

CHIS: A Big Data Infrastructure to Manage Digital Cultural Items

Aniello Castiglione^{a,*}, Francesco Colace^b, Vincenzo Moscato^c, Francesco Palmieri^a

^aUniversity of Salerno, Dipartimento di Informatica, Via Giovanni Paolo II, 132 - 84084 Fisciano (SA), Italy

^bUniversity of Salerno, Dipartimento di Ingegneria Industriale, Via Giovanni Paolo II, 132 - 84084 Fisciano (SA), Italy

^cUniversity of Naples, Dip. di Ingegneria Elettrica e delle Tecnologie dell'Informazione, Via Claudio 21, - 80125 Napoli, Italy

Abstract

In this paper, we describe CHIS (*Cultural Heritage Information System*), a big data infrastructure that can be used to query, browse, analyze and process digital contents related to cultural heritage from a set of heterogeneous and distributed repositories. CHIS is characterized by the following technical features: capability to gather information from distributed and heterogeneous data sources (e.g., Sensor Networks, Social Media Networks, Digital Libraries and Archives, Multimedia Collections, Web Data Services, etc.); advanced data management techniques and technologies; ability to provide useful and personalized data to users based on their preferences and context; advanced information retrieval facilities, data analytics and other utilities/services, according to the SOA paradigm. By means of a set of ad-hoc APIs, and value-added data processing and analytics services, our system can support several applications: mobile multimedia guides for cultural environments, web portals to promote the cultural heritage of a given organization, multimedia recommender and storytelling systems and so on. We discuss the main ideas that characterize the system, showing its use for several applications.

Keywords: Big Data, Cultural Heritage, Resource Management

1. Introduction

Italy boasts one of the largest and most priceless cultural heritage in the world. This invaluable resource can be opportunely protected, preserved and promoted by “embedding” it in the digital ecosystem of a *Smart City*, where economic, tourist, recreational and logistic aspects have to be considered all together.

As highlighted by recent publications and reiterated by European projects discussing Smart Cities [1, 2, 3, 4], the adoption of *Future Internet* technologies, particularly the paradigms of the *Internet of Things* and *Internet of Services*, now represents the “de facto standard” in the design and implementation of IT platforms that can provide effective support to the “smartness” of a city. In such a context, it is possible to design *context-aware* services that take into account both the surrounding environment, whose state is captured by sensors, and the characteristics of the users. These services are then all accessible through a *Cloud Computing* environment [5]. In addition, according to the vision of *participatory sensing*, [6] mobile devices of the latest generations (e.g., smart-phones, tablets, etc.) form an interactive network that allows users to access, analyze and share information and knowledge with an “active” role. Eventually, following the recommendations of the *W3C Semantic Web* framework [7, 8], in order to allow the enormous amount of data collected (*Big Data*[9, 10]) in a smart environment to be used by different applications, the data must be properly processed and stored in the form of *linked open data*, in order to facilitate both access and semantic processing.

A number of proposals, which focus on how the discussed technological solutions should be applied to the cultural domain, has been already presented for the Italian heritage.

*Corresponding author: Aniello Castiglione, Department of Computer Science, University of Salerno, Via Giovanni Paolo II, 132 I-84084 Fisciano (SA), Italy - Phone: +39089969594, Fax: +39089969821, email: castiglione@acm.org, castiglione@ieee.org

Email addresses: castiglione@ieee.org (Aniello Castiglione), fcolace@unisa.it (Francesco Colace), vmoscato@unina.it (Vincenzo Moscato), fpalmieri@unisa.it (Francesco Palmieri)

26 For example, in the Technological District for Cultural Heritage and Activities, Lazio - DTC¹, a platform for the
27 access to cultural heritage has been proposed and is focused on tourists through the realization of digital scenes, virtual
28 reconstructions, augmented reality techniques and mobile applications.

29 Similarly, the Calabria region, in the context of the MESSIAH project², has proposed enabling methods and
30 multi-functional technologies to support cultural heritage identification, monitoring and restoration.

31 Indeed, the problem of evaluation and care of cultural heritage through smart city enabling technologies, is one
32 of the most important issues that has to be taken into consideration, not only at national level, but also within the
33 European scenario.

34 As a first example, Niker³ addresses the problem of the protection of cultural heritage sites from the effects of
35 seismic phenomena, through an integrated approach based on the use of a Knowledge Base; in turn, H-KNOW⁴
36 proposed a similar approach for the restoration of buildings.

37 Other projects, based on *Knowledge Management* and *Knowledge Discovery* techniques aim to promote cultural
38 heritage. Significant examples are: SMARTMUSEUM⁵ which implements a platform for “knowledge exchange”, and
39 PAPHYRUS⁶, which address the problem of “dynamic” knowledge discovery from digital libraries and news archives.

40 Several projects, recently funded by the European community, propose methodologies and the best ways to man-
41 age and organize the knowledge related to cultural heritage.

42 ARIADNE⁷ aims to create an application infrastructure for the management of archaeological data to enable
43 archaeologists and scholars of the ancient world to access the digital archives of European countries online and to be
44 able to use new technologies as an element of the methodology of archaeological research.

45 DC-NET (Digital Cultural heritage NETwork)⁸ aims to develop and strengthen coordination between public re-
46 search programs in the European countries involved in the field of digital cultural heritage, through the development
47 of data infrastructure and services dedicated to virtual community research in digital cultural heritage.

48 MeLa (European Museums and Libraries in/of the Age of Migrations)⁹ aims to develop new strategies to organize
49 multi-inter-trans-cultural conservation, presentation and transmission of knowledge, in ways and forms that reflect the
50 conditions imposed by the migration of people and ideas in the global world. Several other projects dealing with the
51 digitization and use of cultural heritage have been proposed, with their emphasis on digital libraries and museums
52 (EOD, Europeana 1914-1918 HOPE, Linked Heritage, etc..).

53 On the other hand, regarding the challenges in the management of big data, in recent years several projects (EU-
54 DAT, iCordi, EUHIT, smartData, Projectcome, etc..) have been created to propose solutions for the effective and
55 efficient treatment of such data, in particular in the context of scientific data processing. More recent proposals have
56 then investigated the use of big data management techniques and architectures for digital cultural contents, proposing
57 applications [11].

58 Summarizing, the main research problems that can be correlated to the management of Cultural Heritage digital
59 contents in the Smart Cities context are:

- 60 • the adoption of architectural models and standards in the context of Future Internet and Big Data [12, 13];
- 61 • the interfacing and communication with the sensors of a site [14, 15, 16, 17];
- 62 • the access, retrieval, integration and analysis of information from all data sources and the correlation with spatial
63 data [18, 19];
- 64 • the transformation of “captured” data in the form of knowledge and its management [20, 21];
- 65 • the localization and tracking of users on a site [22, 23];

¹<http://www.futouring.it/web/filas/distretto-tecnologico>

²www.culturaeinnovazione.it

³<http://www.niker.eu/>

⁴<http://www.hknow.eu/>

⁵<http://www.smartmuseum.eu/>

⁶<http://www.ict-papyrus.eu/>

⁷<http://www.ariadne-infrastructure.eu/>

⁸<http://www.dc-net.org>

⁹<http://www.mela-project.eu>

- the access to the knowledge based on the user profile, the context and the use of applications [24, 25, 26];
- the analysis of users' opinions from social networks [27].

In this paper, we describe CHIS (*Cultural Heritage Information System*), a system to query, browse and analyze cultural digital contents from a set of distributed and heterogeneous multimedia repositories. In particular, the system prototype has been developed within the DATABENC project ¹⁰.

CHIS is able to manage all the digital contents related to *Cultural Items*. More in details, in our vision each *Cultural Heritage environment* (e.g., museums, archaeological sites, old town centers, etc.) is grounded on a set of cultural *Points of Interest* (PoI), which correspond to one or more cultural items (e.g., specific ruins of an archaeological site, sculptures and/or pictures exhibited within a museum, historical buildings and famous squares in an old town center and so on).

