



University of Dundee

### Hyper-resolution monitoring of urban flooding with social media and crowdsourcing data

Wang, Ruo Qian; Mao, Huina; Wang, Yuan; Rae, Chris; Shaw, Wesley

Published in: Computers and Geosciences

DOI: 10.1016/j.cageo.2017.11.008

Publication date: 2018

Document Version Peer reviewed version

Link to publication in Discovery Research Portal

Citation for published version (APA): Wang, R. Q., Mao, H., Wang, Y., Rae, C., & Shaw, W. (2018). Hyper-resolution monitoring of urban flooding with social media and crowdsourcing data. Computers and Geosciences, 111, 139-147. https://doi.org/10.1016/j.cageo.2017.11.008

#### **General rights**

Copyright and moral rights for the publications made accessible in Discovery Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

• Users may download and print one copy of any publication from Discovery Research Portal for the purpose of private study or research.

You may not further distribute the material or use it for any profit-making activity or commercial gain.
You may freely distribute the URL identifying the publication in the public portal.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

# Hyper-resolution Monitoring of Urban Flooding with Social Media and Crowdsourcing Data

Ruo-Qian Wang<sup>1,2</sup>, Huina Mao<sup>3</sup>, Yuan Wang<sup>4</sup>, Chris Rae<sup>5</sup>, Wesley Shaw<sup>5</sup>
1 Department of Civil and Environmental Engineering, University of California, Berkeley, CA
94720
2 School of Science and Engineering, University of Dundee, Dundee, UK DD2 1BW
3 Oak Ridge National Laboratory, Oak Ridge, TN 37831
4 Department of Civil and Environmental Engineering, Tufts University, Medford, MA 02155
5 Blue Urchin LLC, 419 11<sup>th</sup> Ave E, Seattle, WA 98102

### 10 Abstract

11 Hyper-resolution datasets for urban flooding are rare. This problem prevents detailed flooding risk analysis, urban flooding control, and the validation of hyper-resolution numerical models. 12 13 We employed social media and crowdsourcing data to address this issue. Natural Language 14 Processing and Computer Vision techniques are applied to the data collected from Twitter and 15 MyCoast (a crowdsourcing app). We found these big data based flood monitoring approaches can 16 complement the existing means of flood data collection. The extracted information is validated 17 against precipitation data and road closure reports to examine the data quality. The two data 18 collection approaches are compared and the two data mining methods are discussed. A series of 19 suggestions is given to improve the data collection strategy.

### 20 1. Introduction

21 Urban flooding is a global problem that costs lives and money. In 2010 alone, 178 million people 22 suffered from floods. The total economic losses in 1998 and 2010 both exceeded \$40 billion (Jha 23 et al., 2012). Urban floods can be caused by a variety of reasons, including natural hazards of 24 river overflow, coastal storm surge, sea-level rise, flash floods, groundwater seepage, sewer 25 overflow, lack of permeability, and lack of city management. As urbanization proceeds and 26 climate change intensifies, urban planners and city managers are facing the challenge of preparing 27 for and mitigating flood damage. They need tools to monitor and predict the event for emergency 28 response and development planning.

29 Monitoring and predicting urban floods needs high-resolution data with good coverage. High-

30 resolution data can capture the variation of flood flows among streets or parcels, so that the

31 heterogeneity of flood flows caused by heterogenous urban landscape can be captured. In this

- 32 study, we define data that can reflect the variation on the parcel and street scale as "hyper-
- 33 resolution" data. In addition to resolution, it is important to have a good coverage of flood data to
- 34 obtain complete information.
- 35 The traditional method of obtaining flood related data lack both resolution and coverage. Remote
- 36 sensing is a commonly used data source. Aerial photography, for instance, is being conducted by
- 37 many research teams, engineering companies, emergency response services as well as
- 38 governmental departments, and has demonstrated its value (Marcus and Fonstad, 2008).
- 39 However, the systematic application of aerial photography is limited by vegetation canopies and

40 cloud cover during floods (Hess et al. 1995, 2003). The limitation also exists in radar and satellite

