

WATERSensing: A smart warning system for natural disasters in Spain

José M. Cecilia, Juan-Carlos Cano, Carlos T. Calafate, Pietro Manzoni, Carlos Perinián-Pascual

Universitat Politècnica de València

Francisco Arcas-Túnez, Andrés Muñoz

Universidad Católica de Murcia

Abstract—Floods are expected to increase in the coming years due to global warming. The early identification of water-based disasters can be life-saving, and the challenge is to identify appropriate and timely warning measures. Social-media tools such as Twitter provide citizens with a real-time communication channel for reporting problems related to our environment, which allows humans to act as social sensors. In this paper, we show the main results and lessons learned from the research project WATERoT, funded by the Spanish government. In this project, we designed a social sensing application (called *WATERSensing*) for the prevention and evaluation of water-related disasters with the participation of individuals through social networks. This tool crawls micro-texts from different social networks such as Twitter, RSS feeds or Telegram, which are analyzed with natural language processing (NLP) techniques. A case study of Storm Gloria, a Mediterranean storm that heavily affected eastern Spain last January 2020, is presented to evidence that the system can correlate data from social media with actual events. We demonstrate that the analysis of different sources of information opens up new opportunities in the development of warning-systems for the prevention, early identification, and management of natural disasters.

I. INTRODUCTION

CLIMATE change is one of the main challenges that modern societies have to face. Its impact, often associated with extreme meteorological events, has dramatic consequences all over the world. The risks associated with this type of episodes are expected to continue growing with the continuously increasing global warming. Unfortunately, many regions have already experienced

these consequences in terms of loss of human lives, natural ecosystems and economic losses, including damage to housing, infrastructure, primary production systems and tourism, not to mention food, insecurity, and the significant impact on health. In the context of the Sendai Framework for Disaster Risk Reduction (SFDRR - <https://www.undrr.org>) 2015-2030, preventing and reducing disaster risks involves understanding them, strengthening disaster-management policies, increasing investments to reduce the impacts of these events, and improving the preparedness for effective disaster response [1].

The early identification of water-based disasters can be lifesaving; hence, the challenge is to identify appropriate and timely warning measures in a continuously changing environment [2]. New methodologies for obtaining observed flow data are currently being explored through crowd-sensing tools [3]–[5]. In this regard, mobile crowd-sensing (MCS) is a priority research line, which relies on data collection from many mobile sensing devices to empower ordinary citizens to extract intelligence from crowds and deliver human-centric services thorough the aggregation and fusion of this data in the cloud [6]. MCS is inexpensive, since there is no need for network deployment, and its spatio-temporal coverage is outstanding. Following Guo et al. [7], two data-generation modes in MCS can be distinguished: (1) mobile sensing, which leverages raw data generated from the hardware sensors that are embedded in mobile devices (e.g., accelerometer, GPS, camera, or microphone, among others), and (2) social sensing (or social networking), which leverages user-contributed data from social media. The latter mode considers participants as “social

sensors”, i.e., agents that provide information about their environment through social-media services after the interaction with other agents.

In the context of MCS, this article explores the applicability of social sensing to achieve the early detection of water-related issues through an analysis of people’s posts and opinions in social networks. Indeed, the ubiquitous use of social networks, such as Twitter or Facebook, establishes a generic framework for having a sufficiently large social critical mass to provide citizens with a real-time communication channel for reporting problems related to our environment. Several studies have employed social media to automatically identify natural hazards. Data collected from Twitter have been analyzed for the detection/location of earthquakes and typhoons [8], wildfires [9] and heat waves [10]. Moreover, Spielhofer et al. [11] provided interesting insights into a flooding crisis in UK by leveraging data from Twitter, where they concluded that this type of techniques could be used in real time to provide actionable intelligence for emergency services. Similarly, Shi et al. [12] analyzed and classified a batch of tweets generated during a flooding episode in the Kinu River Basin, Japan. They proved that the information contained in these tweets actually provided valuable information for disaster management in the target area. The above studies used Twitter in an offline batch-based manner, downloading information for *a posteriori* analysis following a streaming approach. With respect to real-time applications to track social-media information for water-related issues, Bruijij et al. [4] developed a database for detecting floods in real time on a global scale using Twitter. Their tools are useful to obtain datasets for an offline analysis but, they do not provide actionable information, i.e., they do not identify issues related to floods.

