# GENERATING TEXT DESCRIPTIONS FOR GEOGRAPHICALLY DISTRIBUTED SENSORS

Martin Molina and Javier Sanchez-Soriano
*Department of Artificial Intelligence, Technical University of Madrid*
*Campus de Montegancedo S/N, 28660, Boadilla del Monte, Madrid, Spain*

## ABSTRACT

Sensor networks, with thousands of geographically distributed sensors and different types of quantitative measures, need software tools to help users understand the meaning of measures. In this paper we pay attention to the problem of automatic generation of geographic descriptions in natural language for geographically distributed sensors. We describe this problem in the context of a web application in the domain of hydrology which is part of a more complex multimedia presentation system that combines text and graphics. We describe the web application and the algorithm that we designed to generate the geographic descriptions for sensors. Besides GIS data files, our method uses two information sources: an online server for geographic names (Geonames) and a specific knowledge base with text patterns that we constructed to process sensor identifiers. The evaluation results confirm that online geographic information resources such as Geonames are useful to generate names for sensors but they need to be combined with other more specific information sources (such as our knowledge base) to obtain good descriptions. We also compare our method with related work and show future lines of work.

## KEYWORDS

Multimedia presentation system, natural language generation, geographic information, sensor network

## 1. INTRODUCTION

Today there is an increasing availability of sensors infrastructures to measure the behavior of dynamic systems such as road networks for traffic surveillance and control or river channels for water management. These infrastructures usually include a large amount of sensors (e.g., thousands of sensors), geographically distributed, that record data on real time periodically (e.g., every hour, every 15 minutes, etc). The measured information is usually quantitative with spatial and temporal references.

To help users analyze this type of information, it is important to have computer systems that interpret and translate the data to more understandable representations. For example, a type of computer system called data-to-text system [Reiter, et al., 2005; Hunter, et al., 2008] translates quantitative data into text in natural language. This system can be part of complex multimedia presentation systems [André, 2000], in combination with graphic visualization tools, to automatically construct effective descriptions that help users to analyze quantitative data.

In the context of sensor networks, with thousands of geographically distributed sensors and different types of measures, an important function of this type of system is to generate text that describes to end users the spatial location of events. For example, in the hydrologic domain, humans describe spatial events using sentences like "there is heavy rain *in the South of Spain*" or the "the water level *in the Ebro River at Ascó* has decreased". The automatic generation of such descriptions is related to what is called referring expression generation. This problem has received attention within the research community of natural language generation and different general approaches have been proposed. However, the domain of geographic expressions presents particular characteristics compared to other more general problems in natural language generation [Turner, et al., 2009]. For example, these solutions use domain knowledge about geography and specific procedures for spatial reasoning to construct appropriate descriptions.

In this paper we pay attention to the problem of automatic generation of geographic descriptions in natural language for geographically distributed sensors. We describe this problem in the context of a web

application in the domain of hydrology. In the paper we describe the web application and the algorithm that we designed construct the geographic descriptions for sensors. We present the results of the evaluation of this method, showing the ability of our method to generate appropriate descriptions and the convenience of using online geographic information servers in combination with other information sources. We also describe related work and future lines of research.



Figure 1. Example presentation generated by the VSAIH system.

## 2. THE VSAIH APPLICATION

VSAIH is a web application [Molina, Flores, 2011] for generating multimedia descriptions that summarize the behavior of hydrologic networks controlled by a national information system in Spain called SAIH. We developed this system to help users who need to interpret and analyze the behavior of rivers and make decisions according to prefixed management goals. The SAIH national system (Spanish acronym for Automatic System Information in Hydrology) [DGA, 2009] includes sensor devices and telecommunications networks in the main river basins to get on real time in control centers hydrologic information about the state of the rivers. The information is recorded periodically (e.g, every hour, 30 minutes or 15 minutes). The SAIH system includes different types of sensors: (1) *pluviometers* that measure the precipitation of rain (in millimeters per hour) at the point they are located, (2) *flow stations* located on the riverbed to measure its flow (in cubic meters per second), (3) *level stations* located at a reservoir or a river to measure the water level (in meters with respect to sea level), (4) *volume stations*, located at the dam of a reservoir to measure the volume (in hectometers) of stored water (for practical reasons, this is considered as a sensor but actually it is deduced locally from the level of the reservoir).

