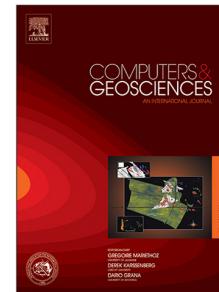


Journal Pre-proof

Real-time social media sentiment analysis for rapid impact assessment of floods

Lydia Bryan-Smith, Jake Godsall, Franky George, Kelly Egode,
Nina Dethlefs, Dan Parsons



PII: S0098-3004(23)00109-7

DOI: <https://doi.org/10.1016/j.cageo.2023.105405>

Reference: CAGEO 105405

To appear in: *Computers and Geosciences*

Received date: 19 October 2022

Revised date: 16 June 2023

Accepted date: 17 June 2023

Please cite this article as: L. Bryan-Smith, J. Godsall, F. George et al., Real-time social media sentiment analysis for rapid impact assessment of floods. *Computers and Geosciences* (2023), doi: <https://doi.org/10.1016/j.cageo.2023.105405>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2023 The Author(s). Published by Elsevier Ltd. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

1 ABSTRACT

1 **Real-time social media sentiment analysis for rapid impact
2 assessment of floods**

3 *Lydia Bryan-Smith*, Jake Godsall*, Franky George*, Kelly Egode*, Nina Dethlefs*, and
4 Dan Parsons**

5 ** University of Hull, Cottingham Road, Hull, HU6 7RX, United Kingdom*

6 **1. Abstract**

7 Traditional approaches to flood modelling mostly rely on hydrodynamic physical simulations. While these
8 simulations can be accurate, they are computationally expensive and prohibitively so when thinking about
9 real-time prediction based on dynamic environmental conditions.

10 Alternatively, social media platforms such as Twitter are often used by people to communicate during a
11 flooding event, but discovering which tweets hold useful information is the key challenge in extracting
12 information from posts in real time.

13 In this article, we present a novel model for flood forecasting and monitoring that makes use of a transformer
14 network that assesses the severity of a flooding situation based on sentiment analysis of the multimodal
15 inputs (text and images). We also present an experimental comparison of a range of state-of-the-art deep
16 learning methods for image processing and natural language processing. Finally, we demonstrate that
17 information induced from tweets can be used effectively to visualise fine-grained geographical flood-related
18 information dynamically and in real-time.

19 **2. Introduction**

20 Natural disasters, such as floods, can occur suddenly and without much warning, forcing people to leave
21 their homes, damaging infrastructure, destroying livelihoods, and having long-term impacts on the health of
22 those affected (FitzGerald et al., 2019; Gould et al., 2020; Khayyam and Noureen, 2020). With increasing
23 occurrence of such events due to climate change and associated phenomena, it is imperative to be able to
24 predict floods in an accurate and timely manner - to give early warnings and avert humanitarian disasters,
25 but also to direct help effectively when a flood has occurred.

26 Traditional approaches rely on hydrodynamic physical simulations. While these simulations can be accurate
27 (Price et al., 2012; Vichiantong et al., 2019), they are computationally expensive and not suited to real-time
28 prediction based on dynamic environmental conditions (Coulthard et al., 2013; Teng et al., 2017).

3 RELATED WORKS

29 Social media platforms like Twitter are often used by people to communicate during a flooding event
 30 (Kongthon et al., 2012) and other natural disasters (Sakaki et al., 2010; Riddell and Fenner, 2021). While
 31 extracting information from such social media posts in real time has the potential to increase situational
 32 awareness during flooding events, the key challenge in achieving this however is discovering which tweets
 33 hold useful information (e.g. "Part of London Road in Carlisle is closed after a building was badly damaged
 34 by #StormFranklin") and which ones do not (e.g. "Never too early for lunch") (Gao et al., 2011).

35 We present a novel model for flood forecasting and monitoring that can simultaneously process and interpret
 36 information from text and images to (a) assess the severity of a flooding situation based on sentiment
 37 analysis of the multimodal inputs, and (b) map the development of floods dynamically using geolocations of
 38 tweets, in combination with the sentiment analysis computed. We present an experimental comparison of a
 39 range of state-of-the-art deep learning methods for image processing and time-series modelling, showing
 40 that models that combine text and images achieve superior performance to unimodal models (e.g. text-only
 41 or images-only) and information induced from tweets can be used effectively to visualise fine-grained
 42 geographical flood-related information dynamically and in real-time.

43 We make the following key contributions in this article:

- 44 • A novel model that uses joint linguistic and visual feature embeddings to create a multimodal
 45 representation of sentiment in flood-related tweets on social media.
- 46 • We show that geographical and sentiment information induced from tweets can model the dynamics
 47 and severity of floods in different geographical areas.
- 48 • A set of benchmarks using state-of-the-art deep learning methodology. All code and data (where we
 49 are able to share) are publicly available.

50 **3. Related Works**

51 **3.1. Flood forecasting**

52 Traditionally, flood forecasting has been approached with physics-based models such as LISFLOOD-FP
 53 (Coulthard et al., 2013), DLEFT3D (Deltas, 2021), ANUGA (Davies and Roberts, 2015), and others (Ming
 54 et al., 2020; Roux et al., 2020; Wu et al., 2020). Taking in current environmental information, these models
 55 run a simulation to calculate a future forecast of environmental conditions. These calculations tend to be
 56 computationally expensive, taking many hours to complete, and must also be manually calibrated - which
 57 can make them prohibitively expensive for real-time use cases.

3 RELATED WORKS

58 GeoAI (Remote sensing and AI) (Janowicz et al., 2020; Li, 2020) can alleviate some of these concerns
 59 (Keung et al., 2018; Furquim et al., 2018). Multiple approaches have been applied here, for example
 60 predicting future sensor values (Le et al., 2019) with an Long Short-Term Memory (LSTM) (Hochreiter and
 61 Schmidhuber, 1997), estimating risk with a Support Vector Machine (SVM) (Cortes and Vapnik, 1995), or
 62 breaking rivers up into a graph of smaller models predicting environmental conditions in the future (Moshe
 63 et al., 2020). These suffer from a number of issues however, from being limited to specific geographical
 64 positions (Le et al., 2019; Moshe et al., 2020), relying on computationally costly simulations (Mojaddadi
 65 et al., 2017), or limited generalisability (Le et al., 2019) to new and previously unseen locations.

66 While physics-based models can predict water levels in specific areas, they aren't fast enough to forecast
 67 and monitor floods dynamically in-situ, e.g. for locations with no historical data available. Also, common to all
 68 approaches reviewed is that they don't take humanitarian needs into account, as these are difficult to infer
 69 from water depth alone.

70 **3.2. Sentiment analysis from social media**

71 In this article, we aim to explore the possibility of using sentiment analysis to assess flood severity and
 72 humanitarian needs in different locations. The idea of applying sentiment analysis to social media data is well
 73 established. Lexicon-based approaches are well explored (Baccianella et al., 2010; Mohammad and Turney,
 74 2013; Vashishtha and Susan, 2019; Rout et al., 2018; Hutto and Gilbert, 2014), but are limited to hard preset
 75 rules. Other linguistically-inspired approaches opt for a stronger grammatical representation of the input. For
 76 example, Fu et al. (2016) use rhetorical structure theory and an LSTM to parse the tweet text - preserving
 77 contextual information in a tree-like form. Some projects have attempted fine-grained classification (e.g. into
 78 6 emotional classes "happy", "sad", "anger", "fear", "surprise", and "disgust") (Purver and Battersby, 2012;
 79 Schoene and Dethlefs, 2016), but challenges remain with respect to accuracy, which is lower than other
 80 binary approaches.

81 Emojis are also often used to express emotions (Rout et al., 2018), but many existing models fail to take
 82 them into account (Sahni et al., 2017; Rout et al., 2018; Kokab et al., 2022). Scope exists to make use of
 83 them in sentiment analysis tasks. For example, Felbo et al. (2017) demonstrates a potential approach to
 84 address this by deriving positive/negative sentiment labels from emojis in tweets, enabling an LSTM-based
 85 model to be trained on a very large dataset of Twitter posts (1.2B tweets). This approach removes the
 86 need for manual and keyword-driven annotation, which reduces manual labour requirements and improves

3 RELATED WORKS

87 representation of the target domain in the training dataset as positive words (e.g. "excited" or "sad") are not
 88 being used as labels.

89 While LSTMs are powerful at modelling natural language (Felbo et al., 2017; Fu et al., 2016), they are not
 90 well suited to being parallelised on a GPU or other parallel computing device, increasing training times
 91 and underutilising equipment (Hochreiter and Schmidhuber, 1997; Vaswani et al., 2017). Transformers
 92 don't have this limitation (Vaswani et al., 2017) so scope exists to apply them to the problem of sentiment
 93 analysis (Agüero-Torales et al., 2021). Zhang et al. (2020) applies Bidirectional Encoder Representations
 94 from Transformers (BERT; a transformer-based sentence embedding model, see section 5.1) (Devlin et al.,
 95 2019), Robustly optimized BERT approach (RoBERTa) (Liu et al., 2019), and XLNet (Yang et al., 2019) to
 96 both manually, automatically (emojis), and crowdsourced (reviews with author-annotated labels) labelled data
 97 to gain performance improvements over baselines, but these models have a large number of parameters
 98 (e.g. 110M for BERT (Devlin et al., 2019)), making them memory and computationally expensive.

99 In comparison to sentiment analysis from other modalities, e.g. news or reviews, social media data faces a
 100 number of challenges. These include non-standard spellings (due to typos or abbreviations), non-standard
 101 use of words and grammar, rapidly evolving vocabulary, mixed languages and images, urls, usernames,
 102 hashtags, etc. Kokab et al. (2022) tries to solve these challenges by splitting words up with BERT as a word
 103 embedding to incorporating out-of-vocabulary words, training an LSTM model to predict positive/negative
 104 sentiment, but strips punctuation and stop words, potentially losing some semantic meaning. To address
 105 this, we propose a transformer-based binary sentiment analysis model using Global Vectors for Word
 106 Representation (GloVe) (Pennington et al., 2014) pretrained on multilingual twitter data. In doing so,
 107 multilingual text, non-standard grammar, and hashtags are included in sentiment predictions. Further,
 108 we use contrastive learning (Radford et al., 2021) to predict sentiment using both text and images at the
 109 same time.

