



**University of Aveiro**  
2015.

Department of Electronics, Telecommunications and  
Informatics

**Vedran  
Semenski**

## **SMARTIE – Secure and Smarter Cities Data Management**

Dissertação apresentada à Universidade de Aveiro para cumprimento dos requisitos necessários à obtenção do grau de Mestre em (designação do Mestrado), realizada sob a orientação científica do Doutor (nome do orientador), Professor (categoria do orientador) do Departamento de (designação do departamento) da Universidade de Aveiro

texto Apoio financeiro do  
(se aplicável)

texto Dedico este trabalho à minha esposa e filho pelo incansável apoio.

(opcional)

## **o júri**

presidente

**Prof. Doutor João Antunes da Silva**  
professor associado da Faculdade de Engenharia da Universidade do Porto

**Prof. Doutor João Antunes da Silva**  
professor associado da Faculdade de Engenharia da Universidade do Porto

**Prof. Doutor João Antunes da Silva**  
professor associado da Faculdade de Engenharia da Universidade do Porto

**Prof. Doutor João Antunes da Silva**  
professor associado da Faculdade de Engenharia da Universidade do Porto

**Prof. Doutor João Antunes da Silva**  
professor associado da Faculdade de Engenharia da Universidade do Porto

**Prof. Doutor João Antunes da Silva**  
professor associado da Faculdade de Engenharia da Universidade do Porto

**Prof. Doutor João Antunes da Silva**  
professor associado da Faculdade de Engenharia da Universidade do Porto

## **agradecimentos**

texto O presente trabalho propõe-se divulgar as mais significativas técnicas de construção existentes em Portugal continental. O livro é composto por uma apresentação dos materiais tradicionais de construção (suas principais características), uma compilação de fichas técnicas (de carácter prático, uma vasta bibliografia comentada e um glossário de termos técnicos. A importante colaboração de diversas personalidades ligadas à área da História da Arquitectura, bem como o levantamento fotográfico realizado contribuem para o conhecimento e valorização de um saber tradicional.

(opcional)

**palavras-chave**

texto livro, arquitectura, história, construção, materiais de construção, saber tradicional.

**resumo**

texto O presente trabalho propõe-se divulgar as mais significativas técnicas de construção existentes em Portugal continental. O livro é composto por uma apresentação dos materiais tradicionais de construção (suas principais características), uma compilação de fichas técnicas (de carácter prático, uma vasta bibliografia comentada e um glossário de termos técnicos.

A importante colaboração de diversas personalidades ligadas à área da História da Arquitectura, bem como o levantamento fotográfico realizado contribuem para o conhecimento e valorização de um saber tradicional.

**keywords**

ABAC, Access Control, IoT, JSON, M2M, NoSQL, SMARTIE, XACML

**abstract**

texto O presente trabalho propõe-se divulgar as mais significativas técnicas de construção existentes em Portugal continental. O livro é composto por uma apresentação dos materiais tradicionais de construção (suas principais características), uma compilação de fichas técnicas (de carácter prático, uma vasta bibliografia comentada e um glossário de termos técnicos.

A importante colaboração de diversas personalidades ligadas à área da História da Arquitectura, bem como o levantamento fotográfico realizado contribuem para o conhecimento e valorização de um saber tradicional.

# Table of content

Table of content .....	I
List of Acronyms .....	I
List of figures .....	III
List of tables .....	IV
1 Introduction .....	1
1.1 Dissertation description text.....	1
1.2 Basic concept .....	1
2 State of the art.....	3
2.1 The Internet of Things (IoT) .....	3
2.2 Machine to Machine communication (M2M) .....	5
2.3 Big Data .....	6
2.4 NoSQL.....	7
2.4.1 Types of NoSQL databases .....	7
2.4.2 General overview of the technologies available at the moment .....	9
2.4.3 Security issues in NoSQL.....	14
2.4.4 Performance studies and comparisons .....	14
2.4.5 NoSQL Conclusion.....	14
2.5 Access Control.....	15
2.5.1 Types of access control.....	15
2.5.2 OASIS, XACML and JSON.....	17
2.5.3 Current Technologies and open source solutions .....	17
2.6 SMARTIE.....	17
2.6.1 Brief Overview .....	18
2.6.2 Conclusion .....	19
3 Project description .....	20
3.1 Introduction.....	20
3.1 General architecture and design choices .....	20
3.1 Development plans and overview .....	21
3.2 Potential improvements/achievements.....	22

3.3	Conclusion .....	23
4	(Main) Development implementation.....	24
4.1	Final structure description in detail.....	24
4.1.1	Global architecture .....	24
4.1.2	Policy database .....	24
4.1.3	Sensor data database.....	24
4.2	Data storage modules .....	24
4.3	Access control module .....	24
4.3.1	PAP.....	24
4.3.2	PIP & PRP .....	24
4.3.3	PEP .....	24
4.3.4	PDP.....	24
5	Testing .....	25
5.1	Test scenatio1 .....	25
5.2	Test scenario2 .....	25
6	Integration in SMARTIE.....	26
6.1	Proposed implementation.....	26
6.2	Challenges.....	26
6.3	Solution.....	26
6.4	Testing results .....	26
7	Conclusion.....	27
8	References .....	28



