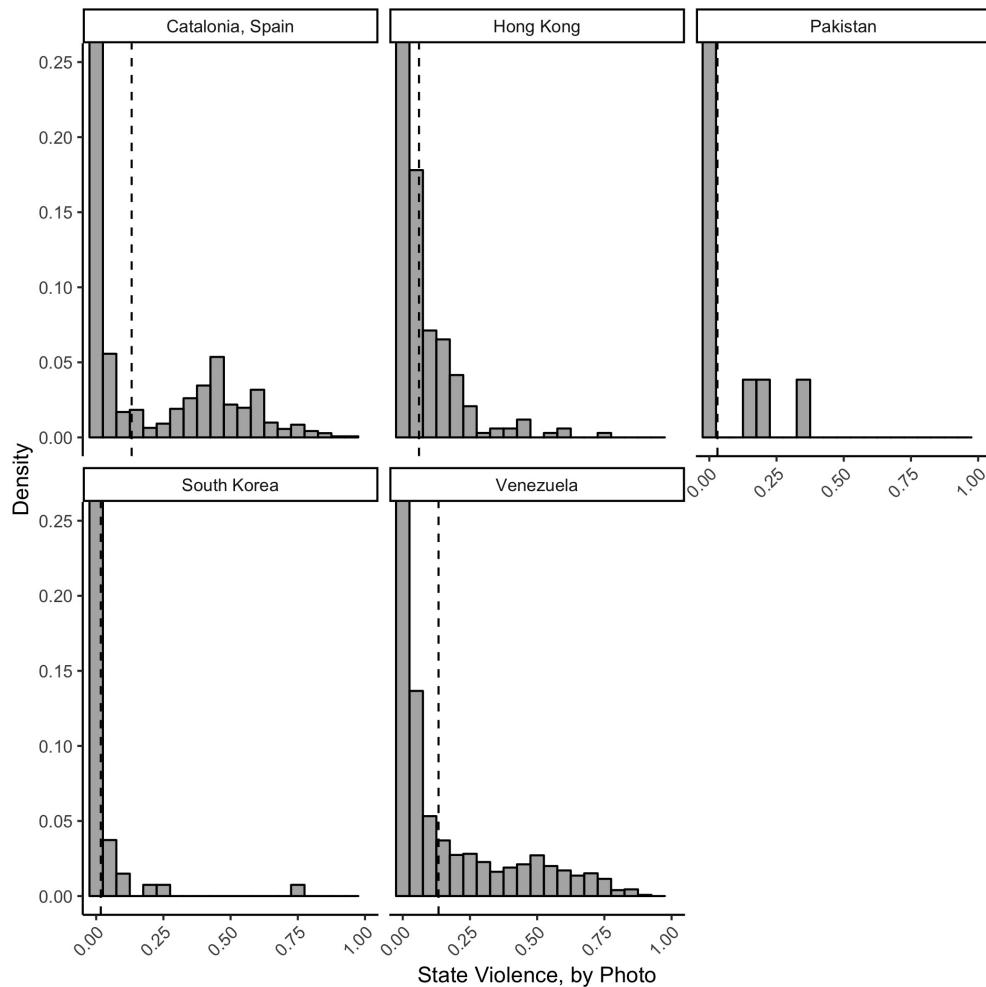
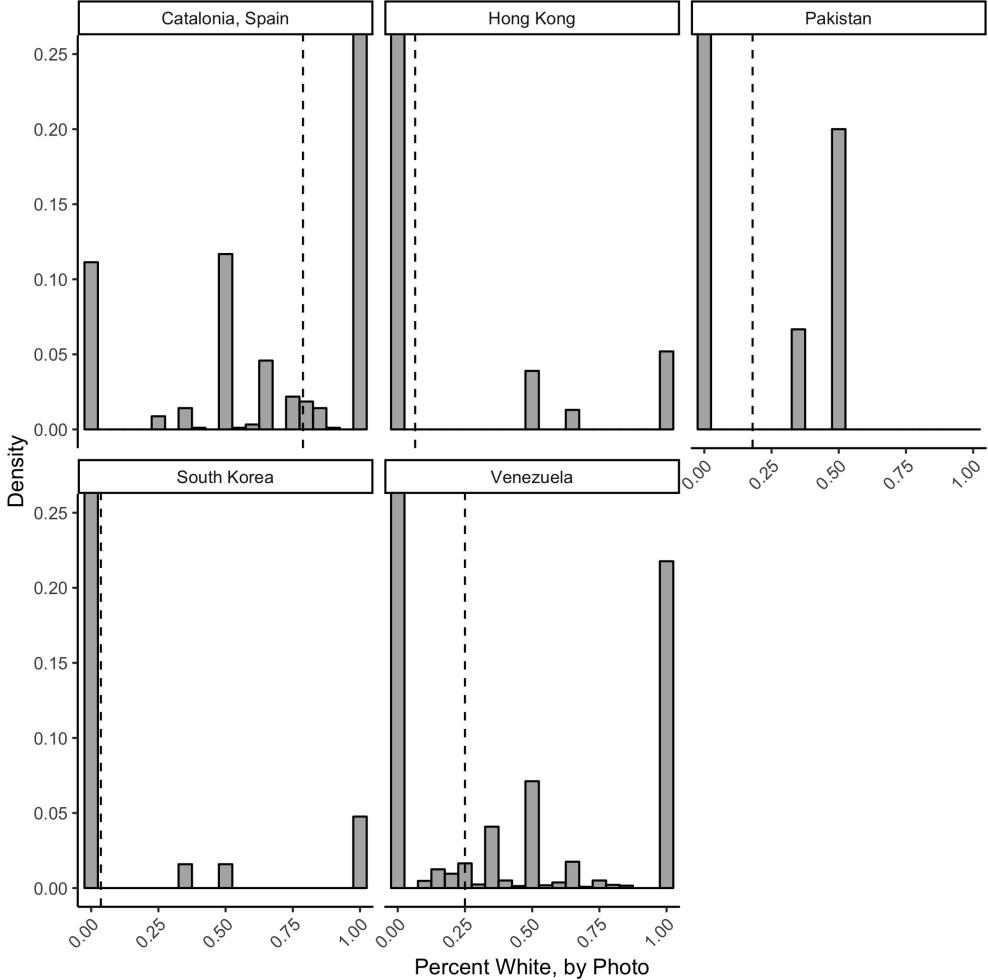


Fig. 4. Distribution of White Faces by Country (vertical line is the mean)



458 Table 3 provides descriptive statistics of these variables. Figure A9 shows the city-day correlation
 459 of these variables, and Figure A10 shows their tweet level correlation. The largest tweet level
 460 correlation, .56, is between $Race Diversity_{i,t-1}$ and $Age Diversity_{i,t-1}$. At the city-day level, the
 461 three diversity variables correlate between .7 and .78. All other correlations are below .32.

462 4.4 Model

463 In addition to the operationalizations detailed in the previous section, we include two control
 464 variables. $Tweets_{i,t-1}$ is the number of protest images per country-day and proxies for the amount
 465 of information available to protesters. This variable captures any effect general knowledge about
 466 a protest will have on protest size (Little, 2015). We also include a lagged dependent variable to

Table 3. Summary Statistics

Statistic	N	Mean	St. Dev.	Min.	Max.
<i>Protest Size</i> _{i,t}	4,401	2.42	14.87	0.00	627
<i>Perceived Protester Violence</i> _{i,t-1}	4,376	0.03	0.12	0.00	1.00
<i>Perceived State Violence</i> _{i,t-1}	4,376	0.02	0.07	0.00	0.94
<i>Police</i> _{i,t-1}	4,376	0.001	0.04	0.00	1.00
<i>Fire</i> _{i,t-1}	4,376	0.07	0.41	0.00	7.00
<i>Gender Diversity</i> _{i,t-1}	4,376	0.02	0.18	0.00	5.00
<i>Race Diversity</i> _{i,t-1}	4,376	0.02	0.11	0.00	0.69
<i>Age Diversity</i> _{i,t-1}	4,376	0.07	0.31	0.00	2.66

467 account for autocorrelation as well as any regression to the mean.

468 We build three models. The first uses only covariates that measure violence, testing H1. The
469 second focuses on the demographic variables. The final model combine the three sets of variables.
470 All independent variables are lagged one day. All models include city fixed effects and city-clustered
471 standard errors, though we run robustness checks with different fixed effects and clustering.

472 To facilitate interpretation, ordinary least squares is the estimator. Since the dependent variable
473 is a logarithm, the interpretation of a coefficient is the percent change in protest size as the result
474 of a one unit increase in the independent variable. Finally, to guard against overfitting, we use
475 five-fold cross-validation: each model is run on five different subsamples of the data and the results
476 are averaged.

477 **5 RESULTS**

478 **5.1 Main Results**

479 Our models most strongly confirm the expectations for protester and state violence. Racial
480 diversity has a signalling effect while gender diversity supports critical mass interpretations.

481 When protesters engage in violence, subsequent protest is smaller. Low amounts of state
482 violence correlate with larger subsequent protests, though severe enough violence will decrease the
483 size of protests. In addition, the more photos that show fire or police at protest, the more people
484 mobilize. Protester violence has a much smaller slope than either state violence variable, with the

485 largest effect occurring when states engage in high levels of violence.

486 We find no statistically significant correlation between racial or gender diversity and subsequent
 487 protest size. In other models, shown soon in the robustness section and in the Supplementary
 488 Materials, gender diversity attains statistical significance with a negative slope and racial diversity
 489 does the same in the opposite direction.

Table 4. Main Result

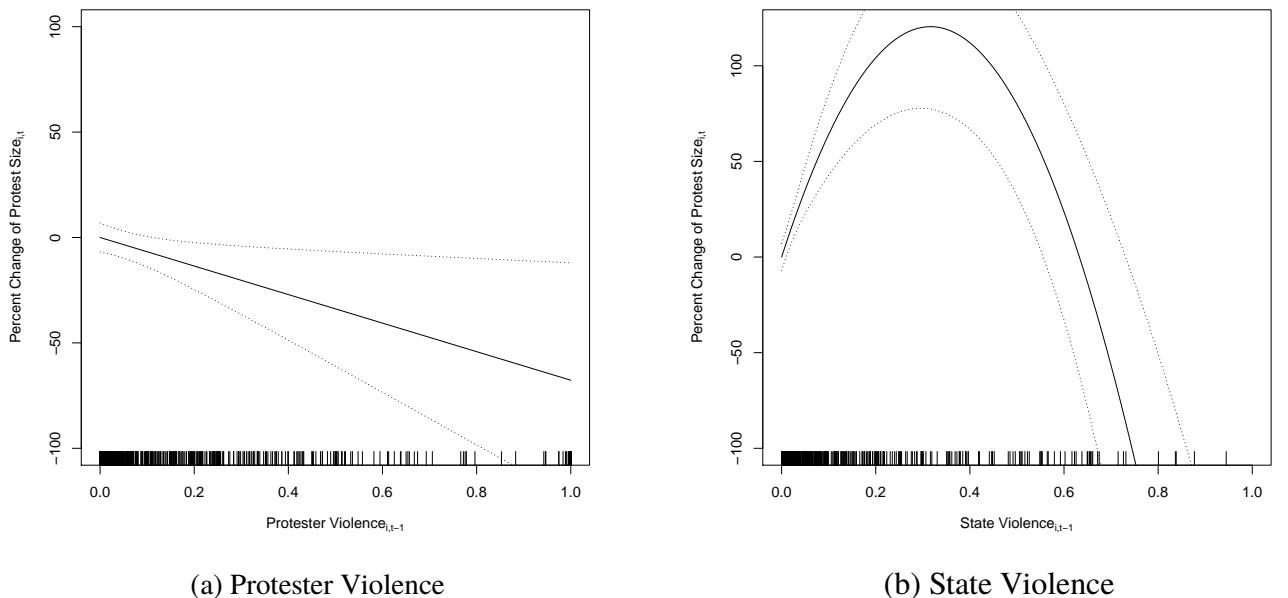
	DV: $\text{Log}_{10}(\text{Sum of Faces})_{i,t}$		
	Violence	Demographics	Combined
	(1)	(2)	(4)
Perceived Prtstr. Violence $_{i,t-1}$	-.1824** (.0714)		-.1674** (.0677)
Perceived Stt. Violence $_{i,t-1}$	1.2850** (.3152)		1.2820*** (.3327)
Perceived Stt. Violence $^2_{i,t-1}$	-2.0956*** (.5820)		-2.1030*** (.6093)
Police $_{i,t-1}$.7566* (.4568)		.7626* (.4493)
Fire $_{i,t-1}$.1099*** (.0203)		.1009*** (.0236)
Gender Diversity $_{i,t-1}$		-.1425 (.0922)	-.1126 (.0939)
Race Diversity $_{i,t-1}$.0955** (.0431)	.0683 (.0440)
Age Diversity $_{i,t-1}$.0233 (.0317)	.0203 (.0289)
Tweets $_{i,t-1}$.0093*** (.0034)	.0110*** (.0039)	.0095*** (.0033)
DV $_{i,t-1}$.1753** (.0722)	.1828** (.0736)	.1578** (.0682)
Intercept	.1227*** (.0173)	.1306*** (.0237)	.1260*** (.0237)
N	4,376	4,376	4,376
City FE	Y	Y	Y
Cluster SE	Y	Y	Y
Adjusted R ²	.2435	.2280	.2450

*p < .1; **p < .05; ***p < .01

City-clustered standard errors shown in parentheses.

490 Figure 5 shows marginal effects of protester and state violence. From values of [0-3), state
 491 violence increases protest, reaching a maximum at .3. At that amount of violence, protest size the
 492 next day is 137% higher than if there was no state violence. Moreover, state repression usually
 493 leads to larger protests: only 77 of 1,467 city-days of protest contain average state violence greater
 494 than .3. Protester violence, on the other hand, monotonically correlates with smaller protest. The
 495 change, however, is much smaller than for state violence: moving from no protester violence to
 496 its mean (.035) decreases protest size by just over 2%, while the difference between state violence
 497 and its mean is an increase of just over 17%. A one standard deviation increase of state violence
 498 from 0 increases protests size by approximately 63%; a one standard deviation increase in protester
 499 violence from the same point decreases protest size by just over 12%.

Fig. 5. Marginal Effects



(a) Protester Violence

(b) State Violence

500 **5.2 Robustness Checks**

501 Three sets of robustness checks confirms the results in Table 4. The first set introduces additional
502 fixed effects and does not cluster standard errors. The second set subsets the raw data to rule out
503 strategic behavior of Twitter users driving the results. The third set removes bots and duplicate
504 images to confirm that results are not due to malfeasance or virality. In all presented results, the
505 first model is the full model from Table 4.

506 Table 5 shows that the main findings are robust to alternate model specifications. A rule of
507 thumb is to not cluster standard errors when there are fewer than 30, and we have 23 ([Cameron et al., 2008; King and Roberts, 2015](#)). Model 2 therefore does not cluster standard errors. Without
508 clustering standard errors, $Gender\ Diversity_{i,t-1}$ and $Race\ Diversity_{i,t-1}$ are statistically signif-
509 icant. Model 3 includes a fixed effect for Saturdays and Sundays, the most popular protest days.
510 $Gender\ Diversity_{i,t-1}$ is once again statistically significant, while $Race\ Diversity_{i,t-1}$ is not. Model
511 4 includes a day of week fixed effect, since some countries (primarily Venezuela and Pakistan)
512 have larger protests outside of the weekend; results match Models 2 and 3. In case unobserved
513 country heterogeneity drives results, we include country fixed effects. Now, $Gender\ Diversity_{i,t-1}$
514 just barely loses statistical significance while $Race\ Diversity_{i,t-1}$ just barely obtains it. In none of
515 the extra checks does inference about perceived protester or state violence change.
516

517 The next set of robustness checks verify that neither strategic behavior nor data pollution drive
518 the results. Table 6 uses different subsets of the raw data to attempt to rule out strategic behavior
519 of individuals driving the results. The patterns hold.

520 One source of bias in newspaper based event data is that newspapers tailor their reporting for
521 their intended audience. The same could be true of Twitter users, and it is more likely to be true
522 the more likely they are to have an audience. Model 2 shows the result when keeping only tweets
523 from accounts that fall within the 25th-75th percentile of their country's follower distribution.
524 The coefficient on $Gender\ Diversity_{i,t-1}$ remains the same, but it is not statistically significant
525 because of a smaller standard error. Other results are not susceptible to keeping users based on
526 their follower count. Model 3 restricts the sample to only users who Twitter has not verified; for

Table 5. Robust to Alternate Specifications

	DV: $\text{Log}_{10}(\text{Sum of Faces})_{i,t}$				
	Original Model	Original Model	Weekend FE	Day-of-week FE	Country FE
	(1)	(2)	(3)	(4)	(5)
Perceived Prtstr. Violence $_{i,t-1}$	−.1674** (.0677)	−.1674*** (.0555)	−.1675*** (.0806)	−.1693** (.0807)	−.1579** (.0793)
Perceived Stt. Violence $_{i,t-1}$	1.2820*** (.3327)	1.2820*** (.2158)	1.2857*** (.3668)	1.3001*** (.3673)	1.3824*** (.3654)
Perceived Stt. Violence $^2_{i,t-1}$	−2.1030*** (.6093)	−2.1030*** (.3583)	−2.1042*** (.6493)	−2.1241*** (.6504)	−2.2609*** (.6353)
Police $_{i,t-1}$.7626* (.4493)	.7626*** (.1489)	.7688*** (.4627)	.7729* (.4613)	.7708* (.4436)
Fire $_{i,t-1}$.1009*** (.0236)	.1009*** (.0161)	.1012*** (.0375)	.1014*** (.0376)	.1093*** (.0384)
Gender Diversity $_{i,t-1}$	−.1126 (.0939)	−.1126*** (.0415)	−.1152*** (.0638)	−.1146* (.0636)	−.1023 (.0650)
Race Diversity $_{i,t-1}$.0683 (.0440)	.0683** (.0280)	.0712 (.0463)	.0718 (.0464)	.0912* (.0468)
Age Diversity $_{i,t-1}$.0203 (.0289)	.0203 (.0230)	.0213 (.0350)	.0206 (.0349)	.0127 (.0357)
Tweets $_{i,t-1}$.0095*** (.0033)	.0095*** (.0009)	.0095*** (.0033)	.0094*** (.0033)	.0093*** (.0033)
DV $_{i,t-1}$.1578*** (.0682)	.1578*** (.0237)	.1559*** (.0405)	.1564*** (.0405)	.1901*** (.0408)
Intercept	.1260*** (.0237)	.1260*** (.0158)	.1186*** (.0186)	.1277*** (.0216)	.0885*** (.0088)
N	4,376	4,376	4,376	4,376	4,376
Adjusted R ²	.2450	.2450	.2459	.2459	.2240
City FE	Y	Y	Y	Y	N
Weekend FE	N	N	Y	N	N
Weekday FE	N	N	N	Y	N
Clustered SE	Y	N	Y	Y	Y

*p < .1; **p < .05; ***p < .01

City-clustered standard errors shown in parentheses for Models 1, 3, and 4. Model 5 uses country-clustered standard errors.

public figures such as celebrities or politicians (and even some academics), Twitter verifies that the account actually belongs to the person it purports to. Since these users should be more likely to engage with Twitter strategically, we drop their tweets from analysis. This process takes away 447 tweets and 5 city-days, which is enough to reduce the coefficient of Police $_{i,t-1}$ below traditional thresholds of statistical significance; all other results match the full model.