110 3.3. Social sensing for emergencies

111 Social sensing approaches such as sentiment analysis/monitoring from social media have been applied to
 112 emergency situations previously. The problem task can be divided into two components:

- 113 1. Identifying new emergency events from social media content, and
- 114 2. Providing real-time intelligence on humanitarian needs during a known emergency event, again from
 115 information provided online.

3 RELATED WORKS

5

116 Different pieces of related research have focused on either of these problems. For example, Arthur et al.
117 (2018) manually labels a dataset of 3879 tweets to train a Naïve Bayes filter to classify tweets by relevance.

118 Relevant tweets are then geocoded using a number of heuristics, before finally performing burst detection to
119 identify flooding events and plotting geocoded tweets on a map. The small dataset limits the transferability of
120 the relevancy classifier (Li et al., 2017), however. Similarly, Smith et al. (2017) streams tweets and filters
121 them using a keyword (relevancy) and geocoding based approach before then detecting bursts of filtered
122 tweets with a simple threshold. When a burst is detected, a hydrodynamic model is run a period of 4 hours
123 with LiDAR and real-time rainfall data as inputs to predict areas which are likely flooded.

124 In this article, we will focus on part of the second problem task, specifically determining locations of
125 humanitarian needs during an emergency situation by quantifying tweets by their sentiment and location, in
126 our case, a flood.

127 An approach focusing on the identification of humanitarian needs is taken by Kankamge et al. (2020),
128 who use frequency analysis and word clustering to discover useful information from tweets during a disaster.
129 This proved effective for their target domain and data, but may not be easily generalisable to other events as
130 tweets were manually-labelled (Li et al., 2017). The difficulty in transferring concepts learnt from a disaster in
131 one place to a disaster in another is highlighted by Li et al. (2017). By using transfer learning performance
132 improvements were made in generalising a model to be effective in multiple disasters. This transfer learning
133 approach however assumes that all the tweets from the target disaster will be available up-front, which may
134 not always be the case. Additionally, lower accuracy is observed when transferring between disasters of
135 different types, and only a small dataset (7000-9000 tweets) is used.

136 Ragini et al. (2018) also classifies tweets using a dictionary: firstly, classifying them by objectivity, secondly
137 categorising by humanitarian need (e.g. water, food, medical emergency, etc), and finally sentiment analysing
138 the subjective tweets with an SVM. The tweets per category evaluated however is again small (2000 for the
139 majority class) and with unbalanced categories the performance (F1: 0.95) isn't directly comparable to other
140 studies.

141 Avvenuti et al. (2014) instead filters tweets by relevancy (i.e. "useful" and "not useful") using a decision
142 tree pretrained on a static dataset, before then performing burst detection to detect events and extracting
143 and geocoding place names to determine where the event happened. The decision tree used though is
144 trained on 1412 manually labelled tweets, which as Li et al. (2017) suggests limits the generalisability of the
145 approach.

¹⁴⁶ Alternative approaches used include clustering and visualisation (Beigi et al., 2016), burst detection (Yin
¹⁴⁷ et al., 2012) analysing images (Ning et al., 2020; de Vitry et al., 2019), hand-labelling training datasets
¹⁴⁸ (Avvenuti et al., 2014), or use keyword-based analysis (Arthur et al., 2018; Ragini et al., 2018; Smith et al.,
¹⁴⁹ 2017), which does not generalise easily to new and sudden events (Li et al., 2017), but despite this wide
¹⁵⁰ range of approaches being taken to the issue, limited attempts to combine text and images have been taken
¹⁵¹ (Wang et al., 2018; Said et al., 2020).

¹⁵² We suggest that scope exists to apply modern machine learning algorithms such as transformers (Vaswani
¹⁵³ et al., 2017) and Contrastive Language-Image Pretraining (CLIP) (Radford et al., 2021) to the problem of
¹⁵⁴ sentiment analysis in flooding situations, in order to create a more generalisable approach - both in terms of
¹⁵⁵ analysing new events, and in terms of handling new situations such as different places or languages.

¹⁵⁶ 4. Data

¹⁵⁷ To collect data for analysis, we identified major flood and extreme weather events in the last ten years and
¹⁵⁸ collected the historical tweets that were related to them based on hashtags and / or keyword search
¹⁵⁹ using Twitter's Academic API. We used the following search terms to source tweets via the Twitter
¹⁶⁰ API: #StormDennis, #StormChristoph, #StormJorge, #HurricaneEta OR #HurricaneEta, #HurricaneBeta
¹⁶¹ OR #HurricaneBeta, #Hurricanelota OR #Hurricanelota, #NSWFloods, #qldfloods, #ColoradoFloods OR
¹⁶² #ColoradoFlooding, #MexicoFloods, flood #snaith, #floods OR #flashfloods, #StormFranklin, finsbury park
¹⁶³ flood, himachalpradesh flash floods, texas floods, #SydneyFloods. The tweets we collected span from
¹⁶⁴ 2007-12-06 to 2020-02-25.

¹⁶⁵ Table 1 details the floods from which data was downloaded. The start and end date columns refer to the
¹⁶⁶ time frame of tweets appearing rather than to the events themselves. To preprocess the data, we excluded
¹⁶⁷ retweets and deduplicated tweets by their IDs (though some duplication is possible if multiple users copy the
¹⁶⁸ same text and post it independently). Replies to tweets matching the search criteria were included.

¹⁶⁹ For all tweets collected, we kept information on: the search term itself, the tweet ID, user name and
¹⁷⁰ description, user location, the number of followers and tweets by a user, the number of retweets and likes
¹⁷¹ of a tweet, any media (e.g. images), and the geolocation of the tweet, if provided. While we downloaded
¹⁷² 13,851 tweets from the hashtag #StormFranklin (including replies), we excluded these from the training
¹⁷³ and validation datasets for later evaluation. This gave us a dataset of 795,065 tweets in total, including the
¹⁷⁴ StormFranklin tweets.

	START DATE	END DATE	REGION	TWEETS
Finsbury park flood	2019-10-08	2021-08-07	London, UK	267
Storm Dennis	2019-12-08	2021-07-22	UK	120 861
Storm Christoph	2021-01-17	2021-07-08	UK	17 085
Storm Jorge	2020-02-27	2021-07-10	UK	25 102
Hurricane Eta	2020-10-30	2021-07-16	Central America	17 265
Hurricane Beta	2020-09-18	2020-11-12	Central America	315
Hurricane Iota	2020-11-09	2022-07-22	Central America	315
New South Wales Floods	2021-02-04	2021-07-23	Australia	15 759
Queensland Floods	2018-12-18	2021-06-28	Australia	1889
Colorado Flooding	2013-08-10	2021-07-23	USA	1420
Mexico Floods	2020-11-07	2021-07-22	Mexico	26
Snaith Floods	2020-02-25	2021-07-13	Snaith, UK	322
Storm Franklin	2022-02-16	2022-03-04	UK	13 851
Himachalpradesh	2021-07-06	2021-07-31	Himachal Pradesh, North India	551
Texas Floods	2021-05-01	2021-05-31	Texas, USA	892
Sydney Floods	2022-02-23	2022-03-21	Sydney, Australia	2553
(floods OR flashfloods)	2007-12-06	2021-07-23	Worldwide	553 218

Table 1

Overview of the floods included in our dataset.

175 Figure 1 shows some examples of tweets with images. Tweets 1a and 1b could be classed as positive ("be
 176 safe", "easy way to do it without getting wet!"), with nobody in immediate distress. Meanwhile, tweet 1c
 177 could be classed as negative ("frustrated and upset"), potentially highlighting an issue that requires human
 178 attention. Below some examples of tweets without images are shown:

179 " @sZL7YcOsTnZhhsf8xFXcA @JF1MQxDGoRT7NsObnSxQpA And it's still bucketing down in
 180 Coffs, landslips at the Big Banana and Thora, Waterfall Way. These hills are supposed to keep
 181 us high and dry but nooooo" –19th March 2021

182 " @fXOQqWAYZqJlpv_r-qly2A @sZL7YcOsTnZhhsf8xFXcA All good here Cows and calves on
 183 the hill paddock All other animals + humans safe" –19th March 2021

184 " @jd2a0Qryu0aGHcz9bj4B2A Extraordinary that Warragamba is full. It's a massive dam. Hope it
 185 goes ok for those downstream" –20th March 2021

186 Tweet 2 here could be classed as positive, whereas tweet 1 would be negative as it indicates someone's
 187 house has been flooded. Tweet 3 could be classed as negative, as it contains information that could have a
 188 significant negative effect on those downstream.

189 Tweets were anonymised using the SHAKE128 hash function (Dworkin, 2015) with a salt. We hashed all
 190 tweet ids, usernames, and conversation ids. We kept geotags and place names extracted by twitter, along
 191 with direct links to associated media.

4 DATA

8

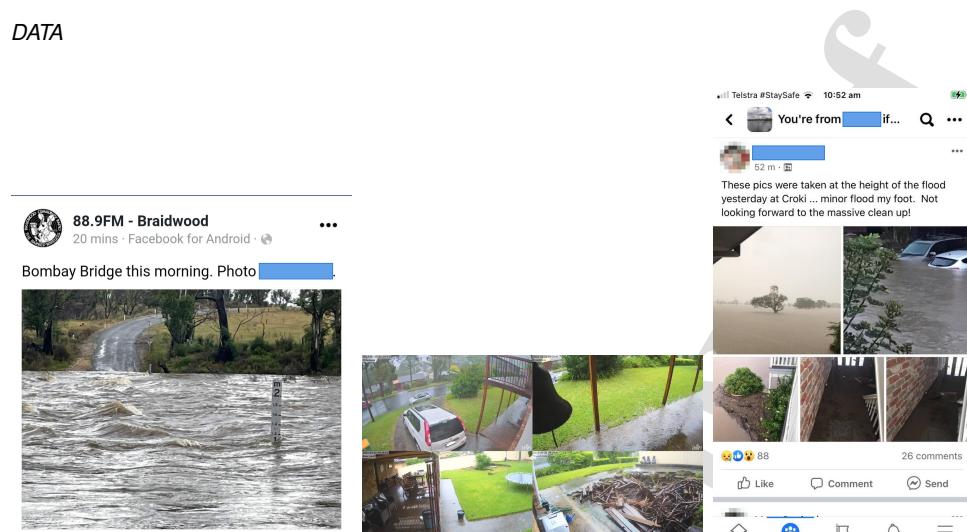


Figure 1: Example tweets with images

A horizontal row of various positive and colorful emojis, including hearts, smiley faces, flowers, and musical notes.

