

A Data Integration Architecture for Smart Cities*

Murilo Borges Ribeiro¹, Kelly Rosa Braghetto¹

¹Department of Computer Science – Institute of Mathematics and Statistics
University of São Paulo (USP)
Rua do Matão 1010, São Paulo, Brazil

{muriloribeiro, kellyrb}@ime.usp.br

Abstract. *The data generated by smart cities have low integration, as the systems that produce them are usually closed and developed for specific needs. Moreover, the large volume of data, and the semantic and structural changes in datasets over time make the use of data to support decision-making even more difficult. In this work, we identify the main requirements of a data integration system to support decision-making in cities, focusing on its challenges. We analyze some existing data integration solutions, to uncover their features and limitations. Based on these results, we propose a new microservice architecture to support the development of software platforms for integrating smart cities' heterogeneous data and a guideline to assess their performance.*

1. Introduction

The increasing availability of electronic devices with sensing capacity and computational power, capable of receiving and sending information, causes a large amount of data from different sources and in different structures to be continuously produced in smart cities. Cities also accumulate data generated by government entities, citizens and systems.

The collection, cleaning, integration, transformation, and analysis of large amounts of data generated by different sources can help us to have a better understanding of the deficiencies of cities. Urban data can be used to assist the evidence-based decision-making and the development of public policies aimed at making the best use of the available resources and the improvement in the quality of life of the population. For example, the crossing of records of the municipal health department with sociodemographic, meteorological, and social networks data can be used to monitor the evolution and prevent of endemic diseases such as Dengue.

The development of solutions to explore urban data faces difficulties due to the lack of accessibility and integration of data [Raghavan et al. 2020]. This happens because many systems are built in silos: they are closed and developed for specific needs. Data integration systems for smart cities have already been the subject of several research works [Psyllidis et al. 2015, Consoli et al. 2015, Cheng et al. 2015, Rathore et al. 2016, Hashem et al. 2016, Costa and Santos 2017, Mehmood et al. 2019]. Despite these works present some approaches for the different stages of data integration (such as ingesting, processing, storing, analyzing, and visualizing data), there are still open issues and place for improvement in such systems. Examples of issues are the insufficient support for

*This research is part of the INCT of the Future Internet for Smart Cities funded by CNPq proc. de Pessoal de Nível Superior – Brasil (CAPES) – Finance Code 001, FAPESP proc. 14/50937-1, and 465446/2014-0, Coordenação de Aperfeiçoamento FAPESP proc. 15/24485-9.

metadata management and the lack of data query facilities for non-specialist users, which makes the discovery and the reuse of urban data more difficult.

In this work, we analyze some of the data integration solutions for smart cities reported in scientific literature, to uncover their features and limitations. Two main research questions have guided our analyses:

1. What are the main challenges and issues identified by the researchers for data integration in smart cities?
2. What are the functional and non-functional requirements of a software platform for data integration in smart cities?

Based on this analysis, we propose a microservices architecture to guide the development of software platforms for integrating smart cities' heterogeneous data and facilitating its use. This architecture was designed to support all the required services (i.e. data ingestion, metadata management, data processing, data analysis, and data visualization) while providing scalability, availability, security, and privacy. We also present a guideline to assess performance of systems that implement the proposed architecture. The guideline follows the Cloud Evaluation Experiment Methodology (CEEM) [Li et al. 2013], used for systematically evaluating cloud services' performance through experiments.

The remainder of this paper is organized as follows. Section 2 analyzes the related works, identifying the requirements of the data integration platforms. Section 3 introduces our microservices architecture, providing details of each one of its services. The guideline for the performance evaluation of implementations of the architecture is presented in Section 4. Finally, Section 5 presents the concluding remarks.

