LEVERAGING SOCIAL MEDIA AS A DATA SOURCE FOR IMPROVED URBAN FLOOD

MONITORING

by

SWAGATO BISWAS ANKON

(Under the Direction of Alysha Helmrich)

ABSTRACT

Urban flooding threatens infrastructure, public safety, and economic stability, with increasing

frequency due to climate change and urbanization. Traditional monitoring methods - sensors, models, and

remote sensing - are effective but limited by cost, time delays, and low spatial resolution. This thesis

explores Twitter as a complementary data source for urban flood monitoring. A framework was developed

to collect, filter, and analyze flood-related tweets using natural language processing, machine learning,

sentiment analysis, and geocoding. Social media data was then integrated with rainfall data to generate near

real-time flood maps. Additionally, a rain-on-mesh simulation using HEC-RAS incorporated terrain, land

cover, and soil data to validate results. Findings show that approximately 75% of flood-affected zones

identified via Twitter matched those from model-generated inundation maps. This demonstrates that social

media can enhance situational awareness and support rapid flood response, making it a valuable tool for

supporting traditional urban flood monitoring systems.

INDEX WORDS:

Crowdsourced Data, Urban Flood Monitoring, Natural Language Processing,

Spatiotemporal Validation, Disaster Resilience

LEVERAGING SOCIAL MEDIA AS A DATA SOURCE FOR IMPROVED URBAN FLOOD MONITORING

by

SWAGATO BISWAS ANKON

BS, Rajshahi University of Engineering & Technology, Bangladesh, 2021

A Thesis Submitted to the Graduate Faculty of The University of Georgia in Partial Fulfillment of the Requirements for the Degree

MASTER OF SCIENCE

ATHENS, GEORGIA

2025

© 2025

Swagato Biswas Ankon

All Rights Reserved

by

SWAGATO BISWAS ANKON

Major Professor: Alysha Helmrich

Committee: Matthew Vernon Bilskie

Linbing Wang Benjamin Rachunok

Electronic Version Approved:

Ron Walcott Vice Provost for Graduate Education and Dean of the Graduate School The University of Georgia May 2025

ACKNOWLEDGEMENTS

I would like to express my deepest gratitude to my advisor Dr. Alysha Helmrich, for mentoring me through this journey, encouraging me to think critically, and inspiring me to bring about the best out of me. Her unwavering support and insightful advice made her contribution invaluable. I am also incredibly grateful to my committee members, Dr. Matthew Vernon Bilskie, Dr. Linbing Wang, and Dr. Benjamin Rachunok, for their support, constructive feedback, and encouragement throughout this process. Their expertise and guidance have strengthened this work in countless ways. This research would not have been possible without the generous support of Funding Organizations. I sincerely appreciate their commitment to advancing knowledge and fostering solutions that address critical challenges. To all those who have supported me – parents, colleagues, friends, and family - thank you for your encouragement and belief in me. This work reflects the collective effort that has shaped my journey.

TABLE OF CONTENTS

| | | Page |
|---------|----------------------------------|------------|
| ACKNO' | WLEDGEMENTS | iv |
| LIST OF | TABLES | vii |
| LIST OF | FIGURES | viii |
| СНАРТЕ | ER . | |
| 1 | INTRODUCTION | 1 |
| 2 | RELEVANT STUDY | 5 |
| 3 | METHODOLOY | 11 |
| | Study Area | 11 |
| | Data Collection | 12 |
| | Data Cleaning and Pre-processing | 13 |
| | Data Analysis | 14 |
| | Rainfall Data | 27 |
| | Satellite Data | 29 |
| | Flood Model Development | 32 |
| | Conceptualization | 36 |
| 4 | Results | 39 |
| | Sentiment Distribution | 39 |
| | Identifying Flood Events | 41 |
| | Removal of False Positives | <i>Δ</i> 1 |

| | Mapping | 43 |
|--------|---------------|----|
| | Validation | 44 |
| 5 | Discussions | 47 |
| | Opportunities | 48 |
| | Limitations | 49 |
| 6 | Conclusion | 51 |
| REFERE | INCES | 52 |

LIST OF TABLES

| | Page |
|--|------|
| Table 1: Tweet filtering using select keyword queries | 13 |
| Table 2: Descriptions (with examples) of the six clusters of tweets within the dataset | 18 |
| Table 3: Comparative performance analysis of machine learning models | 24 |
| Table 4: Details of rain gauge stations across Fulton County | 27 |
| Table 5: Hydrological soil groups and respective soil textures | 31 |
| Table 6: Manning's n values considered for different land class | 32 |
| Table 7: Curve numbers considered for different combinations | 33 |
| Table 8: Validation details of buffer area based on flood depth | 45 |

LIST OF FIGURES

| | Page |
|--|------|
| Figure 1: Annual number of publications. The figure is generated | 7 |
| Figure 2: The study area | 12 |
| Figure 3: Preprocessing steps. First, every word was converted | 13 |
| Figure 4: Top 100 most-weighted words displayed in a word | 15 |
| Figure 5: Silhouette scores with respect to the number of clusters are depicted | 17 |
| Figure 6: Six clusters of tweets as derived from the k-means clustering | 18 |
| Figure 7: Regular expression pattern definition and formation. Six | 21 |
| Figure 8: Expressions and scale of sentiment analysis | 25 |
| Figure 9: Criteria for the fusion of two buffer areas. (a) Fusion criteria | 27 |
| Figure 10: Rain gauges plotted and labeled according to their ID | 28 |
| Figure 11: Active measurement period of each rain station. Eight out of | 28 |
| Figure 12: The terrain map covering the Fulton County boundary is shown | 30 |
| Figure 13: The land cover map covering Fulton County boundary is | 30 |
| Figure 14: The soil layer map covering the Fulton County boundary | 32 |
| Figure 15: Infiltration map covering the study area | 34 |
| Figure 16: (a) Generated mesh at a spacing of 100 ft x 100 ft covering | 35 |
| Figure 17: Satellite precipitation band for January 4, 2023, between | 36 |
| Figure 18: An idealized scenario of a rainfall event leading to flooding and its | 37 |
| Figure 19: Number of tweets in each sentiment category. More than 4800 | 40 |

| Figure 20: Sentiment distribution within Fulton County for the whole | 40 |
|--|----|
| Figure 21: A time series plot of daily rainfall depth and normalized | 41 |
| Figure 22: A time series plot of hourly average satellite rainfall depth | 42 |
| Figure 23: Near real-time mapping of flood-affected zones. Three time | 43 |
| Figure 24: Validation results for 4 January 2023 10:00 – 11:00 AM UTC | 44 |
| Figure 25: Flooded areas that could not be traced with Twitter | 46 |

CHAPTER 1

INTRODUCTION

Flooding is one of the most common natural disasters. It has become more frequent and severe due to a combination of natural and anthropogenic factors. Natural factors include the effects of climate change, such as more intense and frequent rainfall, rising sea levels due to melting glaciers, and shifting weather patterns that prolong or intensify storms (Clarke et al., 2022; Easterling et al., 2012). Other natural drivers are storm surges from hurricanes or typhoons, which overwhelm river systems (Han & Tahvildari, 2024). Anthropogenic factors include urbanization, which has significantly altered natural water recharge systems (Wakode et al., 2018). Expanding impervious surfaces, such as roads, buildings, and parking lots, reduces the land's ability to absorb rainfall. As a result, infiltration deteriorates drastically. Instead of groundwater recharging, the additional water overflows as surface runoff (Arnold & Gibbons, 1996). River alterations, such as channelization and construction in floodplains, increase downstream flood risks (O'Driscoll et al., 2010). Factors like aging dams and insufficient flood planning, further compound these issues. The consequences worsen even more in an urban context.

Urban flooding is critical because it poses severe challenges to the functionality, safety, and resilience of cities worldwide. Flooding in urban areas is pluvial in nature and is mostly driven by rainfall: short intense or prolonged steady. As urban areas expand, they become increasingly vulnerable to flooding due to the prevalence of impermeable surfaces, insufficient drainage infrastructure, and inadequate urban planning (Andreasen et al., 2023; Duy et al., 2018). When cities experience heavy rainfall, stormwater often overwhelms drainage systems, leading to extensive waterlogging in streets, homes, and commercial areas. This can lead to notable economic impacts, including temporary business closures, disruptions to transportation systems, and damage to critical infrastructure such as roads, bridges, and power grids (A. Helmrich et al., 2023). Recovery costs can burden municipal budgets, and property values in flood-prone

areas decline, creating long-term financial instability (BenDor et al., 2020). The consequences of urban flooding extend beyond economic impacts. Public health is affected, as stagnant floodwaters often contain sewage, chemicals, and other pollutants, increasing the risk of waterborne diseases such as cholera and typhoid (Basaria et al., 2023). Breeding mosquitoes in standing water also leads to outbreaks of vector-borne diseases like dengue and malaria (Coalson et al., 2021). Furthermore, the physical toll on affected populations, including injuries and fatalities, can have lasting impacts. Low-income and marginalized communities are particularly at risk, as they are often situated in vulnerable areas and lack access to adequate resources for recovery, exacerbating social inequalities (Moulds et al., 2021). Environmental consequences are also significant, as urban flooding disrupts ecosystems by washing pollutants into nearby rivers, lakes, and wetlands. Changes to land cover, such as soil damage due to extreme weather events like wildfires, can reduce water absorption capacity, leading to increased flood risk (Miller & Hutchins, 2017). These multifaceted impacts underscore the urgent need for sustainable urban planning, improved infrastructure, and proactive policies to mitigate flooding risks. Effective mitigation relies heavily on robust flood monitoring systems that provide real-time insights and predictive capabilities.

Urban flood monitoring is a major concern worldwide, especially in regions prone to severe weather events. It involves observing and analyzing environmental factors such as rainfall, water levels, and drainage systems to assess and mitigate the risks of flooding in cities. It plays a critical role in real-time flood warning systems, helping protect lives, reduce property damage, and enhance community resilience during flood events. Providing reliable flood risk assessment data supports informed decision-making for emergency responders, infrastructure managers, and urban planners. Additionally, in response to climate change impacts, such as intense rainfall and sea level rise, urban flood monitoring is essential. Traditional tools like sensors, satellite imagery, and hydrological models are often used to gather and analyze data for this purpose (Henonin et al., 2013; Tanim et al., 2022). Remote sensing or satellite technology offers good coverage for tracking storm activity and surface water changes (Farhadi et al., 2022; X. Zhu et al., 2024). Hydrological models simulate flood scenarios to evaluate vulnerabilities in flood-prone areas. They are widely adopted in flood management practices across the globe and have proven efficient in modeling flood

behavior and supporting early warning systems. These traditional approaches excel at the physical and environmental dimensions of flooding. For example, quantifying and predicting aspects of flooding like where and when it will occur, the extent of inundation, and the potential damage to infrastructure. However, they often overlook the human aspect – specifically, the public perception of flooding. Flooding oftentimes extends beyond mapped flood zones, and continued development in flood-prone areas leaves residents vulnerable to future flooding events (Hino et al., 2024). Urban flooding is highly localized and cannot be strictly confined to mapped floodplains (Balaian et al., 2024; Son et al., 2023). Public perception of flood risk, driven by personal experiences, community aspects, or cultural contexts, may diverge from scientific assessments (Sawaneh et al., 2024). This gap can worsen the impacts of flooding by influencing preparedness levels, evacuation decisions, or policy priorities (Ahmadi et al., 2022). Besides, traditional methods come with their own limitations. For example, sensors or gauges can be scarce in remote regions, unscaled with limited battery life. Hydrological models can be time-consuming and are mostly used for mitigation purposes. To address these limitations and incorporate public perceptions into flood monitoring, Volunteered Geographic Information (VGI) through crowdsourcing can play a pivotal role. It does not require continuous gauge readings or time-consuming computational models. It can also provide insights into behavioral patterns, risk awareness, and community concerns. By leveraging information shared by individuals through social media, mobile apps, or community reporting platforms, real-time insights into localized flooding conditions and human responses in the form of texts, images, or videos can be gathered. These resources can be analyzed further to derive valuable insights.

However, some concerns have been raised regarding the credibility of social media data. First, it is susceptible to different kinds of biases. There is selection bias, where there is variability of social media usage of different demographic groups (Iacus et al., 2020); self-selection bias, where users engage with content that aligns with their existing beliefs (Persily & Tucker, 2020); and activity bias, where the most active users post more frequently than average users (Baeza-Yates, 2018). These biases can affect the quality and generalizability of social media research outputs. Moreover, lack of fact or cross-checking can alter public perception as misinformation or false data can spread through social media, compromising the

accuracy of research data (Papadopoulos et al., 2016). Besides, there remains the technical challenge of using sophisticated methodologies and analytical tools to process the big, unstructured social media data (Erokhin & Komendantova, 2024).

Building on these considerations, this study evaluates the potential to leverage social media data as a strategic data source for urban flood monitoring despite the limitations. The following three research questions are explored:

- 1. To what extent can integrating social media data contribute to the spatiotemporal resolution of urban flood mapping?
- 2. How can big unstructured data (as found in social media data) be efficiently processed to generate actionable, near-real-time flood maps?
- 3. Does social media data exhibit a spatiotemporal correlation between the user-generated reports and observed flood events?

Chapter 2 provides a review of existing literature that describes the current state of research regarding urban flood monitoring, emphasizing the integration of crowdsourced data. Chapter 3 presents the methodology, comprising the architecture of collection, processing and analysis of Twitter data, rainfall data, satellite data, hydrological model development, and conceptualization of integration between rainfall and Twitter data. Chapter 4 articulates the results and validations, then Chapter 5 discusses the potential opportunities and limitations of crowdsourced data in urban flood monitoring as found in this study. Chapter 6 marks the conclusions of the study.