A horizontal row of various negative and neutral emojis, including sad faces, a skull, a thumbs-down, and a question mark.

Figure 2: The emojis and the categories they were manually assigned to.

192 4.1. Data Labelling

193 The learning models we developed and trained require labels for the training process. To this end we were
194 inspired by Felbo et al. (2017) and extracted a list of all unique emojis from the dataset, which we then
195 manually labelled to be either positive or negative - see Figure 2. By positive and negative, we define positive
196 to be tweets of no concern, and negative to be tweets potentially requiring attention (i.e. where someone
197 may require assistance).

5 APPROACH

198 Then, we extracted all tweets from the data that had at least one (positive or negative) emoji, and automatically
 199 labelled them as either positive or negative based on our categories. Whenever more than one emoji was
 200 present, the majority category was used. If an equal number of positive and negative emojis were present,
 201 the 'positive' category took precedence. Finally, we split the data into two parts, with 80% for training and
 202 20% for validation during the learning process.

203 To evaluate the accuracy of our sentiment labels against a human gold standard, we used Amazon
 204 Mechanical Turk (AMT) to collect human sentiment labels on a representative data sample of 1938 tweets
 205 randomly selected from the #NSWFloods hashtag from 40 different raters. #NSWFloods was chosen as it is
 206 time-limited and has a significant sample size (~15K). Turkers were presented with the tweet text (excluding
 207 images and emojis) and asked to assign a categorical rating from a 1-5 Likert scale, where 1=negative,
 208 2=slightly negative, 3=neutral, 4=slightly positive and 5=positive. We collected fine-grained ratings to allow
 209 a better comparison with models such as Valence Aware Dictionary and sEntiment Reasoner (VADER)
 210 (Hutto and Gilbert, 2014) or RoBERTa who predict such distinctions. Accuracy was low for finer sentiment
 211 distinctions, so 1-2 were collapsed as "negative" and 3-5 were considered "positive". These tweets form the
 212 basis of our experiments. Table 4 provides a comparison of our learning models against the human gold
 213 standard annotations.

214 5. Approach

215 This section will introduce the data representation and learning models for our experiments.

216 5.1. Representation of Inputs

217 In our experiments, we compared the effectiveness of different AI model architectures in predicting the
 218 sentiment of social media posts from Twitter. Specifically, we compare the following architectures:

- 219 • **Pretrained baselines:** VADER, RoBERTa.
- 220 • **Models we trained:** Transformer, LSTM, CLIP, ResNet50 (He et al., 2016).
- 221 When training models on natural language data, a problem is the size of the input data. e.g. if a dataset
 222 contains 15K unique words, using one-hot encoding a vector in the form $[word_0, word_1, word_2, \dots, word_{15\ 000}]$
 223 is required to represent each word, where $word_n \in \{0, 1\}$. This has significant implications on both memory
 224 usage and generalisability - as one-hot encoding does not capture any lexical, semantic, syntactical meaning,

5 APPROACH

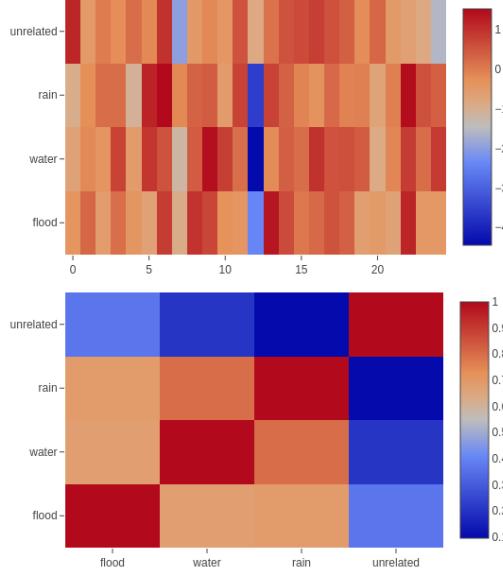


Figure 3: Top: GloVe (trained on 2 billion tweets with vocabulary of 1.2 million words) word embedding vectors with a size (dimension) of 25 visualised as an array of numbers in a heatmap for a selection of words. Bottom: Using cosine difference, two vectors can be compared. `flood` and `rain` have a high similarity, whereas `flood` and `unrelated` have a low similarity. Such semantic comparisons are not possible e.g. using a one-hot encoding model - this enables the AI model that follows the embedding layer to focus on the domain-specific task, rather than having to learn not only the domain-specific task, but also relationships between words. (Pennington et al., 2014)

225 or relationships - so alternate strategies are needed, such as Word2Vec (Mikolov et al., 2013) or GloVe
226 (Pennington et al., 2014).

227 An embedding layer is often the first layer in an AI model, and it encodes the input data - which can in this
228 case be defined as a string of text tokens from a tweet split up into its constituent words - into an array
229 of numerical vectors of a fixed size. Such an embedding layer can be defined by a dictionary in the form

230 $f : \text{word}_{\text{string}}^{1..n} \rightarrow \text{vector}_{\text{float32}}^{\text{dim}}$, which is then executed for every element of the input array of strings, where

231 dim is the size of the resulting vector. Hence, an input to an embedding layer can be defined as $\text{input}_{\text{string}[]}^{\text{words}}$

232 (where words is the number of words in the array of strings), and the output defined as $\text{vector}_{\text{float32}}^{\text{words.dim}}$.

233 Therefore, the memory required to encode the words from the original input is significantly reduced, while
234 also ensuring that semantic meaning is retained by encoding words with a similar semantic meaning to
235 similar numerical values (Koehrsen, 2018) (Figure 3). These embeddings can be trained from the domain
236 corpus, e.g. our tweets, but are typically trained separately from large general-purpose domains. They can

5 APPROACH

237 be seen as a generalised representation of the English language, for example, comprised of different types
 238 of text, domains, and genres, e.g. GloVe (Pennington et al., 2014).

239 **5.2. Learning models**

240 The body of a neural network model follows the embedding layer in the form of a set of hidden layers.

241 Consider an input vector in of inputs that we want to map to a sentiment value out . To do this, we learn
 242 a hidden representation $h(in)$ using a function $f(h, in)$, minimising the loss (error) between a given label
 243 out^{truth} and predicted value $out^{predict}$ - e.g. using cross-entropy loss. We end up with a function that predicts
 244 a sentiment value as $out^{predict} = h(in)$

245 The first model we will train is the transformer network. A full transformer can be applied to e.g. machine
 246 translation tasks by translating one sequence of vectorised inputs into another. Transformers are made up of
 247 two parts: an encoder, which encodes features deemed important by the model into a sequence of vectors
 248 with a lower dimensionality (i.e. a feature map), and a decoder, which converts the feature map into the
 249 desired output.

250 The models used in this paper can be classed more specifically as deep feed forward networks with back
 251 propagation (Goodfellow et al., 2016a) - the process by which these models are trained using pairs of
 252 samples and associated labels. At each step of the training process:

- 253 1. The input sample is put through the model forwards through the directed graph of layers to make a
 254 prediction (feed-forward).
- 255 2. The prediction is compared to the ground truth label, and an error value is calculated using a loss
 256 function, for example mean squared error (i.e. $loss = (actual - predicted)^2$) and cross-entropy loss (e.g.
 257 $loss = -(actual \times \log(predicted) + (1 - actual) \times \log(1 - predicted)))$) (Kaller, 2019; Brownlee, 2019).
- 258 3. Finally, the loss is propagated backwards through the model using an algorithm like gradient descent
 259 (Goodfellow et al., 2016b).

260 We use just the encoder part of a transformer to encode the vectorised input social media post into a feature
 261 map which can be interpreted by later layers of the model. Transformers handle sequences in parallel - as in
 262 $out_i = f(in_i)$. This is achieved by adding a positional embedding signal (explained below) and then dropping
 263 them through layer normalisations and dense layers, and a self-attention layer, which enables the model to
 264 identify which parts of the input are important for making the output prediction.

5 APPROACH

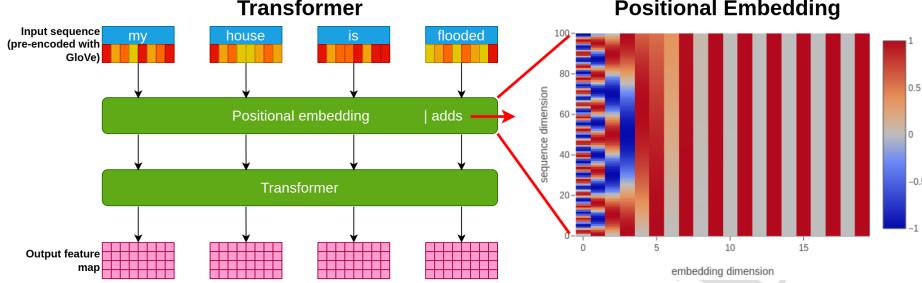


Figure 4: The architecture of a transformer. The input sequence (that has already been encoded with GloVe as described in section 5.1) gets a positional embedding formula added to it (visualised on the right; words in the sequence are along the vertical axis, and the embedding of those words are along the horizontal axis), which gives the information about the ordering of items in the sequence. Finally, the transformer itself processes it in parallel.