In contrast, our application *WATERSensing* takes real-time input from different social-media sources, which is then analyzed with Natural Language Processing (NLP) techniques to identify water-related issues in real time. Several metrics are also used to categorize micro-texts and provide a problem-relatedness perception index (PPI). Moreover, location, person and organization entities are extracted through named-entity recognition (NER) techniques to visualize these issues in real time.

WATERSensing is validated in a real natural disaster scenario, i.e., Storm Gloria, which heavily affected eastern Spain during January 2020. Our results show that the identification of issues performed by our application matches with the timing of a natural disaster, and it is correlated with actual events.

II. SOFTWARE ARCHITECTURE

The overall software architecture of *WATERSensing* is based on three main modules: (1) the crawling or sensing module, which crawls social-media information from different sources, such as Twitter, Telegram and Really Simple Syndication (RSS) messages published by Spanish media, (2) the NLP processing module, that analyzes and categorizes these micro-texts into different topics previously defined by the user, setting a topic score, sentiment score, and combines them into a PPI score, and (3) the visualization module, that shows the processed, geolocated and categorized information, so that it can provide insights for the user. Moreover, the raw and processed information is stored in an Elastic Search database. In the remainder of this section, we introduce the design and implementation details of *WATERSensing*’s main modules.

A. The crawling module

The crawling module is responsible for obtaining a collection of micro-texts from different social networks. Particularly, we have focused on Twitter, Telegram and RSS, based on user-defined settings such as a list of RSS feeds, URLs, Twitter hashtags and/or keywords to be tracked, and Telegram channels/groups. The acquisition of tweets is performed through the Twitter API by setting specific keywords (e.g., *aquifer*, *drought*, *landslide*, or *waste dumping*). A fixed list of RSS feeds is established to track information from representative Spanish media sources. These include national and regional newspapers, official websites of public institutions and government bodies, etc. Finally, the Telegram crawler is linked to the user account of the admin user, filtering private chats. Duplicate feeds are filtered by checking an MD5 hash generated for each micro-text, which was previously stored in the storage module.

B. The NLP processing module

The NLP processing module of the application is based on CASPER [13], a workbench that is aimed at detecting multi-domain problems described in micro-texts. The main tasks in this module are as follows.

1) *Setting topic categories*: One or more categories of user-defined topics (e.g., drought, flooding, etc.) must be previously registered. This category is called “space” in our system, and it is created semi-automatically by selecting significant features, which are relevant words related to the topic of interest. We use seven WordNet semantic relationships to expand the initial keywords, i.e., *x-causes-y*, *x-derived_from-y*, *x-has_hyponym-y*, *x-has_subevent-y*, *x-near_synonym-y*, *x-pertains_to-y* and *x-related_to-y*.

The first step is to select some seed terms that are representative of each category. Currently, the semantic expansion is based on Spanish. In this example, the keyword is “*lluvia*” (i.e., rain), which is the same word used to name the space (in capital letters). Second, the user must also select the relevant meanings to which each keyword is linked, so that the system scans the glosses for definitions related to the given keyword. In this example, the second meaning of “*lluvia*” corresponds to “Atmospheric phenomenon consisting of rain, water falling from the clouds, from the sky”, which is represented by the term *precipitación* (i.e., rainfall) and associated to words such as *aguacero*, *precipitación*, *lluvia*, *pluvioso*, etc. (i.e., downpour, rainfall, rain, pluvial, etc.). Then, each selected keyword is converted into a topical feature, resulting in vector $C_i = (f_{i1}, f_{i2}, \dots, f_{ik})$, where every f_{ij} identifies a feature in the form of a WordNet synset (i.e., set of synonymous words) assigned to the category C_i . Third, the system suggests a set of additional words and phrases from these syn-sets, so that the user can choose them to augment the original list of category descriptors. In this example, some additional words related to the first meaning of “*lluvia*” are *agua dulce* and *bajada descenso* (i.e., fresh water, descent).