In order to meet variety, velocity and volume of the managed information, CHIS is characterized by the following technical features that are typical of a *Big Data* platform [28, 29]:

- capability to gather information from distributed and heterogeneous data sources (e.g., Sensor Networks, Social Media Networks, Digital Libraries and Archives, Multimedia Collections, Web Data Services, etc.);
- advanced data management techniques and technologies;
- ability to provide useful and personalized data to users based on their preferences and context;
- advanced information retrieval services, data analytics and other utilities, according to the SOA paradigm.

Regarding data sources, CHIS can manage the following kinds of data, providing detailed information about the cultural heritage belonging to some specific geographic areas of Campania:

- data coming from sensor networks deployed in a given cultural environment;
- cultural items' descriptions - coming from open web sources (e.g., Wikipedia) or from several digital libraries and archives;
- multimedia data - video, text, image and audio - associated with the various items of interest, retrieved from open (Social Media Networks as Panoramio, Picasa, YouTube, and Flickr) or private collections;
- social data - e.g., user comments and opinions from common on-line social networks like Facebook, Twitter, etc.;
- web service data - any kind of useful data gathered by means of web services.

Thus, CHIS provides all the necessary retrieval and presentation functionalities to search information of interest and present it to the users in a suitable format and according to their needs.

By means of a set of ad-hoc APIs, and exploiting the provided services, our system can support several applications: mobile multimedia guides for cultural environments, web portals to promote the cultural heritage of a given organization, multimedia recommender and storytelling systems and so on.

The paper is organized as follows. Section 2 describes the proposed data model. Section 3 presents the system architecture and related functionalities with several implementation details. Section 4 reports a possible application of our system to support a mobile multimedia guide for the archaeological site of Paestum. Eventually, Section 5 discusses some conclusions and the future work.

¹⁰The High Technology District for Cultural Heritage (DATABENC) management of the Campania Region, in Italy (www.databenc.it).

2. Data Model

As previously described, the introduced data model for the management of cultural contents relies on the concept of “Cultural Item” (*CI*), that is in turn related to a given “Cultural Environment” (e.g., a museum, an archaeological site, an old town center, an historical building etc.). Examples of *CI*s are specific ruins of an archaeological site, sculptures and/or pictures exhibited within a museum, historical buildings and famous squares in an old town center and so on.

In the Cultural Heritage domain, a *CI* can be opportunely described with respect to a variety of annotation schemata, for example the archaeological view, the architectural perspective, the archivist vision, the historical background, etc., that usually exploit different harvesting sets of “metadata” and possibly domain taxonomies or ontologies [30].

In a simplified way, we consider an *ontology* $O = (V, E)$ as a network of concepts belonging to a given domain of interest, where a node $v \in V$ represents a “concept” and an edge $e \in E$ represents a relationship between two concepts¹¹.

Thus, we define an *annotation schema* and a *semantic annotation*[31] for a cultural item as follows.

Definition 2.1 (Annotation Schema). Given a set of ontologies O , an Annotation Schema is a particular tuple $\lambda_O = (A_1, \dots, A_n, B_1, \dots, B_m)$, where A_1, \dots, A_n are attributes for which $\forall i \in [1, n], \exists O = (V, E) \in O$ s.t. $dom(A_i) \subseteq V$ (i.e., A_i assumes values corresponding to nodes of some ontology), and B_1, \dots, B_m are attributes for which $\forall j \in [1, m], \nexists O = (V, E) \in O$ s.t. $dom(B_j) \subseteq V$ (i.e., B_j does not assume values corresponding to nodes of any ontology).

In other words, the attributes A_1, \dots, A_n are *Ontological Attributes* (OAs) and correspond to concepts that are relevant for the specific domain(s) being modeled. In turn, *Non-Ontological Attributes* B_1, \dots, B_m (NOAs) can contain other useful information, such as measures returned by the sensors attached to cultural items, or multimedia objects (e.g., audio, video, images, texts and 3D models, etc.) characterized by a set of low-level features and other metadata.

In addition, a *CI* may be associated with a specific “Point of Interest” (*POI*), defined by a set of geographic coordinates, and corresponding either to a single point or to a set of lines and more complex polygons of the considered environment.

In particular, we can adopt both “literals” or a set of URIs (*Uniform Resource Identifiers*), that allow to access the related cultural information according to the *Linked Data/Linked Open Data* (LD/LOD) paradigms, as values of the annotation attributes.

Definition 2.2 (Semantic Annotation). Given a set of ontologies O , an annotation schema λ_O and cultural item *CI*, a Semantic Annotation of *CI* is a tuple $\lambda_O(CI) = (a_1, \dots, a_n, b_1, \dots, b_m)$, where $\forall i \in [1, n], a_i \in dom(A_i)$ and $\forall j \in [1, m], b_j \in dom(B_j)$

Using various sets of ontologies/taxonomies and semantic annotations, we can thus describe a cultural item from different points of view supporting several applications.

A large set of relationships can also be instantiated among cultural items and the entire system *Knowledge Base* (KB) can be modeled as a particular *graph*.

Definition 2.3 (Knowledge Base). The Knowledge Base is a graph $G = (C, R)$: each node $c \in C$ can be a cultural item or an ontological attribute (concept), while each edge $r \in R$ represents a relationship derived from a semantic annotation or established between two cultural items.

All possible relationships in the model are opportunely defined “a-priori” and the related meaning can be found in a proper thesaurus.

Figure 1 shows how a portion of knowledge related to the Paestum ruins can be easily represented in our model. In particular, we can see some of the several cultural items a tourist in the Paestum ruins can be interested in: the Archaeological Museum of Paestum (containing evidence from the entire archaeological area and related to a specific POI on the cultural environment map), the Temple of Neptune, and the Tomb of the Diver fresco (the famous burial

¹¹Note that a taxonomy is a particular ontology only containing ISA (concept-subconcept) relationships.