2. Software Platforms for Data Integration in Smart Cities

2.1. Related Works

We have searched Google Scholar¹ for works published since 2015 using the search string: (“data integration” or “semantic data integration” or “data warehouse” or “data lake” or “big data”) and “smart cities”). Then, we have analyzed the abstract and keywords of each of the returned papers to filter those that present software architectures to integrate heterogeneous data in smart cities. From this filtering, we reached the works of [Psyllidis et al. 2015], [Consoli et al. 2015], [Cheng et al. 2015], [Rathore et al. 2016], [Hashem et al. 2016], [Costa and Santos 2017], and [Mehmood et al. 2019]. Most of them developed distributed, multi-tiered systems, capable of handling data both in batch and real time, and supporting a large variety of services for applications and final users.

[Psyllidis et al. 2015] developed SocialGlass, a web platform that offers resources for analysis, integration, and visualization of heterogeneous urban data in order to assist in urban planning and decision-making. The SocialGlass architecture is divided into three main modules: ingestion, integration and exploration. The ingestion module refers to the acquisition, cleaning, and processing of social and sensor data. The integration module is responsible for enabling interoperability between different data sources. To achieve that, an ontology-based knowledge representation model was developed, which represents urban systems, the relationships between them, and the corresponding data sources. The

¹<https://scholar.google.com/>

exploration module offers a map-based web interface for data visualization and exploration, making it possible to obtain insights about spatial and temporal parameters of the urban context. Details of the technologies used in the implementation were not provided. The system does not support data processing and access via external platforms.

[Consoli et al. 2015] presented an ontology integration approach using Linked Data (a set of practices for publishing and connecting data on the Web). In their work, each dataset was converted to an RDF (Resource Description Framework) data model using custom processes. With the help of domain experts, ontologies were generated for each dataset, to achieve conceptual interoperability. Data and ontology are accessible by querying the SPARQL API. The ingestion, processing, and visualization of data were out of the scope of their work.

[Cheng et al. 2015] proposed a Big Data architecture integrated with the IoT SmartSantander experimental test environment. This architecture is divided into four main modules: (1) data collection, (2) data storage, (3) data processing and analysis, and (4) API for communicating with external applications. The data collection module is represented by a broker, which is responsible for receiving data from different sources. A NoSQL database is used for data storage. Data processing and analysis are done in batch or stream by using a distributed computing tool. The module for communicating with external applications has a RESTful API to allow external applications to make simple queries, complex queries, and subscriptions. A simple query might request aggregated results about the latest status of all sensors, while a complex query might request aggregated results about historical data within a specified time frame. Subscription is the mechanism used for apps to receive notifications with the latest results, preventing the apps from querying the data all the time. This architecture does not feature a visualization module for stored data and metadata, making it difficult for public managers to use it.

[Rathore et al. 2016] proposed a system for collecting, aggregating, filtering, sorting, pre-processing, computing, and decision-making using the *Data Lake* approach combined with *Data Warehouse*. The proposed system is divided into four layers: (1) data generation and collection, (2) data transmission, (3) data management and processing, and (4) data analysis. The first and second layers are responsible for collecting data using sensors and transferring the data to the storage platform; therefore, they are in a lower level than the services considered in this work. The third layer, data management and processing, use a distributed file system (HDFS) and distributed computing tools for real-time data processing. For historical data, the authors suggested the use of a tool for Data Warehouses (Apache Hive²) and distributed databases (Apache HBase³). The fourth layer is composed of several applications, each one for a different type of planning. The architecture does not support metadata management and data analysis.

Similarly, the smart cities Big Data architecture proposed by [Hashem et al. 2016] is divided into four layers. The first is composed of sources and transferring of data, while the second is responsible for storing the data in a distributed and fault-tolerant database. In the latter, the stored data is processed according to the queries received using a parallel and distributed processing programming model. The third layer, intelligent analytics, was designed to support the use of machine learning and data mining to extract patterns and

²<https://hive.apache.org>

³<https://hbase.apache.org>

knowledge from large amounts of data. The last layer is made up of applications that use the stored data for varied purposes, such as intelligent management of public resources.