VSAIH explains the meaning of measures recorded by the hydrologic sensors in order to make this information more accessible for non-expert users. The information is presented to the user using a journalistic approach, based on the idea of an online virtual newspaper with automatically generated news [Molina, et al., 2011a; Molina, et al., 2011b]. In contrast to a specialized web application that only presents results of analyses to expert hydrologists graphically, VSAIH includes text explanations and it is potentially more useful to a wider range of users (e.g., municipalities, civil protection, engineering consultants, educators,

etc.). VSAIH interprets the quantitative information using knowledge about space, time and hydrology and constructs understandable presentations.

VSAIH generates presentations every hour summarizing about 45,000 measures (sensor readings across Spain). The presentations combine text in journalistic style (headlines and body text organized into coherent discourses), interactive maps with marked locations, interactive temporal series in 2D graphics, and animations (using pictures from meteorological radars and satellite images). For example, Figure 1 shows a headline, a body text (with hyperlinks), and two graphics: an animated illustration (showing the movement of a storm), and an interactive map with the location of relevant sensors. The text summary includes geographic expressions about sensor locations such as the following: "flow above normal *in the Ebro river at Ascó*", "*3 reservoirs in the Ebro River*", "the maximum decrease in volume has occurred *in the Ribarroja reservoir*". In the following sections we describe how we generate geographic descriptions for SAIH sensors, according to the needs of the VSAIH application.

## 3. THE METHOD

The goal of our method is to construct a text description for a sensor or a group of sensors that expresses in natural language their geographic location. This description is used to construct more complex text descriptions explaining different hydrologic situations in combination with graphics (maps, charts, animations, etc).

In order to generate the description, we use as input two characteristics for each sensor: the type of quantity measured (a value of the set of values {*flow, rain, level, volume*}) and the geographic coordinates of the sensor location (latitude and longitude). As a result, the method generates a text description for an individual sensor, such as "the reservoir Ribarroja", "Albalat", "in the river Ebro at Ascó" or for a set of sensors, such as "the majority of reservoirs in the Jucar basin". The method uses different sources of geographic information to construct the description:

- *GIS data files.* This corresponds to hydrologic spatial information as vector data for the representation of rivers (as polylines) and basins (as spatial regions) in Spain.
- *SAIH identifiers for sensors*. This corresponds to alphanumeric codes provided by SAIH experts. They are normally technical identifiers that are understood only by experts (for example, E69_ZAHARA or M06_STE_ARR_VIL). Since there are nine independent SAIH control centers in Spain, each control center uses different criteria to give identifiers for sensors. Some of the identifiers include geographic places, so this information is potentially useful for our method.
- *Geonames server.* Geonames is a name server (*www.geonames.org*) free of charge with Creative Commons license. It includes more than 7.5 million locations and allows, from an application whose parameters are coordinates, provide a suitable name for the specified location. The name is selected based on different features (populated place, road, hydrographic, etc.). Geonames contains place names in numerous languages, integration of multiple data sources for statistics and characteristics such as elevation, population, etc. The coordinates are expressed in the standard WGS84 (World Geodetic System 1984).

Our method works as follows. For a sensor or group of sensors, the method is able to collect the values for the following attributes:

- *River.* If the sensor or sensors are flow stations, the river where they are located is determined. For this purpose the method uses a GIS vector data file with lines corresponding to river segments and applies a spatial procedure to determine the river where each sensor is located. The procedure finds the line segment which is closest (minimum spatial distance) to each sensor location and returns the name of the river associated to this line.
- *Basin.* The method uses a GIS vector data file to find the basin (defined as a spatial region). The method a simple spatial procedure to find the spatial region where the coordinates of the sensors are included and returns the name of this region as the basin's name.
- *Place.* The method uses the Geonames server to find the closest populated place to the sensor using the geographic coordinates. However, if the place can be extracted from the SAIH identifier, then it is preferred instead of the name provided by Geonames, because we assume that the SAIH experts provide better names for places. The place is extracted from the SAIH identifier (when it is present)

using a heuristic procedure using a knowledge base with text patterns (see more details about this procedure below).