- 41 imagery that is based on optical sensors (Wilson et al., 2007). The most realistic and feasible
- 42 approach is considered to be microwave remote sensing, which penetrates cloud cover. However,
- 43 because of the corner reflection principle (Rees, 2001) along with coarse ground resolution, this
- technology is currently unable to extract flood data from urban areas. Another commonly used
- 45 data source is distributed sensors. In the U.S., sensors have been installed in coastal areas by
- 46 National Oceanic and Atmospheric Administration (NOAA) and in rivers and canals by United
- 47 States Geological Survey (USGS), but almost no sensors are distributed on streets that are
- 48 dedicated for urban flood monitoring purpose.
- 49 There are very few existing urban flooding datasets that can be used for detailed model validation
- and urban planning. A dataset about a flood event on January 10, 2005, in the City of Newport
- 51 Beach, California, was built by a collection of 85 digital photographs and eyewitness reports from
- 52 city employees who were dispatched to manage and photo-document the flood. The data
- 53 collection involved interviews and email communications (Gallien et al., 2011). Another example
- 54 in UK showed that contemporary newspaper reports can enrich the evidence of the witness
- reports and photos taken on the day following the storm (Smith et al, 2012). In addition,
- 56 insurance reports can provide some additional information of flooding events such as damage
- 57 evaluation. However, recording of such information is not a priority for civil defense agencies
- during a major coastal disaster, and data collected using this method could be expensive to access,
- 59 incomplete in coverage, and substantially delayed.
- 60 Since urban area usually has a dense population, a big data approach, which relies on volunteered
- 61 data reports from citizens, could be a potential solution to provide primary or complementary
- 62 urban flooding data. The history of using social media data to study natural hazards can be traced
- back to Muralidharan et al. (2011), who compared the difference between non-profit
- organizations and media in the use of Twitter and Facebook during the Haiti earthquake in 2010.
- 65 Since then, a series of studies based on social media data have emerged with a focus on floods.
- 66 Sun et al. (2013) mentioned their effort to use Flickr images to support remote sensing based
- flood maps. Jongman et al. (2015) explored the potential to use Twitter data for early detection of
- flood events in Philippines and Pakistan. Fohringer et al. (2015) was the first to use flooding
- 69 water depth information, which was manually extracted from the photos posted in Twitter.
- Eilander et al. (2016) studied the floods of Jarkata in Indonesia with Twitter data. They
- 71 introduced a probability map to quantify the data uncertainty and found the general accuracy of
- 12 location is around 69% but rises to over 90% if the specific location was mentioned in the content
- of the posts.
- 74 These studies clearly showed the value of data mining in flooding research, but their application
- is still limited due to the poor resolution of data collected. For example, the resolution of
- Jongman et al. (2015) is at city level. To achieve sub-city resolution, researchers heavily relied on
- 77 manual reading as used in Sun et al. (2013) and Fohringer et al. (2015). The best resolution based
- on automatic algorithms was probably obtained by Eilander et al. (2016). Based on a group of
- independent tweets reporting flood depth, they constructed a probability map on the scale of
- 80 urban community with an average size of a few square kilometers. To obtain the location
- 81 information, the authors used manually defined location names and gazetteers to match location

- 82 mentions in tweets. To support disaster monitoring, relief responses, model validation, and
- 83 decision making, we still need higher data quality and accuracy to fully understand the details of
- 84 flooding events.
- 85 The present paper is aimed at demonstrating a new approach to collect and process data for urban
- 86 flooding research using Natural Language Processing and Computer Vision (CV) techniques.
- 87 These techniques are shown promising to extract hyper-resolution data with a wide coverage to
- 88 support urban flooding issues.
- 89 2. Data Source
- 90 2.1 Twitter
- 91 Twitter data was streamed using Twitter API during the period from September 29 to October 28
- 92 2015 and with the filtering keywords of "flood", "inundation", "dam", "dike", and "levee".
- 93 Retweeted posts labeled by Twitter were excluded. There are 7,602 tweets obtained.
- 94 2.2 Crowd-sourcing photos
- 95 We employed the crowd-sourcing platform MyCoast to collect photos about urban flooding.
- 96 MyCoast is a system that, since 2013, has been used by a number of US environmental agencies
- 97 to collect "citizen science" data about various coastal hazards or incidents. The app is built using
- 98 Ionic, a user interface framework based on the Cordova cross-platform app framework. The web
- and data storage component is based upon WordPress and the app communicates with the server
- 100 via a custom REST API. Data is stored in a MySQL back end. The app interface has been
- 101 designed to be intuitive to those unfamiliar with phone apps, as the citizen scientist users often
- skew towards the youth or elderly end of the age spectrum.
- 103 There are two ways to contribute to the database, i.e. via the web and a mobile app. The user
- 104 interface of the app is shown in Figure 2. The system currently contains over 5,000 flood
- 105 photographs, and most of the photos were collected through the mobile app.

= ¢	MyCoast	🗮 🍟 King Tide	🗮 🛛 🖬 King Tide	
1 1	k	14	Estimated water depth	
		Take photo	Ankle (4-6 in)	
North Charleston     Jun 18	North Charleston Jun 18	Select from library	) Vaca (1.0.6)	
High tide at Co at 8×	byote Creek in 3 hours 42 pm (8.9 ft) •	Date of Photo	Knee (1-2 II)	
		16 May 2017	Chest (5+ ft)	
		Time of Photo	Water level appears to be	
	N and There are	06:05 PM	Rising	
© Remind Mel	Pleasty O Hourly	Proximity of photo to worst affected		
		At center of flooding	Retreating	
Share location with a coordinators	state	At edge of flooding	Comments	
• •	dd Report	Page 1 of 2	Page 1 of 2	
		Set Location	Set Location	

#### 107 Figure 1 User interface design of MyCoast.

108 The app uses the phone's sensors to establish location and date/time information. Users are

- 109 prompted to take a photograph and then may optionally add written comments (Figure 2). Users
- are also shown a chart with tide timing so that they can try to optimize the timing of photographs
- 111 with peak tides.