[Costa and Santos 2017] presented an approach to design and implement a Big Data Warehouse in the context of smart cities, with a repository that stores data in raw format. The proposed architecture is divided into four major modules responsible for data collection, preparation and enrichment, storage and access, analysis and visualization. Data can be collected in real-time using a broker (Apache kafka⁴), or in batch using an ETL tool, e.g. Talend⁵ and HDFS Upload. The data collected in batch is directly stored in files on a distributed system, prepared and enriched using a distributed computing tool (Apache Spark⁶), and then stored in a *Data Warehouse* (Apache Hive). Data collected in real time is also stored in a distributed file system, prepared and enriched using a distributed computing tool, and later stored in a distributed database (Apache Cassandra⁷). The stored data can be accessed by a distributed SQL query tool (Presto⁸) and by a data visualization tool. The proposed model was implemented in the SusCity research project and was used to analyze data collected in the city of Lisbon. The solution does not support metadata management and accessing data via external platforms.

Similarly, [Mehmood et al. 2019] proposed an architecture divided into five modules responsible for data collection, ingestion, storage, exploration and analysis, and visualization. For data ingestion, they proposed the use of a stream processing tool (Apache Flume⁹) with storage in a distributed file system (HDFS). Data analysis and exploration were performed using a distributed indexing tool (Apache Solr¹⁰) and distributed computing (Apache Spark). For data visualization, a SQL query web tool (Hue¹¹) and the Matplotlib¹² library were used. The presented metadata management requires data uniformity, making it difficult to analyze the data from the different sectors of the city.

2.2. Requirements for Data Integration Software Platforms

2.2.1. Functional Requirements

The main goal of a platform for data integration in smart cities is to facilitate the development of applications that use data combined from different sources. To this end, most of the analyzed platforms implement requirements for data ingestion, processing, analysis, visualization, and data sharing. Table 1 provides an overview of how the related works cover these functional requirements. Each requirement is described in the sequence.

Data ingestion is the process of importing real-time or batch data into the storage platform. Data can come from different sources, in different formats, such as CSV, TXT, JSON, and others.

Metadata management is the process of collecting and managing information about data stored on the platform. Metadata must contain information about the semantics and structure of data collections. Metadata should also keep information about the mappings needed to standardize data and guarantee backward compatibility, in order to

⁴<https://kafka.apache.org>

⁵<https://www.talend.com>

⁶<https://spark.apache.org>

⁷<https://cassandra.apache.org>

⁸<https://prestodb.io>

⁹<https://flume.apache.org>

¹⁰<https://solr.apache.org>

¹¹<https://gethue.com>

¹²<https://matplotlib.org>

	Ingestion	Metadata	Processing	Machine Learning	Analysis and Visualization	External Access
[Psyllidis et al. 2015]	X	X			X	
[Consoli et al. 2015]		X				X
[Cheng et al. 2015]	X		X			X
[Rathore et al. 2016]	X		X	X		
[Hashem et al. 2016]	X		X	X		X
[Costa and Santos 2017]	X		X		X	
[Mehmood et al. 2019]	X	X	X		X	

Table 1. Functional Requirements

	Scalability	Availability	Security and Privacy
[Psyllidis et al. 2015]	X		
[Consoli et al. 2015]			
[Cheng et al. 2015]	X		
[Rathore et al. 2016]	X		X
[Hashem et al. 2016]		X	
[Costa and Santos 2017]			X
[Mehmood et al. 2019]	X	X	

Table 2. Non-functional Requirements

enable data integration and ease of use.

The data arriving at the platform may be inaccurate, incomplete, inconsistent, or redundant. Additionally, this data may need aggregation, filtering, or analysis before enabling knowledge discovery. Thus, platforms must offer resources for creating and executing **data processing** procedures.

Extracting knowledge and insights of data is of paramount importance to enable better decision-making in cities and support the implementation of efficient public policies. Therefore, data integration platforms must enable the creation, maintenance, and execution of custom **machine learning** models.