# List of Acronyms

ABAC	Attribute-Based Access Control
ACID	Atomicity, Consistency, Isolation and Durability
API	Application Program Interface
BASE	Basically Available, Soft State and Eventually Consistent
BSON	Binary JSON
CQL	Cassandra Query Language / Cypher Query Language
DAC	Discretionary Access Control
DBMS	Database Management System
HDFS	Hadoop Distributed File System
IDE	Integrated Development Environment
IBAC	Identity Based Access Control
IEEE	Institute of Electrical and Electronics Engineers
IoT	Internet of Things
JDBC	Java Database Connectivity
JSON	JavaScript Object Notation
M2M	Machine to machine
MAC	Mandatory Access Control
NFC	Near Field Communication
NoSQL	Not Only SQL
OASIS	Organization for the Advancement of Structured Information Standards
OLAP	Online Analytical Processing
OLTP	Online transaction processing
OWASP	Open Web Application Security Project
PAP	Policy Administration Point
PDP	Policy Decision Point
PEP	Policy Enforcement Point
PIP	Policy Information Point

PII	Personal Identifiable Information
PRP	Policy Retrieval Point
RBAC	Role Based Access Control
REST	Representational State Transfer
RDBMS	Relational Database Management System
RFID	Radio-Frequency Identification
SMARTIE	Smart City
SQL	Structured Query Language
W3C	World Wide Web Consortium
UDAF	User-Defined Aggregate Functions
UDF	User-Defined Functions
UDTF	User-Defined Table Functions
UWB	Ultra-Wide Band
XML	eXtensible Markup Language
XACML	eXtensible Access Control Markup Language
YARN	Yet Another Resource Negotiator

List of figures

Figure 1. High - level overview of the architecture ..... 2

Figure 2. The 5 Vs of Big Data [13]..... 6

Figure 3. Architecture of the proposed solution ..... 21

## List of tables

Nisu pronađeni unosi u tablici slika.

# 1 Introduction

TODO: making the introduction and describing how the document is organized into parts with brief descriptions of each part

## ***1.1 Dissertation description text***

SMARTIE (<http://www.smartie-project.eu/index.html>) is an European project where PT – Inovação e Sistemas is one of the partners. The project is aimed at creating a framework (IoT) to collect and share large volumes of heterogeneous sensor information (Big Data) to be used in smartcities applications. In order to protect sensitive data, a security component is necessary to manage authentication and authorization which will be based on the XACML standard. The original XACML standard is based on XML but in SMARTIE we will try to use JSON.

### **Objectives:**

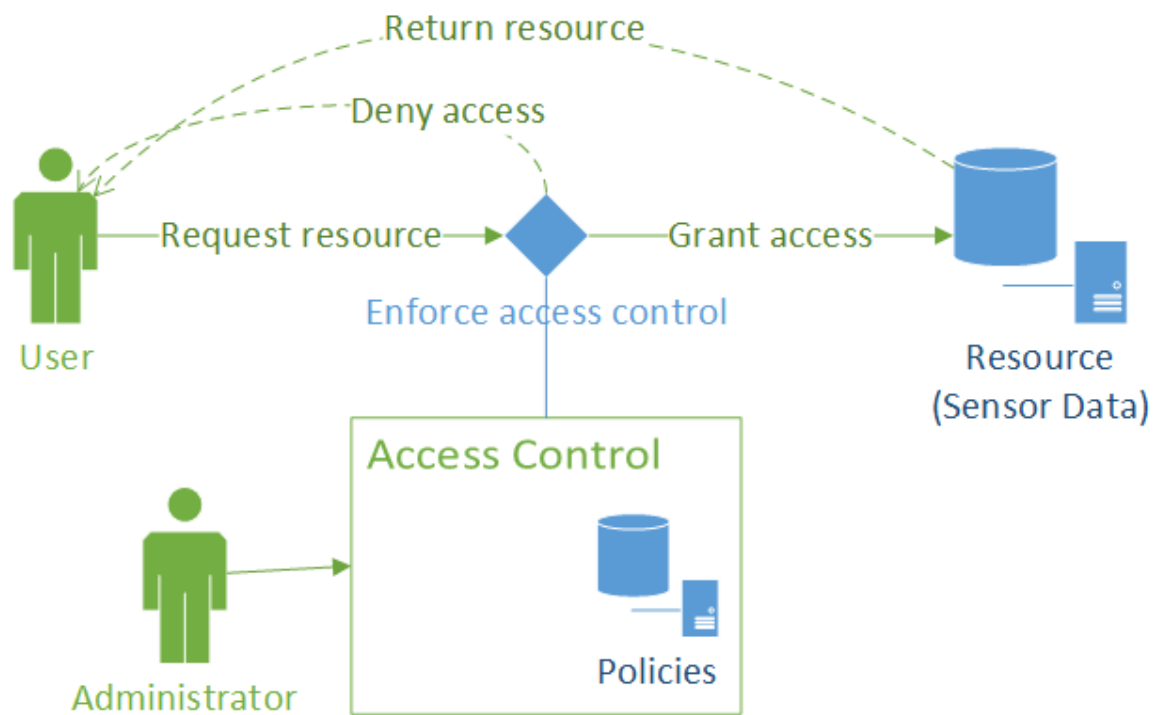
This dissertation is focused on two main development activities: databases and security framework. It comprises the following tasks to be executed:

- Become familiar with SMARTIE;
- Become familiar with IoT;
- Become familiar with ABAC models and particularly with XACML;
- Become familiar with NoSQL databases (Big Data);
- Selection of the NoSQL database to be used to store the collected sensor data;
- Selection of the database to store the security policies to be enforced;
- Conceptual, logical and physical models for both databases;
- Security framework based on XACML.

All the work will be developed at Instituto de Telecomunicações facilities.

## ***1.2 Basic concept***

The main goal is to build an access control framework using ABAC and the OASIS XACML/JSON standard. Along with this one NoSQL storage system will be setup for storing sensor data and another one for storing policies. This will be built in the context of the SMARTIE project but as a separate and independent component that will allow the enforcement of access control in any place needed inside the SmartData platform (that is being developed for the SMARTIE project). The best places to enforce access control at the moment would be when requesting access to sensor data directly.



**Figure 1. High - level overview of the architecture**

A high - level overview diagram of the proposed solution can be seen in Figure 1. It shows how the solution will be integrated to enforce access control and how will it be used by users/system administrators.

## 2 State of the art

This section contains short overviews of areas related to the work done in this dissertation. The overviews contain brief descriptions, analysis of the current state and work related to those areas along with a future oriented perspective.

TODO: elaborate a bit more

### 2.1 *The Internet of Things (IoT)*

In this section the current state of IoT will be presented. The main areas that will be briefly presented and analyzed are: current state and overview, general concerns, recent work and studies, technologies used, implementation examples and expected future developments.

The IoT is a recent paradigm in the area of networks and communication that has had a lot of growth and new developments in recent years. Although there are a lot of definitions of the IoT and the somewhat changing and branching nature of the development trends in this area the basic and simplified idea of this concept is that nowadays everyday objects can be equipped with cheap microcontrollers, sensors, means of connecting to one another and the Internet. The devices can generally be divided into two categories: sensors and actuators, meaning that their purpose is to provide and share data or some kind of readings or to receive commands and react/complete actions accordingly. The overall implementation and use of this could be used in a number of applications including: home automation, industrial automation, medical aids, mobile health care, elderly assistance, intelligent energy management and smart grids, automotive, traffic management and many others [1] [2]. All of these offer beneficial and significant impact on almost all areas of everyday life and have potential for providing advancements in research areas unrelated to networking or computer science. The definition of IoT is not exact as the authors of [1] wrote. They explain that the two main views come from the name "Internet of Things". One vision views the concept from a network perspective concentrating on communication and connection problems while the other is oriented on the "Things" or object perspective concentrating on sensor technologies, new communication technologies like RFID and NFC ,and integration into other devices in a seamless and affordable way. A third view is also present which the authors defined as a "Semantic oriented vision" which is concentrated more on the use, implementation and processing of data.

A number of challenges are in the way of successfully building and utilizing the IoT. Scalability is one of the obvious challenges. Any kind of IoT application that requires a large number of devices commonly face problems with response time, memory, processing and energy constraints. Other challenges include security issues. The data being generated by sensor networks is huge and could over saturate the network, the data could be personal and as such has to be protected from unwanted access. Attacks by hackers, malicious software, viruses and other sources can also disturb the flow and integrity of data/information. As the authors of [3] [4] [5] describe in their work and stress the importance of security for a widely acceptable solution. In their work they describes the IoT divided into 3 or 4 layers and define security requirements for every layer providing the current state of technologies used in these areas. They stress the need for an uniformed and

standardised open architecture and solution. Solutions for many of these problems are already conceptually known and are in some examples implemented but the lack of a open and standardized solution is certainly an issue that would proved to be beneficial if/when solved.

Technological advancements in various sensor modules and cheap and energy efficient microcontrollers along with recent communication technologies like RFID, NFC and network protocols are the main factors that enabled the fast development and spreading the IoT. Because of these advancements the concept became a reality and is gaining importance and more uses and applications. As the means of connection and communication are open and utilizes a number of older and newer technologies the networks and number of devices is fast growing and also brings up the problem of standardisation and implementation of standards in the purpose of uniformly solving known problems with for instance security, integrity and scalability.