265 Alternatively, LSTMs (Hochreiter and Schmidhuber, 1997) or GRUs (Cho et al., 2014) can handle sequenced
 266 data. These are called recurrent models as they have a hidden state, and a system of gates are used to
 267 update the values therein for each element in the sequence, requiring that each element in a sequence is
 268 processed serially - e.g. $out_i = f(out_{i-1}, in_i)$.

269 This is why the transformer architecture enables greater computational parallelisation when training using a
 270 GPU (Vaswani et al., 2017), since without a recurrent element each sequence element can be processed
 271 independently, as in $out_i = f(in_i)$.

272 Transformers instead rely on dense (fully connected) layers. While this makes them more parallelisable
 273 and hence can train more quickly, this also means that they are unable to account for context and relative
 274 positioning in input sequences. To alleviate this weakness, a positional embedding code is added to the input
 275 sequence. If the input sequence (post-embedding layer) is $input^{seq,dim}$ (where seq is the sequence length and
 276 dim is the vector size for each element therein), then the positional embedding can be defined as (Vaswani
 277 et al., 2017):

$$PE_{(i_{seq}, 2i_{dim})} = \sin\left(\frac{i_{seq}}{10000^{2i_{dim}/dim}}\right)$$

$$PE_{(i_{seq}, 2i_{dim}+1)} = \cos\left(\frac{i_{seq}}{10000^{2i_{dim}/dim}}\right),$$

278 ...where i_{seq} is the position in the sequence dimension, i_{dim} is the position in the embedding dimension, and
 279 PE is the positional embedding for a single value in the $input^{seq,dim}$. The embedding dimension alternates
 280 between \sin and \cos for each successive element, as Figure 4 shows.

CLIP tweet text-image pair similarities

Floodwaters have now subsided on the road from Inverell to our Dumaresq substation.\n\nThe site is accessible but remains closed today. Some minor road repairs are required on the access road into the site from the Bruxner highway.

The waters are rising but the sun is out and the skies are blue on the other side of the bridge

Bit of an update on #Moree flooding. 3pm Wed: Gwydir looking to Pally; back to town; town weir to south; looking north with main bridge at bottom of photo. A lot of water coming into Mehi from hills so hard to know what actual height will be

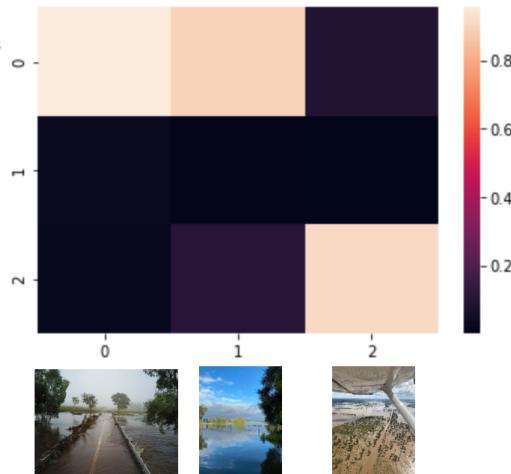


Figure 5: Tweet text-image pairs and their similarities, as calculated by CLIP (higher values mean more similar). Although it is correct most of the time, as shown here sometimes the wording of the tweet confuses CLIP. Where multiple images associated with a tweet, only the first one was chosen.

281 **5.3. Representing text and images jointly**

282 With many users posting many tweets independently, many different topics are discussed in an unstructured
283 manner. In some cases, images associated with tweets contain additional information as illustrated in
284 Figure 1. To explore the effect of images associated with tweets on sentiment analysis tasks, another model
285 architecture we used in our comparison is the Contrastive Language-Image Pretraining (CLIP) architecture.
286 Similarly to the transformer already discussed, CLIP uses an embedding layer as its first layer. However
287 unlike the plain transformer model, CLIP handles not only textual input but images as well. CLIP first has
288 separate encoding layers for textual and image inputs, before later combining them together and training
289 the two encoders to predict which textual and which image inputs were paired with each other. In doing so,
290 CLIP trains to predict how well a textual string and an associated image pair together (Radford et al., 2021)
291 (Figure 5).
292 Figure 6 outlines how CLIP trains and makes predictions by taking a contrastive learning approach. Batches
293 of text-image pairs are compared using cosine similarity, which is used to calculate cross-entropy loss (see
294 Figure 7).

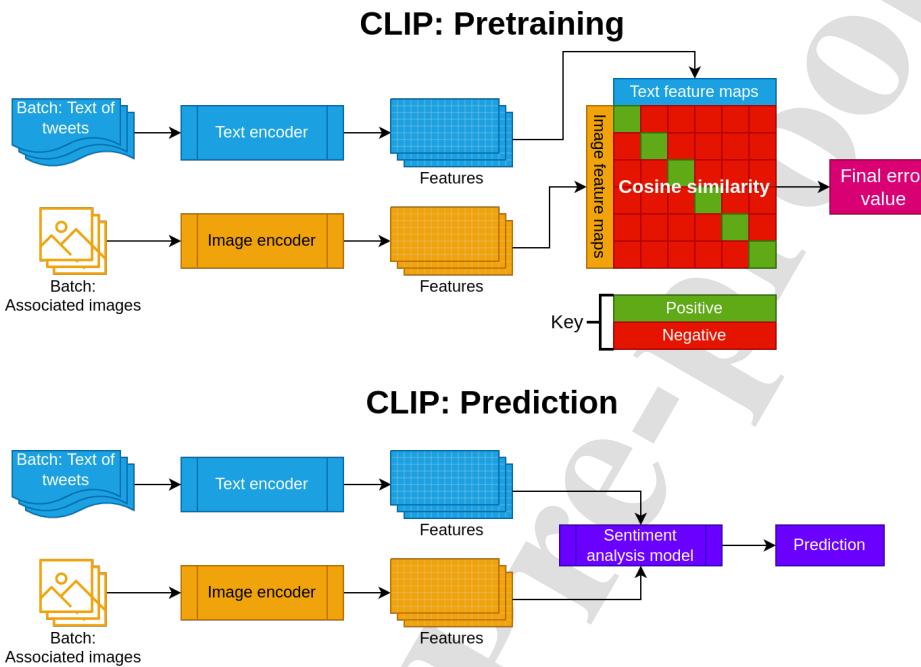


Figure 6: A summary of how the CLIP model operates. A batch of text and image pairs are run through different encoders separately to produce a pair of feature maps. Then, the resulting feature maps are compared using cosine similarity and fed into the loss function which trains the model to learn pairs to be similar to one another (Radford et al., 2021). When making a prediction, the trained encoders can be used to encode new text-image pairs that are fed into a domain-specific model.

295 **6. Experiments**

296 This section presents the experimental setup we adopted for our experiments, presents and discusses
297 results as well as some sample predictions.

298 **6.1. Experimental Setup**

299 We compare a set of different methods for predicting sentiments from tweets given our emoji-based labels
300 from Section 4. Specifically we compare:

- 301 1. **LSTM:** 2 bidirectional layers, 128 units each, batch normalisation.
- 302 2. **Transformer:** 1 transformer encoder, 16 attention heads, 32 units, dropout 0.1, gelu.
- 303 3. **CLIP:** Pretrained, ViT-B/32, followed by 2 x 512 unit dense layers, dropout 0.1.
- 304 4. **ResNet50:** ResNet50 architecture, followed by a softmax dense layer.

Inputs:
$image_embedded^{batch_size, dim_embed}$
$text_embedded^{batch_size, dim_embed}$
$batch_size$
t: Learned temperature parameter
Algorithm:
$logits \leftarrow np.dot(image_embedded, transpose(text_embedded)) \times np.exp(t)$
$labels \leftarrow np.arange(batch_size)$
$loss_image \leftarrow cross_entropy_loss(logits, labels, axis = 0)$
$loss_text \leftarrow cross_entropy_loss(logits, labels, axis = 1)$
$loss = \frac{(loss_image + loss_text)}{2}$

Figure 7: An outline of CLIP’s contrastive learning loss algorithm. Adapted from Radford et al. (2021).

Model	Parameters	Validation Accuracy
Pretrained RoBERTa	345M	n/a
LSTM	731K	81.81%
Transformer	2.6M	80.42%
CLIP (not augmented)	152.1M	89.41%
CLIP (augmented)	152.1M	86.37%
ResNet50	23.6M	73.79%

Table 2

An overview of the models used.

- 305 The hyperparameters of these models were chosen after experimentation with different combinations. In all
 306 cases where we trained a model, we used the Adam optimiser.
 307 We also chose two baseline models for comparison that were shown to perform well on the task of sentiment
 308 analysis in previous work: (1) VADER (Hutto and Gilbert, 2014), a rule-based model developed to predict
 309 the sentiment of social media posts, and (2) RoBERTa (Liu et al., 2019), a generic pre-trained transformer
 310 based model.
 311 We trained our models on various Nvidia GPUs: GeForce 3060, Tesla K40m, Tesla P100, and Nvidia A40,
 312 depending on availability and machine learning library requirements.
 313 Accuracies are reported from models with the same architecture. All models were trained for a total of
 314 50 epochs, and then the checkpoint from the epoch with the highest validation accuracy was chosen. All
 315 models (except CLIP, which has its own inbuilt word embeddings) also used GloVe pretrained on Twitter
 316 data with a dimension of 200 for word embeddings as it is more computationally efficient, although other
 317 word embeddings do exist (Liu et al., 2019; Lewis et al., 2020; Devlin et al., 2019). We test the potential of
 318 this technique against our RoBERTa and VADER baselines.
 319 The transformer encoder (Vaswani et al., 2017) model we trained predicts the positive / negative sentiment of
 320 the tweet text, using the emojis as labels as described in section 4. Although emoji-based labels (see section

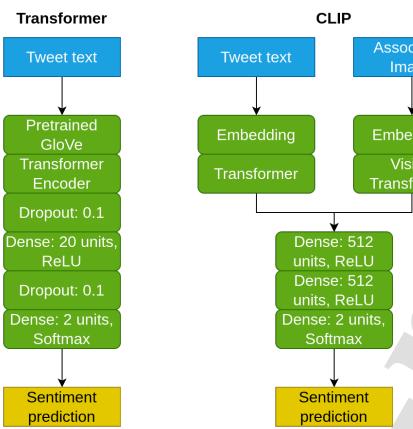


Figure 8: The architectures of our transformer encoder (left) and CLIP (right) models that we trained. The CLIP model concatenates the feature maps from both the image and the text encoders.

