



Multi-perspective Approach for Curating and Exploring the History of Climate Change in Latin America within Digital Newspapers

Genoveva Vargas-Solar, José-Luis Zechinelli-Martini, Javier A Espinosa-Oviedo, Luis M Vilches-Blázquez

► To cite this version:

Genoveva Vargas-Solar, José-Luis Zechinelli-Martini, Javier A Espinosa-Oviedo, Luis M Vilches-Blázquez. Multi-perspective Approach for Curating and Exploring the History of Climate Change in Latin America within Digital Newspapers. Computer Science and Information Systems, 2023, 00, 10.2298/CSIS220110008V . hal-03876742

HAL Id: hal-03876742

<https://hal.science/hal-03876742>

Submitted on 28 Nov 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Multi-perspective Approach for Curating and Exploring the History of Climate Change in Latin America within Digital Newspapers

Genoveva Vargas-Solar¹, José-Luis Zechinelli-Martini², Javier A. Espinosa-Oviedo³,
and Luis M. Vilches-Blázquez⁴

¹ CNRS, Univ Lyon, INSA Lyon, UCBL, LIRIS, UMR5205, F-69622 Villeurbanne, France
genoveva.vargas-solar@cnrs.fr

² Fundación Universidad de las Américas Puebla, 72820 San Andrés Cholula, Mexico
joseluis.zechinelli@udlap.mx

³ Univ Lyon, CPE Lyon, LIRIS, UMR5205, F-69616 Villeurbanne, France
javier.espinosa@cpe.fr

⁴ Centro de Investigación en Computación, IPN, 07738 Mexico City, Mexico
lmvilches@cic.ipn.mx

Abstract. This paper introduces a multi-perspective approach to deal with curation and exploration issues in historical newspapers. It has been implemented in the platform LACLICHEV (Latin American Climate Change Evolution platform).

Exploring the history of climate change through digitalized newspapers published around two centuries ago introduces four challenges: (1) curating content for tracking entries describing meteorological events; (2) processing (digging into) colloquial language (and its geographic variations⁵) for extracting meteorological events; (3) analyzing newspapers to discover meteorological patterns possibly associated with climate change; (4) designing tools for exploring the extracted content.

LACLICHEV provides tools for curating, exploring, and analyzing historical newspaper articles, their description and location, and the vocabularies used for referring to meteorological events. This platform makes it possible to understand and identify possible patterns and models that can build an empirical and social view of the history of climate change in the Latin American region.

Keywords: data curation, metadata extraction, data collections exploration, data analytics.

1. Introduction

Ninety-seven per cent of climate scientists agree that climate-warming trends over the past century are very likely due to human activities⁶. Some observation reports and studies reveal that the planet's average surface temperature has risen about 2.0 degrees Fahrenheit (1.1 degrees Celsius) since the late 19th century. The hypothesis is that this change has been mainly driven by increased carbon dioxide and other human-made atmospheric emissions.

⁵ In Iberoamerica, Spanish has variations in the different countries, even if all Spanish-speaking people can perfectly understand each other.

⁶ <https://climate.nasa.gov/scientific-consensus/>

1 Technological advances have allowed understanding of phenomena and complex sys-
 2 tems by collecting many different types of information. Data collections are exported un-
 3 der different releases with different sizes and formats (e.g., CSV, text, excel), sometimes
 4 with various quality features. Tools helping to understand, consolidate and correlate data
 5 collections are crucial. Even if there is an increasing interest in analysing digital data col-
 6 lections for performing historical studies on climatologic events, the history of climate
 7 behaviour is still an open issue that has not revealed missing knowledge. Long histor-
 8 ical data studies could make it possible to compute more complete models of climatic
 9 phenomena and the conditions in which they emerged. However, meteorology is a young
 10 science that started around the 19th century. It is supported by more or less recent data,
 11 making it challenging to run an analysis that can give more historical pictures of climatic
 12 evolution and its implications using observations instead of extrapolations. Those willing
 13 to promote changes in the behaviour of society and industry to reduce emissions that have
 14 a role in climate change must convince civil society of the importance of the challenges.
 15 For this reason, our work addressed the problem of collecting and analyzing the history
 16 of meteorological events to explore how they were described, lived and perceived by civil
 17 society. In this sense, the digitalization of data collections has an increasingly vital role
 18 in collecting vast amounts of *hidden* data. Thus, considering that digital archives become
 19 more easily accessible every time and contain explicit and implicit spatio-temporal in-
 20 formation, researchers in GIScience [18], are becoming aware of these new data sources
 21 [10], [9], [34], [41]. Moreover, digital data collections make it possible to have an analytic
 22 vision of the evolution of environmental, administrative, economic and social phenomena.
 23 In this context, our work deals with data collections that report the emergence of mete-
 24 orological events (e.g., temperature changes, avalanches, river flow growth, or volcano
 25 eruptions). However, the digitized collections have some implicit issues. They are often
 26 riddled with Optical Character Recognition (OCR) errors that hamper the performance
 27 of information retrieval systems. Therefore, handling OCR errors is one of the two sig-
 28 nificant problems for information retrieval from collections of historical documents. On
 29 the other hand, these sources' problems are related to historical language changes since
 30 digitized texts are written in the language of their origin.

31 This paper proposes an extended description of the Latin American Climate Change
 32 Evolution platform called LACLICHEV [37]. The objective of LACLICHEV is to provide
 33 an integrated platform to expose and study meteorological events described in historical
 34 newspapers that are possibly related to the history of climate change in Latin America.
 35 In this sense, we hypothesize that the history (in Latin America) is contained in newspa-
 36 per articles in digital collections available in national libraries of four countries, namely
 37 Mexico, Colombia, Ecuador, and Uruguay. Considering this starting point, LACLICHEV
 38 addresses the following issues (see Figure 1):

- 39 i First, newspaper archaeology, by chasing articles about climatological events using
 40 specific vocabulary to discover as many articles as possible (see the left side of the
 41 Figure 1). The challenge is choosing adequate vocabulary to increase the chances of
 42 getting articles about climatologic events.
- 43 ii Second, once an article talks about a climatologic event, it is tagged with Geo-Temporal
 44 metadata specifying what happened, where and when it happened, its duration and ge-
 45 ographical extent (see the centre of Figure 1). The objective is to build a climatologic
 46 event history of empirical observations.

- 1 iii Finally, on top of this history, the objective is to run analytics questions and visualize
 2 results in maps given that the content is highly spatial (see the right side of the Figure
 3 1).

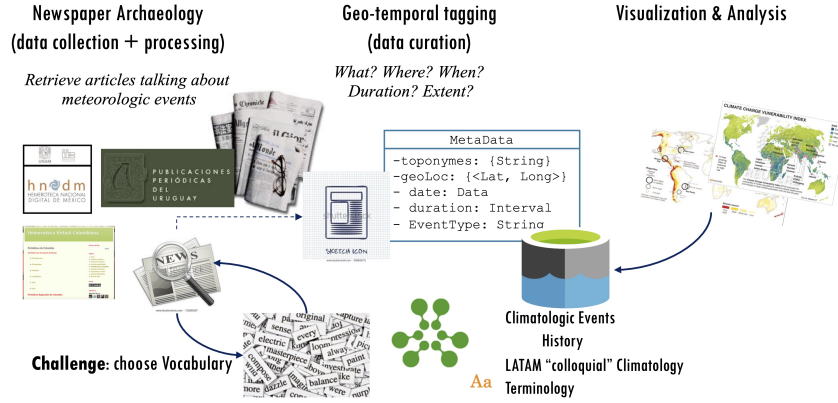


Fig. 1. Problem Statement

4 *Contributions* The main contribution of our work is LACLICHEV. It is a data collections
 5 exploration platform that applies data collections curation and exploration techniques.
 6 These techniques are combined with data retrieval, data analytics, and visualization for
 7 understanding the content of articles that report historical meteorological events. These
 8 data with high geospatial and temporal content can be aggregated into maps. Maps give
 9 a one-shot view of the history of meteorological events observed from the empirical per-
 10 spective of civil society before the emergence of meteorology as a science [39, 6].

11 The second contribution is a meteorological event knowledge model that provides several
 12 perspectives to describe an event. Perspectives organize metadata that represent such an
 13 event is reported through empirical narratives that can appear in newspaper articles, as in
 14 the context of our work.

15 The third contribution is the experimental use of LACLICHEV to build the history of cli-
 16 mate change in Latin America from digital newspapers. To track meteorological events,
 17 we explored newspapers to search articles that could report such events, the conditions in
 18 which they happened, their duration, the places in which they occurred, and their impact
 19 in terms of an approximate number of casualties and the kind of damages, etc. As an ex-
 20 perimental scenario, we chose the XVIII and XIX centuries, which define a golden age
 21 for newspapers in Latin American countries [13], namely, Mexico, Colombia, Ecuador,
 22 and Uruguay.

23 *Organisation of the paper* The remainder of this paper is organized as follows. Section 2
 24 introduces the general architecture of LACLICHEV and its functions implemented by its

1 main modules. Section 3 describes the knowledge model we propose for modelling meteorological events as described in empirical narratives written in natural language. Section 2
 2 meteorological events as described in empirical narratives written in natural language. Section 3
 3 4 describes the general curation and exploration processes implemented by LACLICHEV
 4 to deal with the curation and exploration of historical newspaper articles potentially reporting on climatologic events. It also describes the use cases that we conducted to evaluate it. Section 5 studies approaches that promote datasets exploration for defining the type of analysis possible on top of them. Finally, 6 concludes the paper underlying the contribution and discusses future work.