The results when keeping tweets only in a country's *lingua franca*, shown in Model 4 of Table 6, are particularly interesting. Many Twitter users change their language depending on political context (Metzger et al., 2015), often as a method of attracting foreign audiences (Bruns et al.,

535 2013). For the violence and demographic variables except $Race Diversity_{i,t-1}$, the coefficients are
 536 much larger than the original model. In addition to violence now being estimated to have a stronger
 537 effect, $Gender Diversity_{i,t-1}$ is 80% larger and its standard error is halved, making it statistically
 538 significant. $Race Diversity_{i,t-1}$'s coefficient decreases, but its standard error decreases by even
 539 more, making it statistically significant as well. The most noticeable change is to $Age Diversity_{i,t-1}$,
 540 whose point estimate triples while its standard error remains the same. In only two other models,
 541 shown in the Supplementary Materials, is this variable statistically significant. Overall, restricting
 542 by language produces a model with a 22.5% better fit than the original. This better fit, larger
 543 coefficients, and more precise estimates of those coefficients suggests that language use may be one
 544 of the most common ways users behave strategically on social media.

Table 6. Robust to Strategic Behavior

	DV: $\log_{10}(\text{Sum of Faces})_{i,t}$			
	Original	Normal Users	No Verified Accounts	Country's Language
	(1)	(2)	(3)	(4)
Perceived Prtstr. Violence $_{i,t-1}$	-.1674*** (.0677)	-.1616*** (.0422)	-.1438** (.0670)	-.1899*** (.0425)
Perceived Stt. Violence $_{i,t-1}$	1.2820*** (.3327)	1.2564*** (.3347)	1.3146*** (.3597)	1.4723*** (.4686)
Perceived Stt. Violence $^2_{i,t-1}$	-2.1030*** (.6093)	-2.1084*** (.6093)	-2.0819*** (.6156)	-2.4132*** (.9056)
Police $_{i,t-1}$.7626* (.4493)	.8606** (.3737)	.6409 (.3930)	.4802 (.2929)
Fire $_{i,t-1}$.1009*** (.0236)	.0613* (.0331)	.0876*** (.0234)	.0664*** (.0252)
Gender Diversity $_{i,t-1}$	-.1126 (.0939)	-.1124* (.0662)	-.1121 (.0901)	-.1820*** (.0592)
Race Diversity $_{i,t-1}$.0683 (.0440)	.0316 (.0339)	.0675 (.0411)	.0504* (.0295)
Age Diversity $_{i,t-1}$.0203 (.0289)	.0083 (.0293)	.0149 (.0285)	.0773*** (.0282)
Tweets $_{i,t-1}$.0095*** (.0033)	.0156*** (.0060)	.0124** (.0048)	.0259*** (.0044)
DV $_{i,t-1}$.1578*** (.0682)	.1221** (.0562)	.1412** (.0685)	.0591 (.0706)
Intercept	.1260*** (.0237)	.1283*** (.0156)	.1185*** (.0155)	.0701*** (.0137)
N	4,376	3,715	4,371	3,614
Adjusted R ²	.2450	.1759	.2457	.3002
City FE	Y	Y	Y	Y
Cluster SE	Y	Y	Y	Y

*p < .1; **p < .05; ***p < .01

545 Table 7 presents two attempts to ensure that results are not driven by quirks in the data generation
546 process. The prevalence of bots - social media accounts controlled by computer code - has raised
547 concerns about the veracity of studies relying on social media data (Ferrara and Bessi, 2016).
548 Though other work has found few bots in geolocated tweets (Driscoll and Steinert-Threlkeld,
549 2018), we nonetheless submit every user to the Botometer service and remove tweets with a
550 complete automation probability $\geq .4$, the threshold which has been found to produce the most
551 accurate classification of bots (Varol et al., 2017). Model 2 presents these results, and findings do
552 not change. See Table A11 in the Supplementary Materials for the percent of accounts and tweets
553 that are from bots, by country; no more than 6.5% of tweets in any country are from bots.

554 To confirm that repetition of images do not drive results, we remove duplicate images. Models
555 3 and 4 from Table 7 shows these results. While our data do not contain retweets because Twitter
556 does not assign coordinates to retweets, they do contain replies, and replies contain the image
557 of the original tweet. (Section S8.2 details this methodology, and Table A12 show the percent
558 of tweets per city that are duplicates.) This process removes 2,920 images from the periods in
559 question, and Model 3 presents the results. Results for the violence and demographic variables do
560 not change. Model 4 weights these data by the number of protest tweets per city-day. In this model,
561 the coefficients for the the perceived violence variable are up to twice as large as the original model,
562 though the coefficients for the demographic variables shrink.⁹ Model 4 also produces the best fit
563 of any model we build. The violence and demographic conclusions are not affected by bots or the
564 reproduction of images: our model's original measurements do not appear to measure perception
565 as much as they do actual effects.

566 The Supplementary Materials present seven additional sets of robustness checks in Tables A13
567 through A19. The first set changes the operationalization of the dependent variable. The second
568 uses 15 lags of the dependent variable, as suggested by a partial autocorrelation plot. The third
569 set uses count models, and the fourth weights city-days by their number of tweets. The fifth set
570 increases the probability that tweets are from a protest by discarding those not from mobile devices

⁹Note as well that Race Diversity_{i,t-1} becomes negative, but with a small p-value, in the deduplicated models.

Table 7. Robust to Pollution

	Original	DV: $\text{Log}_{10}(\text{Sum of Faces})_{i,t}$	Deduplicate Images	Deduplicate Images, Weighted
	(1)	(2)	(3)	(4)
Perceived Prtstr. Violence $_{i,t-1}$	-.1674*** (.0677)	-.1593** (.0649)	-.1659* (.0921)	-.3778** (.1513)
Perceived Stt. Violence $_{i,t-1}$	1.2820*** (.3327)	1.2461*** (.3326)	1.2613** (.4913)	2.3127*** (.4773)
Perceived Stt. Violence $^2_{i,t-1}$	-2.1030*** (.6093)	-2.0535*** (.5879)	-2.0571*** (.8778)	-4.1393*** (.8678)
Police $_{i,t-1}$.7626* (.4493)	.6582* (.3806)	.9543* (.5359)	1.4324*** (.2220)
Fire $_{i,t-1}$.1009*** (.0236)	.0923*** (.0268)	.0735 (.0634)	.0178 (.0210)
Gender Diversity $_{i,t-1}$	-.1126 (.0939)	-.1231 (.0862)	.0023 (.0411)	-.0199 (.0618)
Race Diversity $_{i,t-1}$.0683 (.0440)	.0623 (.0390)	-.0115 (.0434)	-.0197 (.0393)
Age Diversity $_{i,t-1}$.0203 (.0289)	.0238 (.0301)	.0215 (.0240)	.0046 (.0341)
Tweets $_{i,t-1}$.0095*** (.0033)	.0127*** (.0048)	.0173** (.0073)	.0067*** (.0008)
DV $_{i,t-1}$.1578*** (.0682)	.1396** (.0641)	.1221 (.0874)	.1471*** (.0482)
Intercept	.1260*** (.0237)	.1214 (.0156)	.1328 (.0169)	.4891*** (.0586)
N	4,376	4,376	2,786	2,786
Adjusted R ²	.2450	.2508	.1912	.4479
City FE	Y	Y	Y	Y
Cluster SE	Y	Y	Y	Y

*p < .1; **p < .05; ***p < .01

571 or from non-protest hours. The sixth runs the full model separately by country, and the seventh
 572 investigates how Police $_{i,t-1}$ and Fire $_{i,t-1}$ correlate with the perceived violence measures.

573 6 DISCUSSION

574 6.1 Images and Measurement

575 Emphasizing the severity of repression during protest policing is not new (Muller, 1985;
 576 Khawaja, 1993); measuring repression as a continuous variable is. For example, the Social Conflict
 577 Analysis Database (SCAD), Urban Social Disorder, and Armed Conflict Locations and Event Data
 578 (ACLED) dataset record repression during an event as occurring or not (Raleigh et al., 2010;

579 [Salehyan et al., 2012](#); [Urdal and Hoelscher, 2012](#)). Repression is sometimes coded as ordinal or
580 nominal as well ([Goldstein, 1992](#); [Stephan and Chenoweth, 2008](#); [Clark and Regan, 2016](#)), and
581 machine-coded event data like the Integrated Conflict Early Warning System use this approach
582 ([Gerner et al., 2002](#); [Boschee et al., 2015](#)).

583 As far as we are aware, all previous approaches generate nominal or ordinal repression variables
584 from primary or secondary sources. This process, completely understandable given how violence
585 is recorded in texts, creates an implicit mapping of a latent quantity onto discrete categories. This
586 mapping is problematic because each researcher has its own mental model, so different studies
587 are likely to map the same latent quantity onto different discrete categories (values of the ordinal
588 variable). Measuring repression on a continuous scale may therefore provide a clearer understanding
589 of how it affects protest dynamics. It also facilitates the inclusion and interpretation of interaction
590 terms for violence, allowing us to test for nonlinear effects ([Moore, 1998](#); [Shellman et al., 2013](#)).

591 The results presented here suggest that measuring violence as a continuous variable may help
592 resolve the repression-dissent puzzle. Mapping violence into discrete bins may be especially
593 pernicious with panel data, explaining why those studies tend to find no correlation between
594 repression and protest. We avoid this pitfall by presenting human coders with over 10,000 pairs
595 of images to label and training a deep learning computer vision model on this training set; the
596 model outputs continuous estimates of protester and state violence, mitigating concerns that a
597 result for repression or protester violence is due to researcher effects. The results presented here
598 are continuous measurements based on primary sources.

599 Using images generated from social media also allows for more precise temporal measurement.
600 A difficulty testing protest dynamics is that action occurs on a timescale difficult to measure with
601 newspaper reports, the primary source of data for these types of studies ([Earl et al., 2004](#)). Most
602 research has therefore analyzed protest dynamics with coarse time scales such as weekly ([Lohmann,
603 1994a](#); [Rasler, 1996](#)) or, usually in the case of surveys, without a time component ([Opp and Gern,
604 1993](#); [Beissinger, 2013](#)). Recent research takes advantage of new datasets, including social media
605 data, to measure protest dynamics at a daily level ([Larson et al., 2016](#); [Ritter and Conrad, 2016](#);

606 Hsuan et al., 2017; Steinert-Threlkeld, 2017). Combining this high level of resolution with the
607 additional information that can be extracted from images has only been attempted twice before (Won
608 et al., 2017; Zhang and Pan, 2019), as far as we are aware, though there is work at scale analyzing
609 how the emotional content of images affects online mobilization (Casas and Webb Williams, 2018)
610 or how different news outlets portray Black Lives Matter protests (Torres, 2018).¹⁰

611 6.2 Ethics

612 The advances in scholarly understanding that the combination of computer vision and social
613 media enable also raises serious ethical concerns. We briefly discuss some and point the reader to
614 (Joo and Steinert-Threlkeld, 2018) for a longer analysis.

615 Like any measurement, the results are only as good as the input data. Many off the shelf
616 computer vision programs reproduce racial biases, and the leading datasets used to train race
617 classifiers have relatively small corpuses of images (Grush, 2015; Lam et al., 2018). The model we
618 use, FairFace, is less biased than other ones, however, because its training data were constructed
619 on racially balanced images whose quality more closely resembles social media photographs than
620 previous datasets'.

621 Treating race as a distinct category around which people may organize is itself problematic.
622 We simply note that in many countries, race encapsulates a multitude of historic power imbalances.
623 While we do not mean to reify race, ignoring it would also do a disservice to its importance in
624 many countries' politics.

625 Protesters may not be as anonymous as they think. Though these data are observational
626 and publicly available, individuals in photographs may not have consented to appear in those
627 photographs. While true of any images of public spaces, the concern is heightened when individuals
628 are engaged in risky behavior. Authorities could monitor images shared on social media to identify
629 people who protested, much as some do with cell phone location data (Davenport, 2014).¹¹ Foreign
630 governments and parts of the United States law enforcement already monitor faces in crowds

¹⁰See Cowart et al. (2016) for an example of manual protest image analysis.

¹¹The logic works in reverse. Shared protest images can be used to identify incriminating state behavior that would otherwise be denied (Lim, 2013).

(Purdy, 2018; Shaban, 2018). This concern about facial recognition also means that individuals who appear in photos but who did not take the photo may not realize they can be implicated in protest.¹² To prevent the identification of individuals in our data, we have chosen not to release the tweet identification number or image URL for the raw data.

6.3 Lower Bounds

This paper's findings are appealing because the precision and resolution of its measures allows for deeper theoretical understanding of protest dynamics. Moreover, the precision of these results should be considered a lower bound, as the number of protest images located to the city level is quite small. Researchers can increase the number of images available, and therefore analyze more events with more precision, using five tactics.

First, the easiest approach would be to accept a less rigorously defined measure of location, users' self-reported location. Twitter profiles contain a location field that individuals can populate with any phrase, e.g. "Los Angeles, CA" or "A Server Somewhere". Approximately 75% of Twitter users have information in this field, but only 8% of users (in the United States) have a string that Google can resolve to a specific location (Mislove et al., 2011). (Globally, $\frac{1}{3}$ of accounts have location or profile information in English (Leetaru et al., 2013).) In the United States, approximately 3.4% of accounts enable GPS coordinates, so using the location field at least doubles the number of images available (Sloan and Morgan, 2015). This increase will be more pronounced the less frequently a country's users geotag: only .3% and .9% of tweets in Korean and Arabic, respectively, contain GPS coordinates (Sloan and Morgan, 2015).

Second, one could purchase tweets from a vendor or download past tweets of users who tweet from protest events.

Third, newspapers and television channels maintain Twitter accounts and share the same stories there that are used in event datasets. While their accounts are not, in our experience, geolocated, it is theoretically feasible to incorporate the articles they share into an event data generating pipeline. Doing so would allow the researcher to determine the sensitivity of event records to source type,

¹²The UCLA IRB approved this study. We emphasize again that we only use publicly available data.

657 more precisely measure bias from sources, and determine if events recorded from traditional media
658 have different effects than those from social media.

659 Two other approaches would move beyond Twitter to collect images. The fourth approach
660 should look to other online platforms, especially Instagram, to collect images. Instagram provides
661 much less data through its application programming interface than Twitter does, so one will have to
662 crawl it. Crawling is a more technically difficult procedure and is actively discouraged. Instagram
663 tends to be used for apolitical postings as well. Flickr has been used to track protests, but it is not
664 a widely used platform (Alanyali et al., 2016). The fifth approach would be to partner with a news
665 images provider, such as The Associated Press or Getty Images.

666 **6.4 Conclusion**

667 We have presented results on a relatively small number of protests, and future work should
668 increase the number of protests analyzed. Doing so will rely on luck, as protests will have to occur
669 in countries use Twitter, or other platforms, heavily. Developing infrastructure to collect more
670 tweets with images will decrease the role luck plays in directing research.

671 Because images contain more information than text, they hold much promise for the study of
672 phenomena of interest to political scientists. For generating event data, images hold particular
673 promise in measuring magnitude, both in terms of crowd size and severity of an event, as well as
674 reducing bias from newspaper data (Sobolev et al., 2019). For explanations of how these data can
675 contribute to subfields like political behavior, communication, or international relations, see Joo
676 and Steinert-Threlkeld (2018). The techniques for generating useful data are similar and different to
677 text-as-data approaches, and this paper demonstrates one area in which computer vision techniques
678 benefit political scientists. Future work should expand on the data and variables introduced in this
679 paper.