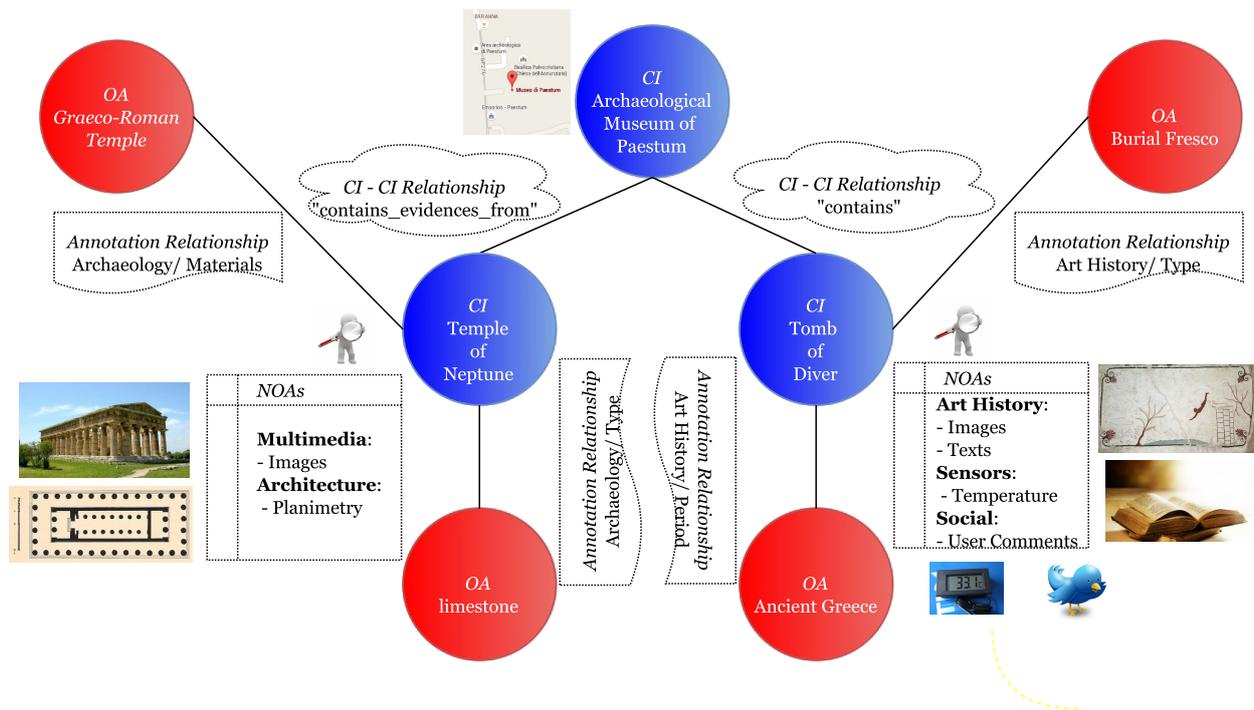


Figure 1: The System Data Model

146 painting belonging to the Greek period in Paestum) that is located within the museum. The non-ontological attributes
 147 of the temple are essentially images related to planimetry (information useful for architects) and some pictures re-
 148 trieved from social media networks; in turn, the fresco's non-ontological attributes represent information related to
 149 Art History (images and text from digital archives), but also comments from on-line social networks and environ-
 150 mental measures collected by the sensors deployed in the museum rooms. Form the other hand, ontological attributes
 151 are then referred to different metadata (type and period for the museum artifact, and type and materials of the ancient
 152 building) whose values can be linked to the nodes of taxonomies and ontologies available for the Cultural Heritage do-
 153 main. Leveraging different annotation schemata and available ontologies, our model allows achieving interoperability
 154 goals[32].

155 The Knowledge Base content can be easily exported in the most used formats (e.g., XML, RDF, OWL) and
 156 according to the most diffused harvesting standards for CH applications (e.g., EDM, Italian ICCD, etc.). On the
 157 other hand, the LD/LOD paradigm permits us to deal with several problems related to data consistency and copyright
 158 constraints in an effective manner: some cultural items descriptions are accessible only using URI, thus the data
 159 management issues are in charge to the related source.

160 3. System Description

161 Our system has to manage a large set of data characterized by significant variety, volume and velocity, representing
 162 the well-known big data features.

163 First of all, we have to deal with the large and heterogeneous amount of information related to the different cultural
 164 items, as an example:

- 165 • annotations and descriptions provided by cultural heritage foundations' archives or by open encyclopedias;
- 166 • multimedia contents – video, text, image and audio – coming from social media networks (e.g., Panoramio,
 167 Picasa, YouTube, and Flickr) and digital libraries;
- 168 • opinions and comments of users from common on-line social networks like Facebook, Twitter, etc.

169 In addition, information about users (in terms of needs, preferences and behaviors) and measures (e.g., humidity
170 and temperature) captured by the different sensors deployed in a cultural environment should be also taken into account
171 together with web service data (any kind of useful data gathered by means of web services such as touristic attractions
172 or accommodations in the same geographic area, or meteorological information and so on) in order to provide “smart”
173 services to final users and applications.

174 The functionalities provided by our system and built on the top of such data are:

- 175 • capability to gather information from distributed and heterogeneous data sources (e.g., Social Media Networks,
176 Digital Libraries and Archives, Multimedia Collections, Web Data Services, Sensor Networks, etc.);
- 177 • advanced data storage and management techniques and technologies;
- 178 • system resources management;
- 179 • advanced information retrieval services (ability to provide useful and personalized data to users/applications
180 based on their preferences and needs), browsing facilities, data analytics and other utilities.

181 By means of a set of ad-hoc APIs, and exploiting the provided services, our system can support several applica-
182 tions: mobile guides for cultural environments, web portals to promote the cultural heritage of a given organization,
183 multimedia recommender and storytelling systems, and so on, also integrated with sensor data, applications for cul-
184 tural heritage preservation and security.

185 Thus, our system provides all the necessary retrieval and presentation functionalities to search information of
186 interest and present it to the users/applications in a suitable format and according to their needs.

187 3.1. Architecture

188 Figure 2 describes at a glance a functional overview of the proposed system that presents a layered architecture
189 typical of a Big Data platform [28, 29], exploiting the related stack of technologies.

190 The *resource management layer* coordinates and optimizes the use of the different resources needed for imple-
191 menting the whole system.

192 In the *data source layer*, each data source is properly “wrapped” in order to extract the information of interest
193 that is then represented as required by the described data model. In particular, each Wrapper is specialized for a
194 particular kind of source (i.e., Social Media Networks, Digital Repositories, Sensor Networks) and must address all the
195 interoperability issues, providing a set of functionalities to access data sources and gather all the desired data, possibly
196 leveraging the available APIs. Data integration problems for heterogeneous data sources are addressed by means of
197 classical schema mapping techniques, record linkage and data fusion techniques [33] or other type of approaches,
198 according to the specific data source. In addition, data stream management problems have to be considered.

199 In the *data storage and management layer*, data are stored in the Knowledge Base in compliance with the above-
200 described data model, and managed also exploiting the LD/LOD paradigm. In addition, specific semantics to be
201 attached to the data is provided using the annotation schemata, including ontologies, vocabularies, taxonomies, etc.
202 related to the Cultural Heritage domain. The KB leverages different Data Repositories realized by advanced data
203 management technologies (e.g., Distributed File Systems, NoSQL and relational systems) and provides a set of basic
204 APIs to read/write data by an Access Method Manager.

205 As a basis for the *data processing layer* our system provides a Query Engine that can be invoked by user appli-
206 cations to search data of interest using information retrieval facilities. In particular, our system supports all the basic
207 functionalities for multimedia and semantic information retrieval by means of proper Information Filters:

- 208 • query by keyword or tags and using specific metadata of the annotation schemata,
- 209 • query by example for cultural items with images,
- 210 • ontology based and query expansion semantic retrieval,
- 211 • browsing of a *CI* collection.

212 The *data analytics layer* is based on different Analytics Services allowing to create personalized “dashboards”
213 for a given cultural environment. In addition, it provides basic data mining, graph analysis and machine learning
214 algorithms useful to infer new knowledge .

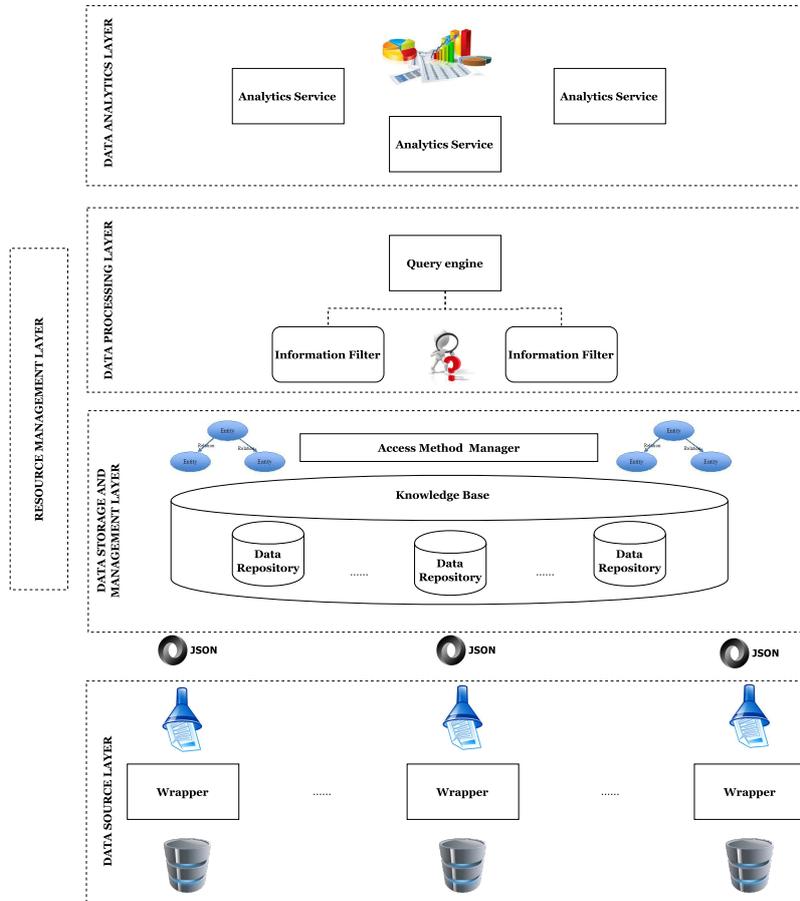


Figure 2: System Architecture.