9 2. LACLICHEV for curating and exploring historical newspapers 10 articles

11 Figure 2 shows the general architecture of LACLICHEV organised into three layers:

- 12 i frontend with an interface providing functions for curating articles and creating events descriptions; and giving access to explore the event history containing curated articles reporting meteorological events;
- 13 ii backend with the meteorological event history stored in a document management system (see number 1 in Figure 2) and modules for curating (pre-processing and tagging the textual content of newspaper articles - number 2 in Figure 2) and exploring events (see number 3 in Figure 2);
- 14 iii external layer connecting to document providers that are available through servers accessible on the Web and APIs exported, for example, by libraries.

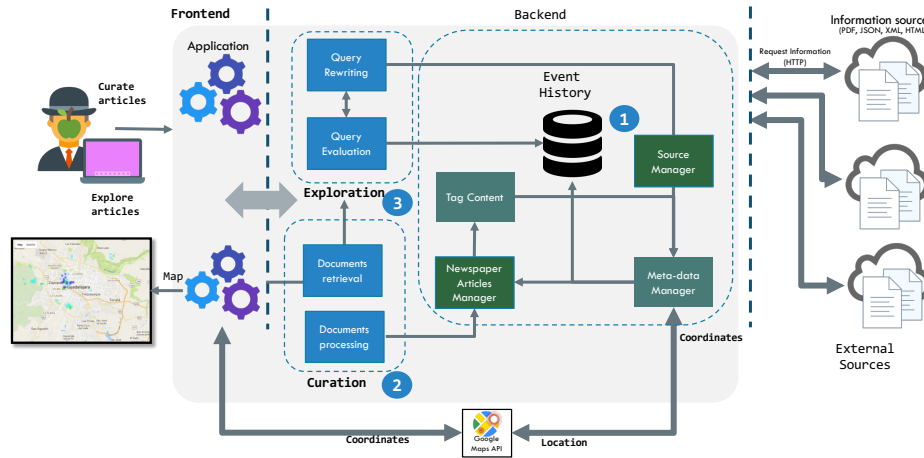


Fig. 2. General functional architecture of LACLICHEV

21 In the following sections, we describe the core layer of LACLICHEV, namely the
 22 backend with its main components, the meteorological event's history, and associated

1 curation modules used to feed the history and exploration modules to process queries to
 2 explore this history.

3 **2.1. Meteorological events history**

4 The event history (storing metadata describing an event from several perspectives, see
 5 number 1 in Figure 2) is based on an event knowledge model that we proposed and that
 6 is described in Section 3. Through this knowledge model, it is possible to represent the
 7 empirical description of a meteorological event with metadata from several perspectives:
 8 descriptive (the vocabulary used for describing a meteorological event and the statistics
 9 of its use); linguistic (the structure of the sentences used in a narrative describing a meteorological event);
 10 the meteorological perspective (represents factual data about an event, location, duration, type, intensity, etc.);
 11 and the domain knowledge perspective (meta-data about empirical and factual observations provided by meteorology experts, e.g., the fact that strong rainfall can correspond to more the 75 mm/hr rain).

12 Metadata is stored in persistence support, a key-value or a document store, depending
 13 on the technology adopted by each library. In contrast, the raw documents remain archived
 14 in a different server or the same store. LACLICHEV uses a document store (i.e., MongoDB⁷) for storing geo-temporally tagged meteorological events. These events' history provides an interface for performing querying and analytics tasks on top of it. The digital collection can be initially queried by filtering the documents by region, country, or year. Digital libraries offer front-ends for performing this classic information retrieval process. For example, select newspapers published in Uruguay (i.e., geographic filter) between 1800-1810 (i.e., temporal filter). It can also perform analytics queries. For example, locate events during the XIX century, enumerate and locate the most famous meteorological events in the region, and create a heat map of the events in Latin America that happened in the last ten years of the XIX century.

26 **2.2. Curating newspapers content modules**

27 The backend of LACLICHEV includes of a set of modules devoted to implement different operations of data curation (see number 2 in Figure 2). The objective of curating (historical) newspaper articles is to build a meteorological events history that newspaper articles reporting events with metadata, providing as much information as possible about the reported event.

28 Figure 3 shows the newspapers curation process that is a semi-automatic process devoted to:

- 29 – find articles reporting this type of events within digital collections available in existing digital libraries repositories;
- 30 – geo-tag interactively and store those articles that actually report such events for building a meteorological event database.

31 LACLICHEV relies on a knowledge graph that integrates a thesaurus classifying meteorological event types, Wordnet and a glossary defining meteorologic characteristics of meteorological events.

⁷ This is a recurrent storage strategy when building databases as a result of processing textual content [32].



Fig. 3. Newspapers curation process

- 1 *Curation process* Curation tasks can be recurrent and include a human-in-the-loop strategy for validating and adjusting results. For example, suppose an event is geo-tagged to
- 2 associate it with a geographic location, and the event is described in an article about Mon-
- 3 tevideo news from Uruguayan newspaper collections. In that case, a human will verify
- 4 that the geographic location refers to Montevideo in Uruguay and not Minnesota (United
- 5 States).
- 6
- 7 During this phase, articles referring to meteorological events are geo-temporally tagged
- 8 to associate them with the region and/or time window in which they happened. The data
- 9 analyst validates tags. Since the result can contain a significant number of articles, the
- 10 user can use three tools to understand the content of the result. The tools let her/him
- 11 manipulate a terms frequency matrix and heat map.
- 12 She/he can also explore the content of the article text using a view that provides in-
- 13 formation about the context in which the terms are potentially describing an event that
- 14 appears in the document. For example, the name of geographic locations in the document
- 15 might refer to the event's location and the region it touched, and a list of geopolitical
- 16 entities (e.g., school, public building, etc.) to determine the damages caused by the event.
- 17
- 18 *Curation functions provided by the backend modules* The data analyst can perform the
- 19 following curation actions:
- 20 - Correct the terms associated with meteorological events that might not be used in such a
- 21 sense in the text. Indeed, some social and political demonstrations are often described as
- 22 meteorological events. For a classic automatic text analysis process, this cannot be easy to
- identify and filter. For example, an article entitled “*Stormy weather within the ails of the*

- 1 *senate in Ecuador*” has nothing to do with the types of events considered but a political
- 2 one.
- 3 - Determine whether personal names correspond to the event’s name (e.g. hurricane or
- 4 storm’s name). If that is the case, this information will be used to insert the event into the
- 5 history.
- 6 - Verify whether the names of cities, regions, and countries correspond to geographic
- 7 entities. The system underlines the names of patronyms, and the data analyst can see the
- 8 location of the possible geographic entities. Thus, the user can also confirm whether the
- 9 article refers to the geographic place that she/he is searching for. For instance, if “Santa
- 10 Clara” is underlined, it can refer to a point of interest, city, or village.
- 11 - Determine the date of the event and its characteristics. The temporal terms and adjectives
- 12 are also underlined to let the data analyst click on those that describe the event.
- 13 - Determine the type of damages caused by the event by exploring those terms that de-
- 14 scribe such information.

15 The previous actions are used to complete the representation of the articles’ content

16 (extracted dynamically) and identify meteorological events more accurately since the data

17 analyst, or domain expert knowledge is used (see Figure 4 showing LACLICHEV inter-

18 faces for curation). Note that one event can be described by several articles. In that case,

19 the information stemming from the different sources is loosely integrated by performing

20 the union of the content by applying some rules. For example, suppose the dates reported

21 in two articles do not entirely correspond (variation of the day or the hour). In that case,

22 the date of the event is modelled as an interval computed by processing the dates. If the

23 dates are too disparate, the system keeps a set of dates. A similar process is done with lo-

24 cations; in this case, the system defines a region. A user can define a threshold of the size

25 of a region associated with an event according to its type. Otherwise, the system keeps a

26 set of geographical points.

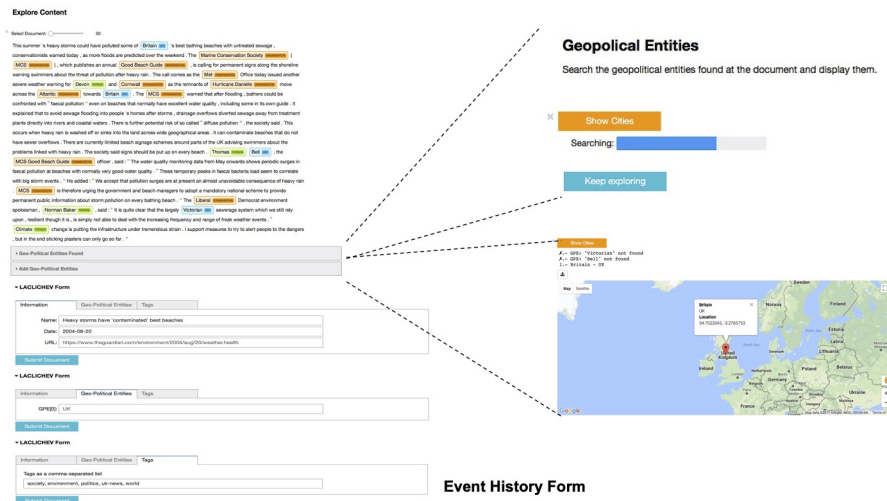


Fig. 4. Event curation process interface for tagging events

2.3. Exploring the collections of digital newspapers

Newspaper articles are explored by conjunctive or disjunctive keyword queries, where keywords can belong to several vocabularies (see number 3 in Figure 2). For example, search articles reporting heavy storms and rivers flooding. The query expressed by a data analyst is automatically completed by using rewriting techniques that consider synonyms, more specific or more general concepts [11]. Thus, three tools can be used for exploring meteorological events depending on expert knowledge of what she/he is looking for.