#### 112 2.3 Authorized Data

- 113 In the scale of the United States, a reliable data related to floods is precipitation statistics. The
- 114 monthly precipitation data of October 2015 was downloaded from Advanced Hydrologic
- 115 Prediction Service (AHPS), a database developed by National Weather Service (NWS)
- 116 (https://water.weather.gov/precip/download.php). The precipitation departure of this month was
- 117 the difference of the observation from the monthly normal precipitation obtained by Parameter-
- elevation Regressions on Independent Slopes Model (PRISM) using the record from 1981 to
- 119 2010. More details can be found in the NWS website.
- 120 Road closure information was collected using local media reports of September 29, 2015. We
- 121 identified two major data sources, including the reports from WCBD online news at
- 122 http://counton2.com/2015/09/29/coastal-flooding-closes-streets-downtown-second-day-in-a-row/
- 123 and the Twitter posts by Ashley Rae Yost, a reporter for WCBD television station in Charleston,
- 124 SC. Integrating the data points, we can list the flooded roads on the day in Table 1.

#### 125 Table 1 Collection of road closure data from media reports.

Name of the road	Starting point	End Point
Lockwood Drive	Gadsden St	N/A
Wentworth Street	Lockwood Blvd.	Gadsden St.
I-26	N/A	N/A
Lockwood Blvd	Beaufain Street	N/A
Broad St.	Barre St.	N/A
Barre St.	Beaufain	Montague
Lockwood St	N/A	N/A
Calhoun St.	Fourth St.	Halsey Blvd.
Market St.	Church St.	East Bay Drive
Long Point Road	Marsh Area	Marsh Area

126

### 127 3. Methods

- 128 3.1 Natural Language Processing methods
- 129 Location information is crucial for mapping floods. However, geotagged tweets are only about
- 130 1% of all the twitter data (Middleton, 2014), which poses a challenge for mapping. Even though
- 131 the tweets are geotagged, the geotag coordinates may not be the same as the location of
- 132 mentioned floods. For accurate location mapping, we aim to extract location mentions within the
- 133 tweets' text. In addition, we extract the quantifier information, e.g. the depth of floods, which can
- help us understand the damage level of flooding. Table 2 shows several examples of tweets
- 135 containing flooding location and depth information.

136 Table 2 Sample flooding tweets that are processed by NLP (location names in **bold** and depth numbers are

#### 137 underlined).

Tweet_ID	Posted_time	Tweet
648973656161394688	2015-09-29 14:32	Roanoke River over <u>12 feet</u> at <b>Walnut St bridge</b> in <b>Roanoke</b> now expected to top 13. Flood stage is 10. #swvawx",
649935518038405120	2015-10-02 06:14	<i>Flood currently on 3 NE Wrightsboro</i> . in New Haven, NC. <u>1-2</u> <u>ft</u> . of water on Tandem CT. #ncwx #flood #flooded #flooding #rain #HurricaneJoaquin
650019989353852928	2015-10-02 11:50	<i>Chicod creek at flood stage in</i> <b>PItt County</b> , <i>Likely to</i> <u>11 ft</u> by tonight. If it hits 12, it's running over the road.
650393844115243008	2015-10-03 12:35	Helped nearby drivers by reporting a flood on <b>PR-1, Ponce</b> on @username - Drive Social.
650403125698891776	2015-10-03 13:12	The Market Street portion of Water Street in Downtown Wilmington is under 9 inches of water

138 139

140 In order to extract location names from tweets, we use the named entity recognition methods.

141 Named Entity Recognition (NER) is a fundamental task in NLP, which aims to classify words

142 into different types of names such as person, organization, location. Stanford's NER

143 (https://nlp.stanford.edu/software/CRF-NER.shtml) is one of the cutting-edge named entity

144 recognition tool available that uses Conditional Random Field (CRF) based classifier (Finkel et

al., 2005). CRF is a conditional sequence model: given observations, it aims to find the highest

146 possible sequence of states. In NER, the observations refer to words, and states refer to named

147 entity tags. The entity tag for each word in a sentence is not predicted independently, but by148 considering the tags of neighboring words.

149

150 The original Stanford NER tool is trained on formal texts such as news data, which is remarkably

151 different from short, informal, and noisy twitter data. So, applying the pre-trained NER model can

152 generate poor performance in analyzing tweets. Lingad et al. (2013) have shown that NER tool

retrained on Twitter data can significantly boost the performance on location detection from

154 tweets. Therefore, we retrained the Stanford NER model using annotated tweets.