215 3.2. Functionalities and Implementation Details

216 3.2.1. Resource Management

217 Resource management is an architectural element of paramount importance when designing big data processing
 218 solutions on a wide area scale. In such a scenario, typically characterized by strong runtime and storage distribution
 219 needs, the right resource management choice can result in significant improvements in both performance and fairness.
 220 Indeed, running complex analytics applications on huge amounts of data, continuously collected from hybrid and
 221 heterogeneous sources, cannot be effectively accomplished neither on a single machine nor by relying on a centralized
 222 data repository.

223 First, we need to break down processing activities into smaller tasks to be distributed and managed in parallel on
 224 multiple machines available on the network. Furthermore, coping with ever increasing data volumes characterized
 225 by continuously changing types and formats/structures, requires extremely flexible and elastic storage management
 226 policies capable to follow the evolving demand and go beyond the vertical scaling limits of the traditional storage
 227 architectures. This implies the adoption of new fully distributed storage paradigms providing horizontal scalability
 228 and hence virtually unlimited storage space available.

229 The *Infrastructure as a Service* (IaaS) facilities provided by modern cloud organizations allow multiple tenants to
 230 dynamically obtain or dispose runtime/storage resources available as *virtual machines* (VMs) or containers according
 231 to a pay-per-use paradigms. Depending on the cloud service providers' commercial strategies and management poli-
 232 cies these object may be provided with fixed or variable amounts of virtual CPUs/cores, available memory and disk
 233 space, so that the involved tenants can modify the number of VMs/containers or their available resource pool when
 234 their resource demands varies over time.

235 Our resource management strategy is based on a multi-tenant cloud model, where tenants can buy multiple VMs
236 or containers, characterized by heterogeneous resource demands (resources of different types, such as storage space
237 or memory), to run their applications. The fundamental resource management tasks include discovery, scheduling,
238 allocation and performance monitoring. The discovery activity refers to the identification of a suitable pool of re-
239 sources that can be used to match the users' demands, whereas the scheduling task aims at selecting the best available
240 options between the matching resources in the above pool, in order to perform provisioning to the requiring runtime
241 objects. This implies allocating runtime or storage resources to specific VMs, Containers and processes, whose usage
242 will be continuously monitored until their final release, after their legitimate use. Several runtime resource manage-
243 ment technologies exploiting multiple levels of parallelism in multi-core, many-core, cluster, grid or cloud computing
244 architectures are available, differentiating in their scalability, performance, cost, usage flexibility, robustness etc. [34].
245 Analogously, many Storage resource management options are available for optimizing the efficiency of shared pools
246 of storage devices distributed throughout the network and support their flexible configuration, together with storage
247 tiering and performance monitoring, by also forecasting future space demands and usage patterns.

248 In particular, YARN [35] is actually considered one of the most effective and promising resource management
249 solutions for data-intensive processing in distributed cloud-empowered scenarios. Its Hadoop-based architecture pro-
250 vides the foundation enabling big data processing applications to share computing clusters and storage space over the
251 cloud while guaranteeing consistent service levels and acceptable response times. We also selected the *Apache Spark*
252 [36] framework, developed at UC Berkley, in order to support in-memory data sharing across complex, multi-step
253 data pipelines implemented by using a direct acyclic graph abstraction, so that multiple different applications can
254 simultaneously work on the same data sets.

255 Spark relies on the *Hadoop Distributed File System* (HDFS) in order to provide advanced storage and access
256 functionalities within a unified and flexible big data management solution that is able to cope with a wide range of
257 use cases and requirements. Essentially Spark consists in an execution engine that is able to operate both in-memory
258 and on stable storage devices. It keeps intermediate processing results in memory rather than on semi-permanent
259 storage, that becomes extremely useful mainly when the same data need to be accessed multiple times by different
260 users. However, Spark is able to effectively perform data processing also when their aggregate volume does not fit into
261 the available memory. In these cases it tries to split the data by putting as much as possible in memory and leaving the
262 remaining part on semi permanent storage devices. Due to its in-memory data storage functions and near real-time
263 processing capabilities, it can obtain the best from the *MapReduce* paradigm, limiting the effect of expensive shuffling
264 operations in processing activities, by also achieving significantly improved performances respect to the other existing
265 big data solutions.

266 Spark provides an abstraction of memorization capabilities based on distributed and fully reliable atomic stor-
267 age elements known as *Resilient Distributed Datasets* (RDDs) that can be resident in memory or on stable storage,
268 depending on their access patterns. Several operations can be accomplished on these element by generating new
269 RDDs:

- 270 • through transformation operators, such as *map*, *filter* and *reduceByKey*, that define the new elements without
271 immediately processing them,
- 272 • as values resulting from action operators like *count*, *collect* and *save*, returning their output to applications or
273 exporting some data sets to stable storage.

274 For efficiency purposes RDDs can be implemented according to multiple caching/storage options such as, for
275 example, MEMORY ONLY, MEMORY and DISK and DISK ONLY.

276 In order to support applications runtime facilities Spark creates a master program known as driver that is in charge
277 for RDDs definition and management activities as well as a set of slaves, known as executors, that perform all the
278 computational tasks. In this scenario, RDD objects can be implemented in several different ways:

- 279 • by using a file within a distributed file system, such as HDFS;
- 280 • by partitioning them into multiple pieces to be delivered on different nodes, through the parallelization of a
281 collection of elements within the driver program;

- 282 • by transforming them through an operation such as Map (passing their components through a user-provided
283 function that maps a set of element A into another set B) and filter (selecting all the elements that match a
284 user-defined predicate);
- 285 • by modifying the persistence characteristics of an already existing RDD, e.g., materializing a block of ephemeral
286 data to be used in the parallelization of some processing operations, or discarding it when no more needed. In
287 detail, there are two specific operations that can be used to cope with data persistence: “save” that writes the
288 dataset to HDFS after evaluating it, and “cache” that suggest to keep a dataset in memory when possible.

289 Spark is also able to perform query processing and evaluation on big data, that can be extremely useful in op-
290 timizing huge data management workflow, by also providing an high level Application Programming Interface that
291 can introduce significant effects on productivity in applications development. Applications can request distributed
292 processing operations such as map, reduce and filter by passing specific closures (i.e., functions) to the Spark runtime
293 framework. Like in traditional functional programming environments, such closures are typically referred to process-
294 ing entities and variables that only exist within the scope in which they have been instantiated, so that, when running
295 a closure on a specific node, such entities have to be copied on the node itself.

296 The overall system architecture is made by using multiple virtual machines rent by a cloud provider, each equipped
297 with Hadoop/YARN and Spark packages, accessing an horizontally scaling dynamic storage space that increases on
298 demand, according to the evolving data memorization needs. It consists in a certain number of slave virtual machines
299 (that can be increased over time in presence of changing demands) and two additional ones that run as masters: one
300 for controlling HDFS and another for resource management.

301 3.2.2. Data Gathering from Big Data Sources

302 One of the most important functionalities provided by our system consists in the capability of gathering the differ-
303 ent kinds of data from the various Big Data sources.

304 In particular, we distinguish four classes of data sources that require different wrapping techniques: *Sensor & User*
305 *Data*, *Social Data*, *Digital Repository Data* and *Web Data*.

306 **Sensor and User Data**

307 The Sensor Data generated by different devices are managed using the “publish-subscribe” paradigm. For such
308 kind of data, we have to deal with data stream management challenges [37]; in addition, event processing techniques
309 can be exploited to reduce the amount of information flow [38].

310 In particular, the wrapper for such kind of data is composed by a set of Gateways and Message Brokers together
311 with a Coordination and Discovery module. Data generated by sensors are captured by gateways and sent to Knowl-
312 edge Base by message brokers; from the Knowledge Base side, a Stream Processing Engine receives messages and
313 convert them in a JSON format in according to our data model.