The rewriting process produces several proposals that the data analyst can adjust and then choose to be evaluated (see details in Section 4). Each chosen query is evaluated using information retrieval techniques, including the article’s text stemming for extracting the terms and constructing a frequency matrix that provides occurrence statistics of the representative terms of the text content within a collection of documents.

In general, information retrieval processes do not exhibit this matrix; it is an internal data structure representing the content of the documents and is used to answer queries. In our approach, this frequency matrix is accessible to the data scientist because it provides an aggregated view of the content of a document collection. Additionally, we compute and exhibit a terms heatmap for a given documents collection to provide a more economical (i.e., consolidated) view of the collection’s content. Our approach provides an interactive interface that lets data scientists manipulate these data structures to define the piece of collections she/he want to explore.

The data scientist can explore them and then decide whether the collection can describe meteorological events and the documents that might be closer to her requirements. She/he can decide eventually to explore some documents directly or reformulate the query. Once a result containing articles that potentially answer the query has been computed, the user can explore the result and validate the selection elements during the next step of the data exploration workflow.

- *Filtering*. Retrieving factual information, for example, filtering events by region, country or year. For example, Uruguay for the country and between 1800–1810 for the temporal filter.

- *Term frequency*. Understanding the content of digital newspaper collection through the vocabulary used in its articles. Therefore, LACLICHEV exposes the terms frequency matrix and a terms heatmap under an interactive interface. The domain expert can see which are, statistically, the terms most used in the articles, group documents according to the terms used, and choose articles using a specific term.

- *Additional information*. Exploring the content of a specific article using a view that provides information about the name of geographic locations in the document. These locations might refer to the event’s location and the region it touched and a list of geospatial features (e.g., school, public building, etc.) to determine, for example, the damages caused by the event.

Exploration process Given a document’s collection and associated data structures describing the content of its articles, the data scientist can explore articles to determine whether they report meteorological events. This phase integrates the human-in-the-loop. The reason is that newspaper articles use colloquial terms that can be tricky and refer to metaphors that might not denote a meteorological event. Language subtleties are not easy

to handle manually, mainly because we are dealing with a language used some centuries ago, which increases the challenge of classifying the content of the articles.

3. Meteorological events knowledge model

We propose a meteorological event knowledge model (see Figure 5) to represent climate event reports in digital documents. The objective is to describe events from different perspectives using the information from the articles and newspapers that report them in an empirical form and complete their description with domain knowledge also described in the model. Newspapers do not describe events scientifically; however, we need to locate and profile them by approximating quantitative characteristics to picture the past climate situation in the region. The different perspectives give context to the quantitative features derived/deduced from the descriptions. As shown in Figure 5, events are associated with the newspaper article(s) that describe them (reading from right to left). Each article can have metadata that curates it, pointing to its “raw” content that has been processed and annotated with linguistic labels.

Classes of documents associated with an event (class *Event* in the figure) contain variables that describe its characteristics, like the date it happened or the geographical scope.

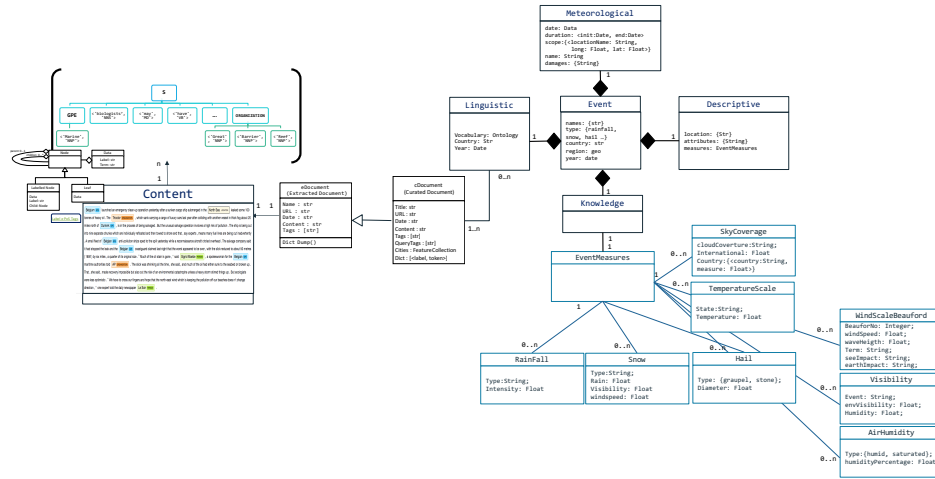


Fig. 5. Event Data Model

According to the perspectives, the event knowledge model provides concepts for representing a *meteorological event*. Each aspect of the knowledge model is implemented using different data structures with associated operations to support exploration actions. The following lines describe the different perspectives of an event and are represented by the event model: descriptive, meteorologic, linguistic and knowledge domain profiles. These perspectives are described in the following sections.

1 3.1. Descriptive profile

2 In newspaper articles, there is no generic list of attributes used for describing a meteorological event. Indeed, meteorological events are described in different ways in historical newspaper articles, depending on the author. However, we can often collect information related to location, date, duration, scope, and damages. Meteorological features (like millimetres of precipitations, wind speed, temperature, pressure, etc.) can be explicitly described in articles or deduced according to the description of the event. For example, an event reported in Montevideo describing an overflow of the river implies winds higher than 100 km/h and rain of more than 10 ml/hour, according to the knowledge provided by meteorologists. This knowledge domain is used to complete the reported event’s meteorologic features.

12 We combine scientific knowledge produced a posteriori with empirical observations reported in colloquial narratives centuries before. This strategy can help to estimate the location of the events. Of course, we could have tried a more appropriate approach correlating the location of the event referred to in the article with ancient distributions and organisation of the territory to have a more precise location of the events. For example, we could have looked for the urban distribution of the city in the publication year of the article. Then, compare this result with contemporary maps and have a more accurate location of the events according to the modern urban distribution of the city. For instance, an event reported in Montevideo city’s “Rambla” sector in 1910 corresponds to a new quarter today. We will develop this approach in our future work.

22 3.2. Linguistic perspective

23 The linguistic perspective gathers the terms used for describing an event in one or several articles belonging to a given newspaper. We propose a tree-based data structure, named *content tree* for representing the content of a historical newspaper article. The tree corresponds to each sentence’s grammatical analysis in the article’s textual content commonly used in Natural Language Processing (NLP) techniques [4]. The **content tree**, as shown below, consists of a set of sentences. A **sentence** is defined as a set of nodes representing grammatical elements of a sentence and leaves representing the terms composing a sentence in a specific article. We use existing classic NLP techniques because we do not aim at contributing to extending or providing novel ways of using them. The objective is to choose adapted methods for processing the meteorologic newspaper texts.

33 In Spanish, we use a simplified grammatical model defined by the following simplified Backus-Naur Form (BNF) specification ⁸. The simplified specification allowed to process the type of articles we explored, of course an extension of the representation in the next versions of LACLICHEV will allow process other texts describing meteorological phenomena for example in historical novels with narratives about major events:

```
38 <sentence> ::= <noun-sentence> | <verb-sentence>
39 <noun-sentence> ::= <named-entity> <conjunction>
40                      <noun-sentence>
41 <noun-sentence> ::= <noun>
```

⁸ We have also used a BNF for English to explore the use of LACLICHEV with other languages. This work is out of the scope of this paper and concerns the next version of LACLICHEV.

```

1 <verb-sentence> ::= <subject> <predicate>
2 <subject> ::= <article> <noun>
3 <predicate> ::= <verb> <direct-object>
4 <direct-object> ::= <article> <noun>
5 <article> ::= EL | LA | UN | UNA
6 <noun> ::= "Spanish nouns"
7 <verb> ::= "Spanish verbs"

```

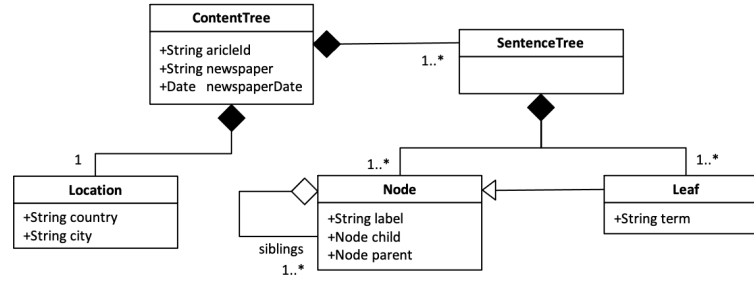


Fig. 6. UML class diagram representing the general structure of a content tree

As shown in Figure 6, a **node** represents a type of grammatical element given in a specific linguistic model defined for a specific language. It is labelled adopting the entity labels produced by classic natural language processing tools known as Part Of Speech (POS) tags. For instance, *noun, proper singular* (NNP), *noun, plural* (NNS), *verb, modal auxiliary* (MD), *Geopolitical entity* (GPE), or Organization. In the case of subjects (NNP), they can be grouped into more general entities that identify geographic locations (GPE), places, names, and organization⁹.