321 4.1) are required during the training process, emojis are not required during inference. Figure 8 shows the
 322 architecture of the transformer encoder model we trained, as well as that of the CLIP-based model.

323 Like the transformer, labels for training the CLIP model came from emojis split into positive / negative
 324 categories. The CLIP model takes both text and images as an input at the same time before then producing
 325 a prediction based on both inputs. Figure 8 shows the architecture of the CLIP model trained. We discovered
 326 that out of the 180K tweets that had an associated image, only 14K tweets also contained an emoji (i.e. an
 327 output label according to our setup).

328 To augment the dataset, we used the CLIP model trained on 14K image-text pairs to annotated each tweet
 329 that had an emoji but no image with a newly associated image that fits the text. Figure 9 shows the algorithm
 330 that we used to augment the data. This augmentation process raised the size of the training dataset to 55K
 331 text-image pairs, and improved the F1 score of the model from **76.3%** to **80.2%** (see table 3). Tweets without
 332 associated images achieved 0.734 F1 (CLIP-augmented) / 0.766 (CLIP-not augmented), and tweets with
 333 images 0.8 F1 (CLIP-augmented) / 0.767 (CLIP-not augmented).

334 6.2. Results

335 *Comparison against emoji labels* To compare the performance of the models, we used the emoji labels, as
 336 we used previously to train the CLIP and transformer models, as a ground truth. Table 3 shows the results of

```

for(tweet of dataset) {
    rankings = [];
    for(image of images) {
        rankings.push(clip.rank_image(image));
    }

    high_rankings = rankings.filter(ranking => ranking > 0.75)
    if(high_rankings.length > 0)
        return random_item(high_rankings);
    else
        return highest_ranking(rankings);
}

```

Figure 9: The algorithm by which we ranked image associations with tweet texts when augmenting the tweets with CLIP.

Model	F1	Recall	Precision	Accuracy	Samples	Truth Pos	Predict Pos	Truth Neg	Predict Neg
VADER	0.339	0.596	0.506	0.471	15485	7183	15224	8302	261
RoBERTa	0.651	0.658	0.651	0.658	15486	7183	6122	8303	9364
Transformer	0.687	0.694	0.693	0.687	15486	7183	8674	8303	6812
LSTM	0.676	0.72	0.695	0.682	15486	7183	10381	8303	5105
CLIP (augmented)	0.802	0.791	0.817	0.841	2491	641	738	1850	1753
CLIP (not augmented)	0.763	0.752	0.814	0.792	2491	641	978	1850	1513
ResNet	0.697	0.691	0.742	0.734	2315	585	918	1730	1397

Table 3

Comparison of the sentiment analysis performance against emojis as a ground-truth label.

Model	F1	Recall	Precision	Accuracy	Samples	Truth Pos	Predict Pos	Truth Neg	Predict Neg
VADER	0.499	0.501	0.501	0.566	1914	1358	1258	556	656
RoBERTa	0.42	0.521	0.521	0.421	1914	1358	517	556	1397
Transformer	0.591	0.595	0.613	0.623	1914	1358	1095	556	819
LSTM	0.589	0.587	0.595	0.645	1914	1358	1263	556	651
CLIP (augmented)	0.512	0.604	0.602	0.512	320	220	96	100	224
CLIP (not augmented)	0.58	0.607	0.624	0.588	320	220	144	100	176
ResNet	0.54	0.559	0.576	0.577	156	116	84	40	72

Table 4

Comparison of the sentiment analysis models we tried against human-labelled tweets.

337 this experiment. We used the 15K tweets from the #NSWFloods hashtag for this emoji-based comparison
 338 and any predictions of neutral were considered positive predictions instead, as detailed in section 4.

339 *Comparison against human ratings* To further explore the comparative performances of these models, we
 340 also used our random human-labelled subset of 1,938 tweets from the #NSWFloods hashtag and analyse
 341 the performance of our models against them - table 4 shows the results of this experiment. As in table 3,
 342 predictions with a class of neutral were considered positive. These results show that - unlike with
 343 the emoji labels in table 3 - our transformer model is the best performing model, with the CLIP and LSTM
 344 models coming in second place. We speculate this could be because emojis and images were not used in
 345 the human labelling process or due to the small sample size.

6 EXPERIMENTS

18

³⁴⁶ *Images-only results* Our ResNet50 model took images associated with the tweets as an input and classified
³⁴⁷ them as positive or negative, using labels predicted by the LSTM model from the associated text as a
³⁴⁸ ground-truth. As this model analyses images rather than text, only a subset of the tweets in tables 3 and 4
³⁴⁹ could be used for the ResNet50 model.

³⁵⁰ **6.3. Discussion**

³⁵¹ Of all the models we tried, under F1 score our CLIP model appears to perform best, followed by our ResNet
³⁵² model when compared against emojis as a ground truth in table 3, while our transformer model comes
³⁵³ in a close third. Linguistic clues can sometimes be quite subtle (e.g. sarcasm), often requiring the reader
³⁵⁴ to infer the correct meaning. Deep learning models find this challenging, so images can help fill this gap.
³⁵⁵ Tweets T1 and T2 in table 5 are examples of this, where the image provides critical context that is otherwise
³⁵⁶ misinterpreted by other models.

³⁵⁷ Our baseline model VADER appears to perform worst. We suggest that this may be because emojis play a
³⁵⁸ key role in identifying the sentiment of tweets and VADER does not adequately take the context in which
³⁵⁹ an emoji is used into account, and its rule based approach is inflexible when compared to models that
³⁶⁰ are trained in a supervised or semi-supervised manner. This is illustrated by tweets T3 and T4 in table 5 -
³⁶¹ which although it is a negative tweet, it is still considered neutral by VADER. This is also backed up by the
³⁶² human-labelled tweets, and an accuracy of 42% (human-labelled) / 33.9% (emoji-labelled).

³⁶³ When compared using human-labelled tweets as a ground truth instead, the story is very different. All scores
³⁶⁴ are generally lower, indicating that the smaller size of the human-labelled dataset may not be completely
³⁶⁵ representative of the entire dataset. Despite this, our Transformer performs best and RoBERTa performs
³⁶⁶ worst, suggesting that the MultiNLI dataset that the RoBERTa model was trained on (Lewis et al., 2020) is
³⁶⁷ not representative of the target domain of social media here.

³⁶⁸ Humans are better at inferring meaning in language than AI, but with visual information not being present for
³⁶⁹ our human raters this makes the task more challenging. Given the human-labelled dataset was small (~2K
³⁷⁰ tweets) and human-labelling tweets is both time consuming and expensive, labelling tweets automatically via
³⁷¹ emojis is significantly more practical.

6 EXPERIMENTS

19

Text	Emoji	Human	Transformer	LSTM	VADER	RoBERTa	ResNet	CLIP
T1  Hastings river port Macquarie #NSWFloods	n/a	neutral	positive	negative	neutral	negative	negative	negative
T2 😊 The local Facebook page is “delivering” today	positive	n/a	negative	negative	neutral	positive	positive	positive
T3 Droughts.. Fires.. Floods.. #Australia #NSWFloods #SydneyFloods Oh and a bit of #COVID19 Wasn’t 2021 meant to be a better year?	n/a	negative	negative	positive	neutral	positive	n/a	negative
T4 Before and after pics of Wauchope railway bridge 😢 #NSWFloods credits to #<name redacted>		negative	negative	negative	negative	neutral	negative	negative

Table 5

Some sample tweets from the #NSWFloods dataset labelled by the various models we tested. The Column CLIP refers to the augmented model.

372 **6.3.1. Sentiment by image**

373 Since images associated with tweets can contain some useful information (Said et al., 2020), we used CLIP
 374 to explore utilising both text and images to predict sentiment in order to understand how images relate to the
 375 overall sentiment of a tweet.

376 With 149.5M more parameters than our Transformer model (approx. 788K of which are in dense layers we
 377 trained after the pretrained CLIP model), it is also significantly more computationally expensive and may
 378 have overfit. The tweet augmentation process is especially computationally expensive, requiring each tweet
 379 to be ranked against every image in the dataset.

380 When compared to emojis as a ground-truth label, our CLIP model easily beats all the other models
 381 that consume only textual data by a significant margin of at least 11%. However, when we compare it
 382 to human-labelled tweets, it does not outperform the text-only transformer even though CLIP also had
 383 associated images as an input.

384 This illustrates that images associated with tweets contain contextual information that was lost to human
 385 raters. Given the small sample size mentioned earlier, this further shows that it is more practical to use
 386 emojis as labels, and to include images for additional visual context.

6 EXPERIMENTS

20

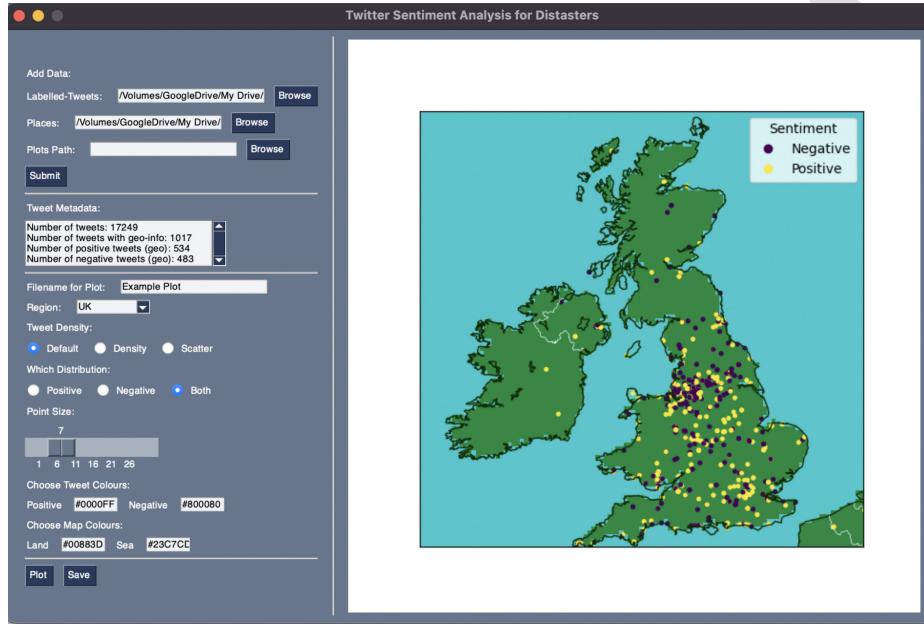


Figure 10: Graphical user interface for the geo-spatial plotting system. Tweets can be optionally grouped by date, allowing for temporal analysis of geospatial trends.