CHAPTER 2

RELEVANT STUDY

Conventional flood monitoring methods play a major role in mitigating flood-related damage and ensuring public safety. Conventional methods mainly involve the usage of water level gauges, in-situ devices, land parcels, survey data, or close-circuit cameras to gather hydrological information. These methods offer several advantages; for example, they provide reliability since ground-based measurements are often highly accurate and trusted by engineers and urban planners (Tao et al., 2024). They also provide direct observation through real-time data, enabling quick decision-making during emergencies. However, they heavily depend on ground-based measurements that are scarce in many regions (Yilmaz et al., 2010). Besides, oftentimes, these gauges face data transmission delays, especially when installed in remote locations. To address these limitations, physical hydrological models (Nkwunonwo et al., 2020), sometimes in combination with machine learning approaches (Koutsovili et al., 2023) have been utilized. In most of these studies, the HEC-HMS (Hydrologic Engineering Center's Hydrologic Modeling System) was used as the baseline physical model, incorporating high-resolution digital elevation, soil maps, and land use data. The integration of machine learning models significantly contributed to flood risk assessment and early warning systems. These models can effectively address the gaps present in conventional methods by integrating large datasets from various sources, such as satellite imagery or remote sensing, to enhance data coverage and reduce dependency on ground-based measurements (Saha & Chandra Pal, 2024). Furthermore, machine learning models can uncover complex patterns in flood behavior that may not be easily identifiable through traditional approaches (Tang et al., 2023). In most cases, the models were trained with time lags to predict critical water depth. However, these models are often time-consuming and struggle to provide short-term forecasts (Liu et al., 2022). Machine learning models need high-quality data and parameter tuning; further, they may face scalability issues when applied to large regions (Koutsovili et al.,

2023). Various sensors (such as optical, pressure, infrared, and ultrasonic) are also used in flood monitoring. In these studies, the research design is based solely on sensor data or hybridized with satellite, field gauges, and/or machine learning models (Finley et al., 2020; Jang & Jung, 2023; Karyotis et al., 2019; Mousa et al., 2016; Sunkpho & Ootamakorn, 2011). The sensor-based methods can provide continuous, highresolution flood monitoring even in remote and inaccessible areas. Integrating multiple sensors improves measurement accuracy if there is a lack of data with sparse, ground-based measurements. But most of the sensors have limited battery lifetime, high installation costs, scalability issues, and are sensitive to weather conditions, causing measurement errors. Besides, some of these sensors are equipped with a global positioning system that requires strong connectivity. In some cases, machine learning models are subjected to proper area-specific training datasets as well as hyperparameter tuning. Even a few of these limitations can affect the quality of the developed study. Some studies have been undergone regarding the usage of unmanned aerial vehicles (UAV) for flood monitoring (Feng et al., 2015; Z. J. Zhu et al., 2017). These UAVs can capture high-resolution special images that can be further processed based on machine learning algorithms for flood mapping. However, UAVs also suffer from limited battery life, regulatory constraints, and specific software requirements. Besides, the captured images may also be compromised by dense vegetation and cloud cover. The applicability of the stated methods in an urban context is also a challenge, as most floods in urban areas are pluvial in nature (Song et al., 2019).

Volunteered reports generated by the massive population living in an urban area can be a solution to go beyond the traditional approaches. In these cases, crowdsourced data through public webcams, social media, or citizen science produces a high volume of fine-scale data that has the potential to be used in disaster management studies (A. M. Helmrich et al., 2021; See, 2019). Social media is perhaps the most common crowdsourcing media. This Volunteered Geographic Information (VGI) is a user generated information stream providing near real-time data that is rich in quality and does not depend on mechanical gauges (Chen et al., 2016; Imran et al., 2014; Karimiziarani et al., 2022). (Muralidharan et al., 2011) was one of the foremost to examine the role of social media – Facebook and Twitter, used by nonprofits and media organizations during the 2010 Haiti earthquake relief efforts. Since 2011, the domain caught an eye

in disaster research and shows a significant rising trend in the number of publications (Figure 1). (Granell & Ostermann, 2016) presented a review of the main categories and different VGI applicability in disaster research. Twitter was found to be the most used among VGIs in all categories. Thus, researchers have leveraged specifically Twitter data to analyze various aspects of disasters from multiple perspectives. For example, situational awareness during a crisis is one of the popular research domains. (Snyder et al., 2020) proposed a methodology to identify relevant tweets using deep learning models during crisis events to support situational awareness from Twitter data. (Q. Huang & Xiao, 2015; Z. Wang & Ye, 2019) examined spatiotemporal situational awareness by separating tweets according to categories at different phases of a crisis using keywords and predefined word lexicons. However, these lacked how it would be beneficial for emergency responders where (Zade et al., 2018; Zhai, 2022) went one step ahead to propose actionability: reach the right information to the right correspondent. (Lachlan et al., 2016) evaluates how well Twitter was used for public communication, especially during the pre-disaster preparedness phase, with the role of hashtags in making actionable information.

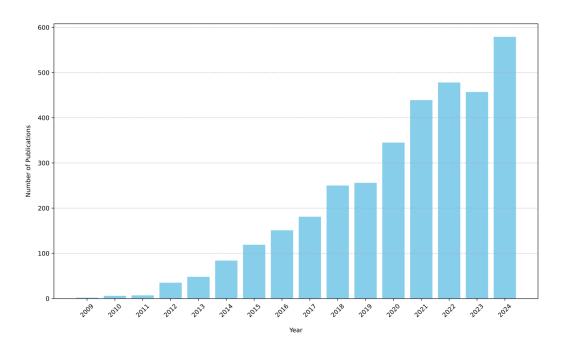


Figure 1: Annual number of publications. The figure is generated from a comprehensive literature search, using search operators such as Booleans, parenthesis, truncation etc. to build the query "("Social Media" OR Twitter) AND (Flood OR "Disaster Management" OR Resilience OR "Crisis Informatic*" OR "Crisis Mapping")" for "Web of Science" document searching. The results show a significant rising trend in the

number of publications since 2011, and it is still growing. This indicates that the emergent field is of growing interest to researchers.

Some studies have been conducted to evaluate disaster resilience. (Crooks et al., 2013) focuses on understanding the spatial and temporal characteristics of Twitter's response to an earthquake using geotagged tweets. (K. Wang et al., 2021, 2023) made sentiment analysis of disaster-related tweets to measure underlying emotion and correlated it to resilience and Twitter indices. (Kryvasheyeu et al., 2016; Mendoza et al., 2019) used Twitter as the medium of rapid damage assessment after a disaster for improved resilience. Some studies focus on using Twitter data to complement existing methodological gaps. (Zou et al., 2018) concluded that the inclusion of social media data improves damage estimation models by offering additional insights. (Panteras & Cervone, 2018) used Twitter data to bridge temporal gaps generated from satellite maps due to infrequent visits and cloud cover. (Cervone et al., 2016) shows that social media can support rapid decision-making in emergency response by identifying critical areas and complementing remote-sensing data gaps when integrated effectively. (Z. Wang & Ye, 2018) conducted a comprehensive review, categorizing existing studies that worked on any of these four key dimensions: spatial, temporal, content, and network. Additionally, they identified studies with hybrid classifications that combined these dimensions among themselves or with remote sensing and census data, providing a nuanced framework for understanding the multifaceted nature of disaster-related research.

(Smith et al., 2017) built a two-dimensional real-time hydrodynamic modelling framework using Twitter data. A standard rainfall hyetograph was applied uniformly across the modelled area, assuming a short-duration, high-intensity rainfall event typical for surface water flooding. The onset of a flood event was detected based on the volume and content of tweets within a defined time window. However, the study is limited by the resolution, scalability, and accuracy of the modelling. (Jongman et al., 2015) focused on analyzing data from the Global Flood Detection System's raster maps and Twitter to enhance the understanding and management of flood events. Twitter signals were available up to two days before floods were reported, showing potential for early flood detection. Even Twitter effectively identified unexpected flood events, such as dam bursts or intentional levee breaches. This information provided qualitative

insights that were not captured by satellite data. The study mostly covered riverine floods, and Twitter data mostly captured responses from urban areas with dense populations. There was still a challenge due to differences in spatial and temporal resolution. A national scale real-time flood monitoring framework was developed by (Barker & Macleod, 2019). They integrated data from sources like the Environment Agency in England/Wales and the Scottish Environmental Protection Agency to obtain up-to-date flood warnings and river levels and fused it with geotagged Twitter data. A Natural Language Processing approach was used: Doc2Vec for feature extraction and Logistic Regression to classify Tweets in "relevant" or "irrelevant" categories based on manually labelled training data. The prototype pipeline successfully retrieved and processed geotagged tweets in real-time, demonstrating its ability to monitor floods at a national scale. However, the number of geotagged Tweets was very few. Also, the study is limited in its applicability to urban areas in cases of pluvial floods. On the other hand, (X. Huang et al., 2018) used remote sensing imagery, stream gauge readings and crowdsourced Twitter data for near real-time inundation probability. The results show that integrating satellite data with social media and gauge readings significantly improves flood extent mapping, providing a higher-resolution and more accurate flood probability distribution. The study achieved a finer scale compared to other models. There still remains the need for an integrated approach that transforms relevant qualitative Twitter texts into quantifiable measures that directly correlate with rainfall events. A framework was conceptualized to address the time lag threshold that can be seen between a rainfall event leading to pluvial flooding and its respective Twitter responses. One major challenge is extracting location information from a large pool of tweets with proper detaining and accuracy. Traditional Named Entity Recognition (NER) tools struggle to accurately identify complex locative references due to the unique characteristics of tweets (Middleton et al., 2020). A variety of approaches, including regular expression pattern formation or regular expression with a combination of other natural language processing techniques, can achieve better accuracy to extract location entities embedded in textual tweets (Subarkah et al., 2023; Yenkar & Sawarkar, 2021; Zhang et al., 1991).

Twitter generally comes with a high-resolution timestamp. Spatial and contextual dimensions are often attributed to questionable reliability and consistency. Proper handling of the content largely determines the

credibility of a study. (Arolfo et al., 2022) conducted quality analysis of Twitter data using four pillars: readability, completeness, usefulness, and trustworthiness. The analysis shows that the average quality of Twitter data streams is higher than anticipated. It varies significantly depending on the filtering techniques and the type of content. Besides, (Fuchs et al., 2013) applied a visual analytics approach to analyze georeferenced Twitter data to investigate how well the tweets reflect disaster events in space and time and found out that filtering tweets based on specific flood-related keywords improved event detection. (Pramanick et al., 2021) shows how language use on social media evolves rapidly during crises, impacting the performance of NLP models. The study states the importance of temporal adaptation during different phases of a crisis. (Paradkar et al., 2022) focuses on analyzing the consistency between the spatially geotagged tweets and the locations mentioned in their content during Hurricane Harvey. Although a small number of tweets are generally geotagged, point locations were the most consistent, followed by area-wise locations. (Gulnerman & Karaman, 2020) evaluated the spatial reliability and accuracy of social media data filtering techniques, with a specific focus on disaster-related events. Filtering techniques that incorporate domain-specific lexicons and machine learning approaches performed better in terms of spatial reliability.

CHAPTER 3

METHODOLOGY

This study employs a multi-step process to integrate near real-time social media with rainfall data. Twitter data related to flooding was collected, preprocessed, and filtered using natural language processing (NLP) and machine learning techniques to identify flood-related tweets. A sentiment analysis was applied to categorize these tweets, focusing on negative sentiments for further geocoding. Geolocation techniques extracted spatial information from tweets to map flood-affected areas. The identified flood events were then cross-referenced with rainfall data from USGS and satellite sources to establish spatiotemporal correlations. A hydrological model using HEC-RAS was developed to simulate flood scenarios for validation. The methodology ensures a comprehensive approach for leveraging social media as a crowdsourced data source for urban flood monitoring.

3.1 Study Area

Situated in the north-central part of Georgia in the southeastern United States, Atlanta is a sprawling metropolitan area consisting of approximately 6.3 million people. It is in the foothills of the Appalachian Mountains and has a humid, subtropical climate. Atlanta experiences annual temperatures between 11.7 and 22.2 ° C and receives a mean 126 cm of annual precipitation (Chang et al., 2021). Extreme rainfall and flooding are exacerbated by urbanization, hilly topography, and floodplain development (Chang et al., 2021). In this study, Fulton County, which covers a major portion of metropolitan Atlanta, is selected as the study area, as shown in Figure 2.

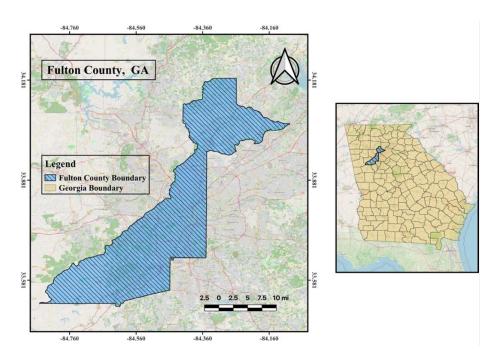


Figure 2: The study area

3.2 Data Collection

Social media data was collected through Twitter, currently known as 'X.' As of 2022, Twitter had more than 368 million monthly active users who generate a large amount of data that can have significant utility (X/Twitter: Number of Users Worldwide 2024 | Statista, 2023). Twitter provides features that support crowdsourcing data, such as geo-tagging, real-time posting, and accessibility (e.g., no user cost, widespread use). Furthermore, until 2023, Twitter hosted a freely available API, which provided a friendly environment to mine data. Crowdsourced data were collected as tweets related to flooding were collected through the Twitter API between September 2021 and March 2022. A 15-mile radius centering Atlanta, GA (33.7490,-84.3880). Three sets of keyword queries were used to filter flooding-related tweets, as described in Table 1. Almost 148,000 tweets were extracted from filtering the keyword queries within the given temporal and spatial constraints. The following attributes were cataloged for each tweet: Tweet ID – the unique identifier of each of the tweet; Timestamp – the time including time zone information when the tweet was posted; Text – the content of the tweet; Search terms – the term(s) used to retrieve the tweets. The tweets were collected anonymously.

Table 1. Tweet filtering using select keyword queries.

| Category | Keyword queries | | | |
|----------|--|--|--|--|
| Rainfall | Raining, rained, pouring, monsoon, rain, rainfall, rainstorm, precipitation, rainwater, | | | |
| | rain shower, | | | |
| Flood | flood, flooded, flooding, flood with, flooded with, flooding with, flood of, inundation, | | | |
| | submerge, immerse | | | |
| Synonyms | Pour, shower, hurricane, storm, soak, soaked, poured, puddle, puddles, sewer overflow, backup, drizzle, ponding, overflow, pond, drown, drowning, drowned, cloudburst, torrent | | | |

3.3 Data Cleaning and Pre-processing

A multi-step process was applied to clean the dataset for analysis. The first cleaning step was to remove the duplicate tweets, including retweets (reposts or forwards of another user's tweet). Retweets were not applicable because dissemination of information was not in the scope of the study. Almost 8% of the tweets were removed through this process – half were duplicate tweets and half were retweets. Next, several preprocessing steps were followed before data analysis to ensure data standardization, consistency, and relevance for analysis, as shown in Figure 3.

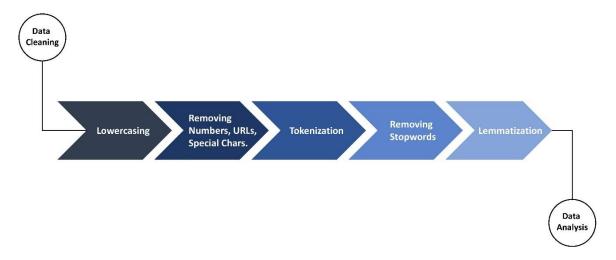


Figure 3: Preprocessing steps. First, every word was converted to lowercase for uniformity. Second, insignificant information such as numbers, hashtags, URLs, and special characters (i.e., emojis and punctuation marks) were removed. Third, the words in tweets were split into comma-separated individual words called tokens. Fourth, all the stopwords were removed in this step. Words that are commonly used in writing but do not contain significant meaning are called stopwords. They include articles, prepositions, pronouns, conjunctions, common verbs, etc. Fifth, every word was transformed to its root form through lemmatization. For example, rained or raining or rains became rain.