- *Quantification*. The geographic reference for sets of sensors is constructed using the most specific area that covers all sensors (river, basin or nation). We quantify the description with three values {*some, majority, all*} to express *some* (number of sensors less than 50% in the area), *majority* (more than 50% but less than 100%), and *all* (100%).

Using these attributes, we construct automatically the descriptions for sensors using text templates as they are shown in Table 1 (for the sake of clarity, the templates are shown in English although the actual templates are in Spanish). In the following section, we describe in more detail how we select the name of the geographic place corresponding to a sensor.

Table 1. Example templates used to construct descriptions for sensors.

| Sensors | Template | Example of generated description |
|---|---|---|
| A volume/level station at a reservoir | reservoir <place> | reservoir Ribarroja |
| A flow/level station at a river | <river> at <place> | the Guadalquivir river at Andújar |
| A pluviometer | <place> | Albalat |
| A set of volume/level stations at reservoirs | <quantification> reservoirs in <basin> | the majority of reservoirs in the Jucar basin |
| A set of flow/level stations in a river | <quantification> sections in <river> | some sections in the Ebro river |
| A set of flow/level stations in a basin | <quantification> rivers in <basin> | the majority of rivers in the Jucar basin |
| A set of pluviometers in a basin | <quantification> points in <basin> | some points in the Jucar basin |

## 3.1 Selecting the Geographic Place

The SAIH identifiers for sensors often include the geographic place. In principle, these identifiers could be used to get automatically the geographic place of sensors. However, there is not a uniform criteria to write the identifiers and, normally, these criteria are different at each SAIH control center (there are nine control centers). Our solution is to use a heuristic method for text processing based on pattern matching that extracts the names of places from these identifiers using a knowledge base with text patterns that represents the criteria followed for the experts in the nine control centers. We constructed this knowledge base by observing the syntactic rules followed by the experts.

For example, experts from the SAIH control center in the Guadalquivir basin use the following identifier for a sensor: "E69_ZAHARA". This sensor corresponds to a volume station located at a reservoir. The strategy for identifiers followed by experts for these sensors is to write an alphanumeric code (E69), the character underscore ("_") and, then, the name of the place (Zahara). Therefore, the pattern in our knowledge base includes the following criteria: (1) the type of sensor is a volume station, (2) the control center is Guadalquivir, and (3) the geographic place is the text after underscore ("_"). Another example is the sensor identifier "AFORO RÍO GUADALHORCE (CÁRTAMA)" from the SAIH control center in the basins of the South of Spain. This name corresponds to a flow station and includes the name of the river (Guadalhorce) and a populated place (Cártama) between parentheses. This criterion is used for all flow stations in the basins of the South of Spain. Therefore, our pattern in our knowledge includes: (1) the type of sensor is a flow station, (2) the control center is South of Spain and (3) the geographic place is the text between parentheses.

We found that with 19 text patterns we were able to model the majority of the strategies followed by the experts in the nine SAIH control centers when the geographic name is used. In addition to this procedure, our method also performs a final text processing procedure to generate the final name of the place. For example, this procedure changes abbreviations by complete words using a dictionary with the typical abbreviations (e.g., *Fte.* = *fuente*, *Sta.* = *santa*, *Pte.* = *puente*, *Af.* = *aforo*, etc.).

This solution works correctly for a number of SAIH identifiers. However, there are cases when it is not possible to extract the name from the identifier. This is, for example, because experts use only alphanumeric codes (e.g., "C001L85PQUIN" in the Ebro basin, and "M06_STE_ARR_VIL at the Guadalquivir basin), unusual abbreviations ("AR31 G.MINCHONES VILL.V" in the Tajo basin), or other exceptions to general rules.