387 To further explore the relationship between images and sentiment, we can look to our image-only ResNet
 388 model, which appears to be the highest performing model in table 3 with respect to both F1 score and
 389 accuracy. This may be due to a small sample size used in the comparison as the ResNet only analyses
 390 tweets containing at least one image (2315 samples vs 15,486 samples for the transformer), and the
 391 relatively unbalanced dataset as compared to the text-only models - suggesting that when people tweet
 392 image(s), these images are more likely to be considered to have a negative sentiment.

393 *6.3.2. Geospatial analysis of tweets*

394 To further improve upon the explainability of the flood sentiment analysis, geo-spatial analysis was performed
 395 using tweets sentiment-analysed by our Transformer model. Social media systems track the locations
 396 of users, thereby making it possible to determine where a person was when they posted on the site or
 397 application.

398 We developed a program to create geo-spatial visualisations from geotagged tweets to explore large
 399 datasets.

6 EXPERIMENTS

400 There are two types of geographical metadata available for tweets from the Twitter API: the exact location of
 401 the origin of the tweet (from GPS-enabled or GeoIP-enabled devices), and the Place object type extracted
 402 from the tweet text content using named-entity-recognition which provides a polygon which bounds the area
 403 from which the user posted the tweet.

404 We used the exact coordinates where possible, and took the central point of the polygon bounding-box
 405 for tweets without precise location data. To account for the overlap of points representing tweets from the
 406 same Place region, two distinct methods were employed: increasing the size of the points as a function
 407 of the number of tweets contained within, and adding Gaussian noise to points superimposed on a single
 408 coordinate.

409 We used tweets from the #StormFranklin hashtag (see also table 1) for our geospatial analysis. This extreme
 410 weather event was chosen as a large number of tweets were available and it was not included in any data
 411 downloaded previously (Hurricane Iota was downloaded separately). Our Storm Franklin dataset has 921
 412 geotagged tweets (7%), and only 28 include precise location information. Sloan and Morgan (2015) showed
 413 that approximately 0.85% of tweets are geotagged, but Twitter made changes in 2019 to reduce the ways by
 414 which precise location geotagging can be achieved, explaining the proportions found (Hu and Wang, 2020).
 415 29% of the tweets were classified as having positive sentiment.

416 The static plot shows that the vast majority of negative tweets are found in high-density clusters, significantly
 417 more so than the positively labelled tweets. This is most obvious for tweets originating in Ireland.

418 Figure 11 shows the temporal-level distribution of tweets for Storm Franklin. Tweet sentiments are visualised
 419 from 2022-02-18 to 2022-02-25, and a time-series of tweet frequency - the majority of tweets were posted
 420 over a two day period: 2022-02-19 to 2022-02-20.

421 Days with most activity correspond with days with the highest number of tweets posted, and is also slightly
 422 before the actual event, which took place on 2022-02-20 to 2022-02-21 (Deltaires, 2022b). This could be
 423 because the Met Office's early warning prompted conversation on Twitter or because of lingering effects
 424 from Storm Eunice, which took place a few days prior (Deltaires, 2022a).

425 By analysing sentiment-analysed tweets geospatially, real-time information on the status of flooding events
 426 can be obtained. By analysing the effectiveness several sentiment analysis models in section 6, we improve
 427 our geospatial analysis of sentiment-analysed tweets. Similarly, by mapping tweets we gain insights into the
 428 effectiveness of our sentiment analysis model.

7 CONCLUSION

22

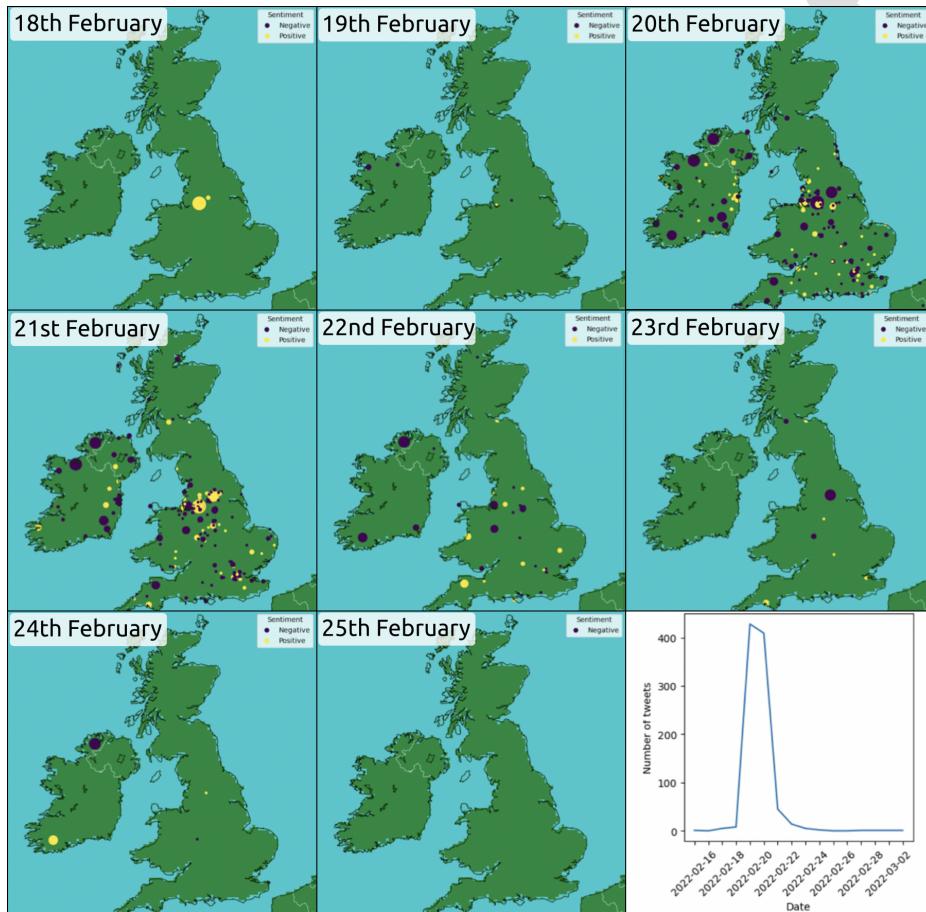


Figure 11: Daily plots of geographically located sentiment-labelled tweets for the #StormFranklin dataset on 18th - 25th February 2022

429 **7. Conclusion**

- 430 We explored using text and images from Twitter for sentiment analysis, comparing VADER, RoBERTa, a
 431 Transformer encoder, and CLIP. Our CLIP-based model takes both text and images into account, which by
 432 taking visual information into account can help close the gap between AI and human raters. We highlighted
 433 the importance of emojis in understanding the sentiment of tweets: both our CLIP and transformer models
 434 outperformed RoBERTa and VADER, neither of which consider emojis.
- 435 Finally, we demonstrate the potential and feasibility of geospatial analysis of sentiment-analysed tweets in a
 436 flooding situation. The scale, frequency and low-latency allow for rapid analysis of disasters in real-time,

7 CONCLUSION

23

Text	Label
P1 Before retiring to bed this evening, folks, don't forget to make sure your beer's properly tied down #StormEunice #StormFranklin #Beer #BeerBods	positive
P2 #StormFranklin on Sunday so don't give the neighbours back fences.wheelie bins and trampolines from #stormeunice just yet"	positive
P3 @ciAD_xsnxg-KiLUkM8tK6g @oB2hntlJShfk6tmZZ7z8Q The next storm will begin with a F. (or is that an F??) They're calling it #StormFranklin I'm looking forward to #StormInATeacup 🤪🤩🤩	positive
P4 Ik ben benieuwd wat Eunice ons gaat brengen! Stay strong! #eunice #storm #StormEunice #StormFranklin	positive
N1 Woohoo! #StormFranklin is due after #StormEunice takes her hook 🤪	negative
N2 How much do these extreme storms make you worry about climate change?#StormEunice #StormDudley #StormFranklin	negative
N3 The storm coming in Sun/Mon #StormFranklin could even more concerning than Eunice. Keep your eyes peeled. Concern is it will arrive only a day or so after Eunice. #windy #storm #StormEunice	negative

Table 6

A sample of tweets from 2022-02-18 for the hashtag #StormFranklin, from before the storm actually hit. In total, there are 11 positive and 24 negative tweets. Tweets classed as positive generally make light of the storm (P1, P2, P3, P4), while negative tweets are either sarcastic (N1) or worrying about causes and effects (N2, N3). Additionally, 5 of the positive tweets are written in the Dutch language (Google Translate used for identification), but were still classified correctly (P4).

437 enhancing situational awareness and allowing a human-lead approach to identification of affected areas in
 438 real time.

- 439 • An assessment of how our approach could be used to practically support disaster decision making
 440 in-situ is required.
- 441 • An alternative approach such as image segmentation (Pally and Samadi, 2022) or captioning is likely
 442 needed to make better use of images in the context of flooding events.
- 443 • Scope exists for future research to investigate models that consider emojis as well as regular text when
 444 predicting sentiment at inference.
- 445 • Correcting for population density, mass evacuations of people, and potential communication disruption
 446 are all challenges. Making use of other data modalities such as satellite data, mobile phone cell tower
 447 information, and traffic data may help here.
- 448 • Screening tweets from irrelevant and automated sources remains challenging.
- 449 • The relationship between social media response and flood severity is difficult to study given no
 450 consistent metric could be found that isn't limited by country borders. For example, media reporting
 451 may have an effect on social media responses.