314 In the current system implementation we deal only with *Wireless Sensor Network* (WSN) data sources. We lever-
315 age the *PerLa* WSN management middleware [39] as gateway, which abstracts a sensor network as a virtual distributed
316 database system providing a SQL-like language for query formulation and data retrieval. Data integration issues are
317 coped by a LAV (*Local As View*) schema-mapping approach.

318 Concerning the other implementation details, we use *Apache Kafka* to realize message brokers and the publish-
319 subscribe communication, *Apache Storm* as stream processing engine (see Figure 3) and *Apache Zookeeper* as coordi-
320 nator.

321 Measures returned by sensors’ stream are finally associated with one or more cultural items in the considered
322 environment¹². An example of resource that could be derived from such data sources is reported in Figure 4.

323 The example shows that in a given indoor location (corresponding to a room of the Paestum museum) a certain
324 sensor (identified by the string “S01”) has performed a temperature measure (whose value is “32°”) in a given temporal
325 instant (“12-Apr-16 14:20”) related to a given *CI* (with identifier “00ayr4hfdD2”)¹³. The meaning of all the metadata

¹²Sensor data can be used both to report environmental conditions (temperature and humidity of indoor or outdoor locations), and especially – for monitoring aims – to support preservation and security of cultural items.

¹³All the information about this measure is accessible via web using the URI www.databenc.it/resource/measure/?id=321

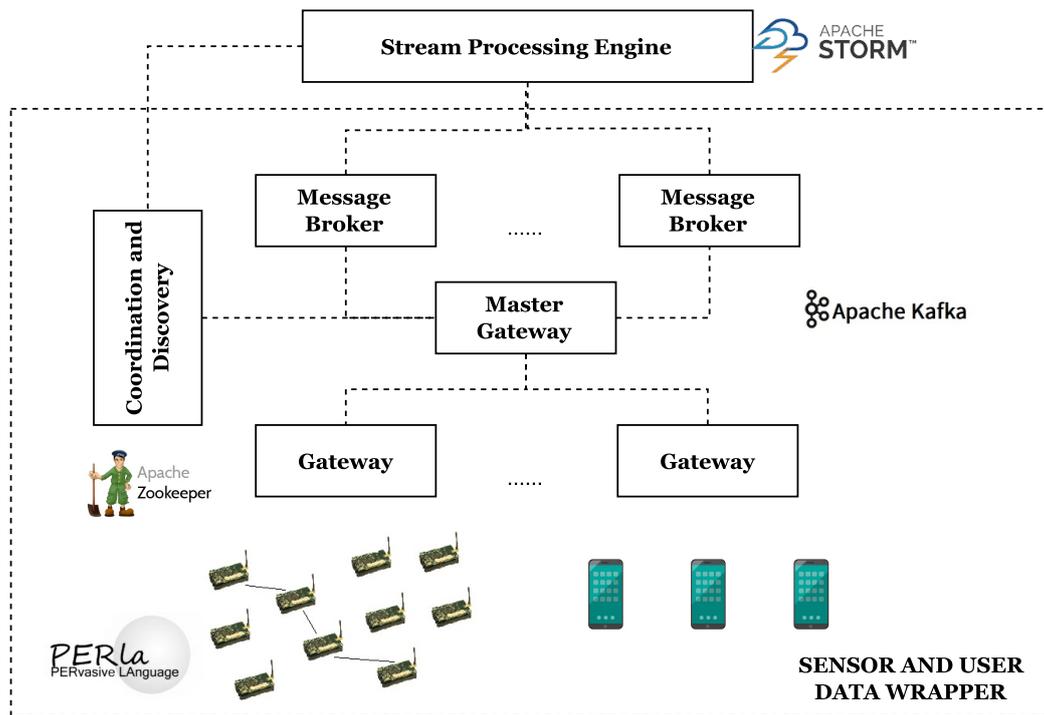


Figure 3: Sensor and User Data Wrapper.

326 describing the measures (e.g., Type, Timestamp, Location, Value, Source) can be retrieved from the annotation schema
 327 *Sensor*.

328 We note that even the user devices, which can communicate the effective users' location by using the GPS tech-
 329 nology or capture their preferences and needs, constitute a different sort of sources, generating User Data that have to
 330 be opportunely wrapped and coded in the described format.

331 For user localization within indoor cultural environments, the applications communicating with our system can
 332 use beacons or NFC technologies.

333 User data may include other static and dynamic information, useful to define the related *profiles*. In particular,
 334 recall that user preferences, needs and behaviors can be acquired either explicitly, by means of suitable questionnaires
 335 provided by the applications, or in an implicit manner using application logs. Static information (e.g., favorite artis-
 336 tic genre) is stored in the KB and possibly updated; in turn, dynamic information (e.g., behaviors) is continuously
 337 updated¹⁴.

338 Also for this kind of data, we use a publish-subscribe based wrapping: the gateway provides a set of basic API
 339 allowing applications to communicate information about position and profiles (see Figure 3).

340 Social and Digital Repositories Data

341 As to *Social Data*, the current prototype considers multimedia information coming from Twitter and different
 342 social media networks (e.g., Flickr and Panoramio). In particular, the wrapper is constituted by a Social Harvesting
 343 module that retrieves user comments and multimedia contents about a given *CI* by exploiting the related APIs and
 344 directly covert them in the described JSON format (see Figure 6). We note that the implementation of the described
 345 module is based on the *Social Harvest* framework. The example in Figure 5 shows that user "Vincenzo.Moscato" has
 346 expressed a comment ("Wonderful") on a post about *The Tomb of the Diver* in a certain date ("14-Apr-16"). The
 347 meaning of all the metadata describing the user activity (e.g., Type, Timestamp, Language, Text, Topic, User, etc.)

¹⁴Such information constitutes the *Personally Identifiable Information* (PII).

```

Sensor_message={
  "URI": "www.databenc.it/resource/measure/?id=321"
  "Annotation Schema": "Sensor"
  "CI" [
    {
      "id": "00ayr4hfdD2"
      "URI": "www.databenc.it/resource/CI/?id=00ayr4hfdD2"
    }
  ]
  "Type": "Temperature"
  "Value": "32°"
  "Timestamp": "12-Apr-16 14:20 "
  "Source": "WSNSensor:S01"
  "Location": "Indoor/Room"
}

```

Figure 4: An example of JSON message related to sensor data

```

Social_message={
  "URI": "www.databenc.it/resource/social/?id=76421"
  "Annotation Schema": "Social Users"
  "CI" [
    {
      "id": "00ayr4hfdD2"
      "URI": "www.databenc.it/resource/CI/?id=00ayr4hfdD2"
    }
  ]
  "Type": "Twitter Comment"
  "Language": "en"
  "Text": "Wonderful"
  "Timestamp": "14-Apr-16 10:40 "
  "Topic": "The Tomb of the Diver"
  "hashtag": "#tombofdiver"
  "User": "foaf:Vincenzo.Moscato"
}

```

Figure 5: An example of JSON message generated by the Twitter Wrapper

348 can be retrieved from the annotation schema *Social Users*. Social data can be used in applications requiring a “social
 349 vision” of cultural items (e.g., general users’ mood of a given social community with respect a given picture).

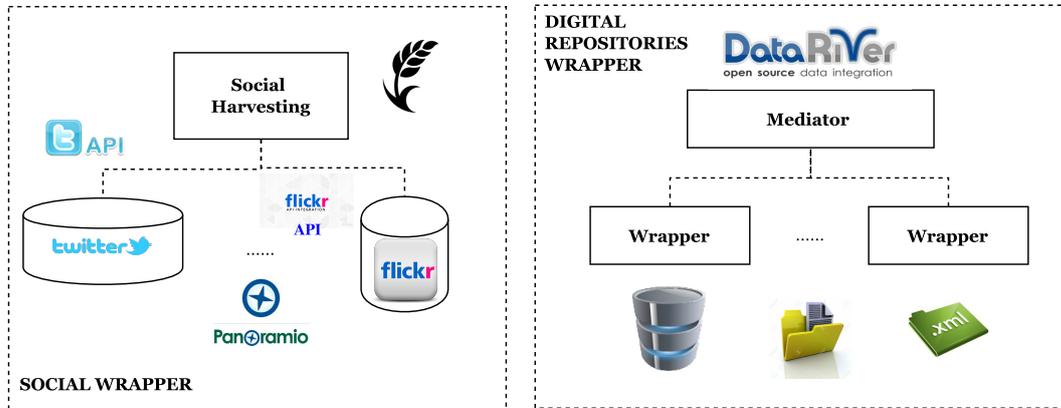


Figure 6: Social and Digital Repositories Wrapper.