When it is not possible to match any text pattern with the SAIH identifier, we use the Geonames server and the sensor's geographic coordinates. We ask to the Geonames server the name of the populated place that is closest to the location defined by the geographic coordinates. For instance, we selected the name of place

"Andújar" (a populated place near the sensor location) using Geonames for a sensor with the SAIH identifier "M10_GLQUVIR_AND" in the Guadalquivir basin.

## 4. EVALUATION

We applied three procedures to evaluate our method. First of all, we made a non-exhaustive evaluation by observing randomly descriptions generated by our method while it was working in the online web application. This evaluation procedure does not allow a complete validation of the system, but allowed us to detect some errors (e.g., names whose extension beyond what is desirable, wrong capital letters, and names poorly expressed). To correct these errors, we calibrated the procedure for spatial reasoning (e.g., minimum distance to rivers or maximum distance to populated places) and modified or extended the content of the knowledge base with text patterns.

We applied a second procedure to evaluate correctness based on human judgment. For this purpose, we obtained a sample of the generated descriptions. We randomly generated 90 descriptions (10 descriptions for each one of the nine basins). In order to consider a result correct, we verified the following requirements for each description:

- There are not syntactic errors (excessive length, capital letters improperly set, abbreviations, etc.).
- The description is easy to read (it is not poorly expressed).
- The geographic place and type of measure (rain, flow, etc.) are correct.

We found that 85 out of 90 descriptions were correct (94.4%) and 5 descriptions were wrong (5.6%). This result shows a high rate of valid descriptions according to the goals of our method. For the wrong cases, we found as problems: (1) unusual abbreviations that are not included in our model, (2) wrong capital letters ("río Turia O Guadalaviar en Tramacastilla"), (3) descriptions that match a text pattern but whose information is not appropriate. Other specific causes of problems are repetitions of determinants ("embalse *de del* Limonero") or problems with roman numerals.

Table 2. Evaluation results for the generation of geographic names

| Number | Algorithm | Correct results |
|--------|-----------|-----------------|
| 1 | *Geonames spot* | 7% |
| 2 | *Geonames hydrographic* | 8% |
| 3 | *Geonames populated place* | 58% |
| 4 | *SAIH identifiers* | 75% |
| 5 | *SAIH identifiers + Geonames populated place* | 89% |

This second evaluation procedure for correctness considers that a geographic place is valid when this place corresponds to a correct geographic place near the sensor. However, in practice, humans select names for geographic places based on several geographic features. To evaluate how our method selects appropriate names of geographic places, we applied a third evaluation procedure. In this case, we used 100 complete descriptions for sensors which were generated manually by hydrologic experts for the SAIH system in two basins. For each sensor, we compared the human authored description and the generated descriptions using different versions of our algorithm: (1) *Geonames spot*, i.e, an algorithm that generates the description using the geographic feature *spot* (a landmark like a historic monument, a commercial place, etc.) provided by Geonames, (2) *Geonames hydrographic*, i.e, an algorithm that generates the description using the feature *hydrographic* provided by Geonames, (3) *Geonames populated place*, i.e., an algorithm that generates the description using the closest populated place with Geonames, (4) *SAIH identifiers*, i.e., an algorithm that uses our knowledge base with text patterns to extract the geographic place from SAIH identifiers, and (5) *SAIH identifiers + Geonames populated place*, i.e., an algorithm that combines algorithms 3 and 4 (as it is explained in Section 3.1). We consider a result correct when the human authored description includes the same geographic reference than the generated description (this is based on the assumption that expert hydrologist select the most appropriate geographic reference).

Table 2 shows the results of this comparison. The results show that the algorithms using the features *spot* and *hydrographic* obtained low values: 7% and 8% respectively (we also evaluated the features provided by Geonames *vegetation*, and *road*/*railroad* but we obtained 0% in these cases). The results show that the

algorithm based on using Geonames with the feature *populated place* obtains a value of 58%, i.e. it generates the same geographic references than the descriptions written by SAIH experts in 58% of the cases. The results also show that the algorithm based only on SAIH identifiers (algorithm number 4) obtains a high value (75%). We obtain the best value (89%) when we combine the extraction of places from SAIH identifiers and Geonames with the feature populated place (algorithm number 5). Therefore, the results confirm that using Geonames is useful as source of information to generate descriptions but it needs to be combined with other sources (e.g., SAIH identifiers) to increase the quality of the descriptions up to an acceptable level of quality.