A **node** has children, where each child can also be a Node or a Leaf, and a set of siblings, which are other nodes. A **leaf** specializes in a node, and it represents a term contained in the article. A term is a string with a parent, a **node** means a POS tag.

According to the model, the **ContentTree** represents a document's content where the vocabulary used is determined by a **Location** in a country and a city. These classes represent that the same language, Spanish, varies among countries and cities. Recall that in different locations, people describe meteorological events using different vocabulary.

Every article in a newspaper is associated with its content tree. A data analyst or expert domain can explore the articles by navigating their content trees without reading the full content. For example, *retrieve articles reporting heavy storms in Uruguay in December 1810*. Nodes are related through two relation types: instance, correlation. The relation of type correlation describes two terms that appear in the same sentence with a given distance given by the number of intermediate terms.

⁹ A full list of POS tags can be found in <https://www.cms.gov/>

1 3.3. Meteorological perspective

2 The meteorological perspective characterizes the event with attributes used to describe it
3 in one or several newspaper articles. Nevertheless, not all the attributes can have an asso-
4 ciated value since there might be no evidence within the articles that report it. Attributes,
5 like the date of the event, its geographical scope, or the location of the damaged regions,
6 are computed by navigating through the *content tree* of every article reporting the event.

```
7 MeteorologicEvent: <date: Data,  
8                     duration: <init:Date, end:Date>  
9                     scope:{<locationName: String,  
10                        long: Float, lat: Float>}  
11                     name: String, damages: {String}>
```

12 3.4. Knowledge Domain Perspective

13 The knowledge domain perspective describes meteorological events using knowledge do-
14 main statements created by experts of the National Library of Uruguay. This knowledge
15 has been associated with events through manual analysis of newspaper collections and
16 meteorologists interacting. This knowledge can help interpret the empirical information
17 reported in the articles and complete the information associated with the event description.
18 For example, if the river was flooding due to a storm, it is possible to estimate the wind
19 speed and the approximate litres of rain. The knowledge domain perspective is modelled
20 as a glossary. Figure 7 shows the intuition of its structure.

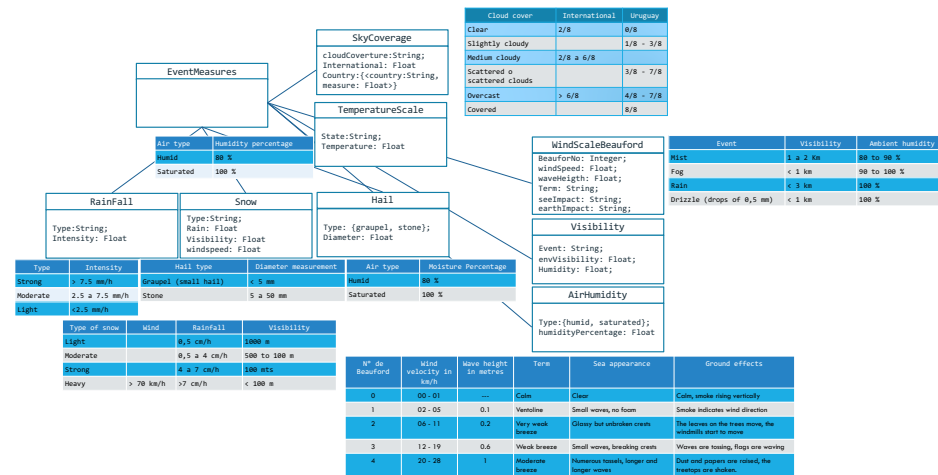


Fig. 7. Climate events glossary

21 Modelling the empirical knowledge about meteorological events is critical when cu-
22 rating newspapers' descriptions. It represents the interpretation of the emerging content

by observing the phenomena and associating it with metering techniques available today. The principle can be stated as follows: *“today, based on the metrology performed during meteorological events, we know that when the river floods, there is an approximate wind speed and more than “x” litres/meter² of rain. So we can estimate the conditions in which the events could have happened in the past.*

4. LACLICHEV in action

LACLICHEV is a client-server system for executing the human-in-the-loop tasks that implement the data exploration process. We have configured LACLICHEV to process historical newspapers of four countries provided by the national libraries of each country. The curated event history has been explored by the librarians of the participating countries. The idea was not to experimentally test the system but to calibrate it according to the characteristics of the digital collections.

4.1. Building a Latin American meteorological event history

We have worked with the national libraries of Mexico, Colombia, Ecuador, and Uruguay to access their newspapers’ digital collections. For our experiments, we worked with the collections of the XVIII and XIX centuries of newspapers written in Spanish with the linguistic variations of those mentioned above Latin American countries. The National Libraries of these countries manage historical newspapers with about 4 to 7 million images of newspapers between the XVIII and XIX centuries, depending on the country. For example, the National Library in Mexico maintains 7 million images of digital national newspaper collections. In Colombia the newspaper library is made up of publications published between the end of the 18th century and the first half of the 20th century, including: “El Papel Periódico Ilustrado”, “Diario Político de Santafé de Bogotá”, “El Alacrán”, “El Mosaico”, “Semanario del Nuevo Reyno de Granada”. It includes newspaper collection from Ecuador and Argentina, namely “La Verdad Desnuda (Guayaquil, Ecuador) and “Vida Intelectual” (Santa Fe, Argentina). The current version of LACLICHEV processed around 19 million images in the newspapers of the fourth countries. The event history has curated 800 different meteorological events.

We curated collections and generated the vocabulary used on articles identified as reporting a meteorological event (see Figure 8). Digital newspaper collections remain in the initial repositories that belong to the libraries. Then, terms and links to the OCR (Optical Character Recognition) archives containing documents with articles reporting meteorological events were stored in distributed histories managed in each country. As shown in Figure 8, the process consists of five steps usually used in natural language processing techniques: sentence segmentation, tokenization, speech tagging, entity and relation detection. LACLICHEV implements these phases in Python, relying on the NLTK library.

The first phase of the pre-processing process of newspapers leads to graphs representing the content of the articles and classic inverse index and frequency matrices used for performing exploration queries.

Besides curating the data collections’ content, we wanted to discover linguistic variations in different Latin American countries to describe meteorological events. People’s

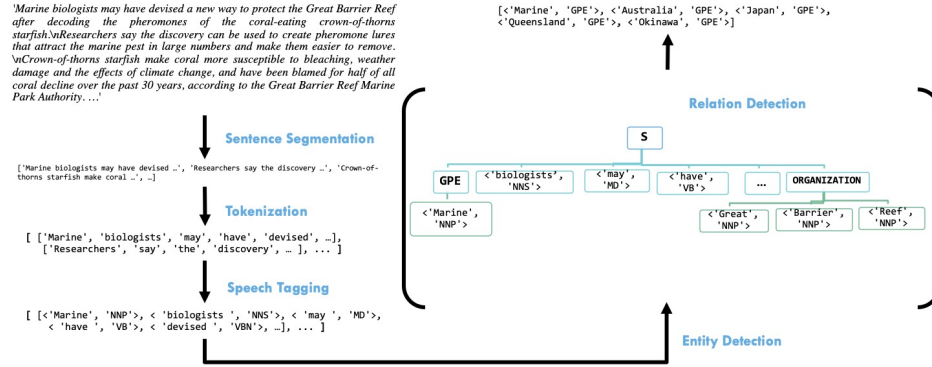


Fig. 8. Pre-processing text pipeline

- 1 language and variations can picture civilians' perception of these events, consequences,
- 2 and associated explanations. Thus, local vocabularies were created out of the terms used in
- 3 newspapers' articles (see Figure 9). For example, referring to a storm as a stormtrooper¹⁰
- 4 Then we updated and enriched through queries, exploration and analytic activities, these
- 5 vocabularies through human-in-the-loop actions. Data analysts tagged "colloquial" terms
- 6 used to describe meteorological events and associated them with more scientific terms.
- 7 These terms can be then used for defining keyword queries for exploring newspaper
- 8 datasets.

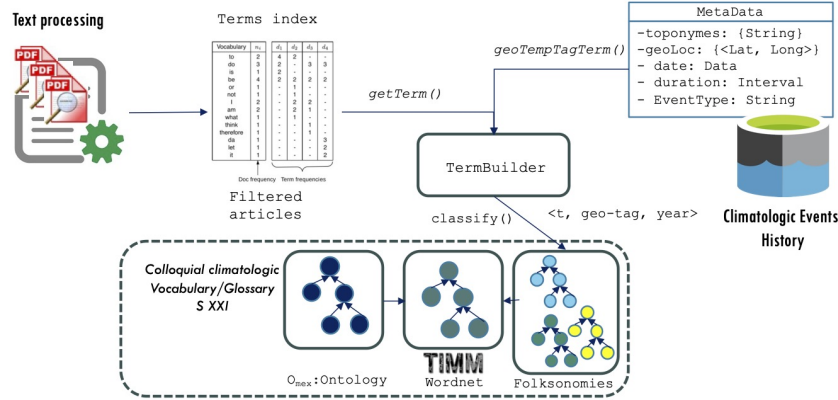


Fig. 9. Collecting colloquial vocabulary

¹⁰ In Mexico, a storm is called a "chaparrón" and in Uruguay, it is called a "chubasco".