3.4 Data analysis

3.4.1 Keywords Visualization

A term frequency-inverse document frequency (TF-IDF) algorithm was applied to determine the most important words in the corpus. TF-IDF algorithms are effective for assessing large datasets, such as Twitter data, due to their fast computation and simplicity. It effectively highlights contextually significant words and filters out common, less informative terms, improving text analysis accuracy. The term frequency (TF) represents the ratio of how often a word appears in a document compared to the total number of words in that document and is represented as:

$$TF(t, D) = \frac{Number\ of\ times\ term\ (t)\ appears\ in\ document\ (D)}{Total\ number\ of\ terms\ in\ document\ (D)}$$

The inverse document frequency (IDF) reflects the significance of a word within a corpus, calculated as the logarithm of the ratio between the total number of documents in the corpus and the number of documents containing the word. It is calculated as:

$$IDF(t, N) = \log \left(\frac{\text{Total number of documents in the entire corpus (N)}}{\text{The number of documents in the corpus that contain the term }(n_t)} \right)$$

Finally, the important keywords are defined by TF-IDF as,

$$\mathsf{TF} - \mathsf{IDF}\left(\mathsf{t}, \mathsf{D}, \mathsf{N}\right) = \, \mathsf{TF}(\mathsf{t}, \mathsf{D}) \, . \, \mathit{IDF}\left(\mathsf{t}, \mathsf{N}\right) = f_{t, D} \, . \, \mathsf{log}\left(\frac{N}{n_t}\right)$$

Where, t is a specific term or word, D is a specific document, N is the total number of documents in the entire corpus, and n_t is the number of documents in the corpus that contain the term t. Words with the highest TF-IDF weight indicate their strong importance within a document, as determined by their

frequency and the overall size of the corpus. A word cloud of the top 100 most-weighted words is shown in Figure 4.



Figure 4: Top 100 most-weighted words displayed in a word cloud. Disclaimer: The word cloud is generated from a collection of tweets and contains taboos (e.g., English swear words). The inclusion of such terms reflects the original content of the tweets and does not endorse or promote any inappropriate language.

3.4.2 Clustering

In large datasets, the texts are generally vast and unstructured. Organizing them is essential to obtain meaningful insights. While each tweet is short and unique, it can contain valuable information, opinions, and news that can be clustered based on similarity. This helps identify patterns, trends, or topics that may not be visible when analyzed individually. K-means clustering was applied in this study on the previously derived TF-IDF vectors using Algorithm 1. K-means clustering is one the most widely adopted clustering algorithms across various domains because of its efficiency and low computational complexity (Ikotun et al., 2023).

Algorithm 1: K-means clustering

```
Input: [T = t_1, t_2, \dots, t_n] #a list of tweets (documents)
          Output: \overline{E} = \overline{e_1}, \overline{e_2}, \dots, \overline{e_m} #cluster of tweets
          #compute TF-IDF
1
          tfidf_matrix \leftarrow TF-IDF(T)
2
          tfidf_matrix_normalised ← normalise (tfidf_matrix, norm = '12')
3
          #initialize k-means
4
          K ← number of clusters (user-defined)
5
          #iterate until convergence
6
          for each iteration i
8
               for each tweet t_1 \in T
                    S(t_i, c_j) = \frac{t_i \cdot c_j}{|t_i| \cdot |c_j|} #compute cosine similarity
9
               c_j \leftarrow t_i with the highest cosine similarity #nearest centroid
10
               Recalculate centroids using c_{j(new)} = \frac{1}{|c_i|} \sum_{t_i \in c_j} t_i #averaging vectors of each assigned cluster
11
               if d (c_{j(old)}, c_{j(new)}) < \text{threshold} \in, break
12
13
          #assign final cluster labels
          \overline{E} \leftarrow [\overline{e_1}, \overline{e_2}, \dots, \overline{e_m}]
14
          #visualize
15
          Reduced_matrix ← PCA(2)
16
          Return \overline{E}
17
```

Here, the algorithm takes a list of tweets as input and aims to produce clusters of tweets based on their content similarity. First, it computes the TF-IDF vectors for the tweets that transform the text into numerical representations (as described in section 3.4.1). It reflects the importance of terms within the dataset. These vectors are then normalized to ensure consistent distance calculations among the post vectors when applying cosine similarity:

$$cos(t_i, c_j) = \frac{t_i \cdot c_j}{|t_i| \cdot |c_j|} = \frac{\sum t_{ik} c_{jk}}{\sqrt{\sum t_{ik}^2 \sum c_{jk}^2}}$$

Where t_i is the tweet vector and c_j is the centroid vector. Next, the K-means clustering is initialized with a user-defined number of clusters. Different values of k are tested using silhouette scores to determine the optimum value, as demonstrated in Figure 2. The optimum value of k was determined to be six. The algorithm iterates until convergence by first calculating the cosine similarity between each tweet and the

centroids of the clusters. The tweet is assigned to the cluster with the highest cosine similarity, aggregating the tweets. After all tweets have been assigned to clusters, the centroids of the clusters are recalculated by averaging the vectors of all the tweets in each cluster. The process repeats, with the centroids being updated at each iteration, until the movement of the centroids falls below a user-defined threshold (10^{-4}) , signaling convergence. Once convergence is achieved, the algorithm assigns the final cluster labels to the tweets. Lastly, the dimensionality of the data is reduced using principal component analysis for visualization purposes, and the resulting clusters are returned. The resulting clusters are presented in Figure 5. Table 2 describes the general themes of the discussions in each of these clusters. The themes are generalized based on the top 30 most-weighted keywords in each cluster.

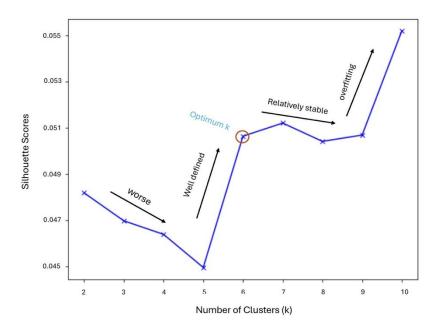


Figure 5: Silhouette scores with respect to the number of clusters are depicted. As k increases from two to five, silhouette scores drop, suggesting that the clustering is worsening. This might happen because the tweets are not being grouped into well-defined clusters. The silhouette score jumps significantly at k=6. This suggests that the data is much better separated into six clusters than five or fewer clusters. At this point, the clusters are well-defined. After k=6, the silhouette score fluctuates slightly with small increases and decreases between k=7 and k=9. This suggests that adding more clusters doesn't dramatically improve the quality of the clusters. At k=10, there's a notable jump in the silhouette score again. The magnitude of the increase indicates that the model is overfitting, as too many clusters may start to capture noise.

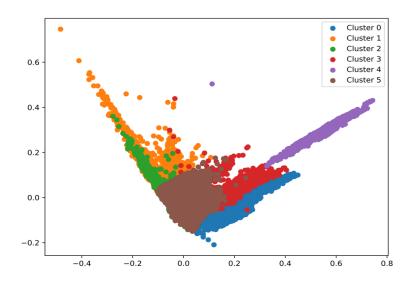


Figure 6: Six clusters of tweets as derived from the k-means clustering (optimum number of k = 6).

| T 11 1 D ' ' | / •.1 1 | \ C.1 | . 1 | C | • .1 • .1 • . |
|-----------------------|--------------------|-----------|----------------|-------------|----------------------|
| Table 2. Descriptions | (with example | oci at th | o six clustors | at tweets w | vithin the dataset |
| Tubic 2. Descriptions | (Willi Cadillipid | cojojiii | c six ciusicis | of incers n | viiiiii iiic aaiasci |

| ID | Description | | | |
|-----------|---|--|--|--|
| Cluster 0 | Personal reflections on life, emotions, and weather | | | |
| | Example 1: I needa hurry up and get out the house before it rain. Glad I did my baby hair | | | |
| | already. | | | |
| | Example 2: Recent extreme weather is a great opportunity to remake the global economy | | | |
| Cluster 1 | Flood warnings and emergency alerts | | | |
| | Example 1: LIX issues Flood Warning for Mississippi River at Red River Landing [LA] | | | |
| | till May 25, 7:00 AM CDT | | | |
| | Example 2: Flood Alerts are active in parts of the Midwest and Northeast as the rain from | | | |
| | Winter Storm #Miles meets snowpack and frozen ground. | | | |
| Cluster 2 | Severe weather (tornado and thunderstorm warnings) | | | |
| | Example 1: TAE issues Severe Thunderstorm Warning [wind: 60 MPH (RADAR | | | |
| | INDICATED), hail: 0.75 IN (RADAR INDICATED)] for Bay, Calhoun, Gulf [FL] till | | | |
| | 12:45 PM CDT | | | |
| | Example 2: Severe thunderstorm may hit today at Fulton County Georgia [wind: 50 | | | |
| | MPH] | | | |
| Cluster 3 | Casual, everyday conversations about weather and life | | | |
| | Example 1: If it's going to rain, could it at least storm and go crazy smh | | | |
| | Example 2: I love this game so much I played in the rain! Thank you! #tennis | | | |
| | #ultimatetennis #tenniscouples #couplesplay #hittingballs @ Decatur, Georgia | | | |
| Cluster 4 | Daily life influenced by weather and astrological signs (sagittarius, aquarius etc.) | | | |
| | Example 1: I hate shooting in the rain. This lightning delay wasn't long enough. | | | |
| | Example 2: If you, Äôre an aquarius i hope you ride your motorcycle in the rain. | | | |
| Cluster 5 | Severe thunderstorm warnings, evacuation and monitoring. | | | |
| | Example 1: A flooded river in Oregon is prompting the evacuation of roughly 50 people | | | |
| | from an RV Park. 10 have already been successfully rescued. | | | |
| | Example 2: Today, we're monitoring a severe weather threat that brings the possibility | | | |
| | of thunderstorms, tornadoes, damaging winds, and flash flooding. | | | |

3.4.3 Location Detection

In order to identify where flooding may be occurring, it was crucial to geolocate the tweets in order to get a spatial visualization. However, less than 1% of the tweets had the built-in geolocation function enabled on Twitter. This number of tweets was insufficient to support an urban flood monitoring system meaningfully. However, many of the tweets had user-generated location information embedded in their text content. Therefore, these text-embedded locations needed to be detected and assigned. Initially, the research team used methods such as Natural Language Toolkit (NLTK) or TextBlob, but the results were unsatisfactory, as both methods detected many false positives and missed numerous location information. This might be due to significant variability in location embeddings in tweets of different clusters. For example, a term may be either a place or a name (e.g., Georgia), cardinal directions being identified as a location, multiple locations listed in one tweet, or shorthand for street names (e.g., Street as St.). For this reason, 500 tweets were randomly selected from each cluster to learn how people reference locations in their tweets. Broadly, the identified content was categorized into three types:

- Type 1: County-wise location (e.g., Fulton County)
- Type 2: Area-wise location, such as districts (e.g., downtown), cities (e.g., Alpharetta), and neighborhoods (e.g., Old Fourth Ward)
- Type 3: Street-wise location (e.g., Peachtree Street)

The randomly selected tweets from each cluster showed that the locations could be identified through patterns to recognize all three types of locations. These patterns could be unique or common among the clusters. So, in the first phase, a regular expression pattern recognition tool was used to detect and assign these locations. As the tool takes user-generated predefined patterns as inputs, these patterns were defined first. Each unit of the pattern was considered and constructed to account for possible variations. Figure 7 shows a detailed overview of the defined patterns for all three types of locations. The patterns were applied to the tweets based on the details that could be extracted using 'if' statements (as indicated in the order column). This process begins with order 1. If a pattern of order 1 successfully extracts a location from a

tweet, the succeeding order patterns will not be applied to that tweet. Conversely, if the order 1 pattern fails to identify a location, the order 2 pattern will be applied, and so forth. Finally, the order 5 pattern, which is the least detailed, will be applied to any remaining tweets.

Definition:

- 1. Place: Must start with a capital letter; optional multiple words; optional punctuation marks between words; optional 'and' between multiple words
- 2. State Abb: Two letter abbreviation of all states; optional brackets; optional dots after each letter
- 3. State: Full names of all states; optional brackets.
- 4. Relay: for, of, near, between, in, at, on, above, under
- 5. County: county, counties, co., co, cos; allows cases
- 6. Street: street, st., st, road, rd., rd; allows plurals; allows cases

| Pattern format | Mostly Applicable to | Order |
|---|---|-------|
| Type 1: County-wise Location | | |
| • General: Relay (optional) Place1, Place2,, PlaceN County Input: A strong thunderstorm will impact portions of Meriwether County and Muscogee County through 145 PM EST [wind: 40 MPH, hail: 0.25 IN] Output: Meriwether County and Muscogee County Input: first thunderstorm of the season here in Fulton Co. The dog and I are at complete opposite ends of the excitement spectrum Output: Fulton Co | Cluster 0, Cluster 1, Cluster 4, Cluster 5 | 5 |
| Type 2: Area-wise Location | | |
| • General: Relay (optional) Place1, Place2,, PlaceN Relay (optional) or comma (optional) State Abb or State Input: Flood warning has been issued in Portage, Stark, Summit [OH] till May 14, 9:45 PM EDT Output: Portage, Stark, Summit [OH] | All Clusters | 2 |
| Input: I get home from having a wonderful day and my neighbors done flooded the apartments at Ballard, KY Output: Ballard, KY Input: Atlanta GA and this damn rain!!! Im about to eat and sleep good.' Yall stay safe | | |
| Output: Altanta GA Nested: Relay Place1, Place2,, PlaceN Relay Place1, Place2,, PlaceN Relay (optional) or comma (optional) State Abb or State Input: LIX extends time of Flood Warning for Mississippi River at Red River Landing [LA] till May 25, 7:00 AM CDT Output: Mississippi River at Red River Landing [LA] Input: Severe thunderstorm warning going on here in Colombus in Muscogee County, GA Output: Colombus in Muscogee County, GA | Cluster 1, Cluster 2, Cluster 5 | 1 |
| Type 3: Street-wise Location | | |
| • General: Relay (optional) Place1, Place2,, PlaceN Street Relay (optional) or comma (optional) State Abb or State Input: This whole Bouldercrest, Candler, Wesley chapel road, GA area all up in there in the rain. Output: Bouldercrest, Candler, Wesley chapel road, GA | Cluster 3, Cluster 4 | 4 |
| • Nested: Relay (optional) Placel, Place2,, PlaceN Street Relay Place1, Place2,, PlaceN Relay (optional) or comma (optional) State Abb or State Input: Floodwater accumulating at Leonard Street between Jabbertown and Firefly in North Carolina. Output: Leonard Street between Jabbertown and Firefly in North Carolina. | Cluster 3, Cluster 4, Cluster 5 | 3 |

Figure 7: Regular expression pattern definition and formation. Six constitutional units were defined first. The units were structured to form patterns for county-wise, area-wise, and street-wise location extraction. The patterns were then applied to the tweets to capture matches. The order of application was based on the level of subtlety that could be extracted for a location. The matches were deemed as location outputs. The extracted locations were further processed for minor modifications to make them suitable for the OpenStreetMaps geolocation application. For example, the relays were replaced with a comma where applicable. All terms within the 'County' and 'Street' definitions were replaced with 'County,' 'Street,' and 'Road,' respectively, where applicable.