## 5. RELATED WORK AND DISCUSSION

The research work presented in this paper is related to what is called referring expressions generation, an area of natural language generation [Reiter, Dale, 2000]. However, as it is shown by RoadSafe [Turner, et al., 2009], the specific case of geographic descriptions makes this problem different from the general task of referring expression generation.

RoadSafe [Turner, et al., 2009] is a data-to-text system that generates natural language descriptions about the weather condition of roads. Roadsafe generates approximate geographic descriptions such as "in some far southern and southwestern places". RoadSafe and our method generate geographic descriptions but RoadSafe is specialized in summarizing groups of locations instead of generating descriptions for individual places. RoadSafe uses four geographic characteristics: altitude, coastal proximity, population and direction. In contrast, our method can generate descriptions for individual places using appropriate proper names. In addition, we generate expressions for groups of sensors using other geographic characteristics (e.g., river and basin).

The generation of spatial descriptions has been analyzed by [Varges, 2005] using the Map Task dialogue corpus. This approach generates referring expressions that distinguish particular points on the map form other points (the spatial reference for each one of these points is prefixed). In contrast to this, the spatial reference for each point in our approach is not prefixed but it is generated by selecting appropriate geographic attributes as they are selected by experts.

An interesting characteristic of our method is its ability to show some degree of independence from the sensor network. As a result of maintenance procedures, every year new sensors are added to the SAIH network and, sometimes, some sensors are removed. Our method accepts automatically small changes and generates appropriate names for the new sensors without human intervention. However, as the evaluation of our method shows, the best results of our method are obtained when it uses the SAIH identifiers and the knowledge base with text patterns (together with Geonames). The knowledge base depends on the criteria used by SAIH experts, which are relatively stable. If the criteria change, then the knowledge base must be updated manually. However, the knowledge base can be updated with an acceptable effort (according to our experience, this can be done in less than one month by a software programmer, non-expert in hydrology).

Our method could be improved to make it completely independent from the SAIH network. For this purpose, instead of using the knowledge base with text patterns that require to be updated manually, we plan to use in the future additional available online information sources such as geographic ontologies or other geographic data sources similar to Geonames such as VGI data (Volunteered Geographic Information) like OpenStreetMap (*www.openstreetmaps.org*).

In general, geographic descriptions can use locative expressions such as "near Madrid" or "west of Denver" [Creary, et al., 1989]. Our method uses a few locative expressions (e.g., "the Ebro river at Ascó"). A potential improvement of our method is to include more complex locative expressions. For example, a framework for generating locative expressions is proposed by [Kelleher, Kruijff, 2005] addressing the issue of combinatorial explosion inherent in the construction of relational context models. A more general solution may require additional functions such as improvements in spatial reasoning (distances and orientations), more sophisticated geographic locations (roads, coasts, etc.), approximate reasoning (e.g. fuzzy logic models), etc.

# 6. CONCLUSION

In this paper we have described our method to generate automatically descriptions in natural language for geographically distributed sensors. We have described this problem in the context of VSAIH, a web application in the domain of hydrology. VSAIH is a multimedia presentation system that explains the meaning of measures recorded by the hydrologic sensors in order to make this information more accessible for non-expert users. The information is presented to the user using a journalistic approach, based on the idea of an online virtual newspaper with automatically generated news.

Our method generates text descriptions for sensors using GIS data files with spatial reasoning procedures together with a general online geographic database (Geonames) and a knowledge base with text patterns to process SAIH identifiers. This knowledge base represents criteria used by different hydrologists to write sensor identifiers and it is useful to extract automatically the geographic place (when it is present) from these identifiers.