2) *Processing natural language*: The main goal of this stage is to perform the aggregation and filtering of the monitored topics. Thus, each micro-text is divided into sentences, each sentence is

tokenized, and then each token is labelled with a part of speech (POS). A tweet is now represented as the vector $T_m = (w_{m1}, w_{m2}, \dots, w_{mp})$, in which w_{mn} is an object for every word that appears in the tweet T_m , being p the number of words of T_m . Each w_{mn} is represented with several attributes, such as the word form, the position in the micro-text, the stem, the POS, the topic, and the sentiment, just to mention a few.

The topic value of each w_{mn} is obtained as follows. The ngrams (i.e., sequences of n contiguous lexical items in a given text) of each w_{mn} in T_m are searched in the vector C_i . If it matches any f_{ij} in C_i , the topic value is set to 1. Calculating the value of sentiment involves detecting significant ngrams with respect to sentiment. These significant ngrams are discovered by making use of two distinct datasets. On the one hand, a polarity lexicon called SENTIMENTS holds the polarity type (i.e., positive or negative) for a number of words, where their POS (i.e., verb, adjective, noun or adverb) is taken into consideration. On the other hand, the user can set any f_{ij} in C_i as a negative word when setting up topic categories. This feature allows users to expand the lexicon of polarity with additional negative words that are related to the particular context domain. Thus, the ngrams of each w_{mn} in T_m are searched in both datasets. The application assigns the values +1 or -1 (for positively and negatively marked ngrams, respectively) to the feeling attribute of each w_{mn} in T_m when a hit is found.

3) *Determining the topic score*: Since tweets and categories are represented as vectors, a measure of similarity can be used to evaluate the degree of relationship between these two vectors. In this context, the “cosine similarity” (or normalized point product) is used as a measure of semantic distance.

4) *Determining the sentiment score*: Our application does not use a straightforward approach to sentiment calculation, i.e., just summing up the sentiment values of the ngrams in the micro-text and determining the polarity of the message by the sign of the final score. Our sentiment measure is grounded on the degree of sentiment in each tweet through a metric originally used to assess political positions in texts [14]. In particular, a scaling procedure changes the counts of sentiment-coded ngrams in the tweet T_m into a point on the

sentiment dimension S . It is important to note that, as we are only concerned with problem detection, the sentiment-relatedness function automatically assigns the value of zero to those tweets that do not express negative polarity.

5) *Detecting the problem*: With the aim of having a final score that could serve to identify if a particular tweet is likely to deal with a given water-related problem, the problem-relatedness perception index (PPI) is eventually computed, resulting from the geometric mean of the topic score and the sentiment score. Therefore, tweets can be sorted according to the PPI, which can be used to establish alert thresholds, so that messages can be sent to the system administrator to report about any potential risk, hazard, or disaster. This alert system is now being implemented through email notifications. Actually, when there are more than N messages with $PPI > 0.5$, the system sends an email to the administrator with a summary of these messages.

6) *Detecting entities*: We are also capable of labelling up to three different types of entities, i.e., LOC for location, PER for person, and ORG for organization, through a named-entity recognition (NER) module. This module is built using the *Spacy* library, which has a pre-trained model for Spanish. However, for a better performance, a custom model has been trained using a corpus automatically extracted from Wikipedia. The framework uses a transition-based prediction model (i.e., embedding plus encoding with a convolutional neural network) capable of understanding the context of each word in a sentence, and using this context as input for the neural network to predict its next move, i.e., whether it tags a word as a named entity or as part of a named entity, skip to the next word, etc.