However, there could be tweets that did not fall under a defined pattern. Also, some people may misspell a location due to various reasons, such as English not being their first language, tweeting in rush due to the severity of an event, etc. To accommodate those locations, a second step was undertaken. In the second phase, each unique pattern-captured location from the first step was used as a reference to find a match from the remaining tweets. In this case, if a match of 80% or higher was found between the reference locations and any of the word(s) in the remaining tweets, it was also extracted and counted as a location. The process identified 12.9% (~19,000) of the tweets as containing location information provided by users. This pool of tweets was named 'GeoTweets Pool.' Finally, the 'GeoTweets Pool' was processed further in the next steps.

3.4.4 Separating Flood-Related Tweets

Filtering flood-related tweets was essential for refining the dataset and identifying tweets pertaining to flood events. The 'GeoTweets Pool,' identified in the previous step, contained writings on a diverse array of topics. Twitter users may tweet about their sentiments regarding activities, experiences while driving, poetry or story about rain, and other subjects during a rainfall event. However, it is important to recognize that not every rainfall event will result in flooding. Therefore, it was necessary to identify the tweets referencing a flooding event to support an urban flood monitoring system. For example, tweets discussing flood warnings, extreme weather alerts, personal experiences during a flood (e.g., reporting an inconvenience), damage reports, emergency evacuation information, road closures, etc. To identify the flood-related tweets, machine learning models were leveraged. Five hundred tweets were randomly selected

from each cluster within the 'GeoTweets Pool.' The tweets were manually inspected to detect if the tweet indicated a flooding event or not. Each tweet was labeled as either not indicative of flooding (0) or indicative of flooding (1).

The tweet text content (X) and corresponding labels (y) are extracted from the manually inspected dataset. The overall dataset is then split into a training set (70%) and a test set (30%). The feature extraction is completed with the TF-IDF vectorizer to convert the tweet text into numerical features representing the word's importance in each tweet. Four classification models (Logistic Regression (LR), Random Forest (RF), Naïve Bayes (NB), and Support Vector Machine (SVM)) were executed to have a comparative analysis of their performance. The vectorizer is fitted on the training data and transformed for the training and test data. Since urban flooding events are scarce, tweets indicating a flood event were relatively low compared to tweets that did not indicate a flooding event. The balancing ratio (flood-indicative to non-flood-indicative ratio) was nearly 0.25, indicating a moderate class imbalance. This could influence the performance of the machine learning models, especially for the minority class if not addressed. Therefore, optimal class weights (LR, RF, SVM) or class priors (NB) were assigned—higher for the minority class and lower for the majority class. Considering weights, each classifier is trained on the training set, and predictions are made on the test set. The classified results fall under the following four categories:

- True positive (TP): the number of tweets that correctly fall under flood tweets
- False positive (FP): the number of tweets that incorrectly fall under flood tweets
- True negative (TN): the number of tweets that correctly fall under non-flood tweets
- False negative (FN): the number of tweets that incorrectly fall under non-flood tweets

The labeling performance is measured in terms of accuracy, precision, recall, and F-1 score. Precision measures the proportion of tweets predicted as flood-related that are actually flood-related as:

$$Precision = \frac{TP}{TP + FP}$$

Recall measures the proportion of actual flood-related tweets that were correctly identified by the model. It is calculated as:

$$Recall = \frac{TP}{TP + FN}$$

F1-Score is the harmonic mean of precision and recall, which balances the two as:

$$F1 Score = 2' \frac{Precision' Recall}{Precision + Recall}$$

Accuracy means the overall proportion of correctly classified tweets (both flood-related and non-flood) and is measured as:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Finally, a comparative performance analysis of the four models predicting flood tweets is shown in Table 3.

Table 3. Comparative performance analysis of machine learning models

| Model | Precision | Recall | F-1 Score | Accuracy |
|------------------------|-----------|--------|-----------|----------|
| Logistic Regression | 0.94 | 0.74 | 0.83 | 0.85 |
| Random Forest | 0.92 | 0.89 | 0.90 | 0.89 |
| Naïve Bayes | 0.90 | 0.83 | 0.86 | 0.88 |
| Support Vector Machine | 0.94 | 0.90 | 0.91 | 0.91 |

Since the support vector machine achieved the best results in all four metrics, it was applied to the remaining unlabeled tweets for the classification task using the same TF-IDF vectorizer. Almost 5,200 tweets were identified as flood-related tweets. This updated pool of tweets was named 'Flood GeoTweets Pool.'

3.4.5 Sentiment Analysis

The 'Flood GeoTweets Pool' was further processed with sentiment analysis, a natural language processing (NLP) technique used to determine the emotional tone expressed in a piece of text. The goal of sentiment analysis is to identify whether the sentiment conveyed in the text is positive, negative, or neutral emotions.

In this study, text was processed with vader, a predefined lexicon-based sentiment analyzer where the words in a document are annotated with semantic scores between -1 and 1 (Hutto & Gilbert, 2014) (Figure 8). Vader can score individual words and sentence fragments (aggregating scores of each word in the sentence). It also accounts for intensifier words (e.g., 'flooding' versus 'severe flooding,' where severe is working as an intensifier.) In this study, the compound score of each tweet classified as:

- Positive: If the compound score falls within the range 0.1 to 1.00
- Neutral: If the compound score falls within the range -0.1 to 0.1
- Negative: If the compond score falls within the range -0.1 to -1.00

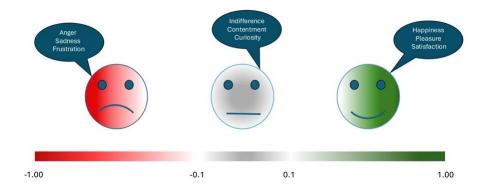
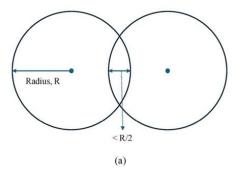


Figure 8: Expressions and scale of sentiment analysis

The 'Flood GeoTweets Pool' was processed through the sentiment analyzer for sentiment classification. The collection of real-time field data on critical aspects such as infrastructure damage, road closures, and public distress was essential for the flood monitoring system. These data provide valuable insights into the severity of flooding and its impact on mobility and community well-being. To ensure that the most relevant information was analyzed, tweets expressing compound negative sentiment were specifically targeted for the study, as firsthand reports of disruptions, safety concerns, and frustration caused by the flood are more likely to be captured by negative sentiments (Karmegam & Mappillairaju, 2020). The updated pool of tweets was named 'Sen Flood GeoTweets Pool' and was processed further in the next step.

3.4.6 Geocoding

OpenStreetMaps (OSM) geocoding feature in QGIS was used to geolocate the 'Sen Flood GeoTweets Pool.' OSM primarily relies on the Nominatim service, which follows a structured process for interpreting and matching location queries. When a user enters a place name or address (extracted locations from tweets described above), the input query is parsed into its component parts (e.g., house number, street, city, postal code, country). The query is structured hierarchically (from specific to broad) to improve matching accuracy. Nominatim searches OpenStreetMap's geospatial database, which contains location data tagged with attributes. The system uses full-text search indexing on pre-processed address components to find relevant entries quickly. If a multiple match is found for a query, relevant results within a county boundary are selected for display. Once a match is found, Nominatim assigns geographic coordinates (latitude & longitude) from OSM's point, line, or polygon data. The exact point coordinates are returned for pointbased features (e.g., buildings, landmarks). For polygon features (e.g., cities, administrative boundaries), a centroid (central coordinate) is computed. In this study, areawise and streetwise locations returned centroid coordinates of respective polygon features. However, since a centroid coordinate does not fully represent a polygon area, an average buffer radius of 0.5 mile was considered around the centroid coordinate. If the 'buffer area' of two points overlapped by more than or equal to half of the radius, they were fused together to a single point, as shown in Figure 9. On the other hand, countywise locations were regarded as administrative boundaries of each county. As the boundaries of counties are fixed, the county-wise locations are not required to have geocoding. Tweets that contain countywise locations are assumed to have an effect on the whole county boundary.



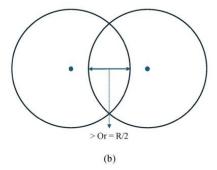


Figure 9: Criteria for the fusion of two buffer areas. (a) Fusion criteria do not apply as the overlapping distance is less than half of the radius. (b) Fusion criteria apply as the overlapping distance is greater or equal to half of the radius.

3.5 Rainfall Data

3.5.1 In-situ Rainfall Data

In-situ rainfall data was collected from USGS archives for the specified timeframe. The rainfall data was collected daily and contained other information like rain gauge station names and coordinates, as mentioned in Table 4. The gauges have been plotted using QGIS (Figure 10). Not all of the gauges were actively measuring rainfall throughout the timeframe of tweet collection. Figure 11 shows an overview of each rain station's active status, which highlights the periods of data availability and gaps in measurement.

Table 4. Details of rain gauge stations across Fulton County.

| ID | Station name | Coordinates |
|----|----------------------------------|-------------------|
| 1 | Palmetto 3.2 NW, GA US | 34.0139, -84.6997 |
| 2 | Atlanta 3.7 N, GA | 33.8169, -84.4174 |
| 3 | Roswell 0.7 SE, GA | 34.0298, -84.3466 |
| 4 | Alpharetta 2.1 NNE, GA | 34.0983, -84.2600 |
| 5 | Alpharetta 1.6 SE, GA | 34.0529, -84.2506 |
| 6 | Alpharetta 4.8 WNW, GA | 34.0903, -84.3515 |
| 7 | Atlanta 3.2 S, GA | 33.7171, -84.4210 |
| 8 | Roswell 5.9 SE, GA | 33.9789, -84.2800 |
| 9 | College Park 1.6 NNW, GA | 33.6620, -84.4641 |
| 10 | Atlanta 2.3 NE, GA | 33.7881, -84.3966 |
| 11 | Atlanta 2.3 SE, GA US | 33.7363, -84.3995 |
| 12 | Atlanta Fulton Co Airport, GA US | 33.7775, -84.5246 |
| 13 | Alpharetta 3.7 ENE, GA US | 34.0962, -84.2168 |
| 14 | Roswell 3.0 ESE, GA US | 34.0139, -84.3119 |
| 15 | Roswell 2.7 NW, GA US | 34.0611, -84.3945 |

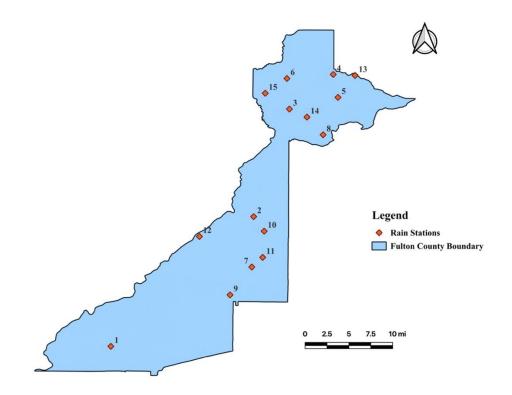


Figure 10: Rain gauges plotted and labeled according to their ID for Fulton County.

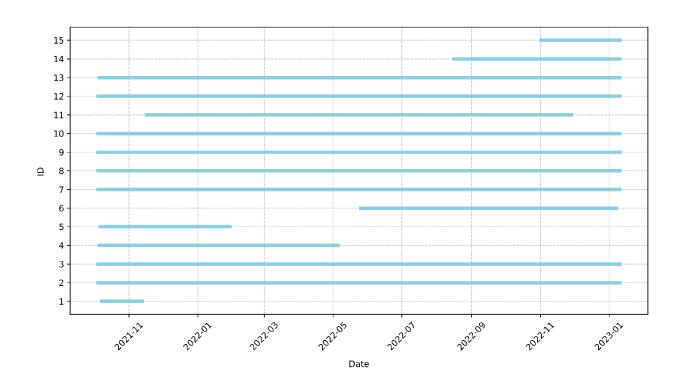


Figure 11: Active measurement period of each rain station. Eight out of 15 rain stations monitored rainfall throughout the timeframe.

3.5.2 Satellite Rainfall Data

Satellite rainfall data was collected from the Center for Hydrometeorology and Remote Sensing (CHRS) of the University of California Irvine (Nguyen et al., 2019). The raster resolution was 0.04 x 0.04 degree. The collected data provided hourly precipitation data throughout the study area, ensuring that real-time, localized variations in rainfall patterns were captured to supplement the coarse in-situ rainfall data.

3.6 Satellite Data

3.6.1 Digital Elevation Model (DEM) Data

A digital elevation is a representation of the Earth's topographic bare surface in a digital format. It is usually created from elevation data collected through remote sensing methods like LiDAR, radar, or satellite imagery. In this study, the DEM raster with a resolution of 30 meters was collected from the NASA/NGA Shuttle Radar Topography Mission (SRTM) data and processed through QGIS. It provides a grid-based format where each cell contains an elevation value, making it possible to generate contour maps, analyze slope gradients, or simulate water flow across landscapes. A detailed terrain map is shown in Figure 12.

3.6.2 Land-Use and Land-Cover (LULC) Data

Land-use and land cover (LULC) data represent the physical characteristics of the Earth's surface, distinguishing between natural features like forests, water bodies, and grasslands, as well as human-made features such as urban areas, agricultural fields, and infrastructure. In this study, Sentinal-2 10-meter land cover raster from 2021 was extracted from the ESRI National Land Cover Dataset (NLCD) and processed through QGIS. A detailed land cover map is shown in Figure 13.

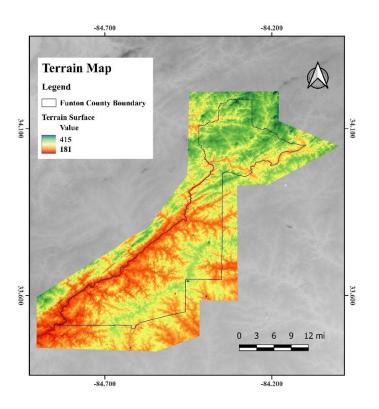


Figure 12: The terrain map covering the Fulton County boundary is shown. The color gradient implies lower (red) to higher (green) elevation.

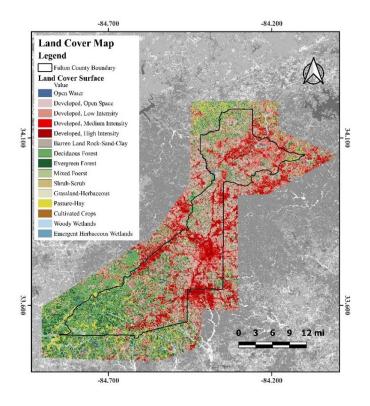


Figure 13: The land cover map covering Fulton County boundary is shown. Fifteen unique color codes demonstrate various land cover types.