The evaluation of our method shows good results for correctness (94%). We also evaluated the quality of our method for selecting appropriate geographic places by comparing the generated descriptions and human authored descriptions. We compared different versions of our algorithm and we found that the best result (89%) was obtained with the algorithm that combines two sources of information: Geonames (with the strategy of the closest populated place) and the SAIH identifiers with our knowledge base with text patterns. This result contrasts to the best value (58%) for an algorithm that only uses Geonames.

These results confirm that an online geographic information resource such as Geonames is useful to generate names for sensors but it is not enough to obtain good descriptions. In order to improve the performance, we need to use additional information sources such as our knowledge base with text patterns and SAIH identifiers.

Our solution was designed for the SAIH sensor network but it shows some degree of independence from the sensor network. For example, the knowledge base with text patterns depends on the SAIH network (although it requires an acceptable effort for maintenance). In order to design a network independent solution, we are interested in online knowledge sharing using other online information resources with additional geographic data. For example, our future work includes using other web servers with geographic information (e.g., VGI data like OpenStreetMap) and information related to the semantic web with ontologies for geographic information and sensor knowledge (e.g., [Compton, et al., 2009]). In addition, we are planning to improve our method by using additional spatial reasoning mechanisms that help to generate descriptions that use more complex locative expressions.

# ACKNOWLEDGEMENT

# REFERENCES

André, E. 2000. The generation of multimedia presentations. In: Dale, R., Moisl, H., Somers, H. (eds.) *A Handbook of Natural Language Processing: Techniques and Applications for the Processing of Language as Text*, pp. 305–327. Marcel Dekker Inc., New York.

Compton, M. et al, 2009. A Survey of the Semantic Specification of Sensors. *Proceedings of 2nd International Workshop on Semantic Sensor Networks, at 8th International Semantic Web Conference*. Washington DC. USA.

Creary, L. G. et al, 1989. Reference to Locations. In *Proceedings of the 27th annual meeting on Association for Computational Linguistics.* Vancouver, Canada.

D.G.A. (Dirección General del Agua), 2009. El programa S.A.I.H.: Descripción y funcionalidad. El presente y el futuro del sistema. Ministerio de Medio Ambiente y Medio Rural y Marino (Spain). http://www.mma.es/portal/secciones/acm/aguas_continent_zonas_asoc/saih/pdf/SAIH_WEB_MMA_V301109.pdf

Hunter, J. et al, 2008. Using natural language generation technology to improve information flows in intensive care units. *In Proceedings of the 5th Conference on Prestigious Applications of Intelligent Systems*. Patras, Greece.

Kelleher, J. D. and Kruijff, G. M., 2005. A context-dependent algorithm for generating locative expressions in physically situated environments. *Proceedings of the European Workshop on Natural Language Generation*. Aberdeen, Scotland.

Molina M., et al, 2011a. Generating automated news to explain the meaning of sensor data. *Proceedings of the Tenth International Symposium on Intelligent Data Analysis*. Lecture Notes in Computer Science LCNS 7014, Springer, pp. 282–293.

Molina M., et al, 2011b. Using the journalistic metaphor to design user interfaces that explain sensor data. *Proceedings of the 13th IFIP TC13 Conference on Human-Computer Interaction*. Lecture Notes in Computer Science LNCS 6948, Springer, pp. 636–643, 2011.

Molina, M. and Flores, V., 2011. Generating multimedia presentations that summarize the behavior of dynamic systems using a model-based approach. *Expert Systems with Applications* (in press) doi:10.1016/j.eswa.2011.08.135

Reiter, E. and Dale, R., 2000. *Building natural language generation systems*. Cambridge University Press.

Reiter, E. et al, 2005. Choosing words in computer-generated weather forecasts. *Artificial Intelligence* 67(1-2), pp. 137–169.

Turner, R. et al, 2009. Generating Approximate Geographic Descriptions. *Proceedings of the 12th European Workshop on Natural Language Generation*. Athens, Greece.

Varges S., 2005. Spatial descriptions as referring expressions in the MapTask domain. *Proceedings of the 10th European Workshop On Natural Language Generation*. Aberdeen, UK.