680

REFERENCES

- 681 Alanyali, M., Preis, T., and Moat, H. S. (2016). “Tracking Protests Using Geotagged Flickr
682 Photographs.” *PLoS ONE*, 11(3), 27–30.
- 683 Asal, V., Legault, R., Szekely, O., and Wilkenfeld, J. (2013). “Gender ideologies and forms of
684 contentious mobilization in the Middle East.” *Journal of Peace Research*, 50(3), 305–318.
- 685 Baltrušaitis, T., Robinson, P., and Morency, L.-P. (2016). “Openface: an open source facial behavior
686 analysis toolkit.” *Applications of Computer Vision (WACV), 2016 IEEE Winter Conference on*,
687 IEEE, 1–10.
- 688 Bashir, N. Y., Lockwood, P., Chasteen, A. L., Nadolny, D., and Noyes, I. (2013). “The ironic impact
689 of activists: Negative stereotypes reduce social change influence.” *European Journal of Social
690 Psychology*, 43(7), 614–626.
- 691 Battaglini, M. (2017). “Public Protests and Policy Making.” *Quarterly Journal of Economics*,
692 132(1), 485–549.
- 693 Baum, M. A. and Zhukov, Y. M. (2015). “Filtering revolution: Reporting bias in international
694 newspaper coverage of the Libyan civil war.” *Journal of Peace Research*, 52(3), 384–400.
- 695 Baum, M. A. and Zhukov, Y. M. (2018). “Media Ownership and News Coverage of International
696 Conflict1.” *Political Communication*, 1–28.
- 697 Beissinger, M. R. (2013). “The Semblance of Democratic Revolution: Coalitions in Ukraine’s
698 Orange Revolution.” *American Political Science Review*, 107(03), 574–592.
- 699 Benford, R. D. and Snow, D. A. (2000). “Framing Processes and Social Movements: An Overview
700 and Assessment.” *Annual Review of Sociology*, 26, 611–639.
- 701 Bergamo, A. and Torresani, L. (2010). “Exploiting weakly-labeled web images to improve ob-
702 ject classification: a domain adaptation approach.” *Advances in neural information processing
703 systems*, 181–189.
- 704 Biggs, M. (2003). “Positive feedback in collective mobilization: The American strike wave of
705 1886.” *Theory and Society*, 32, 217–254.
- 706 Biggs, M. (2016). “Size Matters: Quantifying Protest by Counting Participants.” *Sociological
707 Methods & Research*, 1–33.
- 708 Bojarski, M., Del Testa, D., Dworakowski, D., Firner, B., Flepp, B., Goyal, P., Jackel, L. D.,
709 Monfort, M., Muller, U., Zhang, J., et al. (2016). “End to end learning for self-driving cars.”
710 *arXiv preprint arXiv:1604.07316*.
- 711 Boschee, E., Lautenschlager, J., O’Brien, S., Shellman, S., Starz, J., and Ward, M. (2015). “ICEWS
712 Coded Event Data, <<http://dx.doi.org/10.7910/DVN/28075>>.
- 713 Botta, F., Moat, H. S., and Preis, T. (2015). “Quantifying crowd size with mobile phone and Twitter
714 data.” *Royal Society Open Science*, 2, 150162.

- 715 Bradley, R. A. and Terry, M. E. (1952). "Rank analysis of incomplete block designs: I. the method
716 of paired comparisons." *Biometrika*, 39(3/4), 324–345.
- 717 Bratton, M. and Walle, N. V. D. (1992). "Popular Protest and Political Reform in Africa." *Comparative Politics*, 24(4), 419–442.
- 719 Bruns, A., Highfield, T., and Burgess, J. (2013). "The Arab Spring and Social Media Audiences:
720 English and Arabic Twitter Users and Their Networks." *American Behavioral Scientist*, 57(7),
721 871–898.
- 722 Bueno de Mesquita, B., Smith, A., Siverson, R. M., and Morrow, J. D. (2003). *The Logic of
723 Political Survival*. MIT Press, Cambridge.
- 724 Buolamwini, J. and Gebru, T. (2018). "Gender shades: Intersectional accuracy disparities in
725 commercial gender classification." *Conference on Fairness, Accountability and Transparency*,
726 77–91.
- 727 Cameron, A. C., Gelbach, J. B., and Miller, D. L. (2008). "Bootstrap-Based Improvements for
728 Inference with Clustered Errors." *Review of Economics and Statistics*, 90(3), 414–427.
- 729 Cantu, F. "The Fingerprints of Fraud: Evidence From Mexico's 1988 Presidential Election."
730 *American Political Science Review*, Forthcoming.
- 731 Carey, S. C. (2006). "The Dynamic Relationship Between Protest and Repression." *Political
732 Research Quarterly*, 59(1), 1–11.
- 733 Casas, A. and Webb Williams, N. (2018). "Images That Matter: Online Protests and the Mobilizing
734 Role of Pictures." *Political Research Quarterly*.
- 735 Celestino, M. R. and Gleditsch, K. S. (2013). "Fresh carnations or all thorn, no rose? Nonviolent
736 campaigns and transitions in autocracies." *Journal of Peace Research*, 50(3), 385–400.
- 737 Chenoweth, E. and Schock, K. (2015). "Do Contemporaneous Armed Challenges Affect the
738 Outcomes of Mass Nonviolent Campaigns?" *Mobilization: An International Quarterly*, 20(4),
739 427–451.
- 740 Chenoweth, E. and Stephan, M. J. (2011). *Why Civil Resistance Works*. Columbia University Press,
741 New York City.
- 742 Clark, D. H. and Regan, P. M. (2016). "Mass Mobilization,
743 <<https://www.binghamton.edu/massmobilization/about.html>>.
- 744 Cowart, H. S., Saunders, L. M., and Blackstone, G. E. (2016). "Picture a Protest: Analyzing Media
745 Images Tweeted From Ferguson." *Social Media and Society*, 2(4), 1–9.
- 746 Danneman, N. and Ritter, E. H. (2013). "Contagious Rebellion and Preemptive Repression."
747 *Journal of Conflict Resolution*, 58(2), 254–279.
- 748 Davenport, C. (2007). "State Repression and Political Order." *Annual Review of Political Science*,
749 10(1), 1–23.

- 750 Davenport, C. (2014). “Old Wine in an E-bottle (or, the Text that Mistook Itself for
751 a Tactical Shift), <<http://politicalviolenceataglance.org/2014/01/28/old-wine-in-an-e-bottle-or-the-text-that-mistook-itself-for-a-tactical-shift/>>.
- 752
- 753 Davenport, C. and Armstrong II, D. A. (2004). “Democracy and the Violation of Human Rights: A
754 Statistical Analysis from 1976 to 1996.” *American Journal of Political Science*, 48(3), 538–554.
- 755 Davenport, C. and Soule, S. A. (2009). “Velvet Glove, Iron Fist or Even Hand? Protest Policing in
756 the United States, 1960-1990.” *Mobilization*, 14(1), 1–22.
- 757 D. Della Porta and H. R. Reiter, eds. (1998). *Policing protest: The control of mass demonstrations
758 in Western democracies*. University of Minnesota Press, Minneapolis.
- 759 DeNardo, J. (1985). *Power in Numbers: The Political Strategy of Protest and Rebellion*. Princeton
760 University Press, Princeton.
- 761 Deng, J., Dong, W., Socher, R., Li, L.-j., Li, K., and Fei-fei, L. (2009). “ImageNet : A Large-Scale
762 Hierarchical Image Database.” *IEEE Conference on Computer Vision and Pattern Recognition*,
763 248–255.
- 764 Downs, A. (1957). *An Economic Theory of Democracy*. Harper and Row, New York City.
- 765 Driscoll, J. and Steinert-Threlkeld, Z. C. (2018). “Does Social Media Enable Irredentist Information
766 Warfare?” *Working paper*.
- 767 Earl, J. (2003). “Tanks, Tear Gas, and Taxes : Toward a Theory of Movement Repression.”
768 *Sociological Theory*, 21(1), 44–68.
- 769 Earl, J., Martin, A., McCarthy, J. D., and Soule, S. A. (2004). “The Use of Newspaper Data in the
770 Study of Collective Action.” *Annual Review of Sociology*, 30, 65–80.
- 771 Earl, J., McKee Hurwitz, H., Mejia Mesinas, A., Tolan, M., and Arlotti, A. (2013). “This
772 Protest Will Be Tweeted: Twitter and protest policing during the Pittsburgh G20.” *Information,
773 Communication & Society*, 16(4), 459–478.
- 774 Fassiotto, M. and Soule, S. A. (2017). “Loud and Clear: the Effect of Protest Signals on Congress-
775 sional Attention.” *Mobilization: An International Quarterly*, 22(1), 17–38.
- 776 Feinberga, M., Willer, R., and Kovacheff, C. (2017). “Extreme Protest Tactics Reduce Popular
777 Support for Social Movements.
- 778 Ferrara, E. and Bessi, A. (2016). “Social bots distort the 2016 U.S. Presidential election online
779 discussion.” *First Monday*, 21(11), 1–17.
- 780 Fisher, D. R., Dow, D. M., and Ray, R. (2017). “Intersectionality takes it to the streets : Mobilizing
781 across diverse interests for the Women’ s March.” *Science Advances*, 3, 1–8.
- 782 Francisco, R. A. (1995). “The Relationship between Coercion and Protest: An Empirical Evaluation
783 in Three Coercive States.” *Journal of Conflict Resolution*, 39(2), 263–282.

- 784 Francisco, R. A. (1996). "Coercion and Protest: An Empirical Test in Two Democratic States."
785 *American Journal of Political Science*, 40(4), 1179–1204.
- 786 Francisco, R. A. (2004). "After the Massacre: Mobilization in the Wake of Harsh Repression."
787 *Mobilization: An International Journal*, 9(2), 107–126.
- 788 Gause, L. (2018). "The Advantage of Disadvantage: Legislative Responsiveness to Collective
789 Action by the Politically Marginalized.
- 790 Gerner, D. J., Schrodt, P. A., Abu-Jabr, R., and Yilmaz, O. (2002). "Conflict and Mediation Event
791 Observations (CAMEO): A New Event Data Framework for the Analysis of Foreign Policy
792 Interactions." *Annual Meeting of the International Studies Association*.
- 793 Girshick, R. (2015). "Fast r-cnn." *Proceedings of the IEEE international conference on computer
794 vision*, 1440–1448.
- 795 Girshick, R., Donahue, J., Darrell, T., and Malik, J. (2014). "Rich feature hierarchies for accurate
796 object detection and semantic segmentation." *Proceedings of the IEEE conference on computer
797 vision and pattern recognition*, 580–587.
- 798 Goldstein, J. S. (1992). "A Conflict-Cooperation Scale for WEIS Events Data." *Journal of Conflict
799 Resolution*, 36(2), 369–385.
- 800 Goldstone, J. A. (2001). "Toward a Fourth Generation of Revolutionary Theory." *Annual Review
801 of Political Science*, 4, 139–187.
- 802 Gonzalez, F. (2019). "Collective action in networks: Evidence from the Chilean student movement."
803 *Working paper*.
- 804 Gruber, D. A. (1996). "Say It with Pictures." *The ANNALS of the American Academy of Political
805 and Social Science*, 546(1), 85–96.
- 806 Grush, L. (2015). "Google engineer apologizes after Photos app tags two black people as gorillas."
807 *The Verge*, <<https://www.theverge.com/2015/7/1/8880363/google-apologizes-photos-app-tags-two-black-people-gorillas>>.
- 808 Güler, R. A., Neverova, N., and Kokkinos, I. (2018). "Densepose: Dense human pose estimation
809 in the wild." *arXiv preprint arXiv:1802.00434*.
- 810 Gunitsky, S. (2015). "Corrupting the Cyber-Commons: Social Media as a Tool of Autocratic
811 Stability." *Perspectives on Politics*, 13(01), 42–54.
- 812 Gupta, D. K., Singh, H., and Sprague, T. (1993). "Government Coercion of Dissidents: Deterrence
813 or Provocation?." *Journal of Conflict Resolution*, 37(2), 301–339.
- 814 Gurr, T. (1970). *Why Men Rebel*. Princeton University Press, Princeton.
- 815 Gurr, T. R. and Moore, W. H. (1997). "Ethnopolitical Rebellion: A Cross-Sectional Analysis of
816 the 1980s with Risk Assessments for the 1990s." *American Journal of Political Science*, 41(4),
817 1079–1103.

- 819 He, K., Zhang, X., Ren, S., and Shun, J. (2016a). “Deep Residual Learning for Image Recognition.”
820 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 770–778.
- 821 He, K., Zhang, X., Ren, S., and Sun, J. (2016b). “Deep residual learning for image recognition.”
822 *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- 823 Heaney, M. T. and Rojas, F. (2008). “Coalition Dissolution, Mobilization, and Network Dynamics
824 in the U.S. Antiwar Movement.” *Research in Social Movements, Conflicts, and Change*, 28(08),
825 39–82.
- 826 Hellmeier, S., Weidmann, N. B., and Geelmuyden Rød, E. (2018). “In The Spotlight:Analyzing
827 Sequential Attention Effects in Protest Reporting.” *Political Communication*, 00(00), 1–25.
- 828 Hess, D. and Martin, B. (2006). “Repression, Backfire, and the Theory of Transformative Events.”
829 *Mobilization: An International Journal*, 11(2), 249–267.
- 830 Hollander, E. J. and Byun, C. C. (2015). “Explaining the Intensity of the Arab Spring.” *Digest of*
831 *Middle East Studies*, 24(1), 26–46.
- 832 Hsuan, T., Chen, Y., Zachary, P., and Fariss, C. J. (2017). “Who Protests? Using Social Media
833 Data to Estimate How Social Context Affects Political Behavior.
- 834 Huval, B., Wang, T., Tandon, S., Kiske, J., Song, W., Pazhayampallil, J., Andriluka, M., Rajpurkar,
835 P., Migimatsu, T., Cheng-Yue, R., et al. (2015). “An empirical evaluation of deep learning on
836 highway driving.” *arXiv preprint arXiv:1504.01716*.
- 837 Joo, J., Li, W., Steen, F. F., and Zhu, S.-C. (2014). “Visual persuasion: Inferring communicative
838 intents of images.” *Proceedings of the IEEE conference on computer vision and pattern*
839 *recognition*, 216–223.
- 840 Joo, J., Steen, F. F., and Zhu, S.-C. (2015). “Automated facial trait judgment and election outcome
841 prediction: Social dimensions of face.” *Proceedings of the IEEE international conference on*
842 *computer vision*, 3712–3720.
- 843 Joo, J. and Steinert-Threlkeld, Z. C. (2018). “Image as Data: Automated Visual Content Analysis
844 for Political Science.” *Working paper*.
- 845 Kalyvas, S. N. (2004). “The Urban Bias in Research on Civil Wars.” *Security Studies*, 13(3),
846 160–190.
- 847 Kärkkäinen, K. and Joo, J. (2019). “Fairface: Face attribute dataset for balanced race, gender, and
848 age.” *arXiv preprint arXiv:1908.04913*.
- 849 Kern, H. L. (2011). “Foreign Media and Protest Diffusion in Authoritarian Regimes: The Case of
850 the 1989 East German Revolution.” *Comparative Political Studies*, 44(9), 1179–1205.
- 851 Khawaja, M. (1993). “Repression and Popular Collective Action: Evidence from the West Bank.”
852 *Sociological Forum*, 8(1), 47–71.

- 853 King, G. and Roberts, M. E. (2015). “How Robust Standard Errors Expose Methodological
854 Problems They Do Not Fix, and What to Do About It.” *Political Analysis*, 23(2), 159–179.
- 855 Koopmans, R. (1993). “The Dynamics of Protest Waves: West Germany, 1965 to 1989.” *American
856 Sociological Review*, 58(5), 637–658.
- 857 Kovashka, A., Parikh, D., and Grauman, K. (2012). “Whittlesearch: Image search with relative
858 attribute feedback.” *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference
859 on*, IEEE, 2973–2980.
- 860 Kuran, T. (1989). “Sparks and Prairie Fires: A Theory of Unanticipated Political Revolution.”
861 *Public Choice*, 61(1), 41–74.
- 862 Lam, O., Wojcik, S., Broderick, B., and Hughes, A. (2018). “Gender and Jobs in Online Image
863 Searches.” *Report no.*, Pew Research Center.
- 864 Larson, J. M., Nagler, J., Ronen, J., and Tucker, J. A. (2016). “Social Networks and Protest
865 Participation: Evidence from 130 Million Twitter Users.
- 866 Lawrence, A. K. (2016). “Repression and Activism among the Arab Spring’s First Movers:
867 Evidence from Morocco’s February 20th Movement.” *British Journal of Political Science*, (May),
868 1–20.
- 869 Leetaru, K. H. (2014). “Fulltext Geocoding Versus Spatial Metadata for Large Text Archives:
870 Towards a Geographically Enriched Wikipedia.” *D-Lib Magazine*, 18(9), 1–16.
- 871 Leetaru, K. H., Wang, S., Cao, G., Padmanabhan, A., and Shook, E. (2013). “Mapping the global
872 Twitter heartbeat: The geography of Twitter.” *First Monday*, 18(5-6), 1–33.
- 873 Lim, M. (2013). “Framing Bouazizi: ‘White lies’, hybrid network, and collective/connective action
874 in the 2010-11 Tunisian uprising.” *Journalism*, 14(7), 921–941.
- 875 Little, A. T. (2015). “Communication Technology and Protest.” *Journal of Politics*, 78(1), 152–166.
- 876 Liu, Z., Luo, P., Wang, X., and Tang, X. (2015). “Deep learning face attributes in the wild.”
877 *Proceedings of the IEEE International Conference on Computer Vision*, 3730–3738.
- 878 Lohmann, S. (1993). “A Signaling Model of Informative and Manipulative Political Action.”
879 *American Political Science Review*, 87(2), 319–333.
- 880 Lohmann, S. (1994a). “The Dynamics of Informational Cascades: The Monday Demonstrations
881 in Leipzig, East Germany 1989-91.” *World Politics*, 47(1).
- 882 Lohmann, S. (1994b). “The Dynamics of Informational Cascades: The Monday Demonstrations
883 in Leipzig, East Germany, 1989-91.” *World Politics*, 47(1), 42–101.
- 884 Malik, M. M., Lamba, H., Nakos, C., and Pfeffer, J. (2015). “Population Bias in Geotagged Tweets.”
885 *9th International AAAI Conference on Weblogs and Social Media*, 18–27.