**Data preprocessing**: We pre-processed the data as follows. First, we remove non-english tweets

156 (the python library "guess\_language" is used to identify if a given tweet is English or not.) After

157 filtering, 6,593 English tweets remain. Second, url links and mentioned usernames in tweets are

replaced with unique words, "<URL>" and "@username", respectively. Third, each tweet is

tokenized into as a list of words, which is fed into NER tool for location recognition.

### 160 **NER model training and testing:**

161 The training and testing data is provided by the ALTA 2014 Twitter Location Detection shared

162 task. The original training and testing sets include 2,000 and 1,000 tweet ids, respectively. Since

some tweets have been deleted or become invalid, we have obtained 1,851 in training and 930

164 tweets in testing. The training data is randomly split into 1,600 for training and 251 for validation.

165 The annotated locations include place names and point-of-interests (POIs), in either main text

166 strings or hashtags.

167

- 168 We have compared the performance of the original Stanford NER model and the re-trained model
- based on Twitter data. Results are evaluated in three metrics: precision, recall, and F-score, which
- are defined as:
- 171 Precision = TruePositive / (TruePositive + FalsePositive)
- 172 Recall = TruePositive / (TruePositive + FalseNegative)
- 173  $F-Score = 2 \times Precision \times Recall / (Precision + Recall)$
- 174 So, precision is the number of positive (or correct) results divided by the total number of all
- 175 returned results (i.e. results that are predicted as relevant, which include both true positive and
- 176 false positive samples), while recall is the number of positive (or correct) results divided by the
- 177 total number of the actual relevant results. F-score combines precision and recall.
- 178 As it can be seen in Table 3, F-score is significantly improved after retraining. Even though
- 179 original Stanford NER has a high precision rate, the recall is low (32.07%). In other words, a
- 180 large number of location names are missed by this method. We have found that since the original
- 181 NER model is trained based on the formal news data, it cannot capture many location information
- 182 expressed in tweets, which are considered as informal language (e.g. abbreviations, misspellings,
- 183 hashtags). Therefore, we use the re-trained Stanford NER model to detect locations from our
- 184 flooding tweets.
- 185 Table 3 Results of location detection in ALTA 14 datasets.

Models	Precision	Recall	F-score
Stanford NER	94.51%	32.07%	47.88%
Stanford NER retrained	86.68%	69.72%	77.28%

<sup>186</sup> 

**Geocoding:** From above, we have obtained a list of location names extracted from tweets. At this step, we need to geocode location names into geographical coordinates. Bing Maps Location API is used via Python's geocoding library (*geopy*), which accepts a location string as the input, and

190 generates address name and latitude/longitude coordinates. For example, "*Battleship Rd*"

- identified from a tweet is converted to a standard address as (Battleship Rd, Camden, SC 29020,
- 192 United States, (34.26021, -80.62693)). If more than one result are found by Bing Maps, the one
- 193 with the highest probability is selected.

Flood depth information extraction: in order to capture the water depth mentions, we define
regular expression patterns to capture numbers followed by "ft", "feet", "inch", "in" and their
plural forms.

- 197 3.2 Computer Vision
- 198 The present study employs the technique of Convolutional Neural Networks (CNN) to
- automatically classify the crowdsourcing photos (Zeiler and Fergus, 2014). Originally inspired by
- 200 animal visual perceptron, CNN is often used to detect and recognize objects. A good example of

- 201 CNN algorithm is Clarifai, which has won the competition of ImageNet Large Scale Visual
- 202 Recognition Competition in 2013 (ILSVRC-2013). We applied this code through Clarifai API,
- 203 which can be accessed as a remote web service. This service receives feedings of individual
- 204 photos and feedbacks a list of tags that describe the objects present too busy to clean in the
- 205 photo. For each tag, a probability is assigned to quantify its confidence of the recognition. Three
- 206 examples of the recognition can be found in Figure 2 and the corresponding processed results are
- shown in Table 4.



Latitude: 32.78208 Longitude: -79.9446 Time: 16:38:00 Date: 09/29/2015



Latitude: 42.083821 Longitude: -70.6459 Time: 13:43:00 Date: 10/18/2016

Figure 2 Sample photos of the crowdsourcing data. Each photo is labeled by coordinates and time information

- as shown below the corresponding photo. The first photo shows a flooded street with trapped cars, bus, and
- trucks. The second photo has a water paddle at a port, but no spreading flood was observed. The third photo
- 212 shows a river filled by water, but no flood can be recognized.
- 213 Table 4 Sample results of the computer vision applied to the sample photos shown in Figure 2.