350 We collect *Digital Repositories Data* information describing cultural items from on-line digital repositories (e.g.,
 351 museums, libraries, theatrical foundations digital archives, multimedia collections, etc.).

352 Resources are often described on the base of different annotation schemata and provided in XML or RDF based
 353 formats (as in *DBpedia*).

354 The wrapper for this kind of sources can import such data descriptions and convert them into the described JSON
 355 format according to the system data model using schema-mapping techniques; in particular we use the *DataRiver*
 356 system (see Figure 6).

357 The example in Figure 7 shows a possible description (together with an image) of the *The Tomb of the Diver*. The
 358 meanings of all the metadata describing the painting (e.g., Type, Name, Author, Period, Description, Subject, etc.) can
 359 be retrieved in annotation schema *Art History*. Clearly, textual descriptions of cultural items can be made available in
 360 different languages and users’ applications can choose the desired ones.

361 All the multimedia web resources (e.g., images, texts, video, etc.) related to a given cultural item can be similarly
 362 collected using the wrapper facilities, and associated with a given annotation schema.

363 In particular, the descriptions in terms of basic metadata are captured and stored within our systems, while raw
 364 multimedia data can be opportunely linked and, in other cases, temporarily imported into the system for content-based
 365 analysis [40, 41]¹⁵.

366 Web Data

367 Finally, other information of interest, such as that related to meteorological or traffic conditions, can also be
 368 captured using the available *Web Data* services and represented in the described format.

369 3.2.3. Data Storage and Management

370 The data gathered by the Wrappers are stored and managed by the Knowledge Base. Data are represented accord-
 371 ing to the described model, based on the LD/LOD paradigm. One of the basic functionalities of the KB is to export the
 372 related content into the *Europeana Data Model*¹⁶ (EDM) format, a proposal that aims at bridging the gaps between
 373 different harvesting standards providing a unique solution to metadata specification.

374 In the EDM [30], each *Cultural Heritage Object* (CHO) - the main entity of the model corresponding to a given
 375 *CI* - can be described at various levels of detail using different metadata schemata by means of the *aggregation*

¹⁵Note that multimedia data that are managed by the system are suitably filtered before the storing process. The number and kinds of multimedia data required by the application are tuned by means of configuration parameters.

¹⁶<http://pro.europeana.eu/edm-documentation>

```

Repository_message={
  "URI": "www.databenc.it/resource/datarepository/?id=7502"
  "Annotation Schema": " Art History"
  "CI" [
    {
      "id": "00ayr4hfdD2"
      "URI": "www.databenc.it/resource/CI/?id=00ayr4hfdD2"
    }
  ]
  "Multimedia" [
    {
      "Type": "external/image"
      "URI": "https://it.wikipedia.org/wiki/.../PaestumTaucher.jpg"
    }
  ]
  "Type": " Burial Fresco"
  "Name": "The Tomb of the Diver"
  "Period": "Ancient Greek"
  "Description" [
    {
      "Text": "The painted decoration of the tomb called the diver,
      found in 1968, shows a great moment of Greek painting,
      around 480 BC, characterized by the same spirit of the
      painters of the vascular severe style."
      "Language": "en"
    }
  ]
  "Subject": "Real Life activity"
  "Actual Location": "Archeological Museum of Paestum"
}

```

Figure 7: An example of JSON message generated by the Digital Repository Wrapper

376 construct. Thus, our annotation schemata correspond to specific metadata aggregations: as an example, social and
377 sensor aggregations permit to associate/link interesting environmental measures or users' comments from an Online
378 Social Network to each cultural item. Similarly, the archeological aggregation supports the description of a cultural
379 item by metadata that are useful for the annotation work of an archeologist.

380 Two further EDM basic concepts are considered and mapped in our model: *Place* (where the object is located
381 corresponding to our *PoI* entity) and *Time Span* (allowing to deal with different temporal views of a cultural item).

382 Possibly, a certain number of web resources (e.g., images, videos, texts, etc.) can be linked to each item. Finally,
383 metadata semantics is provided by the set of annotation schemata (in XML, RDF or OWL formats). For instance, the
384 data can be represented as sequences of triples (\langle subject, predicate, object \rangle) according to the described data model.

385 In particular, the KB is based on several technologies (see Figure 8):

- 386 • The data describing basic properties of *CIs* (e.g., name, short description, etc.) and basic information on users
387 profiles are stored into a *key-value data store* (i.e., *Redis*).
- 388 • The complete description in terms of all the metadata of *CIs* using the different annotation schemata are in turn
389 saved using a *wide column data store* (i.e., *Cassandra*). We use a table for each kind of *CIs* having a column
390 for each “metadata family”; column values can be literals or URIs.
- 391 • The *document store* technology (i.e., *MongoDB*) is used to deal with JSON messages, complete user profiles
392 and descriptions of internal resources (multimedia data and textual documents, etc.) associated with a cultural
393 item.
- 394 • All the relationships among cultural items within a cultural environment and interactions with users (behaviors)
395 are managed by means of a *graph database* (i.e., *Titan*).
- 396 • The entire cartography related to a cultural environment together with *PoIs* is managed by a GIS (i.e., *PostGIS*),
397 which provides the functionalities to filter and visualize on a map the geographic area around a given *PoI* ;
- 398 • Multimedia data management have been realized using the *Windsurf* library [41];
- 399 • We exploit an *RDF store* (i.e., different *Allegrograph* instances) to memorize data views in terms of triples
400 related to a given cultural environment and useful for specific applications, providing a SPARQL endpoint for
401 the applications;
- 402 • All system configuration parameters, internal catalogs and thesauri are stored in a relational database (i.e.,
403 *PostgreSQL*)
- 404 • All the basic analytics and query facilities are provided by *Spark* based on *HDFS*;
- 405 • Semantics of data can be specified by linking values of high-level metadata to some available internal (managed
406 by *Sesame*) or external ontological schemata.

407 This heterogeneous Knowledge Base provides basic *Restful APIs* to read/write data and further functionalities for
408 importing/exporting data in the most common diffused Web standards.

409 We chose to adopt such heterogeneous technologies in order to meet the specific requirements of the applications
410 dealing with the huge amount of data at stake. For example, Social Networking applications typically benefit of graph
411 database technologies because of their focus on data relationships. In turn, more efficient technologies (key-value ore
412 wide-column stores) are required by Tourism applications to quickly and easily access the data of interest.

413 3.2.4. Data Processing and Analytics

414 The search of data useful for the applications can be eased by using different Information Filters that implement
415 the right queries to the various databases.

416 The implementation of such filters is based on *Spark*, and we distinguish three four kinds of query on the KB:

- 417 • *query by keywords/tags*: through such a query a user/application can search a set of *CIs* using keywords (as in
418 Google search engine) or specific tags (the query is then “expanded” with similar search terms leveraging the
419 system thesauri);

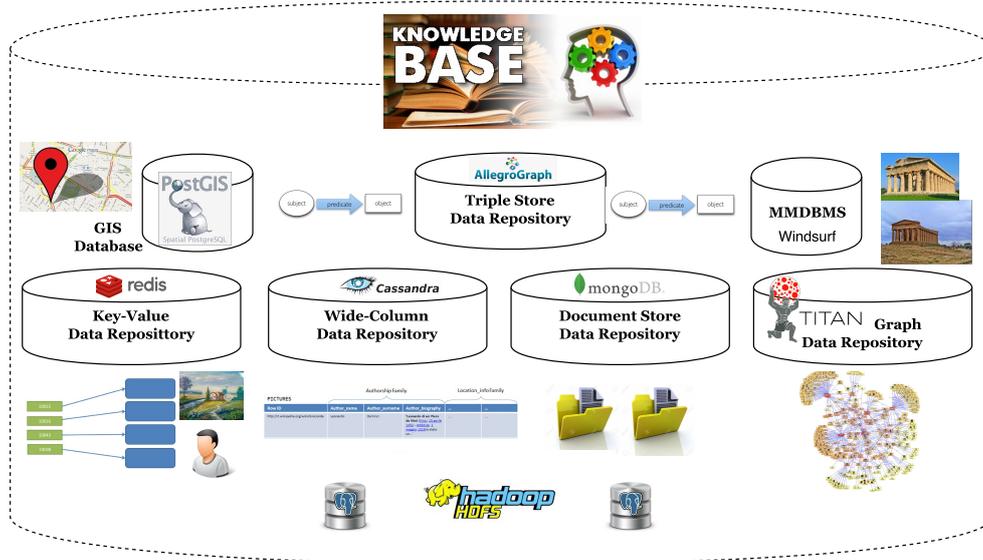


Figure 8: Knowledge Base.