A more future oriented view and analysis is provided by [6]. It describes a more human centric rather than thing centric future of the IoT. Devices being linked to people and monitoring their state and condition. Others are used for monitoring or controlling things in the environment but always in close relation to human needs and/or wants. It describes a need for IoT applications to provide better quality of service and seamless integration into areas of life. The more relevant and faster growing application domains are: Environmental Monitoring, Smart Retail, Smart Agriculture, Smart Energy and Power Grids and Smart Healthcare. The architecture proposed is separated into 5 layers stacked one on top another in this order:

1. Things - containing devices or elements that are data generators and/or consumers of information. The devices could range from small and unintelligent embedded systems to complex devices or virtual entities. The communication would be done different communication technologies and a wide variety of protocols.
2. Network - this layer contains functionality and means of managing a sensor network. Discovery of new devices, maintenance, scalability, universal abstractions of the Inputs and Outputs and general abstraction for the upper and lower layers and enabling plug-n-play like use using already known models like push/pull or publish/subscribe models and REST based protocols.
3. Data Management - it is defined as "Big-Little" Data Management referring to the data usually generated by sensor networks meaning that the data generated by for instance temperature sensors is small in individual reading size large considering a large number of sensors over a period of time. This layer is responsible for categorizing and aggregating data retrieved from the Network layer in order for it to be used and/or by the Analytics layer. Data generated by sensors is often slow changing so this fact should be taken advantage of for more efficient data storage.
4. Analytics - this layer mines/retrieves data from the Data Management layer and performs data processing and analysis depending on the type of data and end user of the result. It provides the applications (which would be located on top) with useful data and information for subsequent use.

There a lot of commercial implementations of IoT on a smaller scale which are mainly focused on personal use in home energy consumption, health monitoring and environment monitoring applications. These



solution usually utilize custom solutions along with custom hardware (regarding the sensor devices). Introducing standards that would be accepted and the creation of publicly available frameworks that utilize those standards would be very beneficial and would allow easier further developments. The lack of said standards and frameworks mean during initial developments of applications many of them face the same problems and implement a custom solution. This leads to incompatibility issues and stand in the way of a truly global IoT. The surveys and proposals given in [5] [3] give a good overview of the current situation in IoT and a give proposals for more uniformed future research and developments.

## ***2.2 Machine to Machine communication (M2M)***

This section will provide a brief overview of the current state, issues and a future developments in M2M. M2M, which means "Machine to machine" refers to the technologies used for communication between devices. It is a broad term that can sometimes be used for terms like "Machine to Mobile" and "Man to Machine" which mainly refer to the implementation in a more precise manner. It is also commonly regarded in the context of IoT because it is essential part if IoT. IoT concentrates more on a higher overview considering network problems and viewing object in a more general way, while M2M deals more with one-to-one communication between devices and the functioning of devices regarding one another, less regarding the overview from a network perspective [7]. The motivation behind M2M comes as stated in [8] mainly comes from two observations:

- 1) a networked machine is more valuable than an isolated one
- 2) when multiple machines are effectively interconnected, more autonomous and intelligent applications can be generated
- 3) smart and ubiquitous services can be enabled by machine-type devices intelligently communicating with other devices at anytime and anywhere

M2M systems are mainly used in the areas of personal health monitoring and smart homes/smart houses. M2M communication can be achieved with a number of technologies but the most significant and most promising ones are wireless technologies. In [8] the main technologies that are mentioned and review are: Zigbee, Bluetooth, UWB, IEEE 802.15.6, Wi-Fi, HomeRF, 60GHz transmission and Visible light communications. They offer easy integration of devices and many options for developing application. Applications utilizing M2M systems are usually focused on improving the quality of life for individuals. The two areas that are most significant and are the main drivers for developments in M2M are "personal health monitoring" and "smart home" applications. Trends in M2M as mentioned in [8] suggest exponential growth in the number of M2M-enabled devices. From 50 million in 2008 and over 200 million in 2014, it is expected to grow up to 50 billion in 2020. Along with communication and integration challenges, security is one of the more important areas that need to be deled with regarding M2M.

The architecture of M2M systems can generally be divided into 3 layers.

- 1) Terminal layer - contains the access gateway , area network and M2M nodes, M2M enabled devices
- 2) Network layer - responsible for transmitting data between the other two layers

- 3) Application layer - contains the application that utilizes the M2M system

Security issues in M2M systems are similar to ones found in providing security and privacy in communication over other networks, the internet... Additional problems that are more specific to M2M occur in the Terminal layer. Devices could potentially be easily accessible to attackers. This means that they could be tampered with or be controlled by attackers and false data could be injected threatening the overall performance of the M2M system and application.

M2M systems in the IoT context provide the means to build applications that can further the developments in many scientific areas along with bettering personal quality of life. Solving security issues and the general standardisation of undefined aspects in communication, as well as the development of quality and publicly available solutions are the future steps that will enable the wide commercial use of these technologies.