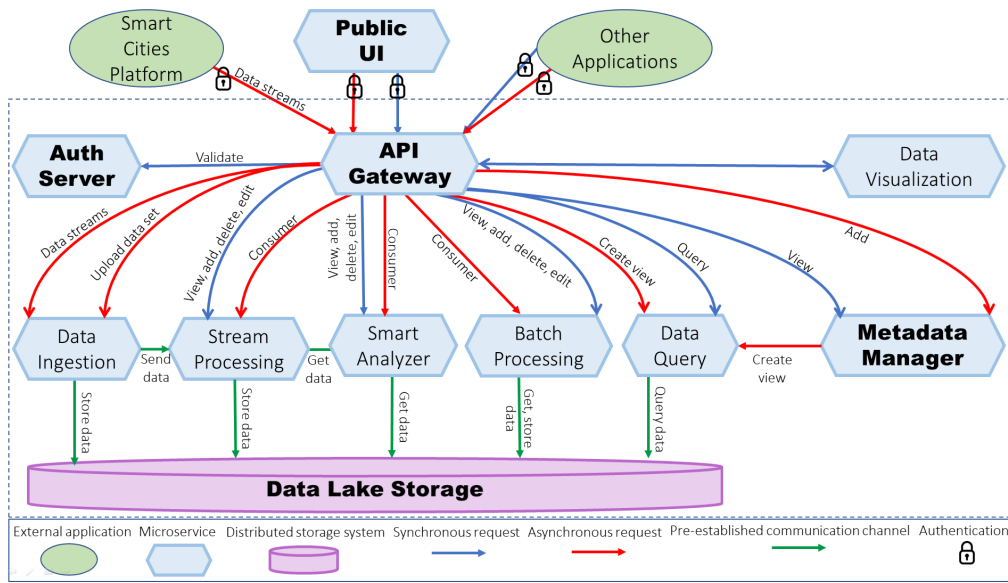
Data analysis and visualization refer to the presentation of data in user-friendly graphical formats, to help users understand the behavior of cities and the use of resources. A data integration platform must offer features to support the creation of custom reports and dashboards for managers to convert data into knowledge.

The **external access** refers to the possibility of external systems consuming the data stored on the platform, allowing the development of new applications to improve the services provided to the population or to optimize the use of available resources. Only authorized users or systems should have access to the data.

2.2.2. Non-functional Requirements

Table 2 shows the non-functional requirements mentioned in the related work of Section 2.1. We describe each one of the requirements in the sequence.

Scalability refers to the ability to increase or decrease computational resources according to the need of the system. Scalability can be vertical, meaning adding (remov-



Microservices with features not provided by the solutions of the related works are highlighted in bold.

Figure 1. The proposed software architecture for data integration in smart cities.

ing) resources to (from) a single node; or horizontal, when adding (removing) nodes to (from) a distributed system.

Availability refers to the ability of the platform to be resilient to hardware, software, and power failures to keep services available for as long as possible.

Security and privacy refer to restrict access to the stored data to authorized users and systems, preventing leakage and misuse of information. It also includes the platform's ability to comply with data protection policies so that sensitive data is properly anonymized and secured.

The implementation of these requirements imposes challenges due to the large volume of data and the heterogeneity of sources and formats. According to [Cheng et al. 2015], scalability is an important issue as the amount of data greatly increases over time, with the availability of new services and technology, and also with the population growth. Security and privacy are important issues as well since there is a lot of sensitive data being collected in smart cities, and the cyber attack attempts become each day more frequent [Cheng et al. 2015, Hashem et al. 2016].

3. A Microservice Architecture for Data Integration in Smart Cities

In this section, we present a new architecture to guide the development of data integration software platforms for smart cities. Figure 1 shows the diagram of the architecture. This architecture was initially derived from the architectures of related works, and then it was significantly improved to address the requirements and issues discussed in Section 2.2.

3.1. Components of the Architecture

The lowest-level component in the proposed architecture is the **Data Lake Storage**, a distributed storage system responsible for storing data as close as possible to its original format, so that end-users do not need to understand details of how data is stored in order

to be able to use it. To prevent the storage system from being a bottleneck for query engines, data must be stored in a standardized, compressed format that facilitates analytical querying and reduces storage size and cost, but information cannot be discarded or lost. The most used open-source distributed file system is HDFS. Data can be standardized into Avro¹³ or Parquet¹⁴ format and compressed using Snappy or Gzip compression.