1 4.2. Curating data collections

2 LACLICHEV proposes functions that data scientists may exploit through diverse func-
 3 tionalities. Next, we present the type of functions of LACLICHEV's API (application
 4 programming interface). The implementation of these functions were adapted to the case
 5 of historical newspaper collections:

6
 7 - **Curating data collections** by exploring and processing their content for building the
 8 history of meteorological events possibly related to climate change in the considered
 9 Latin American countries. The functions for processing texts in Spanish are the core of
 10 LACLICHEV. They were coupled with other functions to extract, derive and associate as
 11 much data as possible to articles describing meteorological events.

12 Curation tasks were performed on a collection of textual digital documents with min-
 13 imum associated metadata, particularly those used by digital libraries that own the col-
 14 lections. Each library adopts its metadata schema, but they generally specify the news-
 15 paper's name, the country, the date and number, the number of pages, and the window
 16 time in which it circulated. Libraries export the metadata schema used to describe these
 17 resources and align them to standards used by digital libraries. For example, the editions
 18 of the collection of Uruguayan newspapers were published during the first 10 years of the
 19 XIX century.

20 The curation process generated data structures that provide an abstract representation of
 21 the content of each article describing an event. A frequency matrix integrated the terms
 22 representing the content of articles extracted from the different libraries' collections. This
 23 matrix was sharded and allocated to the servers devoted to interacting with each library.
 24 This strategy implies having queries evaluated on different servers. This distributed query
 25 evaluation was supported by an inverted index that provided information about the doc-
 26 uments containing specific terms and their location. With the inverted index, the cura-
 27 tion process also created initial vocabularies, classified by location (country and city) and
 28 year. These vocabularies classified the terms used to describe climatologic events in the
 29 different Hispanophone countries in LATAM. The temporal dimension allowed to store
 30 information about their evolution.

31
 32 **Querying the event's history** of already tagged events can be done by keyword ori-
 33 ented queries (e.g., locate the most famous events in Mexico during the XVIII century).
 34 Users decide to use some terms that can belong to any of the vocabularies generated in
 35 the pre-processing phase. LACLICHEV applies query rewriting techniques to extended
 36 user-expressed queries with synonyms, subsuming and general terms. The particular char-
 37 acteristic of this task is that the user (i.e., data analyst) can interact and guide the process
 38 according to her/his knowledge and expectations about what she/he expects to explore
 39 and search. The first result of this process, based on a "queries as answers approach",
 40 is a set of queries that can potentially provide the largest number of results stemming
 41 from the collections of the different libraries. The details of the approach we proposed for
 42 LACLICHEV is detailed in the following section.

43 - **Analytics operations and analysis results** are generally presented within maps (e.g.,
 44 how did rainy periods evolved in the region?). In the current version of LACLICHEV
 45 analytics queries cannot be expressed in the frontend. They are implemented manually
 46 through notebooks running on top of the event history. The analytics queries concerning

1 aggregative queries on the event’s history, for example, the number of events happening in
 2 a country within a specific time window. The average wind speed and millimetres of water
 3 per hour deduced for events regarding rainfalls and hurricanes in Montevideo. Classify-
 4 ing the terms used for describing specific types of events. The event history is a curated
 5 and clean data collection on top of which other analytics models can be applied for dis-
 6 covering knowledge. This characteristics open analytics perspectives for future uses of
 7 LACLICHEV. For example, applying supervised learning for analysing newspapers ar-
 8 ticles and determining whether they describe meteorological events. This can allow to
 9 semi-automatise the curation process and enhance it with a recommendation system.

10
 11 - **Managing vocabularies**, adding terms, guiding their classification and studying the lin-
 12 guistic connections between the terms used in the different countries. The vocabularies in
 13 the current version of LACLICHEV are implemented as RDF ontologies, and it relies on
 14 SPARQL mechanisms for querying them. This version mainly addresses the construction
 15 of vocabularies and their maintenance as new terms are identified in events’ descriptions.

16 Next subsections describe exploration techniques implemented for meteorological
 17 events in the history built through the newspaper articles in the Latin American coun-
 18 tries we used.

19 4.3. Query rewriting

20 *Queries-as-answers exploration.* Data analysts can express queries that can potentially
 21 explore historical newspaper content to find articles describing meteorological events.
 22 The aim is to have a good balance between precision and recall despite the ambiguity of
 23 the language (Spanish variations in naming meteorological events). The domain experts
 24 must express “clever” queries that can exploit the collections to achieve this goal.

25 Queries can be initially conjunctive and disjunctive expressions combining terms cho-
 26 sen from the built-in vocabularies or not. Then, queries are rewritten in an expression tree
 27 where nodes are conjunction and disjunction operators and leaves are terms, according to
 28 an input query expressed as a conjunction and disjunction of terms potentially belonging
 29 to a meteorological vocabulary.

30 Our approach for rewriting queries is based on a “queries-as-answers” process. This
 31 technique rewrites user queries into queries that can produce more precise results accord-
 32 ing to the explored dataset content. Queries as answers proposed by LACLICHEV consist
 33 of a list of frequently used queries. Thus, we focus on the following aspects:

34 *Extending query alternatives using hypernyms and synonyms.* An initial conjunctive or
 35 disjunctive query is rewritten by extending it with general and more specific terms, syn-
 36 onyms, etc. The terms used to express the query are colloquial vocabulary for denoting
 37 meteorological events. The rewriting process can be automatic or interactive, in which
 38 case the system proposes alternatives, and the user can validate the proposed terms. For
 39 example, if the query is “heavy storms”, the query can be completed by adding “heavy
 40 stormtrooper”, “heavy storm dust”. It can also be rewritten with synonyms for the adjec-
 41 tive heavy. In that case, it creates a combinatorial set of rewritten queries.

42 Note that the colloquial vocabulary stems from the articles of the curated newspapers.
 43 As they are curated, the terms used in the articles feed a vocabulary that is first organised
 44 in the frequency matrix produced when texts are processed as part of the curation process.

1 Then we use Wordnet¹¹ to look for associated terms and synonyms that help address
 2 concepts used in different Spanish-speaking countries. We do not translate the query
 3 terms to other languages because our digital data collections contain Spanish newspa-
 4 pers. LACLICHEV allows equivalent terms searching to morph a query. For a new term,
 5 LACLICHEV generates a node with the operator and then connects the initial term with
 6 the equivalent terms in a disjunctive expression subtree. Thereby, more general terms are
 7 collected and related to the initial term with these terms in a conjunctive expression sub-
 8 tree. The result is a new expression tree corresponding to an extended query Q_{ExT} . The
 9 query morphing algorithm behind LACLICHEV is described in [35].

10 *Extending query alternatives using cultural terms.* Use local vocabularies for generat-
 11 ing new query expression trees that substitute the terms used in Q'_{ExTi} with equivalent
 12 terms used in a target country (e.g., blizzard instead of a heavy storm). This will result in
 13 transformed expression trees each one using the terms of a country ($Q''_{ExT1} \dots Q''_{ExTj}$)
 14 [38].

15 We call metaphorically “folksonomies” a series of vocabularies created by process-
 16 ing newspaper articles “local” vocabulary. We make and feed each vocabulary according
 17 to the country of origin of the processed newspaper article. This lets us extract the vo-
 18 cabulary used during the XVIII and XIX centuries for describing meteorological events
 19 in Latin American countries (i.e. Mexico, Colombia, Ecuador, and Uruguay). Using this
 20 information, LACLICHEV can answer the following queries: *How have terms used to*
 21 *describe meteorological events changed between XIX-XX c.? Which are standard terms*
 22 *used to describe meteorological events across Latin American countries? Which is the*
 23 *distance between terms used in XIX-XX c.? Which are the most popular terms used in XIX*
 24 *c. for describing meteorological events?*

25 *Defining filters using knowledge domain.* We also use domain knowledge for rewriting
 26 the queries. We have a knowledge base provided by domain experts that contains some
 27 meteorological event rules. For example, rules state that in the presence of a heavy storm:
 28 R1. the wind speed is higher than 118 km/h; R2. the rivers can grow and produce big
 29 waves; R3. there are rains between 2,5 7,5 mm/h; R4. the range of surface that can be
 30 reached by a 100 km wind speed storm is of 1000 km.

31 Our approach uses this information to generate possible queries that help the domain
 32 expert better precise her/his query or define several queries that can represent what she/he
 33 is looking for. For example, the previous initial query “ Q_1 : heavy storm” is rewritten
 34 into new additional queries: “ Q_{11} : heavy storm *or storm with wind speed* > 100 km”
 35 (using R1). “ Q_{12} : storms with 100 km speed that reached Mexico City” (using R2 and
 36 knowing the initial point and geographic information). “ Q_{13} : storms touching villages
 37 500 km around Mexico city happening in the same period” (R4). Instead of having a long
 38 query expression, our approach proposes queries that the domain expert can choose and
 39 combine. Note that the system first generates queries, not answers. The answer to a query
 40 is a family of possible queries with some associated samples. The user can then choose
 41 those queries that she wants to execute.

42 A climate glossary associates a term referring to a meteorological event with terms of
 43 the LODE ontology¹²). LODE is an ontology for historical publishing events as Linked

¹¹ <http://timmm.uaen.es/recursos/spanish-wordnet-3-0/>

¹² <http://linkedevents.org/ontology/>

1 Data and physical variables describing events. This information generates new queries,
 2 which help users discover more details about historical meteorological events.