- 886 McCammon, H. J., Campbell, K. E., Granberg, E. M., and Mowery, C. (2001). "How Movements
887 Win: Gendered Opportunity Structures and U.S. Women's Suffrage." *American Sociological
888 Review*, 66(1), 49–70.
- 889 McCarthy, J. D., McPhail, C., and Smith, J. (1996). "Images of Protest: Dimensions of Selection
890 Bias in Media Coverage of Washington Demonstrations." *American Sociological Review*, 61(3),
891 478–499.
- 892 Mellon, J. and Prosser, C. (2017). "Twitter and Facebook are not representative of the general
893 population: Political attitudes and demographics of British social media users." *Research &
894 Politics*, 4(3), 205316801772000.
- 895 Metzger, M., Nagler, J., and Tucker, J. a. (2015). "Tweeting Identity? Ukrainian, Russian, and
896 #Euromaidan." *Journal of Comparative Economics*, 44(1), 16–40.
- 897 Mislove, A., Lehmann, S., Ahn, Y.-Y., Onnela, J.-P., and Rosenquist, J. N. (2011). "Understanding
898 the Demographics of Twitter Users." *Proceedings of the Fifth International AAI Conference on
899 the Weblogs and Social Media*, 554–557.
- 900 Moore, W. H. (1995). "Rational Rebels: Overcoming the Free-Rider Problem." *Political Research
901 Quarterly*, 48(2), 417–454.
- 902 Moore, W. H. (1998). "Repression and Dissent: Substitution, Context, and Timing." *American
903 Journal of Political Science*, 42(3), 851–873.
- 904 Moore, W. H. (2000). "The Repression of Dissent: A Substitution Model of Government Coercion."
905 *Journal of Conflict Resolution*, 44(1), 107–127.
- 906 Morstatter, F., Pfeffer, J., Carley, K. M., and Liu, H. (2013). "Is the Sample Good Enough?
907 Comparing Data from Twitter's Streaming API with Twitter's Firehose." *Association for the
908 Advancement of Artificial Intelligence*.
- 909 Muller, E. N. (1985). "Income Inequality, Regime Repressiveness, and Political Violence." *Ameri-
910 can Sociological Review*, 50(1), 47–61.
- 911 Murdie, A. and Purser, C. (2017). "How protest affects opinions of peaceful demonstration and
912 expression rights." *Journal of Human Rights*, 16(3), 351–369.
- 913 Myers, D. J. and Caniglia, B. S. (2004). "All the Rioting That's Fit to Print: Selection Effects in
914 National Newspaper Coverage of Civil Disorders, 1968–1969." *American Sociological Review*,
915 69, 519–543.
- 916 Newsom, V. A. and Lengel, L. (2012). "Arab Women, Social Media, and the Arab Spring:
917 Applying the framework of digital reflexivity to analyze gender and online activism." *Journal of
918 International Women's Studies*, 13(5), 31–45.
- 919 Nordås, R. and Davenport, C. (2013). "Fight the Youth: Youth Bulges and State Repression."
920 *American Journal of Political Science*, 57(4), 926–940.

- 921 Olzak, S., Beasley, M., and Olivier, J. (2003). “The Impact of State Reforms on Protest Against
922 Apartheid in South Africa.” *Mobilization*, 8(1), 27–50.
- 923 Opp, K.-D. and Gern, C. (1993). “Dissident Groups, Personal Networks, and Spontaneous Coop-
924 eration: The East German Revolution of 1989.” *American Sociological Review*, 58(5), 659–680.
- 925 Parkhi, O. M., Vedaldi, A., Zisserman, A., et al. (2015). “Deep face recognition..” *BMVC*, Vol. 1,
926 6.
- 927 Pearlman, W. (2013). “Emotions and the Microfoundations of the Arab Uprisings.” *Perspectives
928 on Politics*, 11(02), 387–409.
- 929 Pfeffer, J. and Mayer, K. (2018). “Tampering with Twitter’s Sample API.” *EPJ Data Science*, 7(50),
930 1–21.
- 931 Purdy, C. (2018). “China is launching a dystopian program to monitor citizens in Beijing,
932 <<https://qz.com/1473966/china-is-starting-a-big-brother-monitoring-program-in-beijing/>>.
- 933 Qin, B., Strömberg, D., and Wu, Y. (2017). “Why Does China Allow Free Social Media? Protests
934 Versus Surveillance and Propaganda.” *Journal of Economic Perspectives*, 31(1), 117–140.
- 935 Raleigh, C., Linke, A., Hegre, H., and Karlsen, J. (2010). “Introducing ACLED: An Armed Conflict
936 Location and Event Dataset: Special Data Feature.” *Journal of Peace Research*, 47(5), 651–660.
- 937 Rasler, K. (1996). “Concessions, Repression, and Political Protest in the Iranian Revolution.”
938 *American Sociological Review*, 61(1), 132–152.
- 939 Redmon, J., Divvala, S., Girshick, R., and Farhadi, A. (2016). “You only look once: Unified,
940 real-time object detection.” *Proceedings of the IEEE conference on computer vision and pattern
941 recognition*, 779–788.
- 942 Ren, S., He, K., Girshick, R., and Sun, J. (2015). “Faster r-cnn: Towards real-time object detection
943 with region proposal networks.” *Advances in neural information processing systems*, 91–99.
- 944 Ritter, E. H. (2013). “Policy Disputes, Political Survival, and the Onset and Severity of State
945 Repression.” *Journal of Conflict Resolution*, 57(1), 1–26.
- 946 Ritter, E. H. and Conrad, C. R. (2016). “Preventing and Responding to Dissent: The Observational
947 Challenges of Explaining Strategic Repression.” *American Political Science Review*, 110(1),
948 85–99.
- 949 Rizzo, H., Price, A. M., and Meyer, K. (2012). “Targeting Cultural Change in Repressive Environ-
950 ments: The Campaign against Sexual Harassment in Egypt.” *Report No. 614*, Egyptian Center
951 for Women’s Rights, Cairo, <<http://ecwronline.org/?p=1579>>.
- 952 Robertson, G. B. (2007). “Strikes and Labor Organization in Hybrid Regimes.” *American Political
953 Science Review*, 101(04), 781–798.
- 954 Robnett, B. (1996). “African-American Women in the Civil Rights Movement, 1954-1965: Gender,
955 Leadership, and Micromobilization.” *American Journal of Sociology*, 101(6), 1661–1693.

- 956 Rosenblatt, F. (1958). "The perceptron: a probabilistic model for information storage and organi-
957 zation in the brain.." *Psychological review*, 65(6), 386.
- 958 Salehyan, I., Hendrix, C., Hammer, J., Case, C., Linebarger, C., Stull, E., and Williams, J. (2012).
959 "Social Conflict in Africa: A New Database." *International Interactions*, 38(4), 503–511.
- 960 Schaftenaar, S. (2017). "How (wo)men rebel: Exploring the effect of gender equality on nonviolent
961 and armed conflict onset." *Journal of Peace Research*, 54(6), 762–776.
- 962 Schroff, F., Kalenichenko, D., and Philbin, J. (2015). "Facenet: A unified embedding for face
963 recognition and clustering." *Proceedings of the IEEE conference on computer vision and pattern
964 recognition*, 815–823.
- 965 Schweingruber, D. and McPhail, C. (1999). "A Method for Systematically Observing and Recording
966 Collective Action." *Sociological Methods & Research*, 27(4), 451–498.
- 967 Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. (2016). "Grad-
968 cam: Visual explanations from deep networks via gradient-based localization." See <https://arxiv.org/abs/1610.02391 v3>, 7(8).
- 970 Shaban, H. (2018). "Amazon employees demand company cut ties with ICE,
971 <[https://www.washingtonpost.com/news/the-switch/wp/2018/06/22/amazon-employees-
972 demand-company-cut-ties-with-ice>">demand-company-cut-ties-with-ice>](https://www.washingtonpost.com/news/the-switch/wp/2018/06/22/amazon-employees-demand-company-cut-ties-with-ice/) (jun).
- 973 Shakir, O. (2014). "All According to Plan: The Rab'a Massacre and Mass Killings of Protesters in
974 Egypt." *Report no.*, Human Rights Watch.
- 975 Shellman, S. M., Levey, B. P., and Young, J. K. (2013). "Shifting sands: Explaining and predicting
976 phase shifts by dissident organizations." *Journal of Peace Research*, 50(3), 319–336.
- 977 Siegel, D. A. (2011). "When Does Repression Work? Collective Action in Social Networks." *The
978 Journal of Politics*, 73(04), 993–1010.
- 979 Sloan, L. and Morgan, J. (2015). "Who tweets with their location? Understanding the relationship
980 between demographic characteristics and the use of geoservices and geotagging on twitter." *PLoS
981 ONE*, 10(11), 1–15.
- 982 Sloan, L., Morgan, J., Housley, W., Williams, M., Edwards, A., Burnap, P., and Rana, O. (2013).
983 "Knowing the Tweeters: Deriving sociologically relevant demographics from Twitter." *Socio-
984 logical Research Online*, 18(3), 1–15.
- 985 Sobolev, A., Joo, J., Chen, K., and Steinert-Threlkeld, Z. C. (2019). "Newspapers and Social Media
986 Accurately Measure Protest Size." *Working paper*.
- 987 Steinert-Threlkeld, Z. C. (2017). "Spontaneous Collective Action: Peripheral Mobilization During
988 the Arab Spring." *American Political Science Review*, 111(02), 379–403.
- 989 Steinert-Threlkeld, Z. C. (2018). *Twitter as Data*. Cambridge University Press.

- 990 Stephan, M. J. and Chenoweth, E. (2008). "Why Civil Resistance Works." *International Security*,
991 33(1), 7–44.
- 992 Sullivan, C. M. (2016). "Political Repression and the Destruction of Dissident Organizations."
993 *World Politics*, 68(4), 645–676.
- 994 Sun, Y., Chen, Y., Wang, X., and Tang, X. (2014). "Deep learning face representation by joint
995 identification-verification." *Advances in neural information processing systems*, 1988–1996.
- 996 Sutton, J., Butcher, C. R., and Svensson, I. (2014). "Explaining political jiu-jitsu: Institution-
997 building and the outcomes of regime violence against unarmed protests." *Journal of Peace
998 Research*, 51(5), 559–573.
- 999 Thomee, B., Shamma, D. A., Friedland, G., Elizalde, B., Ni, K., Poland, D., Borth, D., and Li, L.-J.
1000 (2016). "YFCC100M: The New Data in Multimedia Research." *Communications of the ACM*,
1001 64–73.
- 1002 Tilly, C. and Wood, L. J. (2012). *Social Movements, 1768-2012*. Paradigm Publishers, 3rd edition.
- 1003 Torres, M. (2018). "Give me the full picture: Using computer vision to understand visual frames
1004 and political communication." Working paper.
- 1005 Tucker, J. A. (2007). "Enough! Electoral Fraud, Collective Action Problems, and Post-Communist
1006 Colored Revolutions." *Perspectives on Politics*, 5(03), 535.
- 1007 Tufekci, Z. (2014). "Big Questions for Social Media Big Data: Representativeness, Validity and
1008 Other Methodological Pitfalls Pre-print." *Proceedings of the 8th International AAAI Conference
1009 on Weblogs and Social Media*, Ann Arbor.
- 1010 Urdal, H. (2006). "A Clash of Generations? Youth Bulges and Political Violence." *International
1011 Studies Quarterly*, 50, 607–629.
- 1012 Urdal, H. and Hoelscher, K. (2012). "Explaining Urban Social Disorder and Violence: An Empirical
1013 Study of Event Data from Asian and Sub-Saharan African Cities." *International Interactions*,
1014 38(4), 512–528.
- 1015 Varol, O., Ferrara, E., Davis, C. A., Menczer, F., and Flammini, A. (2017). "Online Human-Bot
1016 Interactions: Detection, Estimation, and Characterization." Working paper.
- 1017 Weidmann, N. B. (2014). "On the Accuracy of Media-based Conflict Event Data." *Journal of
1018 Conflict Resolution*, 59(6), 1129–1149.
- 1019 Weyland, K. (2012). "The Arab Spring: Why the Surprising Similarities with the Revolutionary
1020 Wave of 1848?." *Perspectives on Politics*, 10(04), 917–934.
- 1021 White, M. (2016). *The End of Protest: A New Playbook for Revolution*. Knopf Canada.
- 1022 Wilson, R. E., Gosling, S. D., and Graham, L. T. (2012). "A Review of Facebook Research in the
1023 Social Sciences." *Perspectives on Psychological Science*, 7(3), 203–220.

- 1024 Won, D., Steinert-Threlkeld, Z. C., and Joo, J. (2017). “Protest activity detection and perceived
1025 violence estimation from social media images.” *Proceedings of the 2017 ACM on Multimedia*
1026 Conference, ACM, 786–794.
- 1027 Xu, H., Gao, Y., Yu, F., and Darrell, T. (2017). “End-to-end learning of driving models from
1028 large-scale video datasets.” *arXiv preprint*.
- 1029 Young, L. E. (2019). “The Psychology of State Repression: Fear and Dissent Decisions in
1030 Zimbabwe.” *American Political Science Review*, 113(1), 140–155.
- 1031 Zhang, H. and Pan, J. (2019). “CASM: A Deep-Learning Approach for Identifying Collective
1032 Action Events with Text and Image Data from Social Media.” *Sociological Methodology*, 49,
1033 1–48.
- 1034 Zhao, D. (1998). “Ecologies of Social Movements: Student Mobilization during the 1989
1035 Prodemocracy Movement in Beijing.” *American Journal of Sociology*, 103(6), 1493–1529.

Supplementary Materials for How Violence Affects Protests

APPENDIX S1. DETAIL

S1.1 Convolutional Neural Networks

We use convolutional neural networks (CNN) to identify and analyze protest images. A CNN is a type of artificial neural network, a machine learning algorithm inspired by the human brain (Rosenblatt, 1958), that has gained widespread adoption in the field of computer vision. It has been successful in various applications including face recognition (Sun et al., 2014; Parkhi et al., 2015; Baltrušaitis et al., 2016), object detection (Girshick et al., 2014; Ren et al., 2015; Redmon et al., 2016), and self-driving cars (Huval et al., 2015; Bojarski et al., 2016; Xu et al., 2017). For methodological detail on computer vision for political scientists, see Joo and Steinert-Threlkeld (2018).

A CNN is a function whose outputs are computed through a series of sequential operations from the input values. For example, in image classification, the input is an image (i.e., an array of color intensities at pixels) and the output is the class that the image belongs to, e.g., an object category such as "police" or "male". A CNN transforms the given input through many operations until it reaches the final step which produces the output. Each operation is also called a layer, and a CNN usually has multiple "convolutional" layers. A convolutional layer performs convolution, which consists of an element-wise multiplication between pixel¹³ values and filter values (connection strengths) and a summation over adjacent pixels: this essentially measures how well the appearance of an input image matches the "template" that the model learned in training. A CNN, as well as other artificial neural networks, is trained to minimize a loss function, a measure of difference between the model prediction and ground truth label. This optimization is typically done by stochastic gradient descent.

¹³A CNN contains many layers and the output of a layer becomes the input of the next layer. The input to the first layer is the original input image's pixel intensities. For non-first layers, their inputs are given from nodes on two dimensional grid in the previous layer, not from the image pixels.

1060 Each CNN is defined by its architecture – the structural configuration specifying the number
1061 of layers, the order of their placement, and the types of non-linear transformations used. There
1062 exist many different CNN architectures with different properties. The architecture of our model is a
1063 “Residual Network” (ResNet) (He et al., 2016b) and has 50 convolutional layers. ResNet has been
1064 used in many of the state-of-the-art computer vision applications such as object detection (Ren
1065 et al., 2015) and human pose estimation (Güler et al., 2018). We use a ResNet model pre-trained
1066 on ImageNet data and finetune it with our data.

1067 This paper does not use tweet text because they do not measure violence, identity, or free riding
1068 as precisely as images. Decades of construction of event data, via hand and computer coding, has
1069 not been able to generate a measure of state or protester violence more refined than an ordinal
1070 measure; images allow for violence to be measured as a continuous variable. Measuring cleavages
1071 from text requires knowing the identity of accounts and would require orders of magnitude more
1072 user data; this exercise would not produce time varying measures because they would be about
1073 the account, not other protesters (Mislove et al., 2011; Sloan et al., 2013). Event datasets that
1074 measure cleavage spanning use newspaper text, which often does not report protester demographic
1075 information, and so measures are fixed at the movement level (Heaney and Rojas, 2008; Kern, 2011;
1076 Wilson et al., 2012; Fisher et al., 2017); images permit the measurement of the mass of protesters
1077 and their daily change. Measuring free riding from tweet text would require building a classifier,
1078 for each language in our dataset, for specific statements such as “I am not going to protest because
1079 it will not make a difference”; images that can induce free riding are easier to identify than specific
1080 tweets because visual language is universal (Grabер, 1996).

1081 **S1.2 Classifier Calibration**

1082 For binary variables in our analysis, we need to transform continuous outputs from CNN to
1083 binary values (0 or 1) by choosing a decision threshold such that we can determine if an image
1084 contains the variable of interest. The optimal decision threshold needs to be chosen so that it can
1085 balance good true positive and true negative rates, evaluated on the target data distribution (i.e., not
1086 the distribution in our development set). To this end, we chose 3,000 protest images from additional
1087 random samples from our Twitter pipeline and used Amazon Mechanical Turk to annotate them. We
1088 then generated a precision-recall curve for each attribute and chose the threshold at the minimum
1089 precision of .85.¹⁴ For each image and each attribute, our model therefore produces a probability
1090 estimate (a real number) via the classifier as well as a binary output (0 or 1). Figure A2 shows
1091 the precision-recall curve for each attribute, providing the threshold value for each. The twelve
1092 attributes and their thresholds are shown in Table A9.

¹⁴One could also use another method such as F-measure to choose the optimal decision threshold. In our study, it is more important to maintain the minimum precision (true positive rate) at a high point for every attribute, rather than trying to detect more relevant images while making more mistakes.

Table A8. List of visual attributes.

Attribute	Description	Hypothesis	Expectation
1. Protester Violence	How violent protesters are.	Violence (H1)	-
2. State Violence	How violent the state is.	Violence (H1)	+, -
3. Police	Police or troops are present in the scene.	Violence (H1)	-
4. Fire	There is fire or smoke in the scene.	Violence (H1)	-
5. Gender	Is the face male or female?	Cleavages (H2)	+, -
6. Race	Is the face White, Middle Eastern, East Asian, Southeast Asian, Black, Indian, or Latino?	Cleavages (H2)	+, -
7. Age	0-2, 3-9, 10-19, [...], 70	Cleavages (H2)	+, -
8. Face	Presence of a face.	Protest size	
9. Group 20	There are roughly more than 20 people in the scene.	Future: free riding	
10. Group 100	There are roughly more than 100 people in the scene.	Future: free riding	
## Children	Children are in the scene.	N/A	N/A
## Shout	One or more people shouting.	N/A	N/A
## Photo	Protesters holding signs or a photograph of a person (politicians or celebrities).	N/A	N/A
## Flag	There are flags in the scene.	N/A	N/A
## Night	It is at night.	N/A	N/A
## Sign	Protesters holding a visual sign (on paper, panel, or wood).	N/A	N/A

NB: Attributes without numbers could not be classified precisely enough to be included in regression models.

NB: Attributes in **bold** are generated using the face classifier.