- 420 • *query by metadata*: by this query a user/application can search a set of *CTs* using specific metadata of internal
421 or external annotation schemata (the query for *OAs* is semantically “expanded” with the concepts of managed
422 ontologies that are similar to the target one);
- 423 • *query by example*: through such a query a user/application can search a set of multimedia contents – related to
424 *CTs* – that are similar to a target one (the query processing is base on the *Windsurf multimedia libraries*);
- 425 • *query by user preferences*: user profiles are exploited to find the set of cultural items that are more similar to
426 user preferences using co-clustering techniques [42].

427 The different types of query are combined by the Query Engine that provides a system query endpoint and support
428 different query languages.

429 Eventually, several analytics can be performed on the stored data using *Graph Analysis* and *Machine Learning*
430 library provided by Spark. As an example, the analysis of the graph coding interactions among users, multimedia
431 objects and cultural items can be properly used to support recommendation or influence analysis tasks[43].

432 4. System Running Example

433 In this section we describe a possible application of our system to support the development of a multimedia guide
434 for the *Paestum archeological site*.

435 The archaeological site of Paestum is one of the major Graeco-Roman cities in Southern Italy. Here, a set of
436 ancient buildings is the main cultural attraction for a tourist: three main temples of Doric style (the *First Temple of*
437 *Hera*, also called *Basilica*, the *Second Temple of Hera*, also known as *Temple of Neptune*, and the *Temple of Athena*),
438 the *Roman Forum* with several ruins, and the *Amphitheater*. The buildings are surrounded by the remains of the city
439 walls. In addition, a museum near the ancient city contains many evidences of the Graeco-Roman life (e.g., amphorae,
440 paintings and so on).

441 The ancient buildings, together with the museum and its main artifacts, constitute the set of cultural items for
442 our case study. Tourists, both from their places and while visiting ruins, can browse these cultural items and enjoy a
443 useful multimedia guide describing them, or be recommended other nearby places, comments of other users and other
444 information of interest.

445 As an example, suppose that some users from the nearby seaside resorts would like to know more about the ruins
446 in front of them, beyond their names and a very generic and basic explanation of the main cultural items.

Find a Cultural Item			
<input type="text" value="Cerca"/>			
Cultural Item	Description	Genre	Social
Temple of Neptune ★★★★★		Outdoor/ Ancient Building/ Graeco-Roman Temple	Rate & Comment
Object	Type	Source	Other Info
	Multimedia/ Image 		
	Multimedia/ Image 		
Tempio di Eia <small>Tempio di Eia, sito delle Terme di Eia, presso il centro storico di Eia, in provincia di Imperia. L'edificio è stato restaurato nel 1980. È stato dichiarato Monumento Nazionale nel 1980. Destinazione: Tempio.</small>	Text/ <i>Italian</i> 		
1/145 	More Objects	Type & Language Other Filters	
<i>"Wonderful..."</i>	Social/ Comments		
<i>"A magic place!"</i>	Social/ Comments		
1/573 	More Comments	Source & Language Other Filters	

Figure 9: System Query Interface

447 When users search a specific cultural item, as an example the Temple of Neptune, our system provide a basic
448 description with the related multimedia objects (i.e., audio, images, video and texts) and detailed users' comments.
449 The list of proposed cultural items with descriptions and multimedia objects depends on the user's preferences and
450 system settings: images rather than voice, or expert-level rather than layman-level descriptions of the art pieces;
451 specific metadata and annotation schemata. In addition, *query by example* facilities can be exploited to determine
452 other images that are similar to a given multimedia object. In addition, a *semantic based search* can be performed
453 on specific ontological attributes to find other cultural items of the same type. At the same time, users can choose to
454 retrieve some interesting information, to read comments, opinions and ratings about the visited cultural items and to
455 express their own ratings and opinions that will be posted on social networks.

456 Figure 9 shows a running example (obtained by assembling different screenshots of our system) concerning the
457 search of *CI*s related to the Paestum ruins. Users can browse the data by means of an appropriate GUI; they can filter
458 objects belonging to a given *CI* using different criteria: type of multimedia data, language, size, and so on.

459 Thus, the information can be dynamically tailored on the users depending on their preferences and needs. As
460 an example, young people with high-school education and basic notions of history and art could receive only some
461 general concepts about the cultural items, while CH experts could be interested in more detailed descriptions and
462 richer multimedia data.

463 Eventually, other cultural items of interest can be automatically suggested using recommendation techniques that
464 exploit context information, user preferences and needs and items' features [42].

465 All the information and functionalities (provided as API) can be exploited by third-part applications such as a
466 multimedia guide for the Paestum ruins.

467 The system prototype is currently running on a cloud-based big data processing environment based on Spark . In
468 particular, we exploited 5 computing nodes, each one composed by 8 cores and 15 GB RAM.

469 5. Conclusions and Future Work

470 In this paper we showed an application and some examples of the use of Big Data technologies for the Cultural
471 Heritage domain.

472 In particular, we described CHIS - a scalable prototype for the management and context-driven browsing of cul-
473 tural environments . The system is characterized by several features that are typical of the modern big data systems: i)
474 information gathering from distributed and heterogeneous pervasive data sources (e.g., Sensor Networks, Social Media
475 Networks, Digital Libraries and Archives, Multimedia Collections, Web Data Services, etc.); ii) context-awareness,
476 provisioning of useful and personalized data and services, on the base of users preferences and of the surrounding
477 environment; iii) advanced data management techniques and technologies for dealing with the information variety,
478 velocity and volume; iv) advanced information retrieval facilities, data analytics and other utilities/services.

479 We detailed all the design choices and characteristics of the implemented big data infrastructure, providing a
480 system running example concerning the cultural data management of the archaeological site of Paestum.

481 Future work will be devoted to collect the huge amount of data related to all the different cultural objects of the
482 Campania region and to experiment our system from the efficiency and effectiveness points of view with respect to
483 the information retrieval and filtering tasks providing a comparison with other systems.

484 References

- 485 [1] A. Caragliu, C. Del Bo, P. Nijkamp, Smart cities in europe, *Journal of urban technology* 18 (2) (2011) 65–82.
486 [2] H. Schaffers, N. Komninos, M. Pallot, B. Trousse, M. Nilsson, A. Oliveira, Smart cities and the future internet: Towards cooperation
487 frameworks for open innovation, in: *The Future Internet Assembly*, Springer, 2011, pp. 431–446.
488 [3] J. M. Hernández-Muñoz, J. B. Vercher, L. Muñoz, J. A. Galache, M. Presser, L. A. H. Gómez, J. Pettersson, Smart cities at the forefront of
489 the future internet, in: *The Future Internet Assembly*, Springer, 2011, pp. 447–462.
490 [4] N. Komninos, H. Schaffers, M. Pallot, Developing a policy roadmap for smart cities and the future internet, in: *eChallenges e-2011 Conference*
491 *Proceedings*, IIMC International Information Management Corporation, IMC International Information Management Corporation, 2011.
492 [5] M. Armbrust, A. Fox, R. Griffith, A. D. Joseph, R. Katz, A. Konwinski, G. Lee, D. Patterson, A. Rabkin, I. Stoica, et al., A view of cloud
493 computing, *Communications of the ACM* 53 (4) (2010) 50–58.
494 [6] J. A. Burke, D. Estrin, M. Hansen, A. Parker, N. Ramanathan, S. Reddy, M. B. Srivastava, Participatory sensing, Center for Embedded
495 Network Sensing.