3 Using the climate glossary for transforming Q_{ExT} into queries with terms that can
 4 serve as filters. There are variables concerning meteorology concepts in the glossary, like
 5 wind speed, rain volume/hour, and the water level of seas, rivers, and lakes. Other vari-
 6 ables involve geographic aspects, like the location of an event and the scope of land it
 7 reaches. Finally, other variables concern damages caused by a climate event with specific
 8 physical and geographic characteristics. These different options generate queries combin-
 9 ing variables of the same group and different groups. For example, “heavy storms with
 10 winds higher than 150 km/h”, “heavy storms with rains higher the 10 mm per square me-
 11 tre”, and “heavy storms with rivers’ overflow”. The result is a set of queries $Q'_{ExT1} \dots$
 12 Q'_{ExTj} .

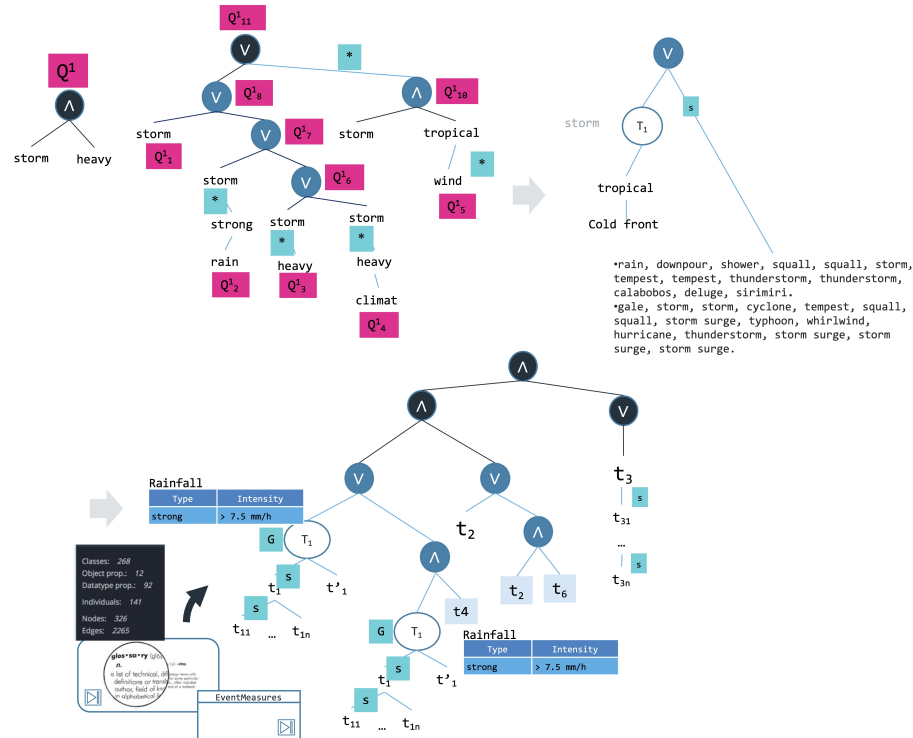


Fig. 10. Queries as answers example

13 Figure 10 shows an example of the general principle of the queries as answers ap-
 14 proach adopted by LACLICHEV. The system rewrites an initial conjunctive query heavy,
 15 storm adding concepts (i.e., terms) related to the terms “storm” and “heavy”. The figure
 16 only shows the rewriting process of the term “storm” for pedagogic reasons. Then in a
 17 second round, the system rewrites the query adding synonyms of the terms, as shown in

the upper right side of the figure. Finally, the query is rewritten according to the rules stated in the glossary. LACLICHEV performs a ranking process for the rewritten queries according to the coverage of their potential answer. The queries with the biggest coverage (those that include the largest subset of events in the history database). The algorithms to estimate the coverage of a documents collection are proposed in [35].

4.4. Evaluating queries

The evaluation process of the query is performed first on top of the curated event store. The result is a set of items (events) that answer, to some extent, the query. We also started to generate maps depicting the events reported in the history [6, 39].

Analytics queries LACLICHEV provides and maintains the meteorological event's history on top of which users can visualize information and perform analytical tasks. For example, LACLICHEV can answer spatio-temporal queries like:

- Q_1 *Locate meteorological events in the XVII century,*
- Q_2 *Enumerate and locate the most famous events in the region or in a specific country,*
and
- Q_3 *Create a heat map of events in Latin America in the last years of the XIX century.*

The objective is to answer analytic queries that imply aggregating information stored in the event's history. For example, *How did rainy time evolve in time in the region?, In which way was climate different between XVII and XIX centuries? How did vocabulary evolve from colloquial to scientific and standardized in the XX century?*

4.5. Scope and limitations

LACLICHEV is running its first version; we expect to enrich the number of digital newspapers digitalised in the libraries. These new items will imply a new curation process that will improve the event history in two directions. First, more articles will describe the already curated events; this will complete the information stored in the history. Second, with more events, we will further test and enhance the analytics queries that require to have a specific volume of data to generate representative maps and analyses about the meteorological events that happened in the past. In future versions, and with more curated events, LACLICHEV is willing to answer prediction queries like *Could it have been possible to predict the evolution of climate behaviour from the data in XVIII and XIX centuries?.* This query requires collecting, curating, and preparing more newspaper articles and other complementary data. However, it concerns future work.

Another limitation of the current LACLICHEV is that it does not provide the adapted mechanisms for exploring the linguistic aspects of the vocabulary. It gathers the terms and organises them in a mesh data structure. Still, it does not provide tools for curating the languages and allowing an analytics exploration of their use of meteorological across countries and time.

1 5. Related work

2 Historical analysis of climate behaviour can explain climatologic phenomena and Earth's
 3 climate behaviour. There exist several scientific efforts to study the history of climate
 4 change. The *Climate of the Past* [1], for example, is an international scientific journal
 5 dedicated to the publication and discussion of research articles, short communications,
 6 and review papers on Earth's climate history. The journal covers all temporal scales of
 7 climate change and variability, from geological time to multidecade studies of the last
 8 century. The Government of Canada provides access to historical observations on climate
 9 in Canada starting from 1840 [2]. However, these data collections are disconnected and
 10 use different reference variables and observation criteria. They are very heterogeneous
 11 and tight to their region. This ad-hoc characteristic is why data curation and exploration
 12 processes are essential to extract knowledge that can be digitally analyzed and correlated.

13 Several domains address aspects that converge in our work, particularly those with
 14 certain originality, like data exploration techniques, geographic information retrieval and
 15 visualization. The following lines summarise the methods and approaches related to those
 16 proposed for LACLICHEV.

17 5.1. Data curation

18 Data curation [14, 33] is the art of processing data to maintain it and improve its interest,
 19 value, and usefulness through its lifecycle, i.e. improving the quality of the data. There-
 20 fore, it implies (i) discovering data collections of interest; (ii) cleaning and transforming
 21 new data; (iii) semantically integrating it with other local data collections; and (iv) dedu-
 22 plicating the resulting composites if required. Data curation provides the methodological
 23 and technological data management support to address data quality issues, maximizing
 24 the usability of the data for analytics and knowledge discovery purposes.

25 Existing commercial and academic systems provide solutions for curating data [36, 29,
 26 40]. They provide operations for modelling and extracting metadata from raw data col-
 27 lections, and they provide tools for exploring them. Prominent commercial examples are
 28 Apache Atlas ¹³ and SolR ¹⁴. Apache Atlas is a framework for governance and manage-
 29 ment of metadata. It offers curation functions for metadata typing and classification, data
 30 lineage, and exploration functions such as data source search. SolR is a document index-
 31 ing system including XML files, comma-separated value (CSV) files, data extracted from
 32 tables in a database, and files in standard file formats such as Microsoft Word or PDF.
 33 Indexing documents can be used as an essential general-purpose curation operation. Its
 34 major exploration features include full-text search, hit highlighting and faceted search.
 35 Other solutions are built on top of these tools for providing end-to-end general-purpose
 36 systems for curating and exploring data, for example, ATLAN ¹⁵.

37 5.2. Data exploration

38 The emergence of the notion of data exploration provides different perspectives of the
 39 data and tools for helping data scientists choose and compound datasets adapted for target

¹³ <https://atlas.apache.org>

¹⁴ <https://solr.apache.org>

¹⁵ <https://atlan.com>

1 experiments [23, 5]. The tools [17] include functions like “data grooming” [27], which
 2 denotes transforming raw data into analyzable data with various data structures. Other
 3 approaches [24] focus on transforming human-readable data into machine-readable data
 4 considering inconsistencies in data formatting given that they are produced under different
 5 conditions. The idea is to exhibit processes, digital spaces, and systems that host datasets
 6 and provide them with access to understand the conditions in which data are processed.

7 *Data grooming* denotes transforming raw data into analysable data with various data
 8 structures. Multi-scale queries propose to split a query into multiple queries executed on
 9 different database fragments and then perform a union of those queries. This allows scal-
 10 ing the query size as the user gets more confident in her query. Result set post-processing
 11 and query morphing go on the premise that the user probably does not need the exact
 12 answer to a query. Result set post-processing assumes an array of simple statistical in-
 13 formation such as min, max, and mean to be more helpful, especially on massive data
 14 sets. Query morphing assumes queries can be wrongly formulated. Query morphing still
 15 focuses on answering the query given by the user but will also use a small portion of
 16 resources in searching data around the original query.