3.6.3 Soil Layer Data

Soils are categorized into hydrologic soil groups (HSG) to represent bare soil's minimum infiltration rate under sustained saturation conditions. A detailed classification of hydrological soil groups and their textures are shown in Table 5.

Table 5. Hydrological soil groups and respective soil textures.

| | 7 |
|-----|---|
| HSG | Texture |
| A | Sand, loamy sand, or sandy loam |
| В | Silt loam or loam |
| С | Sandy clay loam |
| D | Clay loam, silty clay loam, sandy clay, silty clay, or clay |

Group A soils consist of sand or gravel materials with a high water transmission rate (> 0.30 in/hr). These show low runoff potential and high infiltration rates even when thoroughly wetted. Group B mainly consists of drained soils with moderately fine to coarse textures as well as moderate water transmission (0.15 – 0.30 in/hr). These show moderate infiltration rates when thoroughly wetted. On the other hand, Group C soils have moderately fine to fine textures with low water transmission (0.05 - 0.15 in/hr). When thoroughly wetted, these soils show low infiltration rates. Finally, Group D soils consist mostly of clays with high swelling potential. These soils have a very low rate of water transmission (0 - 0.05 in/hr). Moreover, Group D soils have very low infiltration rates, and thus, the runoff potential is maximum. In certain cases, places can have a mixture of multiple soil groups. These are expressed as a combination of existing groups. For example, if an area has both soil Group A and B, it is expressed as Group AB, where A is the dominant group. Impervious developed areas are often referred to as the 'none' category. Figure 14 shows a detailed soil map of the study area extracted from the Soil Survey Geographic Database (SSURGO) and processed through QGIS.

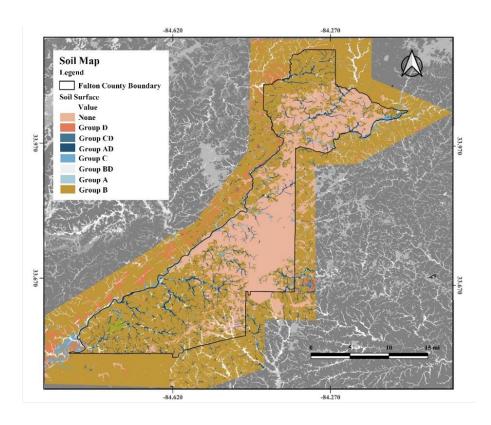


Figure 14: The soil layer map covering the Fulton County boundary is shown. Unique color codes separate different soil groups. As it can be seen, a major part of Fulton County falls under 'none' category as these places are impervious due to developed regions.

3.7 Flood Model Development

3.7.1 Manning's n Value

Manning's n, also known as the Manning's roughness coefficient, is a dimensionless empirical parameter that characterizes the resistance or friction caused by channel surface roughness, which affects the speed and behavior of water flow in natural and artificial channels. It takes into account the effect of channel surface irregularities, vegetation, material composition, and other obstructions on water flow. The following Manning's n values (presented in Table 6) were selected for the NCLD land cover raster.

Table 6. Manning's n values considered for different land classes.

| NLCD Class ID | NLCD Class Name | Manning's n |
|---------------|--------------------------|-------------|
| 0 | No Data | 0.06 |
| 11 | Open Water | 0.04 |
| 21 | Developed, Open Space | 0.05 |
| 22 | Developed, Low Intensity | 0.08 |

| 23 | Developed, Medium Intensity | 0.10 |
|----|---------------------------------------|------|
| 24 | Developed, High Intensity | 0.15 |
| 31 | Barren Land Rock-Sand-Clay | 0.04 |
| 41 | Deciduous Forest | 0.16 |
| 42 | Evergreen Forest | 0.16 |
| 43 | Mixed Forest | 0.16 |
| 52 | Shrub-Scrub | 0.10 |
| 71 | Grassland-Herbaceous | 0.06 |
| 81 | Pasture-Hay | 0.06 |
| 82 | Cultivated Crops | 0.06 |
| 90 | Woody Wetlands | 0.12 |
| 95 | Emergent Herbaceous Wetlands | 0.07 |
| | · · · · · · · · · · · · · · · · · · · | |

3.7.2 Curve Numbers

This study adopted the Soil Conservation Service Curve Number model to estimate runoff from a rainfall event. Developed by the National Resources Conservation Service (NRCS), it is a widely used empirical parameter to evaluate how much rainfall will likely runoff a particular area based on land use, soil type, and moisture conditions. These values are reported in Table 7.

Table 7. Curve numbers considered for different combinations of land cover and soil groups.

| NLCD Class Name | Hydrological Soil Groups | | | | | | | | |
|------------------------------|--------------------------|----|----|----|----|---------|----|----|----|
| | None | A | В | Ċ | D | No Data | AD | BD | CD |
| No Data | 86 | 61 | 74 | 82 | 86 | 86 | 86 | 86 | 86 |
| Open Water | 99 | 99 | 99 | 99 | 99 | 99 | 99 | 99 | 99 |
| Developed, Open Space | 84 | 49 | 69 | 79 | 84 | 84 | 84 | 84 | 84 |
| Developed, Low Intensity | 87 | 61 | 75 | 83 | 87 | 87 | 87 | 87 | 87 |
| Developed, Medium Intensity | 93 | 81 | 88 | 91 | 93 | 93 | 93 | 93 | 93 |
| Developed, High Intensity | 95 | 89 | 92 | 94 | 95 | 95 | 95 | 95 | 95 |
| Barren Land Rock-Sand-Clay | 84 | 49 | 69 | 79 | 84 | 84 | 84 | 84 | 84 |
| Deciduous Forest | 63 | 30 | 48 | 57 | 63 | 63 | 63 | 63 | 63 |
| Evergreen Forest | 80 | 35 | 58 | 73 | 80 | 80 | 80 | 80 | 80 |
| Mixed Forest | 79 | 36 | 60 | 73 | 79 | 79 | 79 | 79 | 79 |
| Shrub-Scrub | 77 | 35 | 56 | 70 | 77 | 77 | 77 | 77 | 77 |
| Grassland-Herbaceous | 89 | 55 | 71 | 81 | 89 | 89 | 89 | 89 | 89 |
| Pasture-Hay | 84 | 49 | 69 | 79 | 84 | 84 | 84 | 84 | 84 |
| Cultivated Crops | 83 | 63 | 73 | 80 | 83 | 83 | 83 | 83 | 83 |
| Woody Wetlands | 93 | 72 | 80 | 87 | 93 | 93 | 93 | 93 | 93 |
| Emergent Herbaceous Wetlands | 93 | 72 | 80 | 87 | 93 | 93 | 93 | 93 | 93 |

3.7.3 Infiltration Layer

The infiltration map in Figure 15 visually represents infiltration rates across the study area by combining hydrological soil group classifications and land cover types from the NLCD dataset. Different colors overlaying the map correspond to distinct levels or classes of infiltration potential.

The abstraction ratio in hydrological modeling refers to the proportion of rainfall that does not contribute directly to runoff but is instead absorbed or retained by the soil, vegetation, and other land features before reaching water bodies. It is a parameter often used in infiltration and runoff models to approximate initial losses due to various interception processes, surface storage, etc. In this study, it is considered that a ratio of 0.2 or 20% of the total precipitation is considered to be lost to initial abstractions before any runoff occurs. The 0.2 factor is a widely used empirical assumption for the SCS Curve Number method to incorporate initial losses.

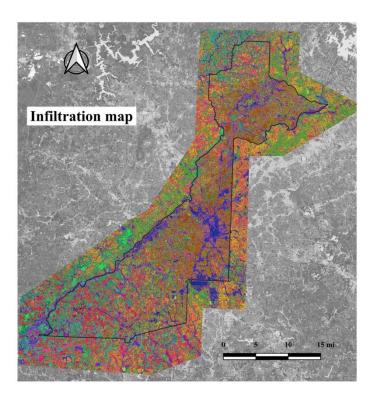


Figure 15: Infiltration map covering the study area.

3.7.4 2D Flow Area Generation

Flow area covering Fulton County for mesh creation and computational points generation was completed for the simulation model in HEC-RAS, as shown in Figure 16.

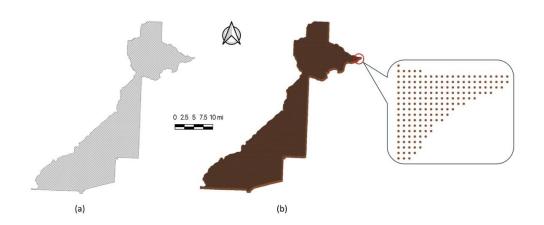


Figure 16: (a) Generated mesh at a spacing of 100 ft x 100 ft covering Fulton County Boundary. 137,929 cells were generated in the mesh. (b) Each cell was used to generate computational points.

3.7.5 Rainfall Event Assignment

Hourly gridded satellite rainfall data, as described in Section 3.5.2, was assigned to the study area. Figure 17 shows an example of a precipitation band over Fulton County for January 4, 2023, between 10:00 and 11:00 AM UTC.

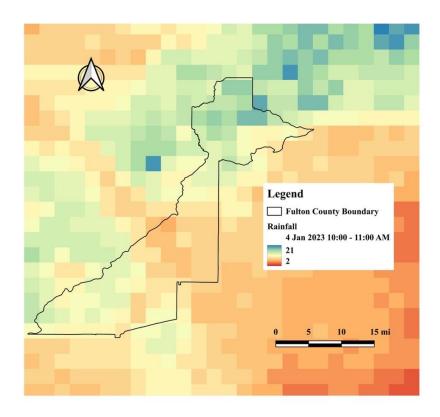


Figure 17: Satellite precipitation band for January 4, 2023, between 10:00 and 11:00 AM UTC. The color contours vary between 2 to 21 millimeters of rainfall per hour over the study area.

3.8 Conceptualization

Extreme natural events generally create a digital social media footprint (Ogie et al., 2022). People talk about the event, express concerns, authorities issue warnings, etc. In this study, public social media sentiment in response to flooding as a direct cause of rainfall has been taken into consideration. The emotional expression regarding flood can be plotted in a time series with the associated rainfall event for a particular spatial zone. An idealized scenario of the time series plot is shown in Figure 18.

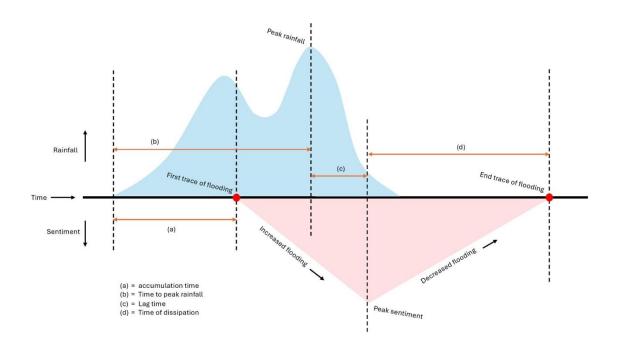


Figure 18: An idealized scenario of a rainfall event leading to flooding and its Twitter response. The blue area represents the magnitude of rainfall over time, and the red area represents the sum of sentiments over time. The time between the start of rainfall and the first trace of flooding reported in social media data is denoted as accumulation time (a). The time to peak rainfall from the start of the event is represented by (b). The peak sentiment value will likely occur after the peak rainfall as it takes time for rainwater to accumulate. The time between peak rainfall and peak sentiment is denoted as lag time (c). When the rainfall decreases and flooding recedes, the perception of flooding will also decrease. The time between the peak sentiment and the end trace of flooding (according to social media) is denoted as the time of dissipation (d).

However, some cases may not always align with the idealized scenario depicted in Figure 18. For example, ideally, the first trace of flooding, according to social media, should fall within the timespan of the rainfall event. However, if the rainfall event starts and ends during overnight hours (e.g., 11:30 PM – 3:00 AM), there may be few tweets until morning when people start their daily activities and notice the flooding. Similarly, peak sentiment should ideally follow peak rainfall with a lag time as rainwater accumulates into flooding. This pattern may shift if peak rainfall begins overnight. Then, the peak in sentiment—driven by user engagement and reactions, which slows in overnight hours—occurs earlier than the peak in rainfall. Moreover, in an ideal case, the tweet sentiment trend is continuous: gradually rising until its peak from the first trace of flooding and then gradually falling until nil. However, Twitter response may be sporadic,

depending upon users recognizing and acknowledging flooding on social media. To counteract these issues, certain assumptions are made:

- 12:00 AM to 06:00 AM, these 6 hours are considered 'off hours' because it falls within typical sleeping hours when most people are inactive. As a result, the number of tweets can become significantly less (Garett et al., 2018). The discontinuity of Twitter responses during the off-hours were disregarded. However, if any social media responses were found during off hours, they were included. It ensured consistency in responses to flooding.
- If discontinuity of Twitter responses for 2 or more consecutive hours other than the 'off-hours' is seen, the later responses after the discontinuation period were disregarded. This served to end the event in regard to real-time tracking. Thus, the false positive flood indications were avoided. In this case, the last tweet before the discontinuity was regarded as the 'end trace of flooding' point.

CHAPTER 4

RESULTS

The results of this study are presented in several key sections. First, sentiment distribution that shows the number of tweets in sentiment categories regarding flood-related tweets and maps their spatial distribution within Fulton County. Next, the correlation between sentiment spikes and rainfall data is explored to identify flood events, followed by removing false positives to refine detection accuracy. Real-time mapping of flood-affected areas highlights impacted zones, and finally, validation using HEC-RAS simulations confirms the reliability of social media-derived flood data.

4.1 Sentiment Distribution

The geographic distribution of sentiments due to flooding is a follow-through of sentiment analysis and geocoding processes, as described in sections 3.4.5 and 3.4.6, respectively. The categorization of the 'Flood GeoTweets Pool' among negative, positive, and neutral sentiments have been shown in Figure 19. As described in section 3.4.5, tweets with negative sentiments named 'Sen Flood GeoTweets Pool,' are processed further for geocoding. The geocoding process showed a wide spatial distribution of tweets throughout the United States and other parts of the world. The spatial distribution of this pool of tweets, limited to only Fulton County during the tweet collection timeframe, is illustrated in Figure 20.

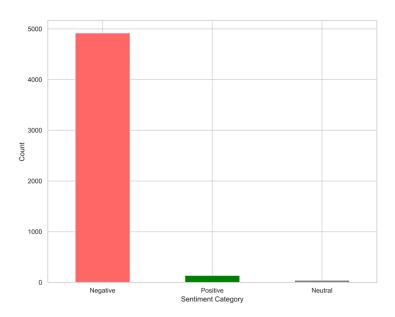


Figure 19: Number of tweets in each sentiment category. More than 4800 tweets were categorized as expressing negative sentiments. The number of tweets expressing positive and neutral sentiments are quite low (Around 300 and 100, respectively) compared to the negative ones.