The **Data Ingestion** microservice is responsible for asynchronously consuming data from the message queues, converting the original data to a standardized and compressed file format, and then sending it to **Data Lake Storage** by using a pre-established communication channel. It is necessary to have at least two message queues, one for real-time data (i.e. from sensors and social networks) and one for batch data. The separation of queues is important to enable the prioritization of the ingestion of data in real-time and to improve the scalability of the system. This microservice can be implemented using Kafka, RabbitMQ¹⁵ and Spark streaming.

The **Metadata Manager** is a metadata catalog with features for registering data sources and information about data schema, data origin, privacy policies, schema versioning, and data mapping rules. It is also responsible for requesting the creation or updating of the data view for the **Data Query** microservice whenever a new version of the metadata is created. The latter are important to make data from different schema versions compatible, so that even legacy data can be easily discovered and used by the authorized persons. The metadata model must follow the World Wide Web Consortium (W3C) standards. We suggest the adoption of standard models such as the RDF vocabulary Data Catalog Vocabulary (DCAT¹⁶). DCAT facilitates the consumption and aggregation of metadata from multiple catalogs. It is integrated with other standards, such as Schema.org and PROV Ontology (PROV-O¹⁷). The catalog is maintained in the microservice's own database (preferably NoSQL, to facilitate the storage of metadata in JSON).

The **Data Query** is the microservice responsible for processing the creation orders of visualization of data stored in **Metadata Manager** and synchronously executing SQL queries on the data in **Data Lake Storage**. Query statistics such as execution frequency and response time should be stored to allow the use of automatic indexing and caching techniques to speed up the access to frequently used data. Softwares like Hive and Spark SQL can be used to implement this microservice.

The analysis and processing of data in real-time are made by the **Stream Processing** microservice. This microservice enables the execution of tasks on data as soon it enters the platform. It provides an interface for developers to create and manage their jobs. Results of the data processing must be stored, enabling their use by other applications. This microservice can be implemented using tools such as Kafka and Spark Streaming.

The **Smart Analyzer** microservice provides tools to support data mining and the creation, management, and execution of machine learning models on the datasets in the **Data Lake Storage**, and data streams provided by the **Stream Processing**. It can generate notification events in a message queue for consumption by other applications. To implement this microservice, Spark ML and Scikit Learn¹⁸ can be used.

¹³<https://avro.apache.org/>

¹⁴<https://parquet.apache.org/>

¹⁵<https://www.rabbitmq.com>

¹⁶<https://www.w3.org/TR/vocab-dcat-3>

¹⁷<https://www.w3.org/TR/prov-o>

¹⁸<https://scikit-learn.org/>

Batch Processing is responsible for enabling the processing of large datasets stored in **Data Lake Storage**, providing an interface for the creation, management, and execution of tasks. The processing results can be saved in **Data Lake Storage** or published in a message queue so that they can be consulted by APIs and visualization tools, or consumed by other applications. For batch data processing, MapReduce can be used.

Data Visualization supports graphical interfaces for presenting and analyzing data stored in **Data Lake Storage** and data generated by **Stream Processing** and **Smart Analyzer**. It allows users to create reports and dashboards. This service can be made available using tools such as Apache Superset¹⁹ or Kibana²⁰.

The **API Gateway** provides for external applications a single access point (with load balancing) to the others microservices. The communication between the **API Gateway** and the external applications must use an encrypted communication channel. Furthermore, every request must consult the **Auth Server** microservice, which authenticates the user and generates a cryptographic token to be used in future requests. After successful authentication, the request is enriched with user information and forwarded to the targeted microservice. The gateway can be implemented using Apache Knox²¹ or Kong²², for example. The **Auth Server** can be implemented using Apache Syncope²³ or Auth0²⁴.

The **Public UI** is responsible for enabling users to access the web interface of other microservices in a single web interface. This way, each microservice offers the type of interface that suits it best. This microservice is public and uses a single access point, the **API Gateway**, with an encrypted communication channel and identified user. Among the technologies that can be used to implement the Web interface and integrate it with the functionalities provided by other microservices are NodeJS²⁵, React²⁶, and HTML5.