Work by the authors of [9] [8] [10] can provide a more detailed overview in M2M.

## 2.3 *Big Data*

Storing, processing, accessing and managing vast amounts of data with unreliable or complex structure all fall in the term Big Data. Although the term refers only to a large amount of data, the challenges that come with handling and using of it also come hand in hand with this term. Big data is often described with the 4 or 5 Vs of Big Data (depending on different sources) as can be seen in the works of [11] [12] and can be seen in Figure 2.



**Figure 2. The 5 Vs of Big Data [13]**

**Volume** refers to the amount of information or data that is being generated and needs to be handled. Because of the sheer amount or just the type, traditional systems (referring to RDBMS) don't offer appropriate solutions and different methods/methods need to be taken into account. The authors of [11] predict that the size of the data being generated to reach the range from petabytes to petabytes.

**Variety** refers to the data's structure. It can vary from a traditional sense of well structured and predictable data to semi-structured or unstructured and/or changing, coming from a variety of resources like Documents, email, Web Pages, Sensor Networks, Social Media etc.

**Velocity** refers to the rate of data flow. This is quite an open definition encompassing both the rate of data coming from different sources, but also rate of data flow in general.

**Variability** refers to the inconsistency of velocity in general. This is an issue that is hard to deal with and is best understood when thinking of the usage of social networking.

**Value** User can run certain queries against the data stored and then user got important results from the filtered data obtained and can also rank it according to the dimensions.

The authors of [11] also added **Complexity** to this as an important and sometimes forgotten aspect referring to the complexity of linking, matching, cleansing and transforming of data which is coming from various sources.

Along with fundamental challenges in Big Data, regarding functional issues, security has to be taken into account if implementations are to be realized. The biggest security challenge in Big Data is as identified in [14] is the protection of user's privacy. Vast amounts of personal identifiable information (PII) that are stored in unstructured databases (NoSQL) means that the location of all information is not always known unambiguously and even direct access is somewhat abstracted. Leaving everything as is and without applying security measures such as Access Control leaves things open for abuse.

## ***2.4 NoSQL***

NoSQL, or otherwise known as "Non only SQL" refers to databases or data management systems which are designed to handle data management problems where conventional RDBMS solutions cannot cope for various reasons [15] [11]. These problems are usually related to: handling large amounts of data and the processing of it, high number of operations, specific types of operations or needs and others. Commonly, NoSQL systems are designed for large scale data storage and parallel processing over a large number of servers. Some systems provide APIs that support SQL or SQL-like languages and convert them into native non-SQL languages and use mechanisms different to ones in RDBMS. Because they provide the ability to handle large amounts of data it usually comes at the price of fully ACID (Atomicity, Consistency, Isolation, Durability) characteristics. This can be explained by the CAP (strong Consistency, high Availability, Partition tolerance) theorem which can be seen in more detail in [16], and because of this most systems can be described as BASE (Basically Available, Soft-state, Eventually consistent).

### **2.4.1 Types of NoSQL databases**

NoSQL databases can be classified in four basic categories as done in [16].

- Key-Value stores
- Document databases (or stores)
- Wide-Column (or Column-Family) stores
- Graph databases

#### ***Key-Value stores***

These storage systems are organized in simple, standalone tables organized like "hash tables". The items stored in tables are key-value pairs where the key is an alpha-numeric identifier and the value is one or a set of values associated with that identifier. As the organisation of table is a "hash table" the limitation that is present is that they are usually limited to only "exact match" type of queries or allowing "<,>" type of operations with a significant reduction in speed. On the other hand read operations are very fast, as to be expected from a hash table data set, and because the keys can also be viewed as the addresses of the value wanted to be retrieved, even data from the same table can be distributed over several locations so these storage systems have linear characteristics regarding scalability.

### ***Document databases***

Document based storage systems, as their name implies, are designed to store data in documents. They use standard data exchange formats such as XML, JSON or BSON to store the data in documents and can distribute these documents on multiple locations. These are considered scheme-less databases because the storage format or storage data structure can be loosely defined. Single columns or single data entries can house hundreds of values and the number or type of values stored can vary from row to row. These are good for storing and managing big collections of documents containing significant amounts of data like text documents, emails, XML documents or objects containing large amounts of values and data. Along with that they are also convenient for storing sparse data collections because of their schema-less data structure. This means that the usual filling out with null values (that is traditionally done in RDBMS) is not necessary and means that the overall amount of space used is correlated to the amount of data stored inside the database. These solutions offer great scalability, unlike key-value based stores allow multiple "< >" types of comparators and both keys and values are fully searchable. Although they can offer MapReduce [17] features they tend to have slow response times to queries. This reason is because fetching data with multiple set parameters for values means reading and parsing data from whole documents.