17 *Query morphing.* Another trend regarding data exploration is to tackle the lack of knowl-
 18 edge a user may have on the dataset. Query morphing and queries as an answer are
 19 rewriting techniques that compute alternative queries (e.g. adding terms) that can poten-
 20 tially better explore a dataset than an initial query. Approaches such as interactive query
 21 expansion (IQE) [30, 8, 19] have shown the importance of data consumers in the data ex-
 22 ploration process. Users’ intention helps navigate the unknown data, formulate queries
 23 and find the desired information. In most occurrences, user feedback acts as vital rele-
 24 vance criteria for the following query search iteration. The key challenge is identifying
 25 bad queries using statistical information or massive scientific databases and identifying
 26 interesting queries to return. Identifying bad queries can be done using a list of frequently
 27 used queries and returning them based on user feedback.

28 SVD/PCA [15] is probably the most known algorithm for exploring data sets. It is used
 29 to reduce high-dimensional data represented as a matrix. From a practical perspective, it
 30 searches for the combination of weighted attributes that expresses the most information,
 31 allowing data analysts to work with the more useful 2 or 3-dimensional graphs. From a
 32 geometric perspective, these techniques search for the vectors with the highest variance
 33 and then express the original matrix according to this new system of dimensions. Using
 34 Eigenvalues makes it possible to estimate the amount of information in each dimension
 35 [12].

36 *Visualization and summarization* are essential to understand the data and maintain it. The
 37 field of visual analytics seeks to provide people with better and more effective ways to un-
 38 derstand and analyze these large datasets while also enabling them to act upon their find-
 39 ings immediately [22]. Visual analytics provides technology [26, 28] that combines hu-
 40 man and electronic data processing strengths. Structured query languages and the graph-
 41 ical interface developed over the top are the standard procedure for accessing data in
 42 a database. Many tools exist to perform data visualization with web visualization tools
 43 such as D3.js or other tools such as Matlab [20] or R programming language [16]. One of
 44 the most critical steps of these tools is to let the data analyst move from confirmatory data

1 analysis (using charts and other visual representations to present results) to exploratory
 2 data analysis (interacting with the data/results). This has led to visual data exploration and
 3 visual data mining [7].

4 **5.3. Text processing in newspaper articles**

5 The discovery of knowledge from large-scale text data or semi-structured data is a dif-
 6 ficult task that can be addressed with text mining techniques. These techniques extract
 7 valuable information to fulfil a user information need. The textual documents available
 8 in unstructured and semi-structured forms can be medical, financial, market, scientific,
 9 and other documents. Text mining applies a quantitative approach to analyse a massive
 10 amount of textual data and tries to solve the information overload problem.

11 The combination of transformers and self-supervised pretraining has been responsible
 12 for a paradigm shift in NLP, information retrieval (IR), and beyond [25]. The approach
 13 in [21] extracts target categories, each including many topics. The method extracts word
 14 tokens referring to topics related to a specific category. The frequency of word tokens in
 15 documents impacts the document's weight calculated using a numerical statistic of term
 16 frequency-inverse document frequency (TF-IDF). The proposed approach uses the title,
 17 abstract, and keywords of the paper and the categories of topics to perform a classification
 18 process. The documents are classified and clustered into the primary categories based on
 19 the highest cosine similarity measure between category weight and documents' weights.

20 The work proposed in [3] discusses the challenge of processing and analysing his-
 21 torical manuscripts. Authors investigate how deep learning models detect and recognise
 22 handwritten words in Spanish American notary records. For dealing with natural language
 23 (ancient Spanish), professional historians prepared a labelled dataset of 26,482 Spanish
 24 words employed in the experiments. The paper [31] proposes a tool that uses raw Spanish
 25 text and Spanish event coders for analysing political news articles. The work combines
 26 natural language processing techniques, including deep learning and encoders, with the
 27 knowledge represented in ontologies to support the automated coding process for Spanish
 28 texts.

29 **5.4. Geographic information retrieval**

30 Within GIScience domains, some approaches have developed. [10] and [9] combined
 31 methods from Geographic Information Retrieval (GIR) and geovisual analytics to obtain
 32 new insights from a digital dictionary about the history of Switzerland. In addition, the au-
 33 thors include sentiment analyses to assess how (historical) places were referred to in texts
 34 over time and provide ways to access and explore spatio-temporal information contained
 35 in many text archives. [34] described a method to supplement existing records of land-
 36 slides in Great Britain by searching an electronic archive of regional newspapers. More-
 37 over, the authors construct a Boolean search criterion by experimenting with landslide
 38 terminology for four training periods. It allowed the discovery of some spatio-temporal
 39 patterns of additional landslides identified in newspaper articles. [41] presented a text-
 40 mining program that extracted keywords related to floods' geographic location, date, and
 41 damages from newspaper analyses of flash floods in Fujairah, UAE, from 2000 to 2018.
 42 Furthermore, this work performed geocoding and validating flood-prone areas generated
 43 through Geographic Information System (GIS) modeling.

1 5.5. Discussion

2 Any query and analysis must be based on a good understanding of the available data
3 collections because the way they are combined and analyzed impacts the quality and
4 accuracy of the results.

5 Existing solutions are not delivered in integrated environments that data analysts can
6 comfortably use to explore data collections. The technical effort is still necessary to com-
7 bine several tools to explore and process datasets and go from raw independent data sets
8 to knowledge, for example, on climate change. Therefore, our research aims to tailor a
9 data exploration environment to help explore digital data collections using a human-in-
10 the-loop approach. In existing solutions, data analysts cannot comfortably explore data
11 collections and design analytics settings, particularly in cases where documents and ques-
12 tions combine scientific observations with empirical observations, like in the case of me-
13 teorological events described empirically in the past.

14 The current version of LACLICHEV did not explore the linguistic aspect, with original
15 or more advanced methods studying texts and combining present and past observations to
16 try to derive conclusions, for example, about climate change.

17 A technical effort is still necessary to combine several tools to explore and process datasets
18 and go from raw independent data sets to knowledge, understanding and prediction, for
19 example, on climate change. Therefore, LACLICHEV aimed to tailor a data exploration
20 environment that could help explore digital datasets using a human-in-the-loop approach.

21 Regarding the qualitative assessment of LACLICHEV, we have not run user experi-
22 ence testing to collect feedback and user experience, and we might perform such testing
23 in the future. For the time being, we focus on the analytics such as correlating different
24 descriptions of the “same” event from articles in various newspapers, the location of me-
25 teorological events in old maps and their correlation with modern maps. We are working
26 on creating historical cartography of meteorological events that can be confronted with
27 contemporary perceptions of such events.

28 6. Conclusion and future work

29 The democratisation of access to data collections opens possibilities for exploring con-
30 tent produced over the years and extracting knowledge that can contribute to understand-
31 ing critical phenomena like climate change. Rather than directly querying collections for
32 searching documents or performing data analytics operations (statistics, correlations), the
33 objective is to let data scientists understand the content of the collections and then decide
34 what kind of queries to ask. Data exploration is a complex and recurrent process that in-
35 cludes calibrating a querying strategy (defining queries as answers) that can increase the
36 scope of content that can be retrieved and possibly analysed to extract evidence around
37 hypotheses or claims. This new paradigm calls for data curation strategies that are well
38 adapted to describe the content of collections with the right metadata and abstractions.

39 Our work contributes to data curation and exploration adapted for Spanish textual
40 content within digital newspaper collections. Using well-known information retrieval and
41 analytics techniques, we developed a data exploration environment named LACLICHEV
42 that provides tools for understanding the content of collections. We used digital news-
43 paper collections for applying such techniques for building and analyzing the history of

1 meteorological events possibly related to climate change in Mexico, Colombia, Ecuador,
 2 and Uruguay. The work reported here is the first step toward this ambitious challenge. We
 3 continue enriching data collections, developing and testing solutions for generating and
 4 sharing step by step this history.