C. Visualising results

The visualization module shows the information collected and processed by the previous modules. We refer the reader to *sensingtools.com* and click WATERSensing for insights. It shows the number of feeds that have been categorized from different sources (e.g., Twitter, RSS and Telegram), along with the processed messages and the meta-information obtained with the NLP module to identify the probability that this text represents a problem (i.e., PPI), the main entities recognized in

the text, and the date when it was published within the selected time frame (i.e., 1 hour, 1 day, 1 week or 1 month). Moreover, the processed messages can not only be filtered by date, but also by relevance, feeling, PPI, etc. There is also a search bar where users can search for particular keywords. To speed up the search process, the entities recognized with the NER are displayed in the search bar, so that the user can select which texts contain a certain entity. The application also shows an image gallery that contains the pictures accompanying the micro-texts that were successfully categorized. Finally, some diagrams are displayed to show the volume of processed messages in this topic, the feeling of these messages, and their localization.

III. RESULTS AND DISCUSSION

Our tool can process information in real time and provide several metrics such as the number of tweets, the PPI for each tweet, the most relevant users, etc. The results shown in Figure 1 and Figure 2 are based on the space *Lluvia* (Rain), with the following keywords associated to the space: torrente, chubasco, aguacero, chaparrón, diluvio, torrencial, anegar, inundar, llover a cántaros, diluviar, diluvial, precipitación, lluvia, pluvioso, lluvioso, llover, precipitar, tromba, and marina (i.e., rainfall, downpour, cloudburst, deluge, waterspout, torrent, rain cats and dogs, rain buckets, showery, rainy etc.), which were carefully selected to avoid false positives, resulting from tweets that include apparently significant words but with senses that are not related to raining events. The search for the information was restricted to Spain and Spanish, so language identification was used to discard messages written in regional languages (e.g., Catalan, Basque or Galician). The crawling was activated on January 19, 2020, when the deadly storm Gloria battered Spain. The consequences of this storm have been dramatic in eastern Spain. For instance, the Balearic Islands and the region around Valencia were the hardest hit on January 20, with heavy flooding and strong winds. More than 30 provinces in Spain were put on high alert.

Figure 2 shows the influx of messages from Twitter, RSS and Telegram for a period of 24 hours following the activation. The blue colour shows the tweets that include some of the above keywords. It

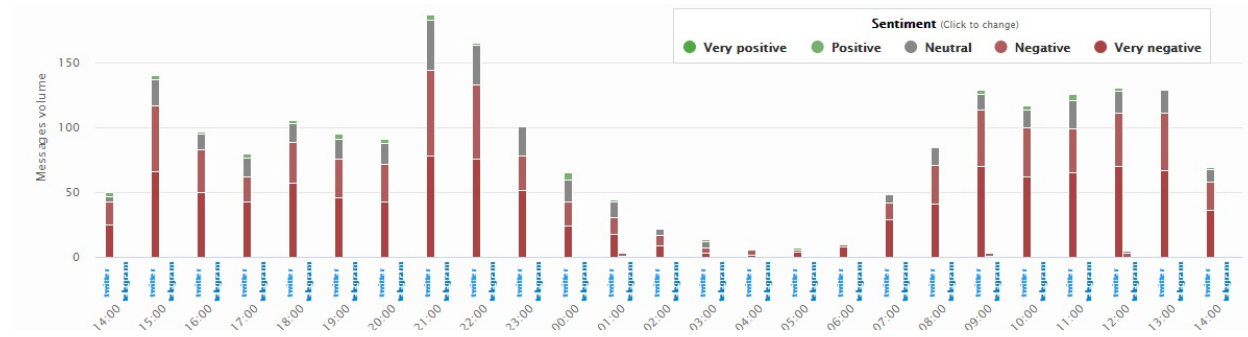


Fig. 1: Number of messages in the Rain space in the last 24 hours according to sentiment analysis from January 19th, 2020 (14:00) to 20th, 2020 (14:00).