Photo 1		Photo 2		Photo 3	
Tag	Probability	Tag	Probability	Tag	Probability
-	(P)	-	(P)	-	(P)
flood	0.9964	no person	0.9678	flood	0.9961
rain	0.9963	wood	0.9656	water	0.9958
water	0.9940	furniture	0.9493	calamity	0.9904
tree	0.9834	chair	0.9381	storm	0.9816
river	0.9804	seat	0.9314	house	0.9796
reflection	0.9787	room	0.8956	building	0.9645
storm	0.9745	luxury	0.8942	river	0.9572
road	0.9696	industry	0.8807	architecture	0.9561
canal	0.9686	house	0.8705	no person	0.9530
no person	0.9652	water	0.8282	seashore	0.9426

214

208

#### 215 Flood detection:

Flood is considered to happen if the flood tag was labeled. Applying this criterion to the photos in

Figure 1, we found that flood was detected correctly in the first photo and no flood was correctly

218 detected in the second, but a flood tag was incorrectly labeled to the third, which is a river in its

219 normal condition. A manual check of 80 photos found that the accuracy is 65%, which is

- 220 comparable to a text-image correlation study of social media (Chen et al., 2013) and an
- 221 application study in earthquake damage estimation (Nguyen et al., 2017).

Visualization: The crowd-sourced data can be shown in Google map using the GPS coordinates contained in the labels of each photo. The road closure data were manually relocated, being

- visualized by a line connecting the starting and end points of the road closure section. Because
- the end points are missing in some road closure records, we only marked the starting points in
- these cases.
- 227

### 4. Results

4.1 Case study 1: the flooding map of the United States

230 The daily volume of flood related Tweets of the United States is shown in Figure 3. The daily volume increases quickly from October 1<sup>st</sup> to October 4<sup>th</sup>. This trend is consistent with the history 231 of the Hurricane Joaquin, which reached the Category 4 major hurricane on October 1<sup>st</sup> and 232 gained its greatest strength on October 3<sup>rd</sup>. Interestingly, the daily volume reached its top on 233 October 4<sup>th</sup>, a day later than the greatest strength of the hurricane. This phenomenon might be 234 attributed to the reason that the flood event needed time to develop before flushing into the 235 236 residential areas. Note that this one-day delay is also reported in the Philippine floods (Jongman 237 et al., 2015). Immediately after the peak, we observe a downturn, which may be partially due to

- the incomplete data streamed from the Twitter API. Despite the incompleteness, this downturn
- 239 was also observed by Jongman et al. (2015) in two of the Philippine floods, but missing in other
- reports, such as Eilander et al. (2016). We suspect that people were too busy cleaning up the
- 241 disaster and the hurricane became of less value in news reports after diminishing.
- 242 The daily volume and the percentage of the hyper-resolution tweets that contained specific
- location are shown in Figure 3. The hyper-resolution tweets generally follow the trend of the total
- flooding related Tweets, and the average percentage is around 51%. This percentage is higher
- than expected and might indicate that most of the hyper-resolution tweets were posted by
- authorized agencies. Note that another hurricane developed on October 24, so the Tweeting
- volume had a second peak.



Figure 3 The daily volume of flooding related tweets and flooding tweets with hyper-resolution location
 information in the United States. The cross marks are the percentage of the hyper-resolution information among
 all the flooding Tweets.

Table 5 shows the number of identified flooded roads in the top 5 states. The top 3 numbers are

253 consistent with the states that were impacted by Hurricane Joaquin. California (CA) and Texas

254 (TX) are listed next to these states, because they are large and populate with high volume of

tweets.

#### 256 Table 5 Number of flooded roads by states: top 5 states.

State code	Number of roads
SC	57
NC	41
FL	21
CA	21
ТХ	17

#### 257

258 The geo-spatial distribution of the hyper-resolution flooding tweets are shown in Figure 4. The

259 identified tweets concentrated at the east coast capturing the extensive floods caused by the two

260 hurricanes. The density of the tweets decreases toward Midwest but increases at the west coast

and Gulf of Mexico.

To examine the reliability of the Twitter based flooding monitoring, we compute the correlation between the precipitation pattern and tweet volume. Using a grid of 50×50 cells, Figure 5

264 compares the geospatial patterns of precipitation observation, the precipitation departure, the

265 percentage of precipitation departure, and the tweet volume. The precipitation data was obtained

266 by averaging the data within the cell. The tweet volume was calculated by counting the tweets in

267 each cell. A common feature among all the patterns is the high concentration around South

268 Carolina, which captures the hurricane caused flooding. The tweet volume pattern has a high

269 value in the stripe of Washington DC – New York – Boston, which is missing in all the other

270 patterns. This mismatch indicates that urban areas that have dense population reported floods

271 relatively more than remote places. This mismatch may bias the Twitter based flooding

272 monitoring. Another mismatch is in Florida, where a high volume of flooding tweets is found