### ***Wide-Column Stores***

Wide-Column (or Column-Family) stores are somewhat in between document based and key-value based storage systems. They have a structure similar to a key-value structure but allow multiple values and require at least one identifier column which fills the role of a primary key. They can form multiple indexes upon other values and allow "equal type" comparisons over the value attributes. Because of the similarity to key-value based systems they share the same faults regarding "< >" types of operations. The similarity to document based systems comes because they are distributed. The key value can be used to distribute data from a table to multiple locations. This offers good characteristics regarding scalability. Reading and writing operations are fast so they are specially suited for MapReduce operations and parallel processing of large amounts of data which is their main purpose and use.

### ***Graph Databases***

Graph databases, by basic concept are relational databases but still are very different from RDBMS in the sense that they also have relations. The relations themselves can be considered as more important

because these are used when we mainly need to store data regarding the relationships and dependencies between objects rather than information about the objects themselves. They store data similar to object-oriented databases as use objects as network nodes which have relationships (edges) and properties or object attributes stored as key-value pairs. The relationships can also have different attributes or properties attached to them. They are suited for storing and visualizing data regarding graphs, networks etc.

## **2.4.2 General overview of the technologies available at the moment**

This section contains an overview of the current "popular choices" in NoSQL data storage and management systems, a brief comparison given from results provided by outside sources and general conclusion.

### ***Hadoop***

This Apache project is a platform for developing open-source software for solving scalability and processing problems in the context of large quantities of data. Data storage supported by Hadoop falls in the wide-column paradigm family of NoSQL systems. The platform has a large developer community and many open projects and because of that it is one of the biggest and fastest evolving platforms. On the Hadoop website it is stated that "The Apache Hadoop project develops open-source software for reliable, scalable, distributed computing" [18]. It is also stated that it is essentially a framework that allows distributed processing across clusters of computers, it is designed for scalability and can easily be scaled from one server to thousands of machines, it handles failures on the application layer and as a result provides a highly-available service on top of clusters, regardless of individual failures on machines inside the cluster.

These are the projects main modules as stated on it's official webpage [18]:

- Hadoop Common: The common utilities that support the other Hadoop modules.
- Hadoop Distributed File System (HDFS™): A distributed file system that provides high-throughput access to application data.
- Hadoop YARN: A framework for job scheduling and cluster resource management.
- Hadoop MapReduce: A YARN-based system for parallel processing of large data sets.

The Hadoop ecosystem of open source is built in a way that allows easy porting and migrating between different storage solutions and also data processing solutions. Therefore a project that is being developed over Hadoop has a bigger value and feature set than the value/feature set than of that single solution. Many additional options are available because of the ecosystem that it is a part of which therefore adds to the value and possible number of implementations. Also, as most Apache projects, it has a large development community so updates and improvements are frequent.

### ***Apache HBase***

Apache HBase is one of the projects built on top of Hadoop, more specifically the HDFS. It falls in the wide-column family of NoSQL database systems and is a distributed database system. As stated on the official pages [19] it is more a "Data Store" than "Data Base" because it lacks many of the features you find

in an RDBMS, such as typed columns, secondary indexes, triggers, and advanced query languages, etc. The way it is built allows it to have linear scalability characteristics and because it is a part of the Hadoop ecosystem of open source projects it is also very modular as the data stored can be used by other data processing systems and it can also be migrated to other solution if there is a need for it.

Most notable HBase features as stated on the official website [19] are:

- Strongly consistent reads/writes: HBase is not an "eventually consistent" DataStore. This makes it very suitable for tasks such as high-speed counter aggregation.
- Automatic sharding: HBase tables are distributed on the cluster via regions, and regions are automatically split and re-distributed as your data grows.
- Automatic RegionServer failover
- Hadoop/HDFS Integration: HBase supports HDFS out of the box as its distributed file system.
- MapReduce: HBase supports massively parallelized processing via MapReduce for using HBase as both source and sink.
- Java Client API: HBase supports an easy to use Java API for programmatic access.
- Thrift/REST API: HBase also supports Thrift and REST for non-Java front-ends.
- Block Cache and Bloom Filters: HBase supports a Block Cache and Bloom Filters for high volume query optimization.
- Operational Management: HBase provides build-in web-pages for operational insight as well as JMX metrics.

As HBase is built on top of HDFS it provides additional functionalities. As HDFS is a general purpose file distribution systems it has its limitations regarding usability ("from a developers perspective"). HBase stores data inside a tables and rows. This is familiar to anyone who worked with standard RDBMSs although this is almost the only similarity.

## ***Apache Hive***

Hive is also one of the open source projects developed on top of Hadoop. It is a mainly a distributed data warehouse system. The list of organisations using Hive as stated on the official website [20] include: eBay, Facebook, LinkedIn, Spotify, Taobao, Tencent, and Yahoo!. As an open source project, Hive has a strong technical development community working with widely located and diverse users and organizations. Hive was originally designed as a translation layer on top of Hadoop MapReduce. As a query language it has its own variant called HiveQL (Hive query language) that will be familiar to anyone already familiar with SQL. It also allows for easy use of the MapReduce framework for more complex analysis and it can be extended with custom scalar, aggregation and table functions (UDFs, UDAFs and UDTFs). It does not offer real time queries or row-level updated and as such is best suited for batch style jobs (over large sets of data). As stated on the official website the main advantages of Hive are [20]:

- scalability (scale out with more machines added dynamically to the Hadoop cluster),
- extensibility (with MapReduce framework and UDF/UDAF/UDTF),
- fault-tolerance
- loose-coupling with its input formats.