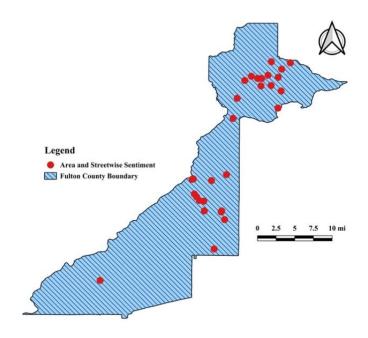


Figure 20: Sentiment distribution within Fulton County for the whole timeframe of tweet collection. As described in section 3.4.6, a buffer radius of 0.5 mile was considered around each coordinate and consequently, fusion rules were applied.

4.2 Identifying Flood Events

The daily tweet sentiments were plotted against the measured daily rainfall data (via USGS gauges) for the whole timeframe of Twitter data extraction. A single or, in most cases, clusters of continuous spikes were found associated with a rainfall event. Single spikes indicate that the flooding event started and ended on the same day. Whereas continuous spikes indicate that the flooding incident continues through multiple days. Figure 21 shows a visualization of rainfall and sentiment spikes.

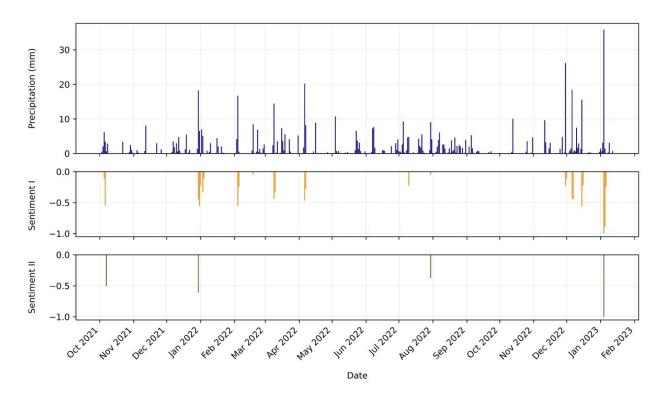


Figure 21: A time series plot of daily rainfall depth and normalized sum of daily sentiment scores from October 2021 to February 2023. Sentiment I is derived from area and streetwise location tweets, and Sentiment II is derived from countywise location tweets. Overall, sentiment spikes are seen for a heavy rainfall event in a single day, resulting in localized floods, or prolonged rainfall events for multiple days, even though the rainfall depth per day is not substantial.

4.3 Removal of False Positives

The hourly average rainfall amount calculated from the satellite rainfall data was plotted against the normalized hourly sum of sentiments. A single rainfall event between 3 and 4 January 2023 was considered for the hourly plot. Figure 22 shows the distribution plot.

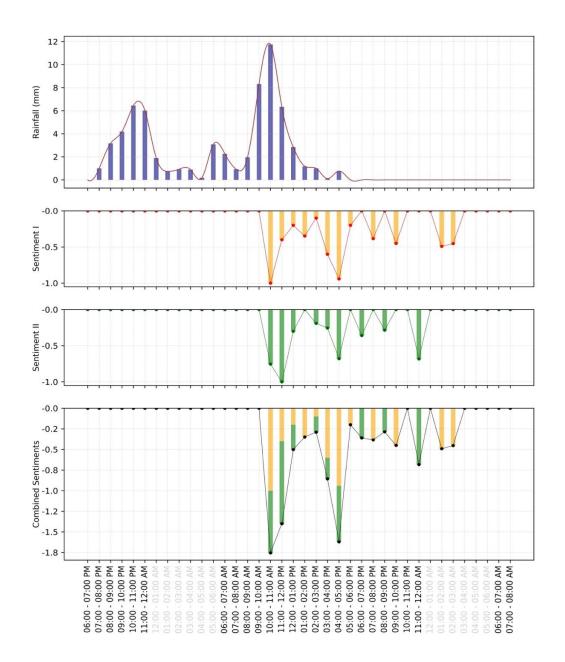


Figure 22: A time series plot of hourly average satellite rainfall depth and normalized sum of hourly sentiment scores for a single rainfall event that started on January 3, 2023, between 6:00 and 7:00 PM UTC and ended on January 4, 2023, between 4:00 and 5:00 PM UTC. Similar to Figure 21, Sentiment I is derived from area and streetwise location, and Sentiment II is derived from countywise location. For both cases, sentiment spikes appear between 10:00 and 11:00 AM on January 4. Combined sentiments represent the summation of Sentiment I and Sentiment II. The 'off hours' have been colored grey in the x-axis. Although there is no sentiment spike recorded between 10:00 and 11:00 PM on January 4, the data collection continues as it ends if a two-hour gap occurs. Two of the six 'off-hours' on January 5 showed sentiment spikes. Therefore, adhering to the conceptualization described in Section 3.8, they have been counted. However, when the 'off hours' ended, there was no spike seen for two consecutive hours on January 5. Hence, 8 AM UTC is considered as the cut-off time for this rainfall event, and any tweet after this time is likely to be a false positive or late post and, therefore, disregarded. Anywhere between 3:00 to 4:00 AM UTC was the end trace of flooding, according to social media.

4.4 Mapping

The tweets that contributed to the 'combined sentiments' in Figure 18 were the real-time user-generated posts. Twenty-six buffer areas were detected throughout Fulton County for the rainfall event between 3 and 4 January 2023. Figure 23 shows a detailed mapping of the flood-affected zones in real time. A flooded zone was defined as an area where the depth of standing water crossed 0.5 ft at minimum.

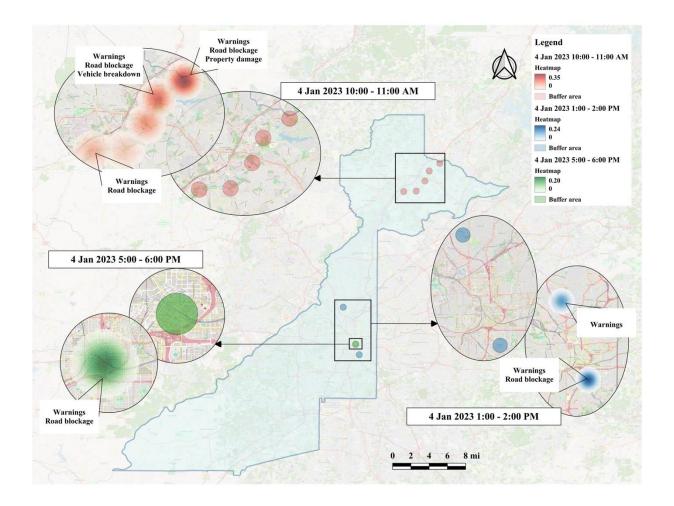


Figure 23: Near real-time mapping of flood-affected zones. Three time-based snapshots have been depicted in this figure for 4 January 2023. The red legend represents flood-affected zones between 10:00 and 11:00 AM UTC, while the blue and green legends represent 1:00 – 2:00 PM UTC and 5:00 – 6:00 PM UTC, respectively. Each time-based snapshot displays two insets. The inner insets reveal the land use, and the outer insets depict a heatmap. The heatmaps focus on the intensity of tweets, where the darker shades indicate comparatively dense Twitter responses. The heatmaps also have callouts that frequently mention emergencies. The red zones include streets, residential zones, and parkside areas near Big Creek vicinity consisting of emergecies such as flood warnings, road blockages, property damage, etc. The blue and green

zones near downtown Atlanta include streets and commercial zones consisting of emergencies such as warnings, road blockages, etc.

4.5 Validation

The hourly gridded rainfall data for 3 to 4 January 2023 was assigned to the HEC-RAS model for the simulation process. The model-generated flood maps for Fulton County were used to validate the Twitter data-generated flood-affected zone map. The spatiotemporal validation results for 10:00 to 11:00 AM UTC are shown in Figure 24.

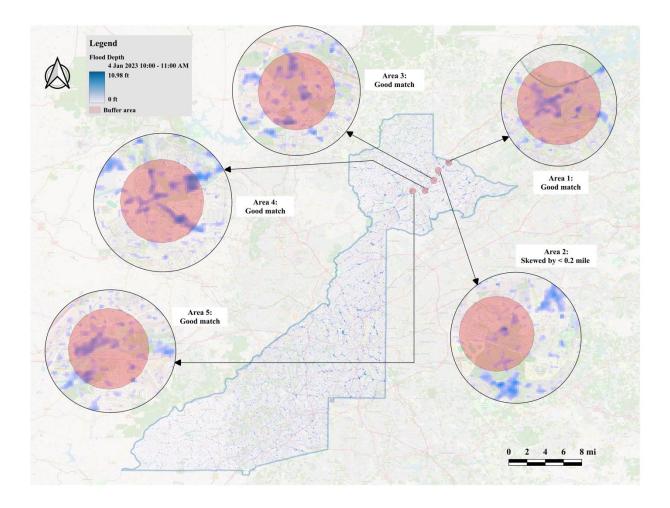


Figure 24: Validation results for 4 January 2023 10:00 – 11:00 AM UTC. The buffer area from Twitter data has been depicted in the red circles from Area 1 to Area 5. The HEC-RAS simulated flood scenario shows a good match where the skewness is less than 0.1 miles for Area 1, Area 3, Area 4, and Area 5 for the timeframe. The flood depth varies between places and reaches up to 4 ft near water bodies in these zones. However, for Area 2, the Twitter-derived area shows a slight skewness of less than 0.2 miles from a major flooded area of at least 0.5 ft in depth from the model-generated results.

Table 8 describes an overview of 26 buffer areas relative to a corresponding or nearby major flooded area with at least 0.5 ft flood depth.

Table 8. Validation details of buffer area based on flood depth.

| Status | Number of buffer areas |
|---------------------------------------|------------------------|
| Good Match | 19 |
| Skewed by ≤ 0.2 miles | 3 |
| Skewed by > 0.2 and ≤ 0.5 miles | 2 |
| Skewed by > 0.5 miles | 2 |

Nearly 75% of the buffer areas show a good match with HEC-RAS simulated results. Only two buffer areas skewed more than 0.5 miles. Although the detected buffer areas showed satisfactory results, many flooded zones that were simulated in HEC-RAS for the timeframe could not be traced in Twitter data. Figure 25 shows some of the undetected areas. Despite having a few undetected flooded zones, the maps provide valuable insights. Many of the flooded zones were located near meandering creeks and water bodies where surging waters overwhelmed their channels. A few occurred due to stormwater accumulation in low-lying areas in cities, streets, parks, or other land use zones.

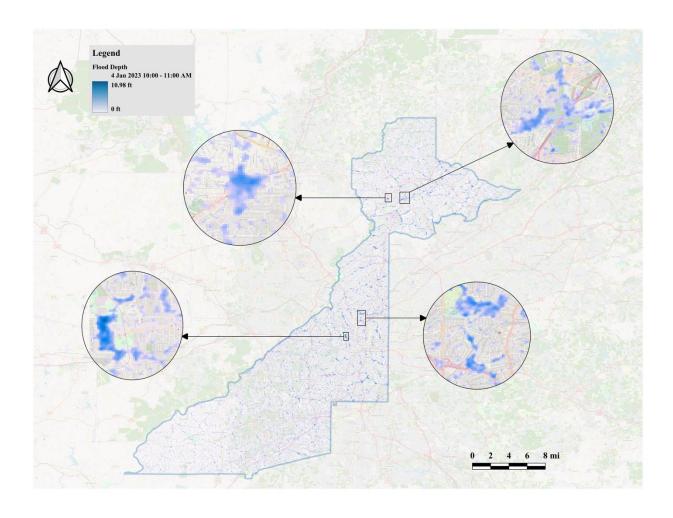


Figure 25: Flooded areas that could not be traced with Twitter data. Although some of these points crossed busy highways, residential areas, commercial areas, or park sides having flood depths ranging from 0 - 5 ft, these areas did not show up in Twitter data.

CHAPTER 5

DISCUSSIONS

The findings of this study demonstrate the potential of social media as a valuable source of crowdsourced data for enhancing urban flood monitoring. By analyzing tweets related to flooding, this study successfully identified flooding events, mapped affected areas in near real-time, and validated these findings against hydrological model simulations. The integration of social media with rainfall data provided deeper insights into both the spatiotemporal distribution of flooding and public sentiment during such events. A major highlight from this study is the extraction of location-based flood information from non-geotagged tweets. Initially, less than 1% of the tweets contained built-in geolocation metadata. However, by leveraging a regular expression-based pattern recognition technique, approximately 13% of the tweets were successfully geolocated based on user-generated location descriptions. This represents a substantial improvement in spatial resolution for flood monitoring as the extracted locations contained more specific to detailed information than tools like Named Entity Recognition for this case. It indicates that user-reported location references can be effectively harnessed for disaster response.

Developing a training dataset by manual labeling for the machine learning classification of flood-related tweets was a crucial step. Support Vector Machine outperformed other classifiers and was selected for the classification task. The temporal analysis of rainfall and sentiment trends provided valuable insights into the lag between precipitation events and public perception of flooding. Although daily rainfall readings from 15 rain stations across Fulton County have been used in this study, there is a discrepancy in daily rainfall amount between the values of gauge rainfall data and satellite rainfall data. It is because not all rain gauges were active throughout the timeframe of tweet mining, as shown in Figure 11. However, the recorded gauge rainfall data still provided valuable insights in identifying flood events, as shown in section 4.2. The flood maps provided a good resolution as the buffer areas were circular, having a radius of 0.5

miles. The radius was set as an optimum value for this study, as a lower radius would increase the redundancy of the flooded zones and a higher radius would lead to less accuracy (Fan et al., 2020). In addition, validation of the Twitter-derived flood maps with this radius of buffer area against HEC-RAS simulated flood scenarios revealed a good degree of agreement. Approximately 75% of the buffer areas closely matched the flood zones identified in the hydrological simulations. Some discrepancies were observed, with a small number of buffer areas showing minor spatial offsets of up to 0.5 miles. Additionally, while social media data effectively captured flooding in urbanized areas, some flooded regions detected by the hydrological model were absent in the Twitter dataset. This could be attributed to several factors, including limited social media activity in certain locations leading to tweets not being posted or some busy highways or city areas may have a very fast water dissipation system that there was no flooding issue in reality (Li et al., 2023). Alternatively, these instances may simply reflect surface runoff rather than true flooding. It is also possible that local authorities could take proactive flood mitigation measures – such as clearing and draining waterlogged areas, releasing water from upstream reservoirs, etc. – before significant social media engagement occurred.

5.1 Opportunities

The study presents several opportunities for enhancing flood monitoring and disaster response. One of the major advantages is the ability to capture near real-time, user-generated data, allowing for immediate situational awareness. Unlike traditional monitoring methods that rely on fixed sensors or satellite imagery with time delays, the study offers instant reporting from individuals directly experiencing the event. Furthermore, the scalability of this approach means that it can be implemented in regions where traditional monitoring infrastructure is limited or non-existent. Additionally, public perception, emotional distress, or urgency can be understood by incorporating sentiment analysis. For example, in Figure 23, specific emergencies – such as road blockages, vehicle breakdowns, and property damage – could be identified.