- there without heavy precipitation or precipitation departure. This finding may be attributed to the
- fact that the Florida floods were caused by tidal and storm surge. A high volume of flooding
- tweets is spotted in California, overlapping the precipitation-departure-percentage pattern. So we
- infer that the Californian floods during that period were caused by departure percentage of
- 277 precipitation and California is more sensitive to the precipitation departure percentage than the
- absolute departure.
- 279 The Pearson correlation coefficient is used to examine the correlation between precipitation
- departure and tweet volume. We use different grids with the sizes of  $100 \times 100$ ,  $50 \times 50$ ,  $25 \times 25$  and
- 281 12×12 to cover the contiguous United States. The Pearson correlation coefficient is calculated at
- different grid resolution in Table 6. The correlation coefficient is in the range of 0.17-0.41, so the
- tweet volume pattern has a weak linear correlation with the precipitation departure. Table 6 also
- shows that the correlation coefficient increases with coarse mesh, which suggests that Twitter
- 285 based flooding monitoring is more reliable with high tweet volume.



286





Figure 4 The spatial distribution of identified Tweets that have hyper-resolution geo-location information. Text
 content of sample posts is shown in the right panel.





Figure 5 Spatial patterns of precipitation and Tweets in October, 2015. (a) The observed precipitation. (b) The

- 295 precipitation departure. (c) The percentage of precipitation departure. (d) The statistics of the Tweet volume.
- 296

297 Table 6 The Pearson correlation coefficients at different resolution of grids.

Grid size	Pearson correlation coefficient
100×100	0.1693
50×50	0.3156
25×25	0.3759
12×12	0.4095

298

The Ripley's K function is used to estimate the spatial pattern of the hyper-resolution flooding
 tweets. The Ripley's K function is a spatial analysis method to describe how point patterns cluster

301 or disperse comparing to randomly distributed points over a given area of interest. The K function

302 is defined (Dixon, 2002),

303 
$$K(t) = A \sum_{i=1}^{n} \sum_{j=1}^{n} I_{t(i,j)} / n^2$$
, (1)

where A is the area covered by the points, i and j are the indices of the points, n is the sample size,

305  $I_t$  is the impact function, which is one if the distance between points *i* and *j* is less than distance *t* 306 and zero otherwise. A more commonly used quantity is the L function defined by

307 
$$L = \sqrt{K(t)/\pi} - t$$
. (2)

308 The positive value of the L function means the data points are more clustered than a random 309 process and the negative value means more dispersed. The L function is plotted against the distance t in Figure 6. The L maximum of around 190 km indicates that the data points are 310 311 clustered the most at this length scale. It is consistent with the length scale of urban areas in the 312 United States, e.g. the length scale of New York metropolitan area is around 185 km (Wikipedia, 313 2017, https://en.wikipedia.org/wiki/New York metropolitan area). This indicates that the twitter 314 pattern follows the population distribution in the US as expected. The L function drops to zero 315 after the length scale of 100 km. So at this length scale the twitter posts distribute close to random

316 pattern.



318 Figure 6 The Ripley's L function against the length scale *t*.

- 319 4.2 Case Study 2: Tidal flood at Charleston, SC
- 320 The most data collected through MyCoast is on September 29, 2015 in Charleston, SC, so
- 321 we used this dataset as a demonstration of our flood monitoring technique. The processed
- 322 crowdsourcing data as well as the road closure data are shown in Figure 7.
- 323 Crowdsourcing data concentrated along the coastline of Charleston and most of the
- 324 places were identified as "no flood" by the CV algorithm. The road closure data are
- 325 distributed relatively inland.



326 327

7 Figure 7 Comparison between crowdsourcing results and road closure data

328 The processed crowdsourcing data can be validated against the road closure data at four spots,

329 because the two datasets roughly overlap in these locations, which are marked in Figure 7. The

figure shows that crowdsourcing data correctly recognized the flood events at Spot 1 and 2, but

made mistakes in Spot 3 and 4. The photos collected at Spot 3 are shown in Figure 8. The first

and third photos show road floods, but the CV algorithm misclassified them. These mistakes

- 333 could be explained by the strong reflection in the first photo and the semi-submerged plants in the
- third photo. These features resemble the natural water bodies and pose difficulties to process. The
- second photo shows no road flooding, but the flood water has reached a bench which should not
- be threatened by tidal floods. Again, the semi-submerged plants might confuse the CV. The last
- 337 photo shows an overflowing drainage inlet. This photo indicates an overwhelmed drainage
- 338 system and a nearby flood, but no road flooding was shown. So the CV code should not be
- blamed. The crowdsourcing photo collected at Spot 4 is shown in Figure 9. The floodwater has a
- regular edge, so the computer vision considered it as a normal water canal in the city. CV
- 341 algorithm is proved useful to extract information, but the comparison shows there is a room to
- 342 improve.



Figure 8 Photos near Road closure #3. From left to right, the photos are corresponding to the four points of Spot
3 in Figure 7 from South to North.