Table A9. Attributes and Thresholds

Attribute	Threshold
Protester Violence	.021
State Violence	.01
Police	.937
Fire	.37
Child	.15
Small Group	.725
Large Group	.509
Shout	.355
Photo	.815
Flag	.187
Night	.359
Sign	.744

Fig. A1. Examples of Our Annotation Interface (in Amazon Mechanical Turk)

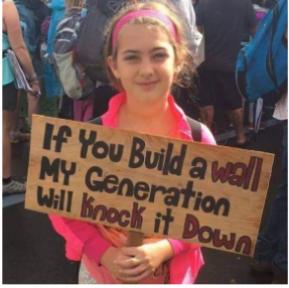
 <p>Q1. Does this image contain a scene of protest?</p> <p><input checked="" type="radio"/> Protest <input type="radio"/> Uncertain <input type="radio"/> Not Protest</p>	 <p>Q2. Does this image contain a scene of protest?</p> <p><input type="radio"/> Protest <input checked="" type="radio"/> Uncertain <input type="radio"/> Not Protest</p>	 <p>Q3. Does this image contain a scene of protest?</p> <p><input type="radio"/> Protest <input checked="" type="radio"/> Uncertain <input type="radio"/> Not Protest</p>
<p>Q0. Please answer all the questions for the below image.</p> <div style="display: flex; align-items: center;">  <div style="margin-left: 10px;"> <p>Protesters (or a protester) holding visual signs (on a paper, panel, or wood).</p> <p><input checked="" type="radio"/> Yes <input type="radio"/> Not sure <input type="radio"/> No</p> <p>There is fire or smoke in the scene.</p> <p><input type="radio"/> Yes <input checked="" type="radio"/> Not sure <input type="radio"/> No</p> <p>Children (or a child) are in the scene.</p> <p><input type="radio"/> Yes <input checked="" type="radio"/> Not sure <input type="radio"/> No</p> <p>There are roughly more than 20 people in the scene.</p> <p><input type="radio"/> Yes <input checked="" type="radio"/> Not sure <input type="radio"/> No</p> <p>It is at night.</p> <p><input type="radio"/> Yes <input checked="" type="radio"/> Not sure <input type="radio"/> No</p> <p>There are flags in the scene.</p> <p><input type="radio"/> Yes <input checked="" type="radio"/> Not sure <input type="radio"/> No</p> <p>There is one or more people shouting.</p> <p><input type="radio"/> Yes <input checked="" type="radio"/> Not sure <input type="radio"/> No</p> </div> </div>		
<p>Q0. Choose the image that you feel is more violent.</p> <p><input checked="" type="radio"/> Image 1</p>  <p><input type="radio"/> Similar</p> <p><input type="radio"/> Image 2</p> 		

Fig. A2. Precision-Recall Curves For Binary Attributes. AP stands for average precision, which is the standard accuracy measure for binary classification. AP is also equal to the area under the precision-recall curve.

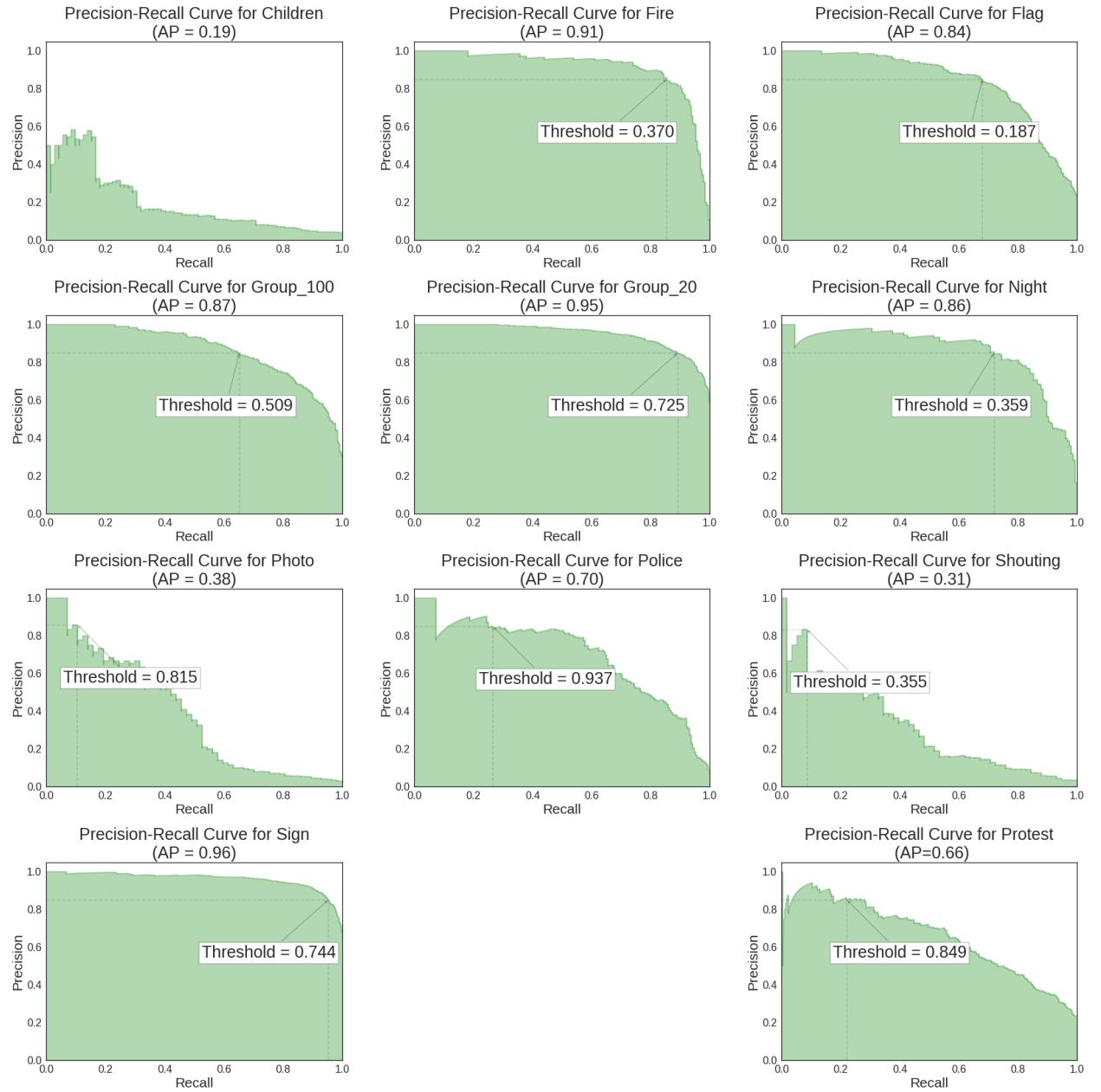


Fig. A3. Example Results of Our Face Model



1093 **S1.3 Evaluating the CNN**

1094 Figure A4 shows the model performance measured on the validation set. The Receiver-Operator
1095 Curve (ROC) documents the relationship between false-positive and true-positives, with a higher
1096 area-under-curve (AUC) corresponding to a better accuracy. Visually, the closer the curve is to the
1097 upper-left corner, the better the classifier for that label.

1098 Figure A5 shows a scatterplot of the classifier's output for violence against the rating recovered
1099 from the Bradley-Terry model. It also shows the ROC curve for protester violence and state violence.
1100 All three subfigures demonstrate strong performance of our classifier's ability to measure perceived
1101 violence.

1102 To intuitively visualize how the classifier works, we use Gradient-weighted Class Activation
1103 Mapping (Grad-CAM) ([Selvaraju et al., 2016](#)). Grad-CAM highlights important regions for classi-
1104 fying the concept in an image. Grad-CAM highlights such regions by tracing back the classification
1105 outcome to the input image through passing gradients. The results are shown in Figure A6, with
1106 red color indicating more important regions. For technical details, see [Selvaraju et al. \(2016\)](#).

1107 Figure A7 arrays images from each category by the classification scores from the CNN. As the
1108 classification score approaches 1 for each category, images more closely exhibit the visual concept.

Fig. A4. Model Performance

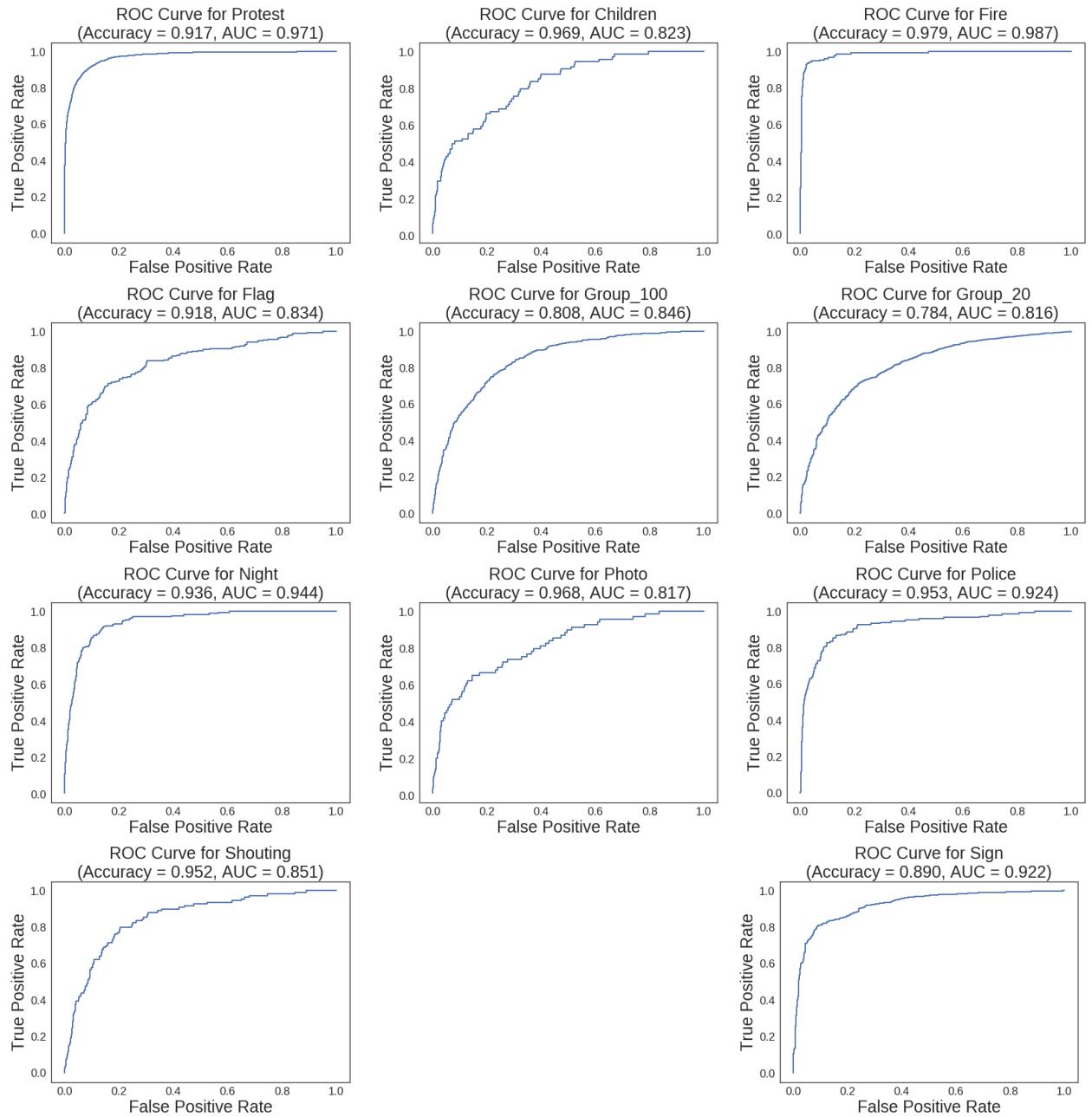


Fig. A5. Validating Violence Measurement

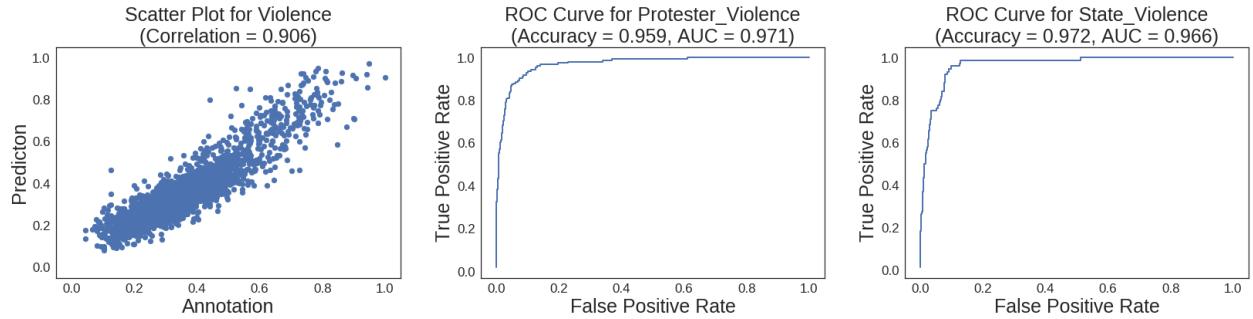


Fig. A6. Visualization of Region Importance in Classification Using Grad-CAM: Important regions which more contribute to the classification for each attribute are highlighted in red.



Fig. A7. Sample Classifier Estimates by Category: Images are ordered by their classification scores. (Blue lines mark the exact classification scores of corresponding images)

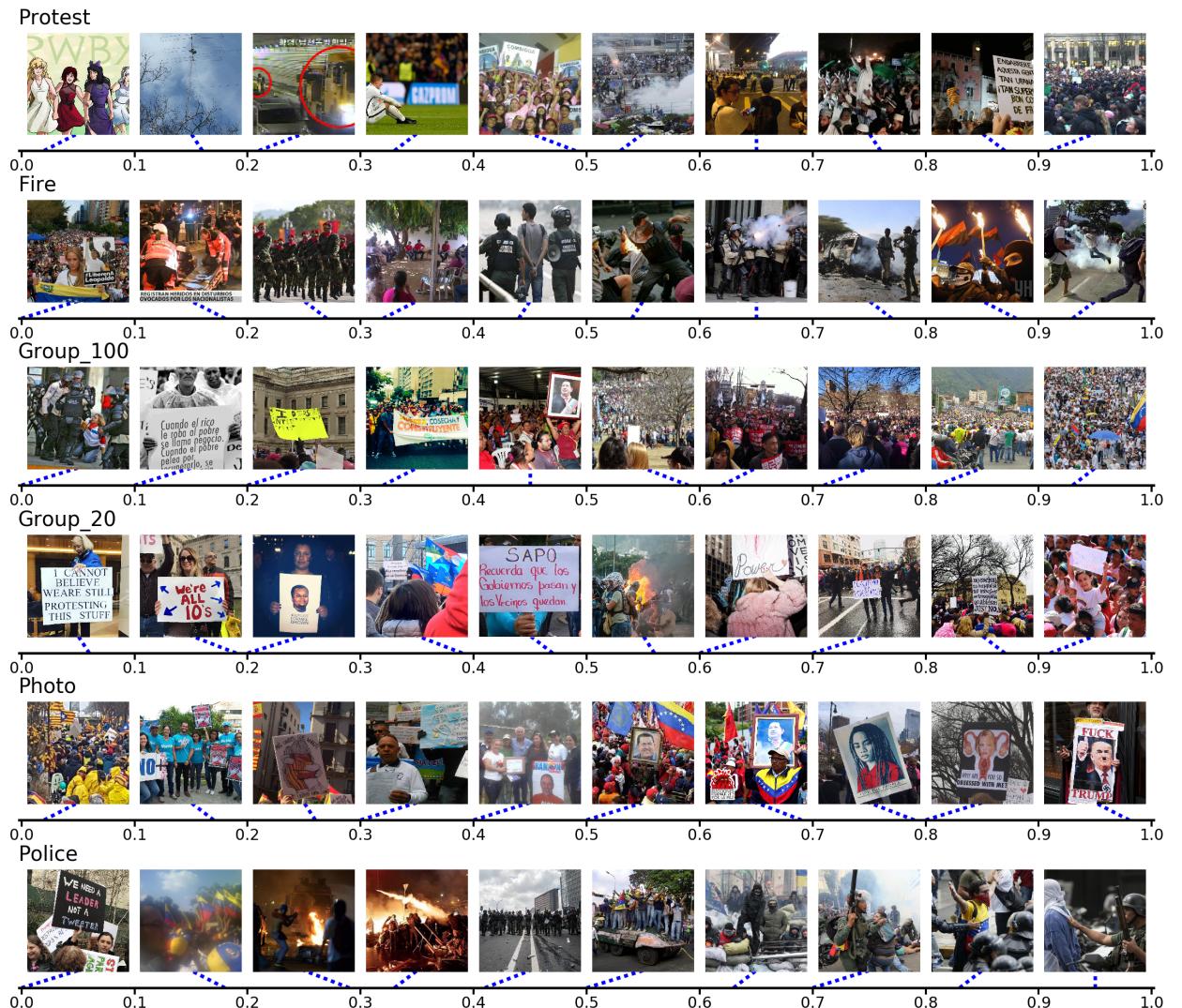
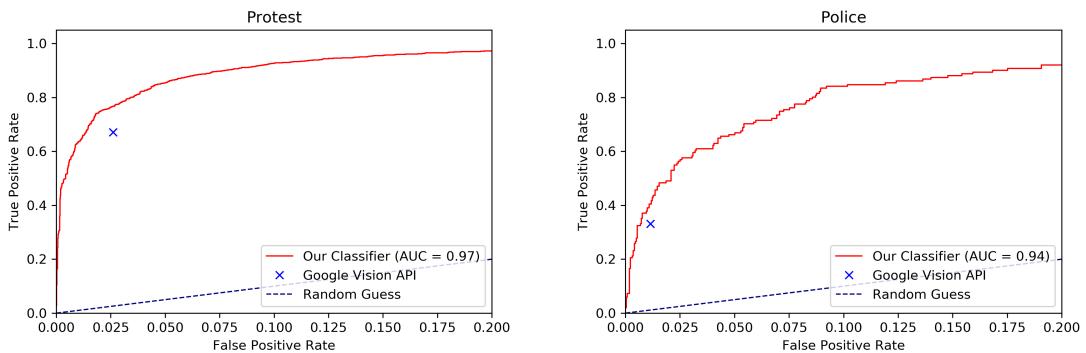


Fig. A7. Sample Classifier Estimates by Category (Continued)



Finally, to compare the classification performance of existing commercial classifiers against our own classifier, we used Google Vision’s label detection on the test set of UCLA Protest Image Database (Won et al., 2017) and measured the classification accuracy. This dataset has 11,000 test images with various labels related to protest activity such as the presence of protesters or police officers in images. Since Google’s label detection automatically identifies visual concepts and objects in many categories, including protest and police, from an input image, we directly compared its accuracy with our model accuracy. As shown in Figure A8, the protest and scene models classified protest and police more accurately than the Google Vision API. The superior result is most likely due to the fact that we specifically collected diverse protest images and hard-negatives (i.e., non-protest images which look like protest) from many sources. The Google Vision API may perform better on general image classification and can be very useful when one does not have any training data.