- 496 [7] C. Bizer, T. Heath, T. Berners-Lee, Linked data-the story so far, *Semantic Services, Interoperability and Web Applications: Emerging*
497 *Concepts* (2009) 205–227.
- 498 [8] T. Berners-Lee, J. Hendler, O. Lassila, et al., The semantic web, *Scientific american* 284 (5) (2001) 28–37.
- 499 [9] S. Lohr, The age of big data, *New York Times* 11.
- 500 [10] C. Bizer, P. Boncz, M. L. Brodie, O. Erling, The meaningful use of big data: four perspectives–four challenges, *ACM SIGMOD Record*
501 40 (4) (2012) 56–60.
- 502 [11] F. Colace, M. D. Santo, V. Moscato, A. Picariello, F. A. Schreiber, L. Tanca, Patch: A portable context-aware atlas for browsing cultural
503 heritage, in: *Data Management in Pervasive Systems*, 2015, pp. 345–361.
- 504 [12] L. Tan, N. Wang, Future internet: The internet of things, in: *2010 3rd International Conference on Advanced Computer Theory and Engi-*
505 *neering (ICACTE)*, Vol. 5, IEEE, 2010, pp. V5–376.
- 506 [13] L. Atzori, A. Iera, G. Morabito, The internet of things: A survey, *Computer networks* 54 (15) (2010) 2787–2805.
- 507 [14] J. Zhang, V. Varadharajan, Wireless sensor network key management survey and taxonomy, *Journal of Network and Computer Applications*
508 33 (2) (2010) 63–75.
- 509 [15] J. J. Rodrigues, P. A. Neves, A survey on ip-based wireless sensor network solutions, *International Journal of Communication Systems* 23 (8)
510 (2010) 963–981.
- 511 [16] N. Mohamed, J. Al-Jaroodi, A survey on service-oriented middleware for wireless sensor networks, *Service Oriented Computing and Appli-*
512 *cations* 5 (2) (2011) 71–85.
- 513 [17] H. Alemdar, C. Ersoy, Wireless sensor networks for healthcare: A survey, *Computer Networks* 54 (15) (2010) 2688–2710.
- 514 [18] S. LaValle, E. Lesser, R. Shockley, M. S. Hopkins, N. Kruschwitz, Big data, analytics and the path from insights to value, *MIT sloan*
515 *management review* 52 (2) (2011) 21.
- 516 [19] P. Russom, et al., Big data analytics, *TDWI Best Practices Report, Fourth Quarter* (2011) 1–35.
- 517 [20] M. Damova, A. Ontotext, B. A. Kiryakov, M. Grinberg, M. K. Bergman, F. Giasson, K. Simov, Creation and integration of reference
518 ontologies for efficient lod management, *Semi-Automatic Ontology Development: Processes and Resources: Processes and Resources* (2012)
519 162.
- 520 [21] R. Ruggles, *Knowledge management tools*, Routledge, 2009.
- 521 [22] H. M. Khoury, V. R. Kamat, Evaluation of position tracking technologies for user localization in indoor construction environments, *Automa-*
522 *tation in Construction* 18 (4) (2009) 444–457.
- 523 [23] L. M. Ni, D. Zhang, M. R. Souryal, Rfid-based localization and tracking technologies, *IEEE Wireless Communications* 18 (2) (2011) 45–51.
- 524 [24] K. Kabassi, Personalisation systems for cultural tourism, in: *Multimedia services in intelligent environments*, Springer, 2013, pp. 101–111.
- 525 [25] F. Ricci, L. Rokach, B. Shapira, *Introduction to recommender systems handbook*, Springer, 2011.
- 526 [26] C. Becker, F. Dürr, On location models for ubiquitous computing, *Personal and Ubiquitous Computing* 9 (1) (2005) 20–31.
- 527 [27] J. Scott, *Social network analysis*, Sage, 2012.
- 528 [28] C. P. Chen, C.-Y. Zhang, Data-intensive applications, challenges, techniques and technologies: A survey on big data, *Information Sciences*
529 275 (2014) 314–347.
- 530 [29] F. Colace, M. De Santo, V. Moscato, A. Picariello, F. A. Schreiber, L. Tanca, Pervasive systems architecture and the main related technologies,
531 in: *Data Management in Pervasive Systems*, Springer, 2015, pp. 19–42.
- 532 [30] M. Doerr, S. Gradmann, S. Hennicke, A. Isaac, C. Meghini, H. van de Sompel, The europeana data model (edm), in: *World Library and*
533 *Information Congress: 76th IFLA general conference and assembly*, 2010, pp. 10–15.
- 534 [31] F. Amato, A. Mazzeo, A. Penta, A. Picariello, Building rdf ontologies from semi-structured legal documents, in: *2008 International Confer-*
535 *ence on Complex, Intelligent and Software Intensive Systems*, 2008.
- 536 [32] F. Amato, A. R. Fasolino, A. Mazzeo, V. Moscato, A. Picariello, S. Romano, P. Tramontana, Ensuring semantic interoperability for e-health
537 applications, in: *Complex, Intelligent and Software Intensive Systems (CISIS)*, 2011 International Conference on, IEEE, 2011, pp. 315–320.
- 538 [33] X. L. Dong, D. Srivastava, Big data integration, in: *Data Engineering (ICDE)*, 2013 IEEE 29th International Conference on, IEEE, 2013, pp.
539 1245–1248.
- 540 [34] J. L. Reyes-Ortiz, L. Oneto, D. Anguita, Big data analytics in the cloud: Spark on hadoop vs mpi/openmp on beowulf, *Procedia Computer*
541 *Science* 53 (2015) 121–130.
- 542 [35] V. K. Vavilapalli, A. C. Murthy, C. Douglas, S. Agarwal, M. Konar, R. Evans, T. Graves, J. Lowe, H. Shah, S. Seth, et al., Apache hadoop
543 yarn: Yet another resource negotiator, in: *Proceedings of the 4th annual Symposium on Cloud Computing*, ACM, 2013, p. 5.
- 544 [36] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, I. Stoica, Spark: cluster computing with working sets., *HotCloud* 10 (2010) 10–10.
- 545 [37] E. Panigati, F. A. Schreiber, C. Zaniolo, Data streams and data stream management systems and languages, in: *Data Management in Pervasive*
546 *Systems*, Springer, 2015, pp. 93–111.
- 547 [38] G. Cugola, A. Margara, The complex event processing paradigm, in: *Data Management in Pervasive Systems*, Springer, 2015, pp. 113–133.
- 548 [39] F. A. Schreiber, R. Camplani, M. Fortunato, M. Marelli, G. Rota, Perla: A language and middleware architecture for data management and
549 integration in pervasive information systems, *Software Engineering, IEEE Transactions on* 38 (2) (2012) 478–496.
- 550 [40] F. Amato, F. Colace, L. Greco, V. Moscato, A. Picariello, Semantic processing of multimedia data for e-government applications, *Journal of*
551 *Visual Languages & Computing* 32 (2016) 35–41.
- 552 [41] I. Bartolini, M. Patella, Multimedia queries in digital libraries, in: *Data Management in Pervasive Systems*, Springer, 2015, pp. 311–325.
- 553 [42] F. Colace, M. D. Santo, L. Greco, V. Moscato, A. Picariello, A collaborative user-centered framework for recommending items in online
554 social networks, *Computers in Human Behavior* 51 (2015) 694–704.
- 555 [43] F. Amato, V. Moscato, A. Picariello, G. Sperli, Multimedia social network modeling: A proposal, in: *Tenth IEEE International Conference*
556 *on Semantic Computing, ICSC 2016, Laguna Hills, CA, USA, February 4-6, 2016*, 2016, pp. 448–453.