345

346 Figure 9 Photo near Spot 4.

### 347 5. Discussion

- 348 Comparing the two big-data approaches, we found they have different characteristics that
- 349 potential users have to be cautious before direct use. Both methods can provide hyper-resolution
- 350 monitoring with the resolution of street and parcel scale. Crowd-sourcing data is more
- advantageous because the GPS location can be up to the accuracy of meters, while the Twitter
- based data is in the scale of street names. In terms of geolocation coverage, Twitter based
- approach might be more favorable, because a much wider space and a worldwide monitoring
- 354 could be feasible with little cost. However, we failed to locate any tweet at Charleston SC on
- 355 September 29, 2015 to be compared with our road closure data. This mismatch indicates that

- 356 Twitter might be more suitable for large-scale monitoring. Crowdsourcing is restricted to the
- 357 number of volunteers and distribution of app users, so its coverage can be specified in particular
- locations, but it might be still inappropriate to conduct large-scale monitoring. Regarding to the
- data volume and data collection speed, Twitter shows its advantage, so it may be more
- 360 appropriate for real-time monitoring and higher volume of data benefits the accuracy of Twitter
- 361 based monitoring.

362 An important issue that limits the big-data platform in practical use is data reliability.

Crowdsourcing provides rich and customized information through its photos and user interface 363 and the mobile phone app has high accuracy GPS sensors, so the collected data may have less 364 365 uncertainty. In comparison, Twitter based approach has significant noise, which needs data 366 cleaning before processing. However, the processing reliability might be opposite, for NLP is 367 more accurate than CV at present. We expect higher reliability in the information extracting 368 process from Twitter and other social media. Considering both of data collection and information 369 extraction, in the case of a large quantity of data and automatic processing algorithm has to be 370 involved, the general data reliability will be case by case, depending on the quality of data and 371 training techniques of the practitioner. If a small amount of data needs to be processed, manual 372 reading of crowdsourcing photos could be conducted, then crowdsourcing approach might be 373 better. This might be the reason that some previous studies used photo data as the ground truth to

374 validate numerical models.

375 On the subject of big data quality improvement, we outline some available measures and point to 376 a series of ongoing research that has a hope to generate promising data quality control tools. For 377 the crowdsourcing based approach, a photo shooting guidance should be provided so that a more 378 complete scene can be captured and more reference objects can be included. A CV code should 379 be retrained using the labeled crowdsourcing data to achieve high reliability. An emerging 380 research trend is to determine water depth from the crowdsourcing photos by comparing to 381 Google Street View. By detecting the difference between normal street view and the flooded 382 streets, water depth can be inferred by calculating how much reference objects are blocked by 383 floodwater. For Twitter based approach, in addition to the development of high accuracy NLP 384 algorithm, a recent trend is to include "citizen science" in a feedback loop, so that a data 385 elicitation can be sent to the original data contributor to request more information about the 386 original report to pinpoint the location and follow the new development of the flood events. 387 Water depth determined from the tweets are shown feasible in the present study, but its value has 388 not been fully explored. In the future, water depth data should be compared against remote 389 sensing data to cross-check the data quality. The photo contained in tweets can also be used to 390 provide information about the flood. Fully exploiting such information can enrich the text-based 391 monitoring effectiveness. A final note is that a good data fusion scheme, an approach to integrate 392 of multiple data and knowledge and resolving data conflicts representing the same real-world 393 object into a consistent, accurate, and useful representation, should be applied before using data 394 to validate numerical models and support city planning and emergency response preparation 395 (Liggins et al., 2017). This method has potential to greatly improve the data quality in general.

396 6. Conclusion

- 397 Urban flooding is difficult to monitor due to various complexities in data collection and
- 398 processing. The present study shows that social media and crowdsourcing can be used to
- 399 complement the datasets developed based on traditional remote sensing and witness reports.
- 400 Applying these methods in two case studies, we found these methods are generally informative in
- 401 flood monitoring. Twitter data is found weakly correlated to precipitation departure. We
- 402 determined a length scale of tweet volume pattern, at which the data points are most clustered.
- 403 The computer vision processed crowdsourcing data is compared against the road closure data.
- 404 The results show that computer vision still has a room to improve, especially in coastal areas.
- 405 These two methods are compared and a series of recommendation is given to improve the big
- 406 data based flood monitoring in the future.

### 408 Acknowledgement

- 409 The first author thanks for the great support from Profs. Mark Stacey and Alexei Pozdnoukhov,
- 410 and contributions from Wei Bai, Diyi Liu, Eason Ruan, Ruonan Ou, Renjie Wu, and Wenjun
- 411 Zhong. The lead author also gratefully acknowledges the support of NSF, grant CBET-1541181.
- 412 Huina Mao would like to thank the financial support for this research from the US government
- 413 for Oak Ridge National Laboratory's Laboratory Directed Research and Development (project
- 414 number LDRD 7677).