## ***Apache Spark***

Spark is also an Apache open source project since 2010. It is a fast processing engine that has some advances over the standard Hadoop MapReduce jobs. It utilizes in memory processing and data caching for faster performance compared to MapReduce. The official site states [21] that it is 100x faster in memory and 10x faster on disk than MapReduce. It is designed to access data from other Apache Hadoop projects like HDFS, Cassandra and HBase and to perform both batch processing (similar to MapReduce) and new workloads like streaming, interactive queries, and machine learning. From the period of January 5<sup>th</sup> 2014. to January 5<sup>th</sup> 2015. it had 6906 commits and 419 contributors [22] making it the most (or one of the most) active among Apache and Big Data open source projects in general. It supports a variety of popular development languages including Java, Python and Scala. This project therefore is already a great solution for many large data processing applications and has great promise for future developments.

## ***Apache Cassandra***

Apache Cassandra Is again a Apache open source project in the Hadoop "ecosystem" of projects. It is mainly a data warehouse and it falls in the row-oriented wide-column family of NoSQL database systems. Same as other projects it has support from other systems in the Hadoop "ecosystem" as the data from Cassandra can be used by, for instance Spark, and processed. Porting/ migrating from other sources should also be easy. It is suited for storing large amounts of data as it also has linear characteristics regarding scalability. Main benefit of using Cassandra is a SQL like query language called CQL ("Cassandra Query Language") and functionalities like column indexes allowing more efficient storage and data manipulations depending on the structure defined by the user. Cassandra is very suitable for environments such as storing sensor data as it is scalable and has proven fault-tolerance on commodity hardware or cloud infrastructure [23]. The distribution of data across nodes is primarily organized according to the value of the primary key. Data which has the same value of the primary key will be located on the same node allowing fast reads. This also means that related data can easily be controlled and stored in one node, therefore allowing for fast access for data processing if needed.

## ***MongoDB***

MongoDB is one of the most popular choices in the document based NoSQL family of database systems. It stores data in a JSON like format called BSON in documents and because of this is very well suited for object mapping and storing unstructured data. Storing data in mongo consists of defining objects with attributes that contain data and unlike other NoSQL systems like most wide-column solutions it doesn't suffer the problem of knowing the structure beforehand. The structure can be changed at any point and as long as the implementation (the application using MongoDB) can cope with this, problems shouldn't occur.

Another good point of MongoDB is that it is well suited for use on a large variety of machines, not necessary on high performance ones. A disadvantage on using MongoDB is that it is somewhat slow regarding reading and modifying the data compared to for instance wide-column solutions but at its main purpose is to be used in implementations where the data structure/schema is frequently changing and/or unknown at the beginning.

Some of the most significant features/characteristics of MongoDB [17] [24]:

- Indexing - primary and secondary indexes are available. The benefits of indexing are the same as those found in traditional RDBMSs.
- Advanced queries - MongoDB supports range search, field search, regular expressions and other functionalities commonly lacking in other scalable NoSQL systems.
- Replication - making replicas increases the availability of data. replicas can either have a primary or a secondary role similar to a master-slave methodology.
- Load balancing and fault tolerant - using sharding based on user defined shard keys, data collections are split into ranges and distributed. The user chooses a shard key, which determines how the data in a collection will be distributed. This can be run over multiple machines and the shards are also replicated and therefore providing fault tolerance.
- File storage - as MongoDB is a document-based NoSQL distributed database system it can also be used as a distributed file storage system. This functionality is called GridFS and it is included in MongoDB.
- Aggregation - MapReduce jobs can be used for aggregation purposes. For instance achieving the usual GROUP BY functionality from SQL in RDBMSs.

## ***Neo4j***

Neo4j is currently one of most popular from the graph-based NoSQL family of database systems. It is mostly useful in applications where relations between entities are the focus of implementation. It is used by a lot of companies and is supported on many platforms and frameworks so easy integration of Neo4j in almost any application in need of a graph database should be easy. This is maybe the biggest benefit of Neo4j over other possible products along with the large set of functionalities and ecosystem of tools and libraries.