Fig. A8. Classification performance comparison between our model and the public model from Google’s Vision API.



1121 **APPENDIX S2. BIAS**

1122 In the United States, Twitter users who geotag are richer, more likely to live in cities, young, and
1123 non-white ([Malik et al., 2015](#)). In the United Kingdom, Twitter users are younger, more educated,
1124 more likely to be male, and more politically engaged (but less likely to vote) than others ([Mellon](#)
1125 and [Prosser, 2017](#)). Once on Twitter, geotagging users are slightly older than non-geotaggers, there
1126 is some difference in rates of geotagging across profession, and there is large variation by tweet
1127 language in the percentage of users who geotag (a low of 0.4% for Arabic accounts to a high of
1128 8.3% for Turkish, with an average of 3.1%) ([Sloan and Morgan, 2015](#)).

1129 Though Twitter users differ from non-Twitter users and those who assign locations to their
1130 tweets differ from those who do not, there is no *a priori* reason to expect that they differ in the type
1131 of protest images they share. Conditional on being at a protest, there is no reason to think that the
1132 contents of a geotagged protest image should systematically differ from a protest image that is not
1133 geotagged. Of anything, the importance of social media for tactical coordination of protests means
1134 that geotagged tweets should be *more* likely to represent a protest than one that is not ([Gunitsky,](#)
1135 2015 ; [Little, 2015](#)).

1136 Section S6 compares users who tweet protest images to those who tweet non-protest images.
1137 More accounts share protest images than non-protest images, and they have fewer followers than
1138 those tweeting non-protest images. There is no statistically significant difference in the account
1139 age or frequency of tweeting. While this comparison does not prove that the protest images are
1140 an unbiased representation of the protest itself, it at least appears to be the case that the accounts
1141 themselves do not appear to be any more biased, and are probably actually more representative,
1142 than the larger Twittersphere. The protest photos appear to come from more “normal” users than
1143 those who normally tweet images.

1144 If bias in protest data from geolocated images shared on Twitter exists, it should nonetheless be
1145 less than that which exists from relying on any text that is not a police archive. The main source
1146 of information for existing event data, newspapers, have large, well-known biases that result from
1147 incentives that are much weaker on social media. Newspapers are much more likely to cover large

events than small ones (McCarthy et al., 1996) as well as events perceived to be of interest to their subscribers (Myers and Caniglia, 2004; Baum and Zhukov, 2015). Events away from urban centers are less likely to receive coverage (Kalyvas, 2004; Weidmann, 2014), as are ones that are parts of a larger wave of events (Hellmeier et al., 2018). Given the increasing consolidation of the newspaper industry, these biases are likely to have become more consistent across sources (Baum and Zhukov, 2018).

These biases exist because newspapers have to maximize readership (advertising, newsstand, and subscription revenue) while constrained by space. This constraint puts an emphasis on reporting novel or unexpected events such as violent attacks or large protests. Even if readership is national - and most newspapers have local or, at best, regional circulation - events are still selected based on their appeal to the readers. The need to daily attract readers means coverage of events quickly tapers regardless of the event duration (Hellmeier et al., 2018). For a more extensive explanation of bias in news coverage, see Earl et al. (2004).

Social media platforms do not face these same pressures. While their business model is more focused on attracting eyeballs than newspapers are, because they do not have subscribers or newsstand sales, there is essentially zero restriction on the space in which to publish.¹⁵ Whether or not the platforms, such as Twitter or Facebook, should be treated as media companies is a separate issue, but one way in which they are not like other media is that they do not employ people to create the information featured on their platform, the way newspapers pay journalists. Given the essentially infinite supply of posts and the lack of control over content providers, there should therefore be much less selection pressure on what appears on social media.¹⁶ Because newspapers face scarcity constraints that social media do not, the latter should be much more likely to provide a less biased account of the world than newspapers. In providing orders of magnitude more posts than newspapers do articles, social media are closer in scope to government archives than they are

¹⁵Each new post imposes a marginal cost - server space and electricity - on the platform that is much smaller than article for a newspaper.

¹⁶Social media platforms increase user engagement by selectively presenting posts to users. While this algorithmic process may present users with biased interpretations of events, that process is not used to decide which tweets to send to the API (Pfeffer and Mayer, 2018).

¹¹⁷² newspapers ([Sullivan, 2016](#)). See [Sobolev et al. \(2019\)](#) for a more extensive comparison of bias in
¹¹⁷³ newspaper and social media event data.

¹¹⁷⁴ **APPENDIX S3. COUNTRY PERIODS ANALYZED**

Table A10. Protest Periods

Country Images	Start Protest Images/Day	End	Issue	Tweets
Belarus	02.18.2017	05.02.2017	Unemployment Tax	2.18
Burundi	04.01.2015	12.01.2015	Elections	.06
Cameroon	11.01.2016	12.01.2017	Bilingualism	.06
Egypt	06.01.2017	06.31.2017	Islands to Saudi Arabia	5.38
Gabon	08.20.2016	09.27.2016	Elections	.235
Hong Kong	2014.09.18	2014.12.23	China Reforms	5.82
Pakistan	11.01.2017	11.30.2017	Blasphemy protests	.941
<i>Russia</i>	<i>03.12.2017</i>	<i>04.26.2017</i>	<i>Corruption</i>	<i>19.3</i>
<i>Catalonia, Spain</i>	<i>2017.09.01</i>	<i>2017.12.31</i>	<i>Secession</i>	<i>31.8</i>
South Korea	2016.10.20	2017.03.14	Anti-incumbency	8.04
Togo	08.01.2017	12.01.2017	Anti-incumbency	.23
Ukraine	11.21.2013	03.21.2014	European Integration	3.32
United States	2017.01.20	2017.01.22	Women's March	9,034.33
Venezuela	2014.03.27	2015.02.08	Grievances	31.20
Venezuela	12.29.2016	12.17.2017	Anti-Maduro	16.40

1175 **APPENDIX S4. CITY-DAY CORRELATION**

1176 Figure A9 shows the correlation between the models' variables. The only variables with
1177 correlation above .8 are the two group variables. These correlations are higher than the per tweet
1178 ones: the correlation comes from aggregating different photos to the city-day level, not from the
1179 classifier producing similar estimates for different labels. The per tweet correlation is shown in
1180 Figure A10.

1181 **APPENDIX S5. TWEET LEVEL CORRELATION**

1182 Figure A10 shows the correlation between variables at the image level. The only correlation
1183 above .8 is male faces and white faces with the number of faces. These correlations are lower than
1184 the country-day correlations, meaning that correlations in the aggregated data come from multiple
1185 mechanisms occurring during a protest, not noise in the classification of individual photographs.

Fig. A9. Covariate Correlation, by City-day

Age Div. _{i,t-1}	0.32	0.05	0.07	0.03	0.26	0.78	0.77	1
Race Div. _{i,t-1}	-0.24	0.02	0.04	0.05	0.16	0.7	1	0.77
Gend. Div. _{i,t-1}	-0.31	0.05	0.06	0.04	0.29	1	0.7	0.78
Fire _{i,t-1}	-0.27	0.54	0.15	0.02	1	0.29	0.16	0.26
Police _{i,t-1}	-0.07	0.07	0.29	1	0.02	0.04	0.05	0.03
Perc. Stt. Violence _{i,t-1}	-0.09	0.45	1	0.29	0.15	0.06	0.04	0.07
Perc. Prstr. Violence _{i,t-1}	-0.09	1	0.45	0.07	0.54	0.05	0.02	0.05
Log(Protest Size _{i,t})	1	0.09	0.09	0.07	0.27	0.31	0.24	0.32
	Log(Protest Size _{i,t})	Perc. Prstr. Violence _{i,t-1}	Perc. Stt. Violence _{i,t-1}	Police _{i,t-1}	Fire _{i,t-1}	Gend. Div. _{i,t-1}	Race Div. _{i,t-1}	Age Div. _{i,t-1}

Fig. A10. Covariate Correlation, by Tweet

	Age Div. _{i,t-1}	Race Div. _{i,t-1}	Gend. Div. _{i,t-1}	Fire _{i,t-1}	Police _{i,t-1}	Perc. Stt. Violence _{i,t-1}	Perc. Prstr. Violence _{i,t-1}	Log(Protest Size _{i,t})
Age Div. _{i,t-1}	0.47	-0.08	-0.02	0.06	-0.06	0.45	0.56	1
Race Div. _{i,t-1}	0.5	-0.11	-0.09	-0.04	-0.06	0.4	1	0.56
Gend. Div. _{i,t-1}	0.3	-0.07	-0.01	0.09	-0.05	1	0.4	0.45
Fire _{i,t-1}	-0.08	0.77	0.04	-0.01	1	-0.05	-0.06	-0.06
Police _{i,t-1}	0.02	0.17	0.47	1	-0.01	0.09	-0.04	0.06
Perc. Stt. Violence _{i,t-1}	-0.08	0.54	1	0.47	0.04	-0.01	-0.09	-0.02
Perc. Prstr. Violence _{i,t-1}	-0.13	1	0.54	0.17	0.77	-0.07	-0.11	-0.08
Log(Protest Size _{i,t})	1	-0.13	-0.08	0.02	-0.08	0.3	0.5	0.47

1186

APPENDIX S6. COMPARING IMAGE SHARERS TO PROTEST IMAGE SHARERS

1187

Interested in whether or not those who tweet images of protests differ from those who tweet images, we analyzed account covariates by country-photo type (protest or not protest). For each set of photos, we kept each user only once, which discarded about 40% of tweets in each set. For users who tweeted the same image type multiple times, we randomly keep one tweet, and an account can appear in both samples. The account characteristics we analyze are the number of followers, following, statuses, and account age (days on Twitter). We also count the number of unique users in each country by image type. The points estimates and 95% confidence intervals are shown in Figure A11.

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

More, less popular people share protest images than non-protest images. In all countries, there are more accounts that tweet protest images than non-protest images, though the difference is only statistically significant in Venezuela. These extra accounts have fewer followers on average than those tweeting non-protest images, though the difference is not statistically significant in Ukraine.

1199

1200

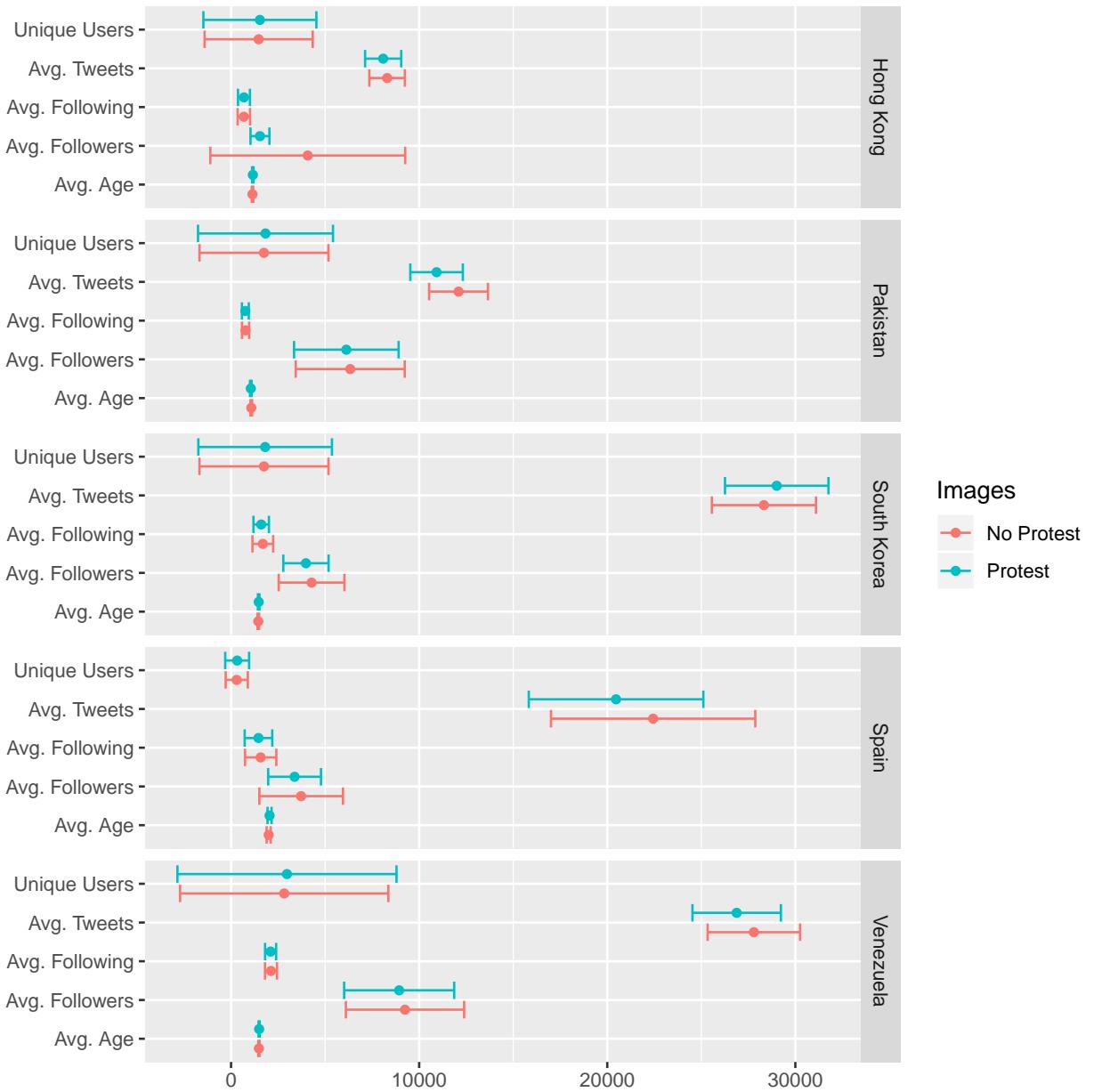
1201

1202

1203

The two sets of users do not differ in how engaged they are on Twitter. In no country do the protest and non-protest accounts follow different numbers of users. They have each been on Twitter the same amount of time. (In Russia only, these accounts have been on Twitter longer, an average of 43.21 days more.) In Venezuela, the protest users are slightly less active on Twitter, but in Russia they are more active; there is no statistically significant difference for Ukraine.

Fig. A11. Users Tweeting Protests vs. Not



1204

APPENDIX S7. VALIDATING THE SCENE AND FACE CLASSIFIERS

1205

Figure A12 shows the distribution of protester violence by country. Figure A13 shows the distribution of race entropy by country. Figure A14 shows the distribution of age entropy by country.