### 415 References

Eilander, D., Trambauer, P., Wagemaker, J. and van Loenen, A., 2016. Harvesting social media for generation of near real-time flood maps. *Procedia Engineering*, *154*, pp.176-183.

- Finkel, J.R., Grenager, T. and Manning, C., 2005, June. Incorporating non-local information into
- 419 information extraction systems by gibbs sampling. In *Proceedings of the 43rd annual meeting on*
- 420 *association for computational linguistics* (pp. 363-370). Association for Computational Linguistics.
- Fohringer, J., Dransch, D., Kreibich, H. and Schröter, K., 2015. Social media as an information
  source for rapid flood inundation mapping. *Natural Hazards and Earth System Sciences*, *15*(12),
  pp.2725-2738.
- Gallien, T.W., Schubert, J.E. and Sanders, B.F., 2011. Predicting tidal flooding of urbanized
  embayments: A modeling framework and data requirements. *Coastal Engineering*, *58*(6), pp.567577.
- 427 Guan, X. and Chen, C., 2014. Using social media data to understand and assess 428 disasters. *Natural hazards*, *74*(2), pp.837-850.
- 429 Hess, L.L., Melack, J.M., Filoso, S. and Wang, Y., 1995. Delineation of inundated area and
- vegetation along the Amazon floodplain with the SIR-C synthetic aperture radar. *IEEE Transactions on Geoscience and Remote Sensing*, 33(4), pp.896-904.
- 432 Hess, L.L., Melack, J.M., Novo, E.M., Barbosa, C.C. and Gastil, M., 2003. Dual-season mapping
- 433 of wetland inundation and vegetation for the central Amazon basin. *Remote sensing of*
- 434 *environment*, 87(4), pp.404-428.

- 435 Jha, A.K., Bloch, R. and Lamond, J., 2012. *Cities and flooding: a guide to integrated urban flood* 436 *risk management for the 21st century.* World Bank Publications.
- 437 Jongman, B., Wagemaker, J., Romero, B.R. and de Perez, E.C., 2015. Early flood detection for 438 rapid humanitarian response: harnessing near real-time satellite and Twitter signals. *ISPRS*
- 439 International Journal of Geo-Information, 4(4), pp.2246-2266.
- 440 Lingad, J., Karimi, S. and Yin, J., 2013, May. Location extraction from disaster-related
- 441 microblogs. In *Proceedings of the 22nd International Conference on World Wide Web* (pp. 1017442 1020). ACM.
- Liggins II, M., Hall, D. and Llinas, J. eds., 2017. *Handbook of multisensor data fusion: theory and practice*. CRC press.
- 445 Marcus, W.A. and Fonstad, M.A., 2008. Optical remote mapping of rivers at sub-meter
- resolutions and watershed extents. *Earth Surface Processes and Landforms*, 33(1), pp.4-24.
- 447 Middleton, S.E., Middleton, L. and Modafferi, S., 2014. Real-time crisis mapping of natural
  448 disasters using social media. *IEEE Intelligent Systems*, *29*(2), pp.9-17.
- 449 Muralidharan, S., Rasmussen, L., Patterson, D. and Shin, J.H., 2011. Hope for Haiti: An analysis
- 450 of Facebook and Twitter usage during the earthquake relief efforts. *Public Relations*
- 451 *Review*, 37(2), pp.175-177.
- 452 Nguyen, D.T., Alam, F., Ofli, F. and Imran, M., 2017. Automatic Image Filtering on Social
- 453 Networks Using Deep Learning and Perceptual Hashing During Crises. *arXiv preprint* 454 *arXiv:1704.02602*.
- Rees, W. G. (2001). Physical principles of remote sensing. Cambridge, UK: Cambridge UniversityPress.
- 457 Smith, R.A., Bates, P.D. and Hayes, C., 2012. Evaluation of a coastal flood inundation model 458 using hard and soft data. *Environmental Modelling & Software*, *30*, pp.35-46.
- 459 Sun, D., Li, S., Zheng, W., Croitoru, A., Stefanidis, A. and Goldberg, M., 2016. Mapping floods
- 460 due to Hurricane Sandy using NPP VIIRS and ATMS data and geotagged Flickr
- 461 imagery. *International Journal of Digital Earth*, 9(5), pp.427-441.
- 462 Wilson, J.N., Gader, P., Lee, W.H., Frigui, H. and Ho, K.C., 2007. A large-scale systematic
- 463 evaluation of algorithms using ground-penetrating radar for landmine detection and
- discrimination. *IEEE Transactions on Geoscience and Remote Sensing*, *45*(8), pp.2560-2572.
- 465 Zeiler, M.D. and Fergus, R., 2014, September. Visualizing and understanding convolutional 466 networks. In *European conference on computer vision* (pp. 818-833). Springer International
- 467 Publishing.
- 468