There are several concerns that an implementation of these microservices have to address. For example, Data Ingestion must be able to handle large data volumes and heterogeneity. Metadata Management must couple with structural and semantic data changes, to provide data compatibilization. Data Query must provide good performance (with automatic indexing, caching, etc.). API Gateway must balance load to handle requests in the most efficient manner. All the microservices must be scalable and fault-tolerant, while ensuring data security and privacy.

To increase the productivity and reliability of the platform's development, we suggest reusing free software in its implementation, mainly those that have an active community of developers; support to stable versions; evolutionary versions; a rich documentation; as well as integration with security and data privacy tools. To facilitate the operation and maintenance of the platform, we suggest the adoption of DevOps best practices, such as continuous integration, continuous delivery, continuous deployment, and monitoring.

3.2. Comparison with Related Architectures

Our architecture supports functionalities for data ingestion and physical integration using approaches similar to those of the works analyzed in Section 2. However, it extends the

¹⁹<https://superset.apache.org>

²⁰<https://www.elastic.co/pt/kibana>

²¹<https://knox.apache.org>

²²<https://konghq.com/kong>

²³<https://syncope.apache.org>

²⁴<https://auth0.com>

²⁵<https://nodejs.org>

²⁶<https://reactjs.org>

related works by supporting some unique features, such as: a single point of access to microservices by external applications; an authentication and access authorization service; a centralizing interface for the services; the creation of new data collections based on existing ones; and compatibilization of data in collections that have suffered structural or semantic changes over time, using the metadata modification history and mapping rules.

[Mehmood et al. 2019] had also proposed using metadata to support the data integration. Their architecture works with data models to provide unified vocabulary among data sources and align syntactic and semantic differences, demanding the definition of data models for each city sector (e.g., environmental, social, and economic data). In our work, we made a more comprehensive proposal for metadata management, using DCAT to describe the data sources. Moreover, we do not enforce uniformization in data ingestion. The data is stored as it comes from the source. Then, it can be compatibilized when it is queried. Views of compatibilized data can be materialized to speed up queries.

4. Performance Evaluation

To evaluate the performance of a system that implements the proposed architecture, the Cloud Evaluation Experiment Methodology (CEEM) [Li et al. 2013] can be used. CEEM is a methodology for systematically evaluating the performance of cloud services by experiments, which can be easily replicated or extended to any environment.

The methodology proposes ten steps to evaluate a service: (1) *Requirement Recognition* – define the problem and objectives of the assessment; (2) *Service Feature Identification* – identify the services and features to be evaluated; (3) *Metrics and Benchmarks Listing* – list the metrics and benchmarks that can be used; (4) *Metrics and Benchmarks Selection* – select the appropriate metrics and benchmarks for evaluation; (5) *Experimental Factors Listing* – list factors that may impact in the experiments' evaluation; (6) *Experimental Factors Selection* – select the factors to be studied and define the acceptance criteria; (7) *Experimental Design* – design the experiments based on the previous steps; (8) *Experimental Implementation* – prepare the test environment and run the designed experiments; (9) *Experimental Analysis* – analyze and statistically interpret the experimental results; and (10) *Conclusion and Reporting*.

In the following, we present a guideline for applying CEEM to evaluate the performance of microservices of our data integration software architecture.

4.1. Requirement Recognition and Service Feature Identification

We want to assess the individual capacity of each microservice to function under both normal and above-normal workload conditions. In particular, we want to evaluate the effectiveness of self-scalability to support an increase in the number of users or simultaneous requests for the microservice's main functionalities, while keeping an acceptable quality of service.

The main functionalities to be evaluated through the experiments are: ingestion of data streams and batches received by the API in the **Data Ingestion** microservice; creating, querying, and compatibilizing metadata in the **Metadata Manager** microservice; and recovery of data from the **Data Lake Storage** in the **Data Query** microservice.