Fig. 2: Influx of messages in the Rain space in the last 24 hours (from January 19th, 2020 (14:00) to 20th, 2020 (14:00)). Blue area shows Twitter messages and the Orange one RSS messages.

shows a sudden increase at the date when the event was consolidated and had greatest impact. Indeed, events like Storm Gloria usually have two critical points. First, authorities warn the population about the potential risks. At this point, if the event is perceived as dangerous by the population, it usually has an impact on social networks and media. This can be seen on the left-hand side of Figure 2. Second, activity in social networks and newspapers is usually even greater when the natural disaster actually occurs and causes damage to the population. This effect is displayed on the right-hand side of Figure 2. Therefore, the course of the activity in social networks and media generated by this event clearly correlates with the usual evolution of natural disasters. This means that people's opinions in social networks are storytelling and the analysis of this information provides a real-time channel of information about the state and consequences of the natural disaster.

Figure 1 shows the number of feeds that were categorized in each source, along with the time when those feeds were published and their impact, being categorized using the PPI metric. This figure shows important information that really matches

the results displayed in Figure 2, where the two control points are also represented. Here the feeling is also introduced, because people can write about a particular topic in a positive or negative way. However, the application is specifically aimed at detecting problems, so we are interested in very negative opinions (i.e., the darkest area in 1).

Finally, Figure 3 shows two maps. The map on the left displays the location of the categorized feeds for the category Rain. They were located by making use of the toponyms mentioned in the micro-texts. The map on the right is the reflectivity map generated by the official State Meteorological Agency (AEMET). Reflectivity images correspond to the lowest elevation of the radar scan (0.5° above the horizontal plane), and the colour scale indicates reflectivity ranges in Z decibels. It should be noted that a correlation exists between these maps, which were generated for the same time frame (i.e., from 21 January 2020 at 00:30 to 21 January 2020 at 12:00). Indeed, as clusters of hazard-related posts correspond to areas with high decibel levels, we conclude that there is a direct correlation between people's opinions and natural events.

IV. CONCLUSION

This paper presented *WATERSensing*, a social sensing application for the early detection of water-related issues through people's opinions in social networks. We introduced the software architecture underlying this application, which is based on three main modules: (1) data collection from different social networks (crawling module), (2) micro-text

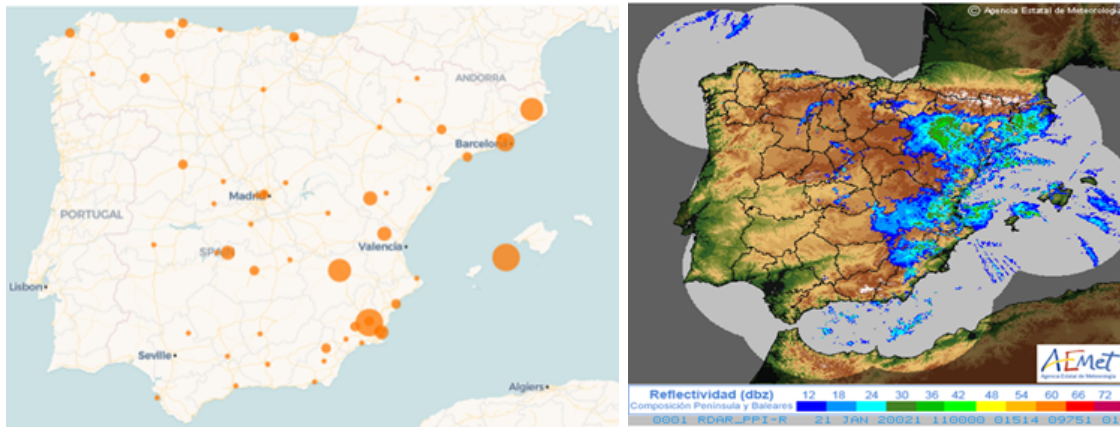


Fig. 3: Map showing the location of categorized feeds for Rain (left) and AEMET's reflectivity map on 21 January 2020 (right).