The findings of this study have significant implications for practice in disaster response and urban planning. Although the depth of the floodwater cannot be directly known from text-based tweets, they provide valuable humanitarian information. Emergency management agencies can integrate social media monitoring into their warning systems to enhance flood response strategies. By identifying near real-time reports of flooding as depicted in Figure 18, it is easier to identify the nature of the emergency at each hotspot. Therefore, first responders can allocate resources more effectively based on the needs of the buffer areas. Additionally, transportation authorities can use this approach to assess road closures and traffic disruptions, improving public safety during flood events. For urban planners, the insights from social media data can inform long-term flood mitigation efforts. By analyzing repeated flooding reports in specific locations, planners can identify vulnerable areas that require improved drainage infrastructure, flood barriers, or other adaptation measures. Policymakers can also use social media sentiment analysis to gauge public concerns and perceptions of flood risk, which can guide the development of more effective communication and preparedness campaigns.

5.2 Limitations

One key challenge in utilizing social media data for flood monitoring is the inherent variability and limitations in data collection, processing, and interpretation. These limitations arise from methodological constraints, biases in data availability, and the evolving nature of digital communication, all of which can impact the accuracy and reliability of the results.

There are methodological factors upon which the results can vary significantly, such as the K-means clustering of tweets, which revealed six themes in this study. Although the silhouette scores imply that the clusters were not mutually distinct (having substantial overlaps), it paved the way to the succeeding steps: identifying patterns for location detecting using regular expressions and developing a training dataset to identify flood relevant tweets. Moreover, robust filtering techniques have been adopted in this study to remove inconsistent, misinformative data. The quality of social media data remains a concern as the local

context may lead to the need for filtering techniques to be adjusted. For example, a city in Kentucky is named Rain. Separating flood-relevant tweets as described in section 3.4.4, requires manual labeling of tweets to develop the training dataset. The task can be painstaking for large datasets. Another limitation of the data quality is the amount of information that can be derived from a social media post. Text-based tweets do not directly provide insights on floodwater depth; future research can be conducted by integrating relevant images and videos gathered from Twitter for a more holistic approach to flood monitoring. For example, in a study by (Couey et al., 2022), camera pictures were used to train and evaluate machine learning models to detect recent moisture in street-level images, aiding in early flood detection. In another instance, (Alizadeh et al., 2022) used crowdsourced street photos to enhance flood mapping and optimize evacuation routes through image processing. Furthermore, social media data is inherently biased towards populations with high internet users and active social media usage. This creates gaps in data availability, particularly in rural or underserved regions (Olteanu et al., 2019). Besides, the accuracy of the regular expression pattern recognition method for location extraction can be compromised for a more diverse or big dataset. It is because as the dataset becomes larger, increasing variations in user-reported locations, ambiguous place names, and informality in language can lead to inconsistency in pattern formation, necessitating continuous improvements in geolocation extraction methods. Besides, two hours of nonactivity in Twitter has been arbitrarily assumed as the cutoff point for the removal of false positives, as described in section 4.3. This is not absolute, as it may not be the case for another rainfall event in another region.

CHAPTER 6

CONCLUSION

This study demonstrates a novel approach of leveraging social media data for near real-time urban flood monitoring. The integration of crowdsourced Twitter data and hourly rainfall data allowed for successfully identifying flood events, mapping affected areas, and validating results against hydrological simulations. The findings highlight the effectiveness of social media in capturing localized flood conditions, particularly in urban areas where conventional monitoring methods face limitations. This approach offers valuable insights for enhancing flood response strategies, optimizing resource allocation, and improving real-time situational awareness. The integration of social media data into flood monitoring frameworks presents a promising avenue for strengthening urban resilience and disaster preparedness.

REFERENCES

- Ahmadi, C., Karampourian, A., & Samarghandi, M. R. (2022). Explain the challenges of evacuation in floods based on the views of citizens and executive managers. *Heliyon*, 8(9), e10759. https://doi.org/10.1016/J.HELIYON.2022.E10759
- Alizadeh, B., Li, D., Hillin, J., Meyer, M. A., Thompson, C. M., Zhang, Z., & Behzadan, A. H. (2022). Human-centered flood mapping and intelligent routing through augmenting flood gauge data with crowdsourced street photos. *Advanced Engineering Informatics*, *54*. https://doi.org/10.1016/j.aei.2022.101730
- Andreasen, M. H., Agergaard, J., Allotey, A. N. M., Møller-Jensen, L., & Oteng-Ababio, M. (2023). Built-in Flood Risk: the Intertwinement of Flood Risk and Unregulated Urban Expansion in African Cities. *Urban Forum*, *34*(3), 385–411. https://doi.org/10.1007/s12132-022-09478-4
- Arnold, C. L., & Gibbons, C. J. (1996). Impervious Surface Coverage: The Emergence of a Key Environmental Indicator. *Journal of the American Planning Association*, 62(2), 243–258. https://doi.org/10.1080/01944369608975688
- Arolfo, F., Rodriguez, K. C., & Vaisman, A. (2022). Analyzing the Quality of Twitter Data Streams. *Information Systems Frontiers*, 24(1), 349–369. https://doi.org/10.1007/s10796-020-10072-x
- Baeza-Yates, R. (2018). Bias on the web. *Communications of the ACM*, 61(6), 54–61. https://doi.org/10.1145/3209581
- Balaian, S. K., Sanders, B. F., & Abdolhosseini Qomi, M. J. (2024). How urban form impacts flooding. *Nature Communications*, *15*(1). https://doi.org/10.1038/s41467-024-50347-4
- Barker, J. L. P., & Macleod, C. J. A. (2019). Development of a national-scale real-time Twitter data mining pipeline for social geodata on the potential impacts of flooding on communities. *Environmental Modelling and Software*, 115, 213–227. https://doi.org/10.1016/j.envsoft.2018.11.013
- Basaria, A. A. A., Ahsan, A., Nadeem, A., Tariq, R., & Raufi, N. (2023). Infectious diseases following hydrometeorological disasters: current scenario, prevention, and control measures. *Annals of Medicine & Surgery*, 85(8), 3778–3782. https://doi.org/10.1097/ms9.0000000000001056
- BenDor, T. K., Salvesen, D., Kamrath, C., & Ganser, B. (2020). Floodplain Buyouts and Municipal Finance. *Natural Hazards Review*, 21(3). https://doi.org/10.1061/(asce)nh.1527-6996.0000380
- Cervone, G., Sava, E., Huang, Q., Schnebele, E., Harrison, J., & Waters, N. (2016). Using Twitter for tasking remote-sensing data collection and damage assessment: 2013 Boulder flood case study.

- *International Journal of Remote Sensing*, *37*(1), 100–124. https://doi.org/10.1080/01431161.2015.1117684
- Chang, H., Pallathadka, A., Sauer, J., Grimm, N. B., Zimmerman, R., Cheng, C., Iwaniec, D. M., Kim, Y., Lloyd, R., McPhearson, T., Rosenzweig, B., Troxler, T., Welty, C., Brenner, R., & Herreros-Cantis, P. (2021). Assessment of urban flood vulnerability using the social-ecological-technological systems framework in six US cities. *Sustainable Cities and Society*, 68. https://doi.org/10.1016/j.scs.2021.102786
- Chen, X., Elmes, G., Ye, X., & Chang, J. (2016). Implementing a real-time Twitter-based system for resource dispatch in disaster management. *GeoJournal*, 81(6), 863–873. https://doi.org/10.1007/s10708-016-9745-8
- Clarke, B., Otto, F., Stuart-Smith, R., & Harrington, L. (2022). Extreme weather impacts of climate change: an attribution perspective. *Environmental Research: Climate*, *1*(1), 012001. https://doi.org/10.1088/2752-5295/ac6e7d
- Coalson, J. E., Anderson, E. J., Santos, E. M., Garcia, V. M., Romine, J. K., Dominguez, B., Richard, D. M., Little, A. C., Hayden, M. H., & Ernst, K. C. (2021). The complex epidemiological relationship between flooding events and human outbreaks of mosquito-borne diseases: A scoping review. In *Environmental Health Perspectives* (Vol. 129, Issue 9). Public Health Services, US Dept of Health and Human Services. https://doi.org/10.1289/EHP8887
- Couey, B., Doerry, E., Gowanlock, M., & Hocking, T. (2022). *AN EVALUATION OF MACHINE LEARNING TO DETECT RECENT MOISTURE IN STREET LEVEL IMAGES*.
- Crooks, A., Croitoru, A., Stefanidis, A., & Radzikowski, J. (2013). #Earthquake: Twitter as a Distributed Sensor System. *Transactions in GIS*, 17(1), 124–147. https://doi.org/10.1111/j.1467-9671.2012.01359.x
- Duy, P. N., Chapman, L., Tight, M., Linh, P. N., & Thuong, L. V. (2018). Increasing vulnerability to floods in new development areas: evidence from Ho Chi Minh City. *International Journal of Climate Change Strategies and Management*, 10(1), 197–212. https://doi.org/10.1108/IJCCSM-12-2016-0169
- Easterling, D., Rusticucci, M., Semenov, V., Alexander, L. V, Allen, S., Benito, G., Cavazos, T., Nicholls, N., Easterling, D., Goodess, C., Kanae, S., Kossin, J., Luo, Y., Marengo, J., McInnes, K., Rahimi, M., Reichstein, M., Sorteberg, A., Vera, C., ... Midgley, P. (2012). 3 Changes in Climate Extremes and their Impacts on the Natural Physical Environment. Cambridge University Press.
- Erokhin, D., & Komendantova, N. (2024). Social media data for disaster risk management and research. *International Journal of Disaster Risk Reduction*, 114. https://doi.org/10.1016/j.ijdrr.2024.104980
- Fan, C., Wu, F., & Mostafavi, A. (2020). A Hybrid Machine Learning Pipeline for Automated Mapping of Events and Locations from Social Media in Disasters. *IEEE Access*, 8, 10478–10490. https://doi.org/10.1109/ACCESS.2020.2965550

- Farhadi, H., Esmaeily, A., & Najafzadeh, M. (2022). Flood monitoring by integration of Remote Sensing technique and Multi-Criteria Decision Making method. *Computers & Geosciences*, *160*, 105045. https://doi.org/10.1016/J.CAGEO.2022.105045
- Feng, Q., Liu, J., & Gong, J. (2015). Urban flood mapping based on unmanned aerial vehicle remote sensing and random forest classifier-A case of yuyao, China. *Water (Switzerland)*, 7(4), 1437–1455. https://doi.org/10.3390/w7041437
- Finley, P., Gatti, G., Goodall, J., Nelson, M., Nicholson, K., & Shah, K. (2020). Flood Monitoring and Mitigation Strategies for Flood-Prone Urban Areas. 2020 Systems and Information Engineering Design Symposium, SIEDS 2020. https://doi.org/10.1109/SIEDS49339.2020.9106583
- Fuchs, G., Andrienko, N., Andrienko, G., Bothe, S., & Stange, H. (2013). Tracing the German centennial flood in the stream of tweets: First lessons learned. *GEOCROWD 2013 Proceedings of the 2nd ACM SIGSPATIAL International Workshop on Crowdsourced and Volunteered Geographic Information*, 31–38. https://doi.org/10.1145/2534732.2534741
- Garett, R., Liu, S., & Young, S. D. (2018). The relationship between social media use and sleep quality among undergraduate students. *Information Communication and Society*, 21(2), 163–173. https://doi.org/10.1080/1369118X.2016.1266374
- Granell, C., & Ostermann, F. O. (2016). Beyond data collection: Objectives and methods of research using VGI and geo-social media for disaster management. *Computers, Environment and Urban Systems*, 59, 231–243. https://doi.org/10.1016/j.compenvurbsys.2016.01.006
- Gulnerman, A. G., & Karaman, H. (2020). Spatial reliability assessment of social media mining techniques with regard to disaster domain-based filtering. *ISPRS International Journal of Geo-Information*, 9(4). https://doi.org/10.3390/ijgi9040245
- Han, S., & Tahvildari, N. (2024). Compound Flooding Hazards Due To Storm Surge and Pluvial Flow in a Low-Gradient Coastal Region. *Water Resources Research*, 60(11), e2023WR037014. https://doi.org/10.1029/2023WR037014
- Helmrich, A., Kuhn, A., Roque, A., Santibanez, A., Kim, Y., Grimm, N. B., & Chester, M. (2023). Interdependence of social-ecological-technological systems in Phoenix, Arizona: consequences of an extreme precipitation event. *Journal of Infrastructure Preservation and Resilience*, 4(1). https://doi.org/10.1186/s43065-023-00085-6
- Helmrich, A. M., Ruddell, B. L., Bessem, K., Chester, M. V., Chohan, N., Doerry, E., Eppinger, J., Garcia, M., Goodall, J. L., Lowry, C., & Zahura, F. T. (2021). Opportunities for crowdsourcing in urban flood monitoring. In *Environmental Modelling and Software* (Vol. 143). Elsevier Ltd. https://doi.org/10.1016/j.envsoft.2021.105124
- Henonin, J., Russo, B., Mark, O., & Gourbesville, P. (2013). Real-time urban flood forecasting and modelling A state of the art. *Journal of Hydroinformatics*, 15(3), 717–736. https://doi.org/10.2166/hydro.2013.132