1206

1207

Fig. A12. Distribution of Protester Violence by Country

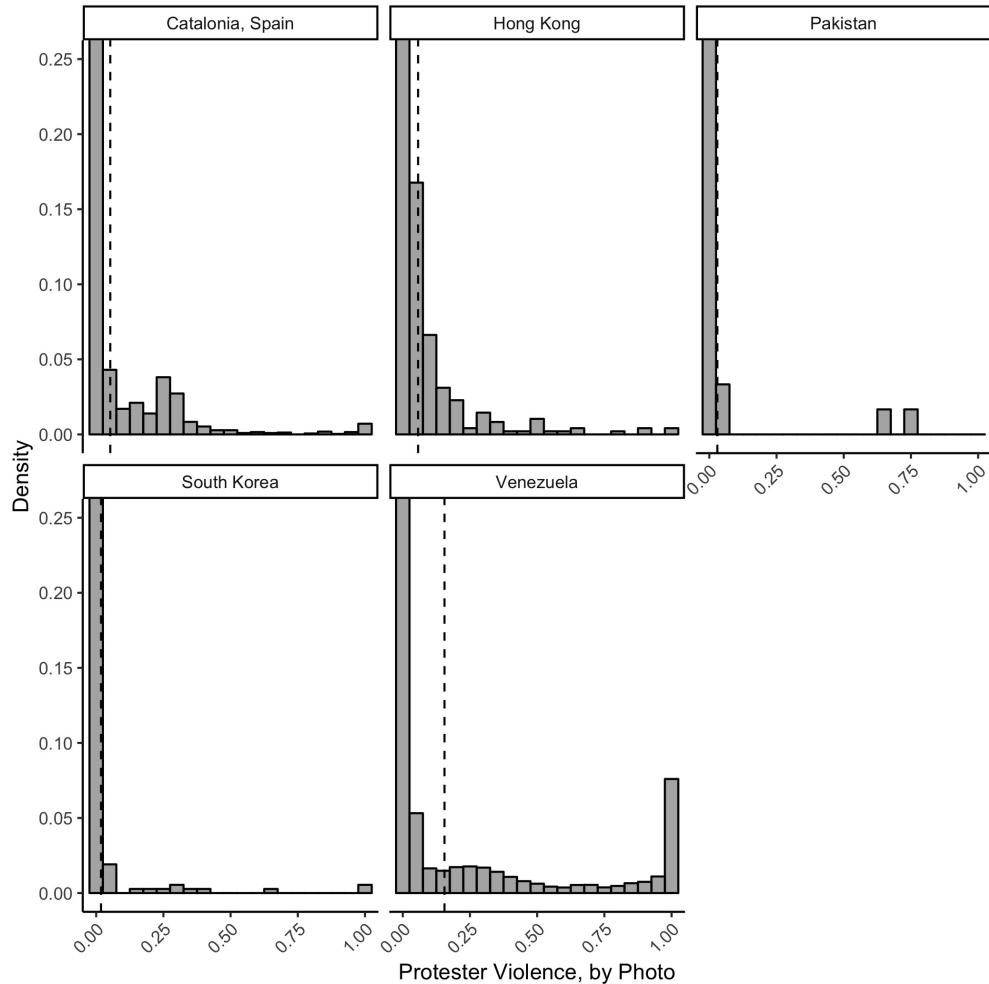


Fig. A13. Distribution of Race Entropy by Country

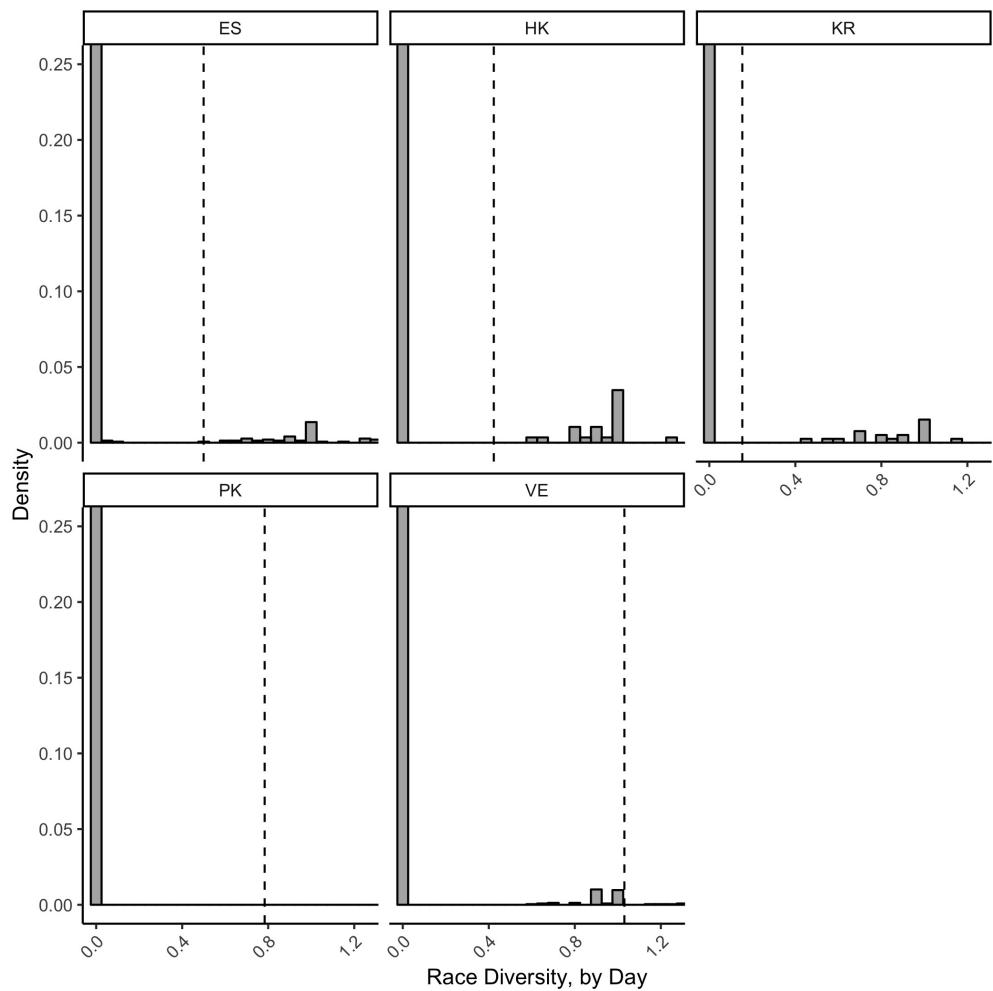
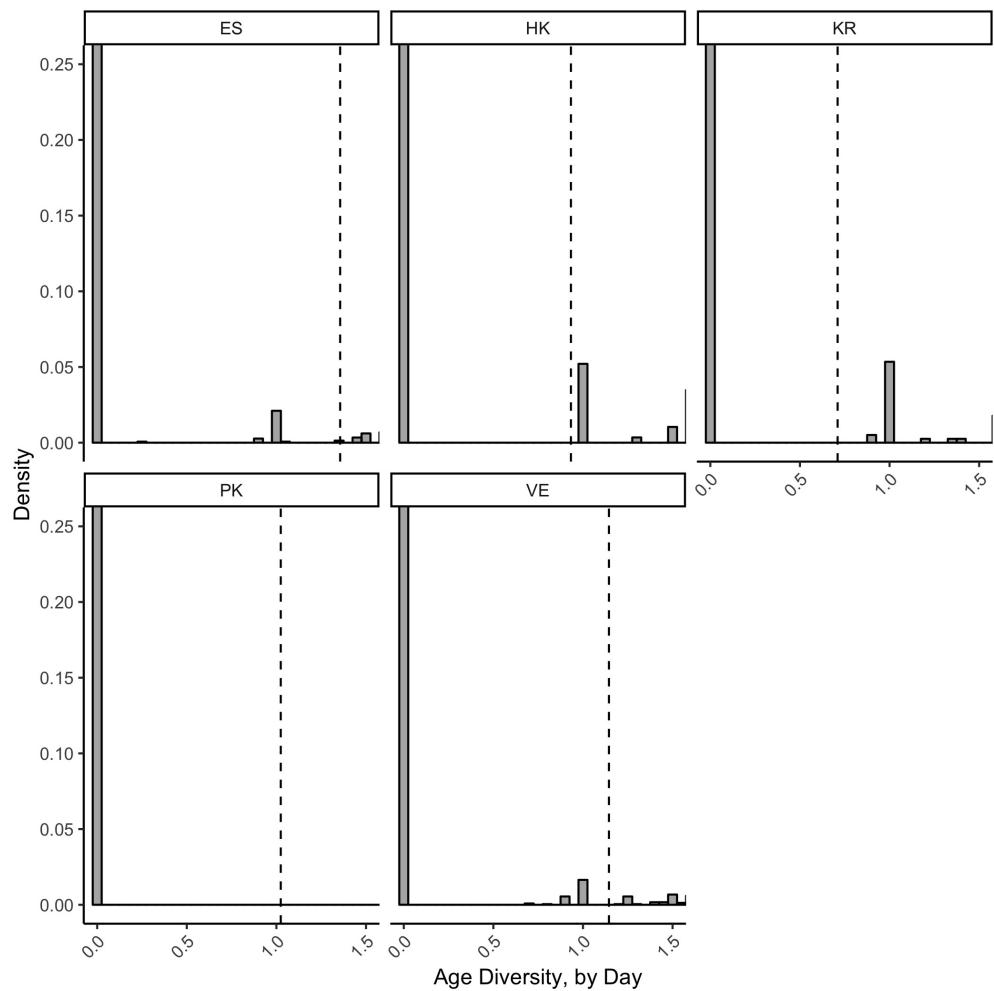


Fig. A14. Distribution of Age Entropy by Country



1208 APPENDIX S8. ADDITIONAL ROBUSTNESS CHECKS

1209 S8.1 Bot Distribution

1210 Table A11 complements Table 7. No more than 6.7% of accounts are from bots, and no more
1211 than 6.5% of tweets.

Table A11. Distribution of Bots by City

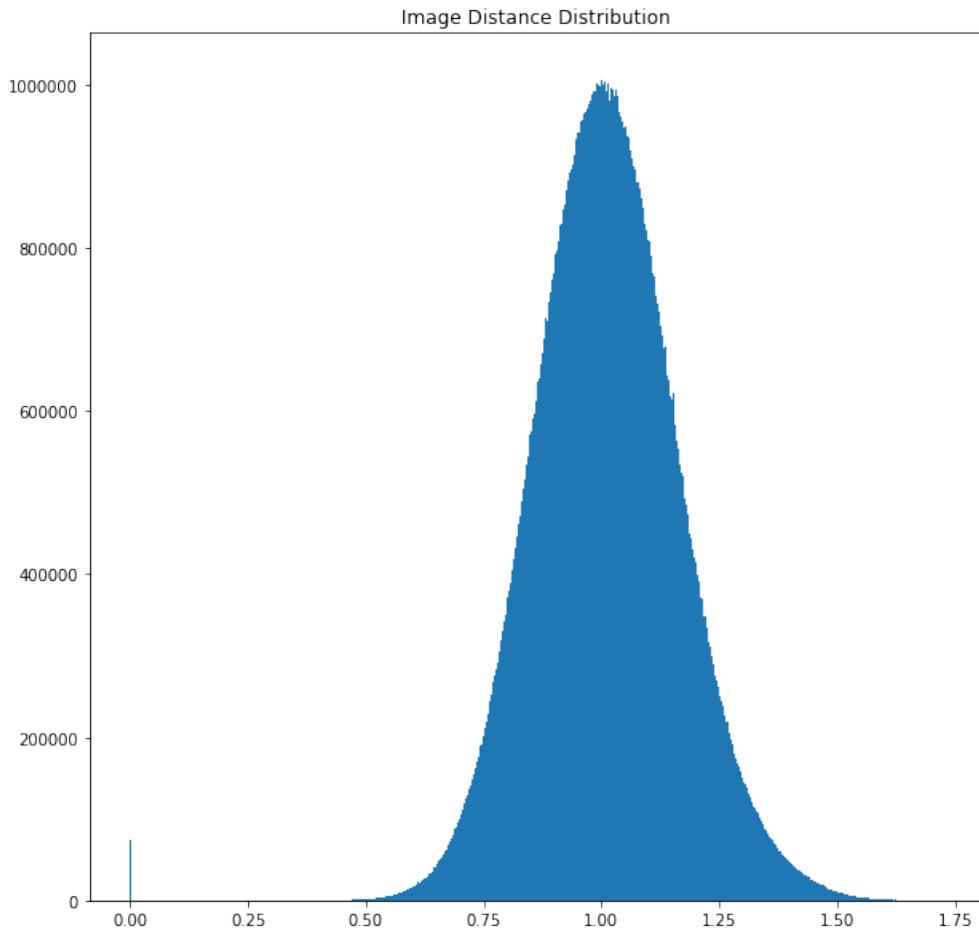
City	Avg. Bot Score	Max. Bot Score	SD Bot Score	Percent Tweets from Bots	Percent Accounts of Bots
Ciutat Vella	0.131	0.948	0.297	0.108	0.040
Lahore	0.067	0.559	0.173	0.100	0.143
Sant Salvador de Guardiola	0.066	0.611	0.155	0.083	0.094
Granera	0.053	0.812	0.156	0.054	0.065
Valencia	0.063	0.921	0.145	0.052	0.058
Tarragona	0.058	0.829	0.172	0.052	0.079
Central	0.050	0.845	0.156	0.051	0.133
Maracaibo	0.050	0.685	0.115	0.039	0.037
Seoul	0.145	0.905	0.170	0.035	0.056
Barcelona	0.029	0.939	0.110	0.024	0.022
Caucagua	0.054	0.905	0.118	0.020	0.027
Boca del Rio	0.032	0.829	0.117	0.020	0.047
Girona	0.034	0.637	0.088	0.019	0.038
Sant Feliu de Pallerols	0.040	0.661	0.086	0.018	0.048
Caracas	0.043	0.942	0.103	0.010	0.025
Granollers	0.011	0.084	0.019	0	0
Kimhae	0.052	0.054	0.009	0	0
Kowloon	0.021	0.355	0.062	0	0
Lleida	0.018	0.355	0.058	0	0
Mataró	0.007	0.054	0.009	0	0
Reus	0.018	0.297	0.051	0	0
Sabadell	0.018	0.355	0.044	0	0
Sant Cugat del Vallès	0.015	0.270	0.047	0	0
Terrassa	0.005	0.030	0.006	0	0

1212 S8.2 Deduplicating Images

1213 To deduplicate images, we extracted 1,000 features from a pre-trained ResNet50 model (He
1214 et al., 2016a). Conventional image preprocessing methods for deep learning models were used.
1215 Each image was resized to 256 x 256 pixels. Then, a center-crop of 224 x 224 pixels was performed.
1216 Finally, the cropped images were normalized to the mean and standard deviation of the ImageNet
1217 dataset (Deng et al., 2009). The 1,000 feature vector of each sample was normalized to unit norm.
1218 The L2 distance among the normalized data is computed, and images are considered matches if the

₁₂₁₉ distance is less than a threshold of 0.2. The histogram of the distribution of distances is shown in
₁₂₂₀ Figure A15.

Fig. A15. Distribution of Pairwise Image Distances



₁₂₂₁ Two manual checks verify these results. The largest 90 clusters were manually inspected and
₁₂₂₂ no images were misidentified as duplicates. The 220 most common images identified as duplicates,
₁₂₂₃ shared 2,500 times, were inspected, and none were misidentified as duplicates.

₁₂₂₄ Table A12 shows the percentage of tweets per city that are duplicates.

₁₂₂₅ Figure A16 shows the number of times each image appears in the dataset.

₁₂₂₆ Figure A17 shows the distribution of the percentage of images per city that are duplicates.

₁₂₂₇ Except for Kimhae, the spike at 1 is cities with 1 image.

Table A12. Duplicate Images

City	Percentage Duplicates	Total Tweets
Kimhae	0.957	46
Sant Feliu de Pallerols	0.621	58
Girona	0.500	108
Caracas	0.463	2,105
Mataró	0.425	40
Sant Cugat del Vallès	0.424	33
Caucagua	0.281	224
Tarragona	0.274	62
Valencia	0.269	167
Sant Salvador de Guardiola	0.250	40
Maracaibo	0.243	152
Terrassa	0.237	59
Boca del Rio	0.227	66
Sabadell	0.187	75
Reus	0.184	38
Barcelona	0.182	1,338
Granera	0.154	65
Granollers	0.150	20
Lleida	0.116	43
Ciutat Vella	0.100	40
Seoul	0.052	326
Central	0	41
Kowloon	0	66
Lahore	0	5

Fig. A16. Distribution of Number of Duplicates

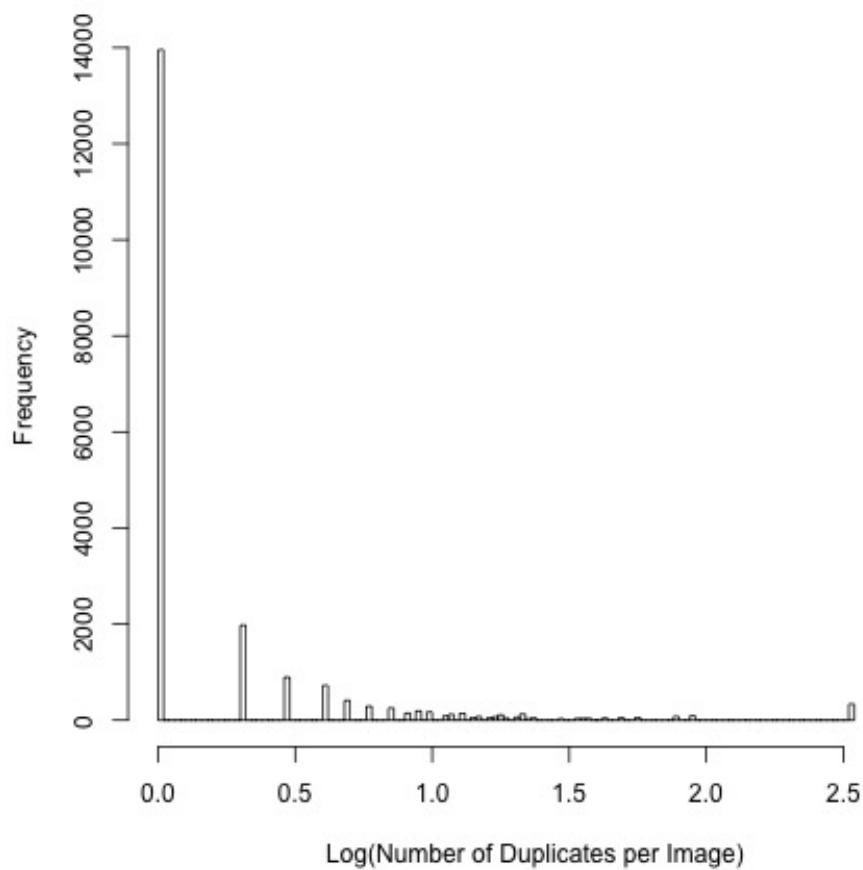
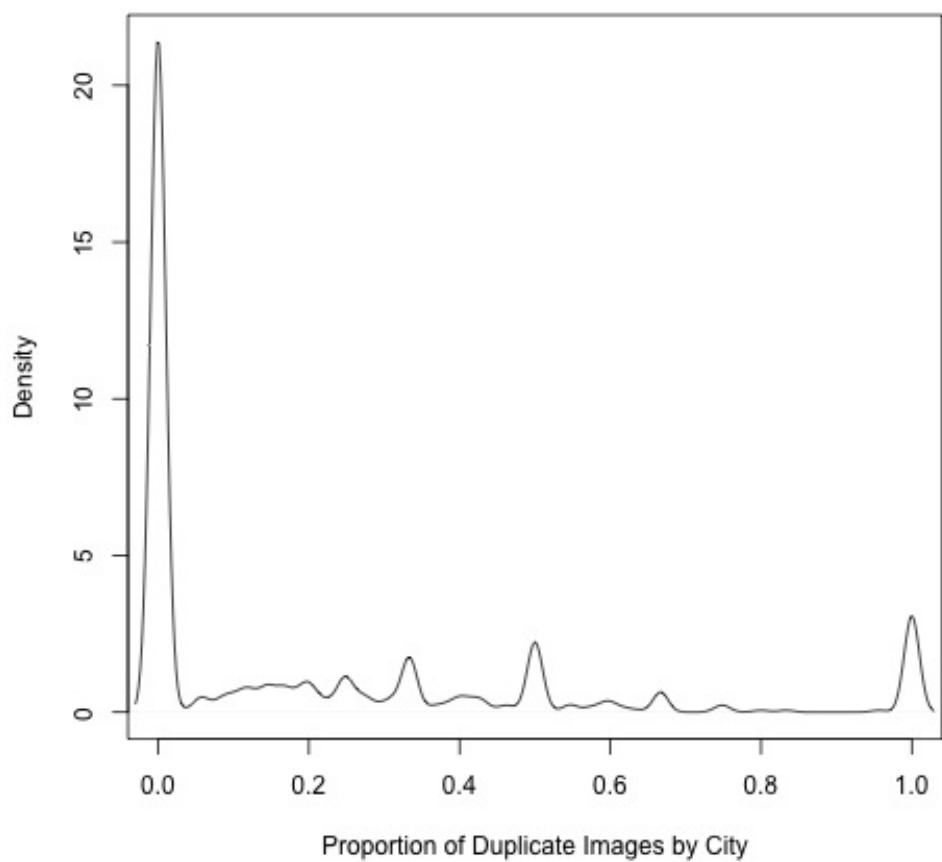


Fig. A17. Percentage of Duplicates by City



1228 **S8.3 Different Dependent Variable**

1229 Table A13 repeats the main model using three different operationalizations of the dependent
1230 variable. The first, shown in Model 2, does not log-transform the number of protesters. The results
1231 for the violence and free-riding variables are the same, except the number of photos containing
1232 police is now statistically significant; its sign matches the original model. Race Diversity_{i,t-1} no
1233 longer correlates with subsequent changes in protest size.