processing to semantically analyze the information retrieved by our system (NLP module), and (3) visual feedback to the user to look for insights (visualization module). The NLP module establishes several problem-identification metrics that are used to classify tweets based on how likely they reflect a given problem. Moreover, we designed an NER module that recognizes several named entities such as locations, persons, and organizations. Location entities were disambiguated at the visualization stage by pinpointing locations on a map.

To assess our approach, we analyzed Storm Gloria, which occurred on the eastern Mediterranean coast from January 19 to 22, 2020. It should be noted that both the alerts issued by the Spanish government before the event and the catastrophic events that occurred had an impact on people's opinions on social networks. The results conclude that the perception of the problem in social networks and the media perfectly correlates with the evolution of the natural disaster. We can anticipate that the benefits of our approach would be relevant to other natural disasters, and the combination with other sources of information, such as YouTube, Facebook or Instagram, could improve the tool to help government authorities and emergency responders in their decision-making processes. Moreover, although our application was specifically designed to address

water issues in Spain, the methodology described in this article is also valid for other hazards, disasters and crisis situations, such as earthquakes, wildfires or even pandemics (e.g. COVID). Indeed, this is a line of research that is currently being explored. It should also be noted that we aim to adapt *WATERSensing* to deal with other languages, providing that we can make use of two types of language-dependent resources, i.e. text-processing resources (e.g., POS tagger) and lexical resources (e.g., polarity lexicon). In this regard, since the former are readily available for many European languages, the main effort is expected to be placed on the latter. Nevertheless, since the construction of semantic spaces when setting topic categories is based on WordNet synsets, further lexical resources are reduced to a few datasets. Indeed, *WATERSensing* has already been expanded to deal with English micro-texts.

ACKNOWLEDGMENT

This work has been partially supported by the Spanish Ministry of Science and Innovation under Grants No. RYC2018-025580-I, RTI2018-096384-B-I00, RTC-2017-6389-5 and RTC2019-007159-5, by the Fundación Séneca del Centro de Coordinación de la Investigación de la Región de Murcia under Project 20813/PI/18, and by the European

Union's Horizon 2020 research and innovation programme under grant agreement No 101017861.