5 References

- 6 1. Climate of the past: An interactive open-access journal of the european geosciences union.
 7 <http://www.climate-of-the-past.net>, european Geosciences Union, Accessed: 2021-04-23
- 8 2. Historical climate data. <http://climate.weather.gc.ca>, government of Canada, Accessed: 2021-
 9 04-23
- 10 3. Alrasheed, N., Prasanna, S., Rowland, R., Rao, P., Grieco, V., Wasserman, M.: Evaluation of
 11 deep learning techniques for content extraction in spanish colonial notary records. In: Proceed-
 12 ings of the 3rd Workshop on Structuring and Understanding of Multimedia heritAge Contents.
 13 pp. 23–30 (2021)
- 14 4. Amavi, J., Ferrari, M.H., Hiot, N.: Natural language querying system through entity enrich-
 15 ment. In: ADBIS, TPD and EDA 2020 Common Workshops and Doctoral Consortium. pp.
 16 36–48. Springer (2020)
- 17 5. Amer-Yahia, S., Koutrika, G., Bastian, F., Belmpas, T., Braschler, M., Brunner, U., Calvanese,
 18 D., Fabricius, M., Gkini, O., Kosten, C., et al.: Inode: building an end-to-end data exploration
 19 system in practice [extended vision]. arXiv preprint arXiv:2104.04194 (2021)
- 20 6. Ballari, D.: Visualizar datos georreferenciados, <http://rpubs.com/daniballari/ipgh-visgeoref>
- 21 7. Battle, L., Stonebraker, M., Chang, R.: Dynamic reduction of query result sets for interactive
 22 visualizaton. In: 2013 IEEE International Conference on Big Data. pp. 1–8 (Oct 2013)
- 23 8. Belkin, N.J.: Some (what) grand challenges for information retrieval. In: ACM SIGIR Forum.
 24 vol. 42, pp. 47–54. ACM New York, NY, USA (2008)
- 25 9. Bruggmann, A., Fabrikant, S.I.: How does giscience support spatio-temporal information
 26 search in the humanities? *Spatial Cognition & Computation* 16(4), 255–271 (2016)
- 27 10. Bruggmann, A., Fabrikant, S.I., Janowicz, K., Adams, B., McKenzie, G., Kauppinen, T.: Spa-
 28 tializing a digital text archive about history. In: CEUR Workshop Proceedings. pp. 6–14. No.
 29 1273, CEUR-WS (2014)
- 30 11. Carvalho, D.A.S., Souza Neto, P.A., Ghedira-Guegan, C., Bennani, N., Vargas-Solar, G.:
 31 Rhone: A quality-based query rewriting algorithm for data integration. In: New Trends in
 32 Databases and Information Systems. pp. 80–87. Springer International Publishing, Cham
 33 (2016)
- 34 12. Chawla, S., Zheng, Y., Hu, J.: Inferring the root cause in road traffic anomalies. In: 2012 IEEE
 35 12th International Conference on Data Mining. pp. 141–150 (Dec 2012)
- 36 13. Comesaña, D., Vilches-Blázquez, L.M.: A study of the latin american newspapers from xix-xx
 37 centuries with a focus on meteorological events. *Revista de historia de América* (156), 29–59
 38 (2019)
- 39 14. Curry, E., Freitas, A., O’Riáin, S.: The role of community-driven data curation for enterprises.
 40 In: Linking enterprise data, pp. 25–47. Springer (2010)
- 41 15. Feldman, D., Schmidt, M., Sohler, C.: Turning big data into tiny data: Constant-size coresets
 42 for k-means, pca and projective clustering. In: Proceedings of the twenty-fourth annual ACM-
 43 SIAM symposium on Discrete algorithms. pp. 1434–1453. SIAM (2013)
- 44 16. Foundation, T.R.: The r project for statistical computing. (2018), <https://www.r-project.org>
- 45 17. Gewers, F.L., Ferreira, G.R., Arruda, H.F.D., Silva, F.N., Comin, C.H., Amancio, D.R., Costa,
 46 L.d.F.: Principal component analysis: A natural approach to data exploration. *ACM Computing*
 47 *Surveys (CSUR)* 54(4), 1–34 (2021)

- 1 18. Goodchild, M.F.: Giscience, geography, form, and process. *Annals of the Association of Amer-*
2 *ican Geographers* 94(4), 709–714 (2004)
- 3 19. Goswami, P., Gaussier, E., Amini, M.R.: Exploring the space of information retrieval term
4 scoring functions. *Information Processing & Management* 53(2), 454–472 (2017)
- 5 20. Inc., T.M.: Matlab (2018), <https://www.mathworks.com/products/matlab.html>
- 6 21. Jalal, A.A., Ali, B.H.: Text documents clustering using data mining techniques. *International*
7 *Journal of Electrical & Computer Engineering* (2088-8708) 11(1) (2021)
- 8 22. Keim, D.A.: Visual exploration of large data sets. *Commun. ACM* 44(8), 38–44 (Aug 2001),
9 <http://doi.acm.org/10.1145/381641.381656>
- 10 23. Kersten, M.L., Idreos, S., Manegold, S., Liarou, E., et al.: The researcher’s guide to the data
11 deluge: Querying a scientific database in just a few seconds. *PVLDB Challenges and Visions*
12 3(3) (2011)
- 13 24. Kumar, M.S., Rajeshwari, J., Rajasekhar, N.: Exploration on content-based image retrieval
14 methods. In: *Pervasive Computing and Social Networking*, pp. 51–62. Springer (2022)
- 15 25. Lin, J., Nogueira, R., Yates, A.: Pretrained transformers for text ranking: Bert and beyond.
16 *Synthesis Lectures on Human Language Technologies* 14(4), 1–325 (2021)
- 17 26. Liu, X., Alharbi, M., Best, J., Chen, J., Diehl, A., Firat, E., Rees, D., Wang, Q., Laramée, R.S.:
18 Visualization resources: A starting point. In: *2021 25th International Conference Information*
19 *Visualisation (IV)*. pp. 160–169. IEEE (2021)
- 20 27. Liu, Y.: Exploring a corpus-based approach to assessing interpreting quality. In: *Testing and*
21 *Assessment of Interpreting*, pp. 159–178. Springer (2021)
- 22 28. Mohammed, L.T., AlHabshy, A.A., ElDahshan, K.A.: Big data visualization: A survey. In:
23 *2022 International Congress on Human-Computer Interaction, Optimization and Robotic Ap-*
24 *plications (HORA)*. pp. 1–12. IEEE (2022)
- 25 29. Rehm, G., Bourgonje, P., Hegele, S., Kintzel, F., Schneider, J.M., Ostendorff, M., Zaczynska,
26 K., Berger, A., Grill, S., Räuchle, S., et al.: Qurator: innovative technologies for content and
27 data curation. *arXiv preprint arXiv:2004.12195* (2020)
- 28 30. Ruthven, I.: Re-examining the potential effectiveness of interactive query expansion. In: *Pro-*
29 *ceedings of the 26th annual international ACM SIGIR conference on Research and develop-*
30 *ment in informaion retrieval*. pp. 213–220 (2003)
- 31 31. Salam, S., Khan, L., El-Ghamry, A., Brandt, P., Holmes, J., D’Orazio, V., Osorio, J.: Auto-
32 matic event coding framework for spanish political news articles. In: *2020 IEEE 6th Intl Con-*
33 *ference on Big Data Security on Cloud (BigDataSecurity), IEEE Intl Conference on High Per-*
34 *formance and Smart Computing,(HPSC) and IEEE Intl Conference on Intelligent Data and*
35 *Security (IDS)*. pp. 246–253. IEEE (2020)
- 36 32. Sodré, A.P., Floriano, L.E.M., Magalhaes, D., Aguiar, C.D., Pozo, A., Hara, C.S.: Compar-
37 ing alternative storage models for words extracted from legal texts. In: *Anais Estendidos do*
38 *XXXVI Simpósio Brasileiro de Bancos de Dados*. pp. 36–42. SBC (2021)
- 39 33. Stonebraker, M., Bruckner, D., Ilyas, I.F., Beskales, G., Cherniack, M., Zdonik, S.B., Pagan,
40 A., Xu, S.: Data curation at scale: the data tamer system. In: *Cidr*. vol. 2013 (2013)
- 41 34. Taylor, F.E., Malamud, B.D., Freeborough, K., Demeritt, D.: Enriching great britain’s national
42 landslide database by searching newspaper archives. *Geomorphology* 249, 52–68 (2015)
- 43 35. Vargas-Solar, G., Farokhnejad, M., Espinosa-Oviedo, J.: Towards human-in-the-loop based
44 query rewriting for exploring datasets. In: *Proceedings of the Workshops of the EDBT/ICDT*
45 *2021 Joint Conference* (2021)
- 46 36. Vargas-Solar, G., Kemp, G., Hernández-Gallegos, I., Espinosa-Oviedo, J., Da Silva, C.F.,
47 Ghodous, P.: Demonstrating data collections curation and exploration with curare. In:
48 *EDBT/ICDT Conference 2019*. p. 4 (2019)
- 49 37. Vargas-Solar, G., Zechinelli-Martini, J., Espinosa-Oviedo, J.A., Vilches-Blázquez, L.M.:
50 LACLICHEV: exploring the history of climate change in latin america within newspapers
51 digital collections. In: *Bellatreche, L., Dumas, M., Karras, P., Matulevicius, R., Awad, A.,*

- 1 Weidlich, M., Ivanovic, M., Hartig, O. (eds.) New Trends in Database and Information Sys-
2 tems - ADBIS 2021 Short Papers, Doctoral Consortium and Workshops: DOING, SIMPDA,
3 MADEISD, MegaData, CAoNS, Tartu, Estonia, August 24-26, 2021, Proceedings. Commu-
4 nications in Computer and Information Science, vol. 1450, pp. 121–132. Springer (2021),
5 https://doi.org/10.1007/978-3-030-85082-1_11
- 6 38. Vargas-Solar, G., Zechinelli-Martini, J.L., Espinosa-Oviedo, J.A.: Computing query sets for
7 better exploring raw data collections. In: 2018 13th International Workshop on Semantic and
8 Social Media Adaptation and Personalization (SMAP). pp. 99–104. IEEE (2018)
- 9 39. Vilches-Blázquez, L.M., Ballari, D.: Unveiling the diversity of spatial data infrastructures in
10 latin america: evidence from an exploratory inquiry. *Cartography and Geographic Information*
11 *Science* 47(6), 508–523 (2020)
- 12 40. Visengeriyeva, L., Abedjan, Z.: Anatomy of metadata for data curation. *Journal of Data and*
13 *Information Quality (JDIQ)* 12(3), 1–30 (2020)
- 14 41. Yagoub, M., Alsereidi, A.A., Mohamed, E.A., Periyasamy, P., Alameri, R., Aldarmaki, S., Al-
15 hashmi, Y.: Newspapers as a validation proxy for gis modeling in fujairah, united arab emirates:
16 identifying flood-prone areas. *Natural Hazards* 104(1), 111–141 (2020)