8 CREDIT AUTHOR STATEMENT

452 8. CRediT Author Statement

453 **Lydia Bryan-Smith:** Conceptualisation, methodology, data curation, software, investigation, and writing
 454 - original draft. **Jake Godsall:** geospatial tweet analysis methodology, analysis, and visualisation.
 455 **Franky George:** RoBERTa methodology and software. **Kelly Egode:** VADER methodology and software.
 456 **Nina Dethlefs:** Primary Supervision, writing - review & editing, and conceptualisation. **Dan Parsons:**
 457 Conceptualization, Secondary supervision. **VIPER high-performance computing team:** Computing
 458 resources.

459 9. Acknowledgments

460 Lydia Bryan-Smith is funded by a PhD stipend from the University of Hull. We acknowledge the VIPER
 461 high-performance computing facility of the University of Hull and its support team. Some of the results
 462 presented in this article were developed during a Hackathon on Sustainable AI, hosted at the University of
 463 Hull, and funded by a NERC Discipline Hopping grant.

464 A. Appendix A: Code availability

465 The code written in support of this paper has been published on GitHub. The following repositories contain
 466 the code in question:

- 467 • <https://github.com/sbirl/twitter-academic-downloader> (Mozilla Public Licence 2.0):
 468 The command line program written to download the tweets from Twitter, using Twitter's Academic API.
 - 469 • <https://github.com/sbirl/research-smflooding> (GNU Public Licence 3.0): The code
 470 written to train and interact with the AI models tested in this paper.
 - 471 • <https://github.com/jakegodsall/twitter-floods> (GNU Public Licence 3.0): The code
 472 written to geolocate and plot the sentiment of tweets on a map.
 - 473 • <https://huggingface.co/siebert/sentiment-roberta-large-english>: The code used
 474 for sentiment analysis with RoBERTa.
- 475 Please note that all of these code repositories use other external open-source libraries to provide some
 476 functionality. For example, Tensorflow (TensorFlow Contributors, 2019) is used as a machine learning
 477 framework. All open source libraries used are open-source, defined in either requirements.txt (Python)

REFERENCES

478 or package.json (Node.js) freely downloadable from either PyPi (Python Software Foundation, 2022)
 479 (Python libraries) or npm (npm, Inc., 2020) (Javascript).
 480 Python was used as the main programming language. Javascript (Node.js (OpenJS Foundation, 2020)) was
 481 also used to initially download the tweets and to manipulate the data. Bash (shell scripting) was used in the
 482 analysis of the data.

483 References

- 484 Agüero-Torales, M.M., Salas, J.I.A., López-Herrera, A.G., 2021. Deep learning and multilingual sentiment analysis on social media
 485 data: An overview. *Applied Soft Computing* 107, 107373.
 486 Arthur, R., Boulton, C.A., Shotton, H., Williams, H.T.P., 2018. Social sensing of floods in the uk. *PLoS ONE* 13.
 487 Avvenuti, M., Cresci, S., Polla, M.N.L., Marchetti, A., Tesconi, M., 2014. Earthquake emergency management by social sensing. 2014
 488 IEEE International Conference on Pervasive Computing and Communication Workshops (PERCOM WORKSHOPS) , 587–592.
 489 Baccianella, S., Esuli, A., Sebastiani, F., 2010. Sentiwordnet 3.0: An enhanced lexical resource for sentiment analysis and opinion
 490 mining, in: LREC.
 491 Beigi, G., Hu, X., Maciejewski, R., Liu, H., 2016. An overview of sentiment analysis in social media and its applications in disaster relief.
 492 Sentiment analysis and ontology engineering , 313–340.
 493 Brownlee, J., 2019. A gentle introduction to cross-entropy for machine learning. Available online: <https://machinelearningmastery.com/cross-entropy-for-machine-learning/> [Accessed 12/10/2020].
 494 Cho, K., van Merriënboer, B., Çaglar Gülcühre, Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y., 2014. Learning phrase
 495 representations using rnn encoder–decoder for statistical machine translation, in: EMNLP.
 496 Cortes, C., Vapnik, V.N., 1995. Support-vector networks. *Machine Learning* 20, 273–297.
 497 Coulthard, T.J., Neal, J.C., Bates, P.D., Ramirez, J., de Almeida, G.A.M., Hancock, G.R., 2013. Integrating the lisflood-fp 2d hydrodynamic
 498 model with the caesar model: implications for modelling landscape evolution. *Earth Surface Processes and Landforms* 38, 1897–1906.
 499 Davies, G., Roberts, S.G., 2015. Open source flood simulation with a 2d discontinuous-elevation hydrodynamic model.
 500 Deltaires, 2021. Home - delft3d. Available online: <https://oss.deltares.nl/web/delft3d> [Accessed 12/07/2022].
 501 Deltaires, 2022a. Red weather warning issued for storm eunice - met office. Available online:
 502 <https://www.metoffice.gov.uk/about-us/press-office/news/weather-and-climate/2022/red-weather-warning-issued-for-storm-eunice> [Accessed 18/07/2022].
 503 Deltaires, 2022b. Storm franklin named - met ofice. Available online: <https://www.metoffice.gov.uk/about-us/press-office/news/weather-and-climate/2022/storm-franklin-named> [Accessed 18/07/2022].
 504 Devlin, J., Chang, M.W., Lee, K., Toutanova, K., 2019. Bert: Pre-training of deep bidirectional transformers for language understanding.
 505 ArXiv abs/1810.04805.
 506 Dworkin, M., 2015. Sha-3 standard: Permutation-based hash and extendable-output functions.

REFERENCES

- 510 Felbo, B., Mislove, A., Søgaard, A., Rahwan, I., Lehmann, S., 2017. Using millions of emoji occurrences to learn any-domain
511 representations for detecting sentiment, emotion and sarcasm, in: EMNLP.
- 512 FitzGerald, G., Toloo, G.S., Baniahmadi, S., Crompton, D., Tong, S., 2019. Long-term consequences of flooding: a case study of the
513 2011 queensland floods. *The Australian journal of emergency management* 34, 35–40.
- 514 Fu, X., Liu, W., Xu, Y., Yu, C., Wang, T., 2016. Long short-term memory network over rhetorical structure theory for sentence-level
515 sentiment analysis, in: ACML.
- 516 Furquim, G., Filho, G.P.R., Jalali, R., Pessin, G., Pazzi, R.W., Ueyama, J., 2018. How to improve fault tolerance in disaster predictions:
517 A case study about flash floods using iot, ml and real data. *Sensors* 18, 907.
- 518 Gao, H., Barbier, G., Goolsby, R., 2011. Harnessing the crowdsourcing power of social media for disaster relief. *IEEE Intelligent
519 Systems* 26, 10–14. doi:10.1109/MIS.2011.52.
- 520 Goodfellow, I., Bengio, Y., Courville, A., 2016a. Deep Learning. MIT Press. Available online: <http://www.deeplearningbook.org>
521 [Accessed 09/10/2020].
- 522 Goodfellow, I., Bengio, Y., Courville, A., 2016b. Deep Learning. MIT Press. Available online: <http://www.deeplearningbook.org>
523 [Accessed 09/10/2020].
- 524 Gould, I., Wright, I., Collison, M., Ruto, E., Bosworth, G., Pearson, S., 2020. The impact of coastal flooding on agriculture: A case-study
525 of lincolnshire, united kingdom. *Land Degradation & Development* 31, 1545 – 1559.
- 526 He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: 2016 IEEE Conference on Computer Vision
527 and Pattern Recognition (CVPR), pp. 770–778.
- 528 Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural Computation* 9, 1735–1780.
- 529 Hu, Y., Wang, R.Q., 2020. Understanding the removal of precise geotagging in tweets. *Nature Human Behaviour* 4, 1219–1221.
- 530 Hutto, C.J., Gilbert, E., 2014. Vader: A parsimonious rule-based model for sentiment analysis of social media text, in: ICWSM.
- 531 Janowicz, K., Gao, S., McKenzie, G., Hu, Y., Bhaduri, B., 2020. Geoai: spatially explicit artificial intelligence tech-
532 niques for geographic knowledge discovery and beyond. *International Journal of Geographical Information Science*
533 34, 625–636. URL: <https://doi.org/10.1080/13658816.2019.1684500>, doi:10.1080/13658816.2019.1684500,
534 arXiv:<https://doi.org/10.1080/13658816.2019.1684500>.
- 535 Kaller, J., 2019. Loss functions in machine learning for beginners | by john kaller | ai³ | theory, prac-
536 tice, business | medium. Available online: <https://medium.com/ai3-theory-practice-business/loss-functions-in-machine-learning-for-beginners-fastai-lesson-9-homework-2-2121954f1f77>
537 [Accessed 09/10/2020].
- 538 Kankanamge, N., Yigitcanlar, T., Goonetilleke, A., Kamruzzaman, M., 2020. Determining disaster severity through social media analysis:
539 Testing the methodology with south east queensland flood tweets. *International journal of disaster risk reduction* 42, 101360.
- 540 Keung, K.L., Lee, C.K.M., Ng, K.K.H., Yeung, C.K., 2018. Smart city application and analysis: Real-time urban drainage monitoring
541 by iot sensors: A case study of hong kong, in: 2018 IEEE International Conference on Industrial Engineering and Engineering
542 Management (IEEM), pp. 521–525.