The number of users and requests per time interval to be supported by a microservice instance must be defined according to the smart city platform to which the system

is coupled. These values will be used as parameters for the execution and analysis of the experiments. Therefore, they must be defined for each microservice to be analyzed.

This guideline does not include experimental scenarios for the **API Gateway**, **Auth Server**, **Stream Processing**, **Batch Processing**, **Smart Analyzer** and **Data Visualization** microservices because we assume that open-source tools with assessed good performance will be used in their implementation. The **Public UI** will not be considered either because it is a front-end application, thus it impacts the system's general performance less than the other components.

4.2. Metrics and Benchmarks Listing and Selection

In this analysis, we will consider the catalog of metrics presented by [Li et al. 2012b]. Four metrics of the catalog are particularly appropriate to assess the performance of the microservices: CPU utilization, RAM utilization, latency, and the number of requests processed per time interval. These metrics enable us to validate whether the developed system is capable of serving the expected number of users and requests with low latency.

4.3. Experimental Factors Listing and Selection

[Li et al. 2012a] point out the operating system and container manager versions as experimental factors to be considered. CPU clock speed and number of cores, type and capacity of RAM memory, and storage capacity are important factors as well, since both software and hardware changes can affect the performance results. To avoid the microservices' instances competing for computing resources, we suggest running them in containers with limited resources.

4.4. Experimental Design and Implementation

This section describes three experiments designed to evaluate the microservices' performance. Each one of the experiments must be run with three different configurations:

Configuration 1 Execution with a single instance of the microservice with auto-scaling disabled, and a workload that gradually increases over time (both in the number of users and in the number of concurrent requests per time unit), until reaching the maximum size expected for the system. The goal in this configuration is to measure the microservice's performance under normal workloads.

Configuration 2 Execution with a single instance of the microservice with auto-scaling disabled, and a workload that gradually increases over time (both in the number of users and in the number of concurrent requests per time unit), until CPU or RAM usage is close to 100%. In this configuration, the goal is to identify the maximum number of users and concurrent requests that a single instance can handle.

Configuration 3 Execution of the microservice with auto-scaling enabled, and a workload varying between the lower and the upper limit values supported by a given fixed number of instances. In this configuration, the goal is to assess the capacity of the microservice to self-adjust to the current demand, increasing or decreasing the number of instances according to the variations in the workload.

During the execution of the experiments, information about the CPU and RAM usage, request processing time, and number of messages processed per time unit (throughput) must be collected and recorded. The number of replications of each experiment

should be defined taking into account the resources available, and the desired sensitivity and confidence of the performance indexes to be obtained from the measurements. Generally, the sensitivity increases with the number of replications of the experiment.

Experiment 1 – Data Intake Latency Start a single instance of the **Data Ingestion** microservice and trigger real-time and batch data provisioning, varying the number of the simulated users and requests over time according to the configuration being executed.

Experiment 2 – Response Time For Operations on Metadata Start a single instance of the **Metadata Manager** microservice and simulate the simultaneous execution of metadata creation, query, and compatibilization requests for different data collections, varying the number of simulated users and requests over time according to the experiment configuration being executed.

Experiment 3 – Response Time of Data Queries To evaluate the **Data Query**, first initialize the **Data Lake Storage** with data from different collections and sources until reaching a considerable percentage of use of its storage capacity, to make possible the evaluation of the response time of queries over a large volume of data. Then, for each data collection stored, create a view in the **Data Query**. After this initialization, launch an instance of **Data Query**. Use a simulator to generate and execute queries with filters and random aggregations over the pre-existing views, varying the number of simulated users and requests over time according to the experiment configuration being executed.

4.5. Experimental Analysis

To analyze data collected in the experiments and drawn conclusions, it is strongly recommended the use of statistical methods to ensure robustness [Li et al. 2013]. Nevertheless, even simple graphical tools (e.g. dot plots, histograms, and box plots), showing the response time, throughput, CPU usage, and RAM usage over time, may help to visualize how well a microservice self-adjust to workload variations.

5. Conclusion

The contribution of this work is twofold. First, it provides a panorama of the requirements of data integration for smart cities and state-of-the-art solutions. Second, it presents an architecture to help researchers and developers to approach the implementation of these requirements and also guidance to assess its performance.