- Hino, M., BenDor, T. K., Branham, J., Kaza, N., Sebastian, A., & Sweeney, S. (2024). Growing Safely or Building Risk?: Floodplain Management in North Carolina. *Journal of the American Planning Association*, 90(1), 50–62. https://doi.org/10.1080/01944363.2022.2141821
- Huang, Q., & Xiao, Y. (2015). Geographic situational awareness: Mining tweets for disaster preparedness, emergency response, impact, and recovery. *ISPRS International Journal of Geo-Information*, 4(3), 1549–1568. https://doi.org/10.3390/ijgi4031549
- Huang, X., Wang, C., & Li, Z. (2018). Reconstructing flood inundation probability by enhancing near real-time imagery with real-time gauges and tweets. *IEEE Transactions on Geoscience and Remote Sensing*, 56(8), 4691–4701. https://doi.org/10.1109/TGRS.2018.2835306
- Hutto, C. J., & Gilbert, E. (2014). VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text. Proceedings of the International AAAI Conference on Web and Social Media, 8(1), 216–225. https://doi.org/10.1609/ICWSM.V8I1.14550
- Iacus, S. M., Porro, G., Salini, S., & Siletti, E. (2020). Controlling for Selection Bias in Social Media Indicators through Official Statistics: A Proposal. *Journal of Official Statistics*, *36*(2), 315–338. https://doi.org/10.2478/jos-2020-0017
- Ikotun, A. M., Ezugwu, A. E., Abualigah, L., Abuhaija, B., & Heming, J. (2023). K-means clustering algorithms: A comprehensive review, variants analysis, and advances in the era of big data. *Information Sciences*, 622, 178–210. https://doi.org/10.1016/J.INS.2022.11.139
- Imran, M., Castillo, C., Lucas, J., Meier, P., & Vieweg, S. (2014). AIDR: Artificial intelligence for disaster response. *WWW 2014 Companion Proceedings of the 23rd International Conference on World Wide Web*, 159–162. https://doi.org/10.1145/2567948.2577034
- Jang, B. J., & Jung, I. (2023). Development of High-Precision Urban Flood-Monitoring Technology for Sustainable Smart Cities. *Sensors*, 23(22). https://doi.org/10.3390/s23229167
- Jongman, B., Wagemaker, J., Revilla Romero, B., & Coughlan De Perez, E. (2015). Early flood detection for rapid humanitarian response: Harnessing near real-time satellite and twitter signals. *ISPRS International Journal of Geo-Information*, 4(4), 2246–2266. https://doi.org/10.3390/ijgi4042246
- Karimiziarani, M., Jafarzadegan, K., Abbaszadeh, P., Shao, W., & Moradkhani, H. (2022). Hazard risk awareness and disaster management: Extracting the information content of twitter data. *Sustainable Cities and Society*, 77. https://doi.org/10.1016/j.scs.2021.103577
- Karmegam, D., & Mappillairaju, B. (2020). Spatiooral distribution of negative emotions on Twitter during floods in Chennai, India, in 2015: A post hoc analysis. *International Journal of Health Geographics*, 19(1), 1–13. https://doi.org/10.1186/S12942-020-00214-4/FIGURES/4
- Karyotis, C., Maniak, T., Doctor, F., Iqbal, R., Palade, V., & Tang, R. (2019). Deep learning for flood forecasting and monitoring in urban environments. *Proceedings 18th IEEE International*

- Conference on Machine Learning and Applications, ICMLA 2019, 1392–1397. https://doi.org/10.1109/ICMLA.2019.00227
- Koutsovili, E. I., Tzoraki, O., Theodossiou, N., & Tsekouras, G. E. (2023). Early Flood Monitoring and Forecasting System Using a Hybrid Machine Learning-Based Approach. *ISPRS International Journal of Geo-Information*, *12*(11). https://doi.org/10.3390/ijgi12110464
- Kryvasheyeu, Y., Chen, H., Obradovich, N., Moro, E., Van Hentenryck, P., Fowler, J., & Cebrian, M. (2016). Rapid assessment of disaster damage using social media activity. *Science Advances*, 2(3). https://doi.org/10.1126/sciadv.1500779
- Lachlan, K. A., Spence, P. R., Lin, X., Najarian, K., & Del Greco, M. (2016). Social media and crisis management: CERC, search strategies, and Twitter content. *Computers in Human Behavior*, *54*, 647–652. https://doi.org/10.1016/j.chb.2015.05.027
- Li, Y., Osei, F. B., Hu, T., & Stein, A. (2023). Urban flood susceptibility mapping based on social media data in Chengdu city, China. *Sustainable Cities and Society*, 88. https://doi.org/10.1016/j.scs.2022.104307
- Liu, J., Cho, H. S., Osman, S., Jeong, H. G., & Lee, K. (2022). Review of the status of urban flood monitoring and forecasting in TC region. *Tropical Cyclone Research and Review*, 11(2), 103–119. https://doi.org/10.1016/j.tcrr.2022.07.001
- Mendoza, M., Poblete, B., & Valderrama, I. (2019). Nowcasting earthquake damages with Twitter. *EPJ Data Science*, 8(1). https://doi.org/10.1140/epjds/s13688-019-0181-0
- Middleton, S. E., Middleton, L., & Modafferi, S. (2020). Real-time crisis mapping of natural disasters using social media. *IEEE Intelligent Systems*, 29(2), 9–17. https://doi.org/10.1109/MIS.2013.126
- Miller, J. D., & Hutchins, M. (2017). The impacts of urbanisation and climate change on urban flooding and urban water quality: A review of the evidence concerning the United Kingdom. In *Journal of Hydrology: Regional Studies* (Vol. 12, pp. 345–362). Elsevier B.V. https://doi.org/10.1016/j.ejrh.2017.06.006
- Moulds, S., Buytaert, W., Templeton, M. R., & Kanu, I. (2021). Modeling the Impacts of Urban Flood Risk Management on Social Inequality. *Water Resources Research*, *57*(6). https://doi.org/10.1029/2020WR029024
- Mousa, M., Zhang, X., & Claudel, C. (2016). Flash Flood Detection in Urban Cities Using Ultrasonic and Infrared Sensors. *IEEE Sensors Journal*, *16*(19), 7204–7216. https://doi.org/10.1109/JSEN.2016.2592359
- Muralidharan, S., Rasmussen, L., Patterson, D., & Shin, J. H. (2011). Hope for Haiti: An analysis of Facebook and Twitter usage during the earthquake relief efforts. *Public Relations Review*, *37*(2), 175–177. https://doi.org/10.1016/j.pubrev.2011.01.010

- Nguyen, P., Shearer, E. J., Tran, H., Ombadi, M., Hayatbini, N., Palacios, T., Huynh, P., Braithwaite, D., Updegraff, G., Hsu, K., Kuligowski, B., Logan, W. S., & Sorooshian, S. (2019). The CHRS Data Portal, an easily accessible public repository for PERSIANN global satellite precipitation data. *Scientific Data* 2019 6:1, 6(1), 1–10. https://doi.org/10.1038/sdata.2018.296
- Nkwunonwo, U. C., Whitworth, M., & Baily, B. (2020). A review of the current status of flood modelling for urban flood risk management in the developing countries. In *Scientific African* (Vol. 7). Elsevier B.V. https://doi.org/10.1016/j.sciaf.2020.e00269
- O'Driscoll, M., Clinton, S., Jefferson, A., Manda, A., & McMillan, S. (2010). Urbanization effects on watershed hydrology and in-stream processes in the southern United States. In *Water (Switzerland)* (Vol. 2, Issue 3, pp. 605–648). MDPI AG. https://doi.org/10.3390/w2030605
- Ogie, R. I., James, S., Moore, A., Dilworth, T., Amirghasemi, M., & Whittaker, J. (2022). Social media use in disaster recovery: A systematic literature review. In *International Journal of Disaster Risk Reduction* (Vol. 70). Elsevier Ltd. https://doi.org/10.1016/j.ijdrr.2022.102783
- Olteanu, A., Castillo, C., Diaz, F., & Kıcıman, E. (2019). Social Data: Biases, Methodological Pitfalls, and Ethical Boundaries. In *Frontiers in Big Data* (Vol. 2). Frontiers Media S.A. https://doi.org/10.3389/fdata.2019.00013
- Panteras, G., & Cervone, G. (2018). Enhancing the temporal resolution of satellite-based flood extent generation using crowdsourced data for disaster monitoring. *International Journal of Remote Sensing*, *39*(5), 1459–1474. https://doi.org/10.1080/01431161.2017.1400193
- Papadopoulos, S., Bontcheva, K., Jaho, E., Lupu, M., & Castillo, C. (2016). Overview of the special issue on trust and veracity of information in social media. In *ACM Transactions on Information Systems* (Vol. 34, Issue 3). Association for Computing Machinery. https://doi.org/10.1145/2870630
- Paradkar, A. S., Zhang, C., Yuan, F., & Mostafavi, A. (2022). Examining the consistency between geocoordinates and content-mentioned locations in tweets for disaster situational awareness: A Hurricane Harvey study. *International Journal of Disaster Risk Reduction*, 73. https://doi.org/10.1016/j.ijdrr.2022.102878
- Persily, N., & Tucker, J. A. (2020). Social media and democracy: The state of the field, prospects for reform. In *Social Media and Democracy: The State of the Field, Prospects for Reform*. Cambridge University Press. https://doi.org/10.1017/9781108890960
- Pramanick, A., Beck, T., Stowe, K., & Gurevych, I. (2021). *The challenges of temporal alignment on Twitter during crises*. http://arxiv.org/abs/2104.08535
- Saha, A., & Chandra Pal, S. (2024). Application of machine learning and emerging remote sensing techniques in hydrology: A state-of-the-art review and current research trends. In *Journal of Hydrology* (Vol. 632). Elsevier B.V. https://doi.org/10.1016/j.jhydrol.2024.130907

- Sawaneh, I. A., Fan, L., & Sesay, B. (2024). Investigating the influence of residents' attitudes, perceptions of risk, and subjective norms on their willingness to engage in flood prevention efforts in Freetown, Sierra Leone. *Nature-Based Solutions*, 6, 100143. https://doi.org/10.1016/j.nbsj.2024.100143
- See, L. (2019). A review of citizen science and crowdsourcing in applications of pluvial flooding. In *Frontiers in Earth Science* (Vol. 7). Frontiers Media S.A. https://doi.org/10.3389/feart.2019.00044
- Smith, L., Liang, Q., James, P., & Lin, W. (2017). Assessing the utility of social media as a data source for flood risk management using a real-time modelling framework. *Journal of Flood Risk Management*, 10(3), 370–380. https://doi.org/10.1111/jfr3.12154
- Snyder, L. S., Lin, Y. S., Karimzadeh, M., Goldwasser, D., & Ebert, D. S. (2020). Interactive Learning for Identifying Relevant Tweets to Support Real-time Situational Awareness. *IEEE Transactions on Visualization and Computer Graphics*, 26(1), 558–568. https://doi.org/10.1109/TVCG.2019.2934614
- Son, C. H., Lee, C. H., & Ban, Y. U. (2023). Analysis of the impact and moderating effect of high-density development on urban flooding. *Heliyon*, *9*(12). https://doi.org/10.1016/j.heliyon.2023.e22695
- Song, J., Chang, Z., Li, W., Feng, Z., Wu, J., Cao, Q., & Liu, J. (2019). Resilience-vulnerability balance to urban flooding: A case study in a densely populated coastal city in China. *Cities*, 95, 102381. https://doi.org/10.1016/J.CITIES.2019.06.012
- Subarkah, A., Kusnanto, G., Permai, S. D., Ohyver, M., & Arifin, S. (2023). Tweet congestion locations identification using natural language processing. *AIP Conference Proceedings*, 2733(1). https://doi.org/10.1063/5.0140149
- Sunkpho, J., & Ootamakorn, C. (2011). Real-time flood monitoring and warning system. In *Songklanakarin J. Sci. Technol* (Vol. 33, Issue 2). http://www.sjst.psu.ac.th
- Tang, Y., Sun, Y., Han, Z., Soomro, S. e. hyder, Wu, Q., Tan, B., & Hu, C. (2023). flood forecasting based on machine learning pattern recognition and dynamic migration of parameters. *Journal of Hydrology: Regional Studies*, 47. https://doi.org/10.1016/j.ejrh.2023.101406
- Tanim, A. H., McRae, C. B., Tavakol-davani, H., & Goharian, E. (2022). Flood Detection in Urban Areas Using Satellite Imagery and Machine Learning. *Water 2022, Vol. 14, Page 1140, 14*(7), 1140. https://doi.org/10.3390/W14071140
- Tao, Y., Tian, B., Adhikari, B. R., Zuo, Q., Luo, X., & Di, B. (2024). A Review of Cutting-Edge Sensor Technologies for Improved Flood Monitoring and Damage Assessment. In *Sensors* (Vol. 24, Issue 21). Multidisciplinary Digital Publishing Institute (MDPI). https://doi.org/10.3390/s24217090
- Wakode, H. B., Baier, K., Jha, R., & Azzam, R. (2018). Impact of urbanization on groundwater recharge and urban water balance for the city of Hyderabad, India. *International Soil and Water Conservation Research*, 6(1), 51–62. https://doi.org/10.1016/j.iswcr.2017.10.003

- Wang, K., Lam, N. S. N., & Mihunov, V. (2023). Correlating Twitter Use with Disaster Resilience at Two Spatial Scales: A Case Study of Hurricane Sandy. *Annals of GIS*, 29(1), 1–20. https://doi.org/10.1080/19475683.2023.2165545
- Wang, K., Lam, N. S. N., Zou, L., & Mihunov, V. (2021). Twitter use in hurricane isaac and its implications for disaster resilience. *ISPRS International Journal of Geo-Information*, 10(3). https://doi.org/10.3390/ijgi10030116
- Wang, Z., & Ye, X. (2018). Social media analytics for natural disaster management. In *International Journal of Geographical Information Science* (Vol. 32, Issue 1, pp. 49–72). Taylor and Francis Ltd. https://doi.org/10.1080/13658816.2017.1367003
- Wang, Z., & Ye, X. (2019). Space, time, and situational awareness in natural hazards: a case study of Hurricane Sandy with social media data. *Cartography and Geographic Information Science*, 46(4), 334–346. https://doi.org/10.1080/15230406.2018.1483740
- *X/Twitter: number of users worldwide 2024 | Statista.* (2023). Retrieved March 7, 2025, from https://www.statista.com/statistics/303681/twitter-users-worldwide/
- Yenkar, P., & Sawarkar, S. D. (2021). Gazetteer based unsupervised learning approach for location extraction from complaint tweets. *IOP Conference Series: Materials Science and Engineering*, 1049(1), 012009. https://doi.org/10.1088/1757-899x/1049/1/012009
- Yilmaz, K. K., Adler, R. F., Tian, Y., Hong, Y., & Pierce, H. F. (2010). Evaluation of a satellite-based global flood monitoring system. *International Journal of Remote Sensing*, *31*(14), 3763–3782. https://doi.org/10.1080/01431161.2010.483489
- Zade, H., Shah, K., Rangarajan, V., Kshirsagar, P., Imran, M., & Starbird, K. (2018). From situational awareness to actionability: Towards improving the utility of social media data for crisis response. *Proceedings of the ACM on Human-Computer Interaction*, 2(CSCW). https://doi.org/10.1145/3274464
- Zhai, W. (2022). A multi-level analytic framework for disaster situational awareness using Twitter data. *Computational Urban Science*, 2(1). https://doi.org/10.1007/s43762-022-00052-z
- Zhang, S., He, L., Vucetic, S., & Dragut, E. C. (1991). *Regular Expression Guided Entity Mention Mining from Noisy Web Data*. Association for Computational Linguistics.
- Zhu, X., Guo, H., & Huang, J. J. (2024). Urban flood susceptibility mapping using remote sensing, social sensing and an ensemble machine learning model. *Sustainable Cities and Society*, *108*, 105508. https://doi.org/10.1016/J.SCS.2024.105508
- Zhu, Z. J., Jiang, A. Z., Lai, J., Xiang, Y., Baird, B., & McBean, E. (2017). Towards Efficient Use of an Unmanned Aerial Vehicle for Urban Flood Monitoring. *Journal of Water Management Modeling*. https://doi.org/10.14796/JWMM.C433

Zou, L., Lam, N. S. N., Cai, H., & Qiang, Y. (2018). Mining Twitter Data for Improved Understanding of Disaster Resilience. *Annals of the American Association of Geographers*, *108*(5), 1422–1441. https://doi.org/10.1080/24694452.2017.1421897