REFERENCES

- 544 Khayyam, U., Noureen, S., 2020. Assessing the adverse effects of flooding for the livelihood of the poor and the level of external
545 response: a case study of hazara division, pakistan. Environmental Science and Pollution Research 27, 19638–19649.
- 546 Koehrsen, W., 2018. Neural network embeddings explained | by will koehrsen | towards data science. Available online: <https://towardsdatascience.com/neural-network-embeddings-explained-4d028e6f0526> [Accessed 05/05/2022].
- 548 Kokab, S.T., Asghar, S., Naz, S., 2022. Transformer-based deep learning models for the sentiment analysis of social media data. Array .
- 549 Kongthon, A., Haruechaiyasak, C., Pailai, J., Kongyoung, S., 2012. The role of twitter during a natural disaster: Case study of 2011 thai
550 flood, in: 2012 Proceedings of PICMET '12: Technology Management for Emerging Technologies, pp. 2227–2232.
- 551 Le, X.H., Ho, H.V., Lee, G., Jung, S., 2019. Application of long short-term memory (lstm) neural network for flood forecasting. Water 11,
552 1387.
- 553 Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., Zettlemoyer, L., 2020. Bart: Denoising
554 sequence-to-sequence pre-training for natural language generation, translation, and comprehension, in: ACL.
- 555 Li, H., Caragea, D., Caragea, C., Herndon, N., 2017. Disaster response aided by tweet classification with a domain adaptation approach.
556 Journal of Contingencies and Crisis Management 26, 16–27. URL: <https://doi.org/10.1111/1468-5973.12194>,
557 doi:10.1111/1468-5973.12194.
- 558 Li, W., 2020. Geoai: Where machine learning and big data converge in giscience. J. Spatial Inf. Sci. 20, 71–77.
- 559 Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V., 2019. Roberta: A robustly
560 optimized bert pretraining approach. arXiv preprint arXiv:1907.11692 .
- 561 Mikolov, T., Chen, K., Corrado, G.S., Dean, J., 2013. Efficient estimation of word representations in vector space, in: ICLR.
- 562 Ming, X., Liang, Q., Xia, X., Li, D., Fowler, H.J., 2020. Real-time flood forecasting based on a high-performance 2-d hydrodynamic
563 model and numerical weather predictions. Water Resources Research 56.
- 564 Mohammad, S.M., Turney, P.D., 2013. Crowdsourcing a word–emotion association lexicon. Computational Intelligence 29.
- 565 Mojaddadi, H., Pradhan, B., Nampak, H., Ahmad, N., Ghazali, A.H., 2017. Ensemble machine-learning-based geospatial approach for
566 flood risk assessment using multi-sensor remote-sensing data and gis. Geomatics, Natural Hazards and Risk 8, 1080 – 1102.
- 567 Moshe, Z., Metzger, A., Elidan, G., Kratzert, F., Nevo, S., El-Yaniv, R., 2020. Hydronets: Leveraging river structure for hydrologic
568 modeling. ArXiv abs/2007.00595.
- 569 Ning, H., Li, Z., Hodgson, M.E., Wang, C., 2020. Prototyping a social media flooding photo screening system based on deep learning.
570 ISPRS international journal of geo-information 9, 104.
- 571 npm, Inc., 2020. npm. Available online: <https://npmjs.org/> [Accessed 18/10/2021].
- 572 OpenJS Foundation, 2020. Node.js. Available online: <https://nodejs.org/> [Accessed 22/10/2020].
- 573 Pally, R., Samadi, S., 2022. Application of image processing and convolutional neural networks for flood image classification and
574 semantic segmentation. Environ. Model. Softw. 148, 105285.
- 575 Pennington, J., Socher, R., Manning, C., 2014. Glove: Global vectors for word representation, in: Proceedings of the 2014 Conference
576 on Empirical Methods in Natural Language Processing (EMNLP), pp. 1532–1543.

REFERENCES

- 577 Price, D.H., Hudson, K.L., Boyce, G., Schellekens, J., Moore, R.J., Clark, P.A., Harrison, T.G., Connolly, E., Pilling, C., 2012. Operational
578 use of a grid-based model for flood forecasting.
- 579 Purver, M., Battersby, S.A., 2012. Experimenting with distant supervision for emotion classification, in: EACL.
- 580 Python Software Foundation, 2022. Pypi • the python package index. Available online: <https://pypi.org/> [Accessed 16/05/2022].
- 581 Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever,
582 I., 2021. Learning transferable visual models from natural language supervision, in: ICML 2021: 38th International Conference on
583 Machine Learning, pp. 8748–8763.
- 584 Ragini, J.R., Anand, P.R., Bhaskar, V., 2018. Big data analytics for disaster response and recovery through sentiment analysis.
585 International Journal of Information Management 42, 13–24. URL: <https://www.sciencedirect.com/science/article/pii/S0268401217307843>, doi:<https://doi.org/10.1016/j.ijinfomgt.2018.05.004>.
- 586 Riddell, H., Fenner, C., 2021. User-generated crisis communication: Exploring crisis frames on twitter during hurricane harvey. Southern
587 Communication Journal 86, 31 – 45.
- 588 Rout, J.K., Choo, K.K.R., Dash, A.K., Bakshi, S., Jena, S.K., Williams, K.L., 2018. A model for sentiment and emotion analysis of
589 unstructured social media text. Electronic Commerce Research 18, 181–199.
- 590 Roux, H., Amengual, A., Romero, R., Bladé, E., Sanz-Ramos, M., 2020. Evaluation of two hydrometeorological ensemble strategies for
591 flash-flood forecasting over a catchment of the eastern pyrenees. Natural Hazards and Earth System Sciences 20, 425–450.
- 592 Sahni, T., Chandak, C., Chedeti, N.R., Singh, M., 2017. Efficient twitter sentiment classification using subjective distant supervision.
593 2017 9th International Conference on Communication Systems and Networks (COMSNETS) , 548–553.
- 594 Said, N., Ahmad, K., Gul, A., Ahmad, N., Al-Fuqaha, A.I., 2020. Floods detection in twitter text and images. MediaEval .
- 595 Sakaki, T., Okazaki, M., Matsuo, Y., 2010. Earthquake shakes twitter users: real-time event detection by social sensors, in: The Web
596 Conference.
- 597 Schoene, A.M., Dethlefs, N., 2016. Automatic identification of suicide notes from linguistic and sentiment features, in: LaTeCH@ACL.
- 598 Sloan, L., Morgan, J., 2015. Who tweets with their location? understanding the relationship between demographic characteristics and
599 the use of geoservices and geotagging on twitter. PloS one 10, e0142209.
- 600 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
601 real-time modelling framework. Journal of Flood Risk Management 10, 370–380.
- 602 Teng, J., Jakeman, A., Vaze, J., Croke, B., Dutta, D., Kim, S., 2017. Flood inundation modelling. Environmental Modelling and Software
603 90, 201–216.
- 604 TensorFlow Contributors, 2019. Tensorflow. Available online: <https://www.tensorflow.org/> [Accessed 06/01/2020].
- 605 Vashishtha, S., Susan, S., 2019. Fuzzy rule based unsupervised sentiment analysis from social media posts. Expert Systems with
606 Applications 138, 112834.
- 607 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I., 2017. Attention is all you need,
608 in: Proceedings of the 31st International Conference on Neural Information Processing Systems, pp. 5998–6008.

REFERENCES

- 610 Vichiantong, S., Pongsanguansin, T., Maleewong, M., 2019. Flood simulation by a well-balanced finite volume method in tapi river
611 basin, thailand, 2017. Modelling and Simulation in Engineering .
- 612 de Vitry, M.M., Kramer, S., Wegner, J.D., Leitão, J.P., 2019. Scalable flood level trend monitoring with surveillance cameras using a
613 deep convolutional neural network. Hydrology and Earth System Sciences 23, 4621–4634.
- 614 Wang, R.Q., Mao, H., Wang, Y., Rae, C., Shaw, W., 2018. Hyper-resolution monitoring of urban flooding with social media and
615 crowdsourcing data. Computers & Geosciences 111, 139–147. URL: <https://doi.org/10.1016/j.cageo.2017.11.008>,
616 doi:10.1016/j.cageo.2017.11.008.
- 617 Wu, W., Emerton, R.E., Duan, Q., Wood, A.W., Wetterhall, F., Robertson, D.E., 2020. Ensemble flood forecasting: Current status and
618 future opportunities. Wiley Interdisciplinary Reviews: Water 7.
- 619 Yang, Z., Dai, Z., Yang, Y., Carbonell, J.G., Salakhutdinov, R., Le, Q.V., 2019. Xlnet: Generalized autoregressive pretraining for language
620 understanding, in: NeurIPS.
- 621 Yin, J., Lampert, A., Cameron, M., Robinson, B., Power, R., 2012. Using social media to enhance emergency situation awareness.
622 IEEE Intelligent Systems 27, 52–59. URL: <https://doi.org/10.1109%2Fmis.2012.6>, doi:10.1109/mis.2012.6.
- 623 Zhang, T., Xu, B., Thung, F., Haryono, S.A., Lo, D., Jiang, L., 2020. Sentiment analysis for software engineering: How far can pre-trained
624 transformer models go? 2020 IEEE International Conference on Software Maintenance and Evolution (ICSME) , 70–80.

- Multimodal AI models are proposed for sentiment analysis of flooding tweets & images
- Using emoji labels for training AI models increases performance on unseen datasets
- Sentiment-analysed tweets plotted on a map enables real-time geospatial analysis

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Lydia Bryan-Smith reports financial support was provided by University of Hull. Jake Godsall, Franky George, and Kelly Egode reports financial support was provided by Natural Environment Research Council. Lydia Bryan-Smith reports equipment, drugs, or supplies and statistical analysis were provided by University of Hull Viper High Performance Computing facility. Lydia Bryan-Smith reports a relationship with University of Hull that includes: employment, funding grants, and non-financial support. Nina Dethlefs reports a relationship with University of Hull that includes: employment. Dan Parsons reports a relationship with University of Hull that includes: employment. Dan Parsons reports a relationship with Loughborough University that includes: employment. Lydia Bryan-Smith receives a PHD scholarship from University of Hull. Jake Godsall, Franky George, and Kelly Egode are students at the University of Hull. The NERC discipline-hopping grant declared was for a hackathon they attended. Dan Parsons was employed by the University of Hull, but has recently moved to Loughborough University.