We have identified functional and non-functional requirements of data integration in smart cities and the challenges involved in their implementation by analyzing the related literature. Then, we have proposed a microservices architecture for a data integration platform that meets these requirements. We have also specified the software components needed to implement them. Our software architecture extends those of related work by providing a solution for metadata management that keeps the history of changes in the structure and semantics of attributes, to enable the compatibilization of data in queries. Another differentiated feature of the architecture is the API Gateway, which provides a secure data access point for external applications (through encrypted and authenticated communication channels). We are currently working on an implementation²⁷ of this architecture on top of InterSCity²⁸, an open-source platform for smart cities.

²⁷<https://gitlab.com/intercity/data-integration.git>

²⁸<https://gitlab.com/intercity/intercity-platform>

Following the CEEM methodology, we have designed a set of experiments that can be used to evaluate the performance of the microservices of the architecture under both normal and above-normal workload conditions. With these experiments, one can assess the effectiveness of self-scalability to keep acceptable quality of service while the number of users and requests vary over time.

References

- Cheng, B., Longo, S., Cirillo, F., Bauer, M., and Kovacs, E. (2015). Building a big data platform for smart cities: Experience and lessons from Santander. In *2015 IEEE International Congress on Big Data, IEEE BigData 2015*, pages 592–599.
- Consoli, S., Mongiovi, M., Nuzzolese, A. G., Peroni, S., Presutti, V., Reforgiato Recupero, D., and Spampinato, D. (2015). A smart city data model based on semantics best practice and principles. In *24th Intl. Conference on World Wide Web, WWW'15*.
- Costa, C. and Santos, M. Y. (2017). The SusCity big data warehousing approach for smart cities. In *Proceedings of the 21st International Database Engineering & Applications Symposium, IDEAS 2017*, page 264–273. ACM.
- Hashem, I. A. T., Chang, V., Anuar, N. B., Adewole, K., Yaqoob, I., Gani, A., Ahmed, E., and Chiroma, H. (2016). The role of big data in smart city. *International Journal of Information Management*, 36(5):748 – 758.
- Li, Z., O'Brien, L., and Zhang, H. (2013). CEEM: A practical methodology for cloud services evaluation. In *IEEE 9th World Congress on Services*, pages 44–51.
- Li, Z., O'Brien, L., Zhang, H., and Cai, R. (2012a). A factor framework for experimental design for performance evaluation of commercial cloud services. In *4th IEEE Intl. Conf. on Cloud Computing Technology and Science Proceedings*, pages 169–176.
- Li, Z., O'Brien, L., Zhang, H., and Cai, R. (2012b). On a catalogue of metrics for evaluating commercial cloud services. In *2012 ACM/IEEE 13th International Conference on Grid Computing*, pages 164–173.
- Mehmood, H., Gilman, E., Cortes, M., Kostakos, P., Byrne, A., Valta, K., Tekes, S., and Riekki, J. (2019). Implementing big data lake for heterogeneous data sources. In *IEEE 35th Intl. Conference on Data Engineering Workshops (ICDEW 2019)*, pages 37–44.
- Psyllidis, A., Bozzon, A., Bocconi, S., and Titos Bolivar, C. (2015). A platform for urban analytics and semantic data integration in city planning. In *Computer-Aided Architectural Design Futures. The Next City - New Technologies and the Future of the Built Environment, CAAD Futures 2015*, pages 21–36. Springer Berlin Heidelberg.
- Raghavan, S., Simon, B. Y. L., Lee, Y. L., Tan, W. L., and Kee, K. K. (2020). Data integration for smart cities: Opportunities and challenges. In Alfred, R., Lim, Y., Haviluddin, H., and On, C. K., editors, *Computational Science and Technology*, pages 393–403. Springer Singapore.
- Rathore, M. M., Ahmad, A., Paul, A., and Rho, S. (2016). Urban planning and building smart cities based on the internet of things using big data analytics. *Computer Networks*, 101:63 – 80.