REFERENCES

- [1] I. Kelman, "Climate change and the sendai framework for disaster risk reduction," *International Journal of Disaster Risk Science*, vol. 6, no. 2, pp. 117–127, 2015.
 - [2] J. Fell, J. Peard, and K. Winter, "Low-cost flow sensors: Making smart water monitoring technology affordable," *IEEE Consumer Electronics Magazine*, vol. 8, no. 1, pp. 72–77, 2018.
 - [3] D. Sabella, A. Vaillant, P. Kuure, U. Rauschenbach, and F. Giust, "Mobile-edge computing architecture: The role of mec in the internet of things," *IEEE Consumer Electronics Magazine*, vol. 5, no. 4, pp. 84–91, 2016.
 - [4] J. A. de Bruijn, H. de Moel, B. Jongman, M. C. de Ruiter, J. Wagemaker, and J. C. Aerts, "A global database of historic and real-time flood events based on social media," *Scientific Data*, vol. 6, no. 1, pp. 1–12, 2019.
 - [5] P. Chaudhary, S. D'Arconco, M. Moy de Vitry, J. P. Leitão, and J. D. Wegner, "Flood-water level estimation from social media images," *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. 4, no. 2/W5, pp. 5–12, 2019.
 - [6] M. Liggins II, D. Hall, and J. Llinas, *Handbook of multisensor data fusion: theory and practice*. CRC press, 2017.
 - [7] B. Guo, Z. Wang, Z. Yu, Y. Wang, N. Y. Yen, R. Huang, and X. Zhou, "Mobile crowd sensing and computing: The review of an emerging human-powered sensing paradigm," *ACM computing surveys (CSUR)*, vol. 48, no. 1, pp. 1–31, 2015.
 - [8] B. Poblete, J. Guzmán, J. Maldonado, and F. Tobar, "Robust detection of extreme events using twitter: worldwide earthquake monitoring," *IEEE Transactions on Multimedia*, vol. 20, no. 10, pp. 2551–2561, 2018.
 - [9] C. A. Boulton, H. Shotton, and H. T. Williams, "Using social media to detect and locate wildfires," in *Tenth International AAAI Conference on Web and Social Media*. AAAI, 2016.
 - [10] A. P. Kirilenko, T. Molodtsova, and S. O. Stepchenkova, "People as sensors: Mass media and local temperature influence climate change discussion on twitter," *Global Environmental Change*, vol. 30, pp. 92–100, 2015.
 - [11] T. Spielhofer, R. Greenlaw, D. Markham, and A. Hahne, "Data mining twitter during the uk floods: Investigating the potential use of social media in emergency management," in *2016 3rd International Conference on Information and Communication Technologies for Disaster Management (ICT-DM)*. IEEE, 2016, pp. 1–6.
 - [12] Y. Shi, T. Sayama, K. Takara, and K. Ohtake, "Detecting flood inundation information through twitter: The 2015 kinu river flood disaster in japan," *Journal of Natural Disaster Science*, vol. 40, no. 1, pp. 1–13, 2019.
 - [13] C. Perrián-Pascual and F. Arcas-Túnez, "Detecting environmentally-related problems on twitter," *Biosystems Engineering*, vol. 177, pp. 31–48, 2019.
 - [14] W. Lowe, K. Benoit, S. Mikhaylov, and M. Laver, "Scaling policy preferences from coded political texts," *Legislative Studies Quarterly*, vol. 36, no. 1, pp. 123–155, 2011.
- José M. Cecilia** is a Ramón y Cajal research fellow (Associate Professor Tenure track) at the Computer Engineering Department, UPV (Spain). His research interest includes HPC, IoT, AI and social sensing. Contact him at jmcecilia@disca.upv.es.
- Juan-Carlos Cano** is a full Professor at the Department of Computer Engineering, UPV (Spain) and Senior Member of IEEE. His current research interests include vehicular networks, mobile ad hoc networks, and pervasive computing. Contact him at jucano@disca.upv.es.
- Carlos T. Calafate** is a full Professor at the Department of Computer Engineering, UPV (Spain). His research interests include ad-hoc and vehicular networks, UAVs, Smart Cities & IoT, QoS, network protocols, video streaming, and network security. Contact him at calafate@disca.upv.es.
- Pietro Manzoni** is a full Professor at the Department of Computer Engineering, UPV (Spain) and senior member of IEEE. His research activity focuses on Mobile Wireless Systems, IoT for Smart Cities and Rural Areas, LPWAN-based networks, and Pub/Sub systems. Contact him at pmanzoni@disca.upv.es.
- Carlos Perrián-Pascual** is an associate professor in the Applied Linguistics Department at UPV (Spain). His main research interests include natural language processing, knowledge engineering, computational linguistics and text mining. Contact him at jopepas3@upv.es.
- Francisco Arcas-Túnez** is an associate professor at UCAM (Spain). His main research interests include natural language processing, knowledge engineering, computational linguistics and text mining. Contact him at farcas@ucam.edu.
- Andrés Muñoz** is an associate professor at UCAM (Spain). His main research interests include Semantic Web technologies, Ambient Intelligence and Intelligent Environments. Contact him at amunoz@ucam.edu.