1234 The third and fourth models measure the size of protest using the number of users who share
1235 a protest photo per city-day. The third uses the raw count, the fourth logged. This quantity is
1236 smaller than the number of protest photos per day because users often share multiple photos.
1237 Results for violence and racial diversity are the same when not taking the logarithm, though they
1238 lose significance when log-transformed. Large Group_{i,t-1} switches signs and is significant in both
1239 models, while Large Group_{i,t-1}² supports the same inference in Model 3 but not Model 4. Of all
1240 robustness checks in the manuscript and supplementary materials, these two differ the most from
1241 the original model.

1242 These results differ the most from the rest of the paper's for two reasons. Most importantly, they
1243 embody a different data generating process than the other operationalizations of protest size. They
1244 do so because counting individual images provides less information about the size of a protest than
1245 counting the faces in an image. It provides less information because the photos are equivalent to
1246 randomly sampling a protest space and the surrounding protesters, akin to the leading methodology
1247 of in-person protest size measurement ([Schweingruber and McPhail, 1999](#)). Other work has shown
1248 that counting the number of protest photos less accurately recovers true protest size than summing
1249 the number of faces in those photos ([Sobolev et al., 2019](#)). Second, there is much less variation in
1250 this measure than in the sum of faces. The maximum value is 158, third quartile 1; when restricted
1251 to days with protest photos, the third quartile is 3.

1252 **S8.4 Accounting for Autocorrelation**

1253 Table A14 shows that the inclusion of lagged dependent variables up to 15 days old do not
1254 change the results for the violence or free riding variables already significant in the original model

Table A13. Different Measures of Protest Size

	Original (1)	No Log (2)	Number Users (3)	$\log_{10}(NumberUsers)$ (4)
Perceived Prtstr. Violence $_{i,t-1}$	-.1674** (.0677)	-9.8923*** (2.8933)	-3.4252** (1.7108)	-.1724*** (.0552)
Perceived Stt. Violence $_{i,t-1}$	1.2820*** (.3327)	86.4146** (37.2147)	24.2096* (13.1782)	.4030* (.2269)
Perceived Stt. Violence $^2_{i,t-1}$	-2.1030*** (.6093)	-168.2910** (79.2780)	-45.4816* (25.5827)	-.6696* (.4007)
Police $_{i,t-1}$.7626* (.4493)	125.1103* (74.2665)	31.9542 (19.9307)	.2401 (.2776)
Fire $_{i,t-1}$.1009*** (.0236)	3.5316*** (.5783)	1.2606* (.7297)	.0690*** (.0242)
Gender Diversity $_{i,t-1}$	-.1126 (.0939)	-9.8651 (7.3636)	-2.0632*** (.7160)	-.0972*** (.0350)
Race Diversity $_{i,t-1}$.0683 (.0440)	7.8310*** (2.6819)	1.0215* (.5586)	.0518 (.0348)
Age Diversity $_{i,t-1}$.0203 (.0289)	1.6471 (1.9407)	.6294 (.5215)	.0209 (.0212)
Tweets $_{i,t-1}$.0095*** (.0033)	.2254* (.1340)	.0112 (.0432)	.0040** (.0021)
DV $_{i,t-1}$.1578** (.0682)	.0910*** (.0348)	.2729*** (.0310)	.3725*** (.0973)
Intercept	.1260*** (.0237)	3.5931*** (.6761)	1.5215*** (.1740)	.1063*** (.0106)
N	4,376	4,376	4,376	4,376
Adjusted R ²	.2450	.2051	.2741	.3710
City FE	Y	Y	Y	Y

*p < .1; **p < .05; ***p < .01

1255 (Model 1). (A partial autocorrelation plot suggested serial dependence for up to fifteen days.)
 1256 The presence of police is now positively correlated with subsequent protest size. More noticeably,
 1257 neither Race Diversity $_{i,t-1}$ nor Any Child $_{i,t-1}$ remain significant. This updated result supports
 1258 the interpretation provided in the main paper of the cleavage variables: they are endogenous
 1259 to the protests themselves, so controlling for enough previous protests removes those variables'
 1260 significance.

Table A14. Robust to Additional Lagged Dependent Variables

	DV: $\log_{10}(\text{Sum of Faces})_{i,t}$	
	Original Model	15 Lags
	(1)	(2)
Perceived Prtstr. Violence $_{i,t-1}$	-.1674** (.0677)	-.0750 (.0788)
Perceived Stt. Violence $_{i,t-1}$	1.2820*** (.3327)	.7844** (.3186)
Perceived Stt. Violence $^2_{i,t-1}$	-2.1030*** (.6093)	-1.3370** (.5303)
Police $_{i,t-1}$.7626* (.4493)	.6663* (.3987)
Fire $_{i,t-1}$.1009*** (.0236)	.0365 (.0332)
Gender Diversity $_{i,t-1}$	-.1126* (.0939)	-.0592 (.0616)
Race Diversity $_{i,t-1}$.0683 (.0440)	.0279 (.0430)
Age Diversity $_{i,t-1}$.0203 (.0289)	.0130 (.0338)
Tweets $_{i,t-1}$.0095*** (.0033)	.0062*** (.0022)
DV $_{i,t-1}$.1578*** (.0682)	.0588 (.0360)
Intercept	.1260*** (.0237)	.0314** (.0129)
N	4,376	4,033
Adjusted R ²	.2450	.3220
15 Lags of DV	N	Y
City FE	Y	Y
Country Fe	N	N

*p < .1; **p < .05; ***p < .01

1261 S8.5 Accounting for Days with No Protests

1262 Table A15 shows attempts to account for days with no protest. Model 2 drops all days with
 1263 no protest images. Model 3 is a Poisson model. Model 4 is a negative binomial, and Model 5
 1264 is a zero-inflated negative binomial model. To converge, Model 5 excludes city fixed effects and
 1265 clustered standard errors; it does use country fixed effects.

Table A15. Count Models

	DV: $\text{Log}_{10}(\text{Sum of Faces})_{i,t}$ Original	DV: $\text{Sum of Faces}_{i,t}$ No Zero Days	Poisson	Negative Binomial	Zero-inflated Negative Binomial
	(1)	(2)	(3)	(4)	(5)
Perceived Prtstr. Violence _{i,t-1}	-0.167** (.0677)	-0.209* (0.121)	-0.798*** (0.120)	0.082 (0.641)	-0.775 (0.504)
Perceived Stt. Violence _{i,t-1}	1.282*** (.3327)	1.175*** (0.447)	9.587*** (0.343)	10.271*** (2.569)	7.575*** (2.089)
Perceived Stt. Violence _{i,t-1} ²	-2.103*** (.6093)	-1.834** (0.768)	-17.328*** (0.811)	-15.779*** (4.462)	-12.593*** (3.629)
Police _{i,t-1}	0.763* (.4493)	1.519*** (0.345)	2.631*** (0.057)	1.843 (1.799)	2.121* (1.122)
Fire _{i,t-1}	0.101*** (.0236)	0.065** (0.025)	0.124*** (0.011)	0.404** (0.181)	0.272** (0.118)
Gender Diversity _{i,t-1}	-0.113 (.0939)	-0.135* (0.078)	-1.315*** (0.045)	-0.285 (0.479)	-0.792*** (0.300)
Race Diversity _{i,t-1}	0.068 (.0440)	0.073 (0.050)	0.513*** (0.027)	0.168 (0.321)	0.684*** (0.206)
Age Diversity _{i,t-1}	0.020 (.0289)	0.025 (0.045)	0.666*** (0.024)	0.345 (0.244)	0.094 (0.170)
Tweets _{i,t-1}	0.009*** (.0033)	0.006*** (0.001)	0.009*** (0.001)		0.012 (0.010)
DV _{i,t-1}	0.158*** (.0682)	0.126*** (0.046)	0.002*** (0.0003)	0.017*** (0.004)	0.002 (0.003)
Intercept	0.126*** (.0237)	0.507*** (0.045)	1.221*** (0.024)	1.018*** (0.181)	2.150*** (0.117)
N	4,376	1,442	4,376	4,376	4,376
City FE	Y	Y	Y	Y	N
Cluster SE	Y	Y	Y	Y	N
Adjusted R ²	0.245	0.199			
Log Likelihood			-19,104.240	-4,387.215	-4,220.237

*p < .1; **p < .05; ***p < .01

1266

S8.6 Weighted Results

1267

Weighting each city-day observation by the number of protest photos shared from it strengthens the paper's results. Race and gender support critical mass theory. The free riding dynamics are more pronounced. The violence coefficients are much larger than the unweighted models, and model fit is almost 50% better than the paper's main models.

1268

1269

1270

1271

S8.7 Most Likely Protest Tweets

1272

1273

The last two models select tweets based on features that increase the probability they come from a protest. Model 5 restricts tweets to those only from mobile devices, based on the source

Table A16. Results Weighted by Protest Tweets per City

	Original (1)	Violence (2)	.outcome Demographics (3)	Combined (4)
Perceived Prtstr. Violence _{i,t-1}	-.1674** (.0677)	-.3383** (.1379)		-.3214** (.1357)
Perceived Stt. Violence _{i,t-1}	1.2820*** (.3327)	2.1774*** (.3740)		2.4331*** (.3725)
Perceived Stt. Violence ² _{i,t-1}	-2.1030*** (.6093)	-3.9986*** (.6880)		-4.4522*** (.6827)
Police _{i,t-1}	.7626* (.4493)	1.3091*** (.1440)		1.2754*** (.1438)
Fire _{i,t-1}	.1009*** (.0236)	.0198*** (.0140)		.0108*** (.0139)
Gender Diversity _{i,t-1}	-.1126 (.0939)		-.2087*** (.0517)	-.2538*** (.0499)
Race Diversity _{i,t-1}	.0683 (.0440)		.1434*** (.0326)	.1289*** (.0318)
Age Diversity _{i,t-1}	.0203 (.0289)		-.0626** (.0311)	-.0388 (.0298)
Tweets _{i,t-1}	.0095*** (.0033)	.0026*** (.0004)	.0027*** (.0004)	.0030*** (.0004)
DV _{i,t-1}	.1578** (.0682)	.2564*** (.0276)	.2986*** (.0332)	.2791*** (.0328)
Intercept	.1260*** (.0237)	.3940*** (.0474)	.5721*** (.0494)	.4906*** (.0490)
N	4,376	4,376	4,376	4,376
City FE	Y	Y	Y	Y
Cluster SE	Y	Y	Y	Y
Adjusted R ²	.2450	.5533	.5253	.5675

*p < .1; **p < .05; ***p < .01

field Twitter provides with each tweet. If that field contains “Twitter Web Client” or “Hootsuite”, the tweet is discarded; this paring leaves 3,743 tweets and 3,129 city days. The results for state violence and large groups match the full model, though with less statistical significance; the other covariates of interest lose statistical significance. The mobile model also fits the data less than half as well as the full model. Finally, we keep only tweets issued between 10 a.m. and 10 p.m., the most likely protest windows. Model 6 shows the results from these 4,664 tweets and 3,134 city-days. The result is a mixture of Models 4 and 5: the violence variables are larger and more precisely estimated, but none of the social cleavage variables are statistically significant, and the results for free-riding do not change.

Table A17. Most Likely Protest Tweets

	DV: $\log_{10}(\text{Sum of Faces})_{i,t}$		
	Original	Source Mobile	Protest Time
	(1)	(2)	(3)
Perceived Prtstr. Violence _{i,t-1}	-.1674** (.0677)	-.0438 (.1005)	-.2242*** (.0586)
Perceived Stt. Violence _{i,t-1}	1.2820** (.3327)	.6405** (.3113)	1.3883*** (.4637)
Perceived Stt. Violence _{i,t-1} ²	-2.1030*** (.6093)	-1.0014** (.4587)	-2.1913*** (.8201)
Police _{i,t-1}	.7626* (.4493)	-.0026 (.1366)	.6971* (.4005)
Fire _{i,t-1}	.1009*** (.0236)	.0578*** (.0214)	.1162*** (.0278)
Gender Diversity _{i,t-1}	-.1126 (.0939)	.0168 (.0466)	-.0933 (.0800)
Race Diversity _{i,t-1}	.0683 (.0440)	-.0126 (.0338)	.1058* (.0603)
Age Diversity _{i,t-1}	.0203 (.0289)	-.0066 (.0246)	.0088 (.0285)
Tweets _{i,t-1}	.0095*** (.0033)	.0097*** (.0015)	.0100*** (.0033)
DV _{i,t-1}	.1578** (.0682)	.1121** (.0534)	.1078 (.0661)
Intercept	.1260*** (.0237)	.1374*** (.0144)	.1381*** (.0161)
N	3,164	3,063	3,067
City FE	Y	Y	Y
Cluster SE	Y	Y	Y
Adjusted R ²	.2450	.1091	.2101

*p < .1; **p < .05; ***p < .01

¹²⁸³

S8.8 Models by Country

¹²⁸⁴

Pakistan is not included because it has too few observations.

¹²⁸⁵

S8.9 Investigating Fire, Police Variables

Table A18. Tables by Country

	DV: $\log_{10}(Sum\ of\ Faces)_{i,t}$				
	Spain	Hong Kong	South Korea 2014-2015	Venezuela 2017	Venezuela
	(1)	(2)	(3)	(4)	(5)
Perceived Prtstr. Violence _{i,t-1}	.3920** (.1941)	-.4626* (.2711)	-.6093* (.3254)	.1124 (.1231)	-.0955 (.0681)
Perceived Stt. Violence _{i,t-1}	.7175** (.3446)	3.7867*** (1.4235)	-1.3734 (1.4230)	.7740 (.5308)	.6684* (.3534)
Perceived Stt. Violence _{i,t-1} ²	-1.6441*** (.5217)	-8.2571** (4.1648)	2.6850 (3.3816)	-1.1302 (.8877)	-.9714 (.6119)
Police _{i,t-1}	1.2854*** (.1970)			.0876 (.4707)	-.1225 (.2208)
Fire _{i,t-1}	.0178 (.0496)	.0715 (.0948)	.3637** (.1457)	-.0451* (.0242)	.0355 (.0317)
Gender Diversity _{i,t-1}	.0642 (.0777)	-.0201 (.1480)	-.0101 (.1195)	-.1568* (.0866)	.0172 (.0855)
Race Diversity _{i,t-1}	-.0781 (.0607)	.1521 (.1221)	-.1716 (.1198)	-.0885 (.0612)	.0340 (.0600)
Age Diversity _{i,t-1}	.0032 (.0372)	-.1156 (.0908)	.0316 (.0665)	.1542*** (.0525)	-.0456 (.0490)
Tweets _{i,t-1}	.0047*** (.0011)	.0098 (.0072)	-.0003 (.0035)	.0153*** (.0033)	.0171*** (.0038)
DV _{i,t-1}	.0680 (.0430)	.0763 (.1200)	-.0795 (.0751)	.0771 (.0583)	.0190 (.0353)
Intercept	.5900*** (.0369)	.1034*** (.0261)	.0004 (.0333)	.0322 (.0357)	.0164 (.0168)
N	1,412	257	365	573	1,752
Adjusted R ²	.2861	.1083	.1185	.6110	.0962
City FE	Y	Y	Y	Y	Y
Cluset SE	N	N	N	N	N

*p < .1; **p < .05; ***p < .01

Table A19. Fire and Police Variables on Their Own

	DV: $\log_{10}(Sum\ of\ Faces)_{i,t}$			
	(1)	(2)	(3)	(4)
Police _{i,t-1}	.4894 (.5121)	.4932 (.5071)		
Fire _{i,t-1}	.1011*** (.0184)	.0945*** (.0206)		
Gender Diversity _{i,t-1}		-.1148 (.0875)		-.1358 (.0978)
Race Diversity _{i,t-1}		.0708* (.0410)		.0968** (.0449)
Age Diversity _{i,t-1}		.0217 (.0288)		.0208 (.0315)
Perceived Prtstr. Violence _{i,t-1}			.0237 (.0627)	.0227 (.0649)
Perceived Stt. Violence _{i,t-1}			1.0602*** (.3149)	1.0529*** (.3342)
Perceived Stt. Violence _{i,t-1} ²			-1.6560*** (.5213)	-1.6396*** (.5458)
Tweets _{i,t-1}	.0095 (.0035)	.0096*** (.0034)	.0107*** (.0040)	.0108*** (.0037)
DV _{i,t-1}	.1926 (.0743)	.1728** (.0704)	.1954** (.0805)	.1684** (.0711)
Intercept	.1239 (.0158)	.1275 (.0158)	.1224 (.0159)	.1273 (.0159)
N	4,376	4,376	4,376	4,376
Adjusted R ²	.2373	.2389	.2304	.2334

*p < .1; **p < .05; ***p < .01

1286 **APPENDIX S9. INTER-CODER RELIABILITY**

1287 We used Fleiss' Kappa to measure the inter-coder reliability of our training image annotations.
1288 In many cases, the inter-coder reliability is typically measured on the coding data on which the
1289 actual analysis is conducted. In our study, the manual coding was performed on the training data,
1290 and the reliability was measured for the annotations to ensure that the models are trained in a
1291 consistent manner. Table A20 shows the estimated reliability statistics.

Table A20. Inter-coder reliability

Label	Kappa
Perceived Violence	.316
Perceived Protester Violence	.566
Perceived State Violence	.473
Large Group	.434
Small Group	.388
Police	.564
Fire	.702
Child	.457