Neo4j most significant features [25] [26]:

- It has a SQL like query language called CQL ("Cypher Query Language")
- It follows Property Graph Data Model
- Advanced features like: Indexing, UNIQUE constraints, transactions, ACID compliant
- It uses Native graph on disk storage with Native GPE(Graph Processing Engine)
- It supports exporting of query data to JSON and XLS format
- Support for many platforms, languages, simple to use APIs...
- High-speed traversals in the node space



Disadvantages Neo4j are mostly regarding selling strategies, expenses, the lack of features in the free "Community" version and issues with dual licenses. Regarding technical disadvantages there aren't many. Compared to other types of databases like wide-column it is less scalable and not suited for storing large amounts of data which is to be expected considering all of the features and that it's a graph-based NoSQL database system.

## ***OrientDB***

OrientDB is a newer NoSQL system (since 2012). This data management system is a mixture of document-based and graph-based NoSQL families. It stores data similar to document-based systems like MongoDB and combines it with the relations found in graph databases like Neo4j. On the official site of OrientDB [27] they have direct comparisons to MongoDB and Neo4j and state clear advantages in comparison to both alternatives. It combines advantages from both graph and document based solutions solving most of the problems by adding relations and fast traversal over relations to classic document-based solution and flexibility, scalability and sharding to a graph-based solutions. Therefore this solution has much promises and can be applied to a wide variety of implementation problems and seems like a perfect blend between MongoDB and Neo4j and a good alternative to both of those and traditional RDBMSs.

Most significant features include [27]:

- Combines best features from document and graph based systems
- OrientDB SQL - supports SQL with some extensions for handling trees and graphs
- Multi-Master Replication - sharding, fault tolerance, availability, scalability...
- ACID Compliance - fully ACID transactions across distributed data storage system
- Community Edition is free - free even for commercial use and contains all of the most significant features as the Enterprise edition except customer support, backups, profiling and similar features

## ***Redis***

Redis is a key-value based NoSQL data store. Unlike other solutions it stores data directly in memory and uses the disk's storage only for persistence, so reading and writing is extremely fast but storage space is limited to memory. It has master-slave replication and can support any number of slaves [28]. One of the more significant advantages over other systems is the supported rich set of data types including: strings, lists, sets, hash tables, sorted sets and others.

Other features that are stated on the official website include [29]:

- Transactions
- Pub/Sub
- Lua scripting
- Keys with a limited time-to-live
- LRU eviction of keys
- Automatic failover

Redis is a simple to integrate solution not only for specific implementation situations and use-cases but can be viewed as a way to speed up fetching and modifying small frequently used parts of data. Problems regarding the storage and usage of data can usually be separated into 2 categories. In one category there would be a large quantity of data that is rarely queried or data that needs to be processed and the other, small quantity of data that needs to be queried almost always before other queries are to be executed like for instance authentication information.

### **2.4.3 Security issues in NoSQL**

NoSQL databases solve many functional problems that conventional RDBMSs have providing greater scalability, flexibility allowing applications to better scale-out and/or utilize different kinds of data. The distribution of data is commonly done by utilizing sharding mechanisms and the usual RDBMS ACID properties are exchanged for BASE properties. While this provides good and usable features and characteristics, it brings to the surface a lot of security issues in need of dealing with. The authors of [30] [31] [32] analyzed a few of the most popular NoSQL choices and made security analysis' of them. They evaluated security features that they use and mention features that are still lacking. Areas that were looked at can be divided into: authentication, access control, data encryption, secure configurations and auditing. The database systems that were analyzed include: MongoDB, CouchDB, Redis, Hadoop, HBase, Cassandra, Redis and Couchbase Server. Although every one implements some features, all of them lack a complete set of security features. MongoDB seems to have the strongest security feature list and Redis has almost no security features. It is easy to see that the main focus of these solutions is performance, leaving security issues to be solved on the application and maintained and configured by system administrators.

### **2.4.4 Performance studies and comparisons**

Comparisons and performance evaluations given by the authors of [16] [33] [12] [34] [35] give a good picture of the current state and capabilities of the most used NoSQL database systems. It has to be noted that the NoSQL databases taken into account were mainly either document or wide column based and these are the best suited for storing and handling large amounts of data. The performance test included measuring latency in scenarios with different ratios of read/write/update operations. Compared to RDBMSs the results indicate that NoSQL systems have greater capacity and can cope with more operations which is to be expected. After a certain number of operations latency in RDBMS increased significantly (exponentially) while in NoSQL systems suffered little. Results indicate similar performance inside the different NoSQL paradigm families. MongoDB (document-based) was shown to be somewhat slower than Cassandra or HBase (wide-column based) with bigger workloads which is to be expected because of the differences between their storage systems and functionalities they offer.

### **2.4.5 NoSQL Conclusion**

NoSQL systems provide solutions in the areas of IoT and Big Data. Same as in other areas, the interest of large Internet companies such as Google, Amazon and Facebook pushed developments, and

because of this there a lot of good solutions available. NoSQL solves problems that conventional RDBMS cannot cope with or are not flexible enough or adapted for. Evaluations regarding data storage, handling and additional features of these solutions suggest that security is still an issue and may be the next significant improvement. Along with this new combined solutions similar to the document-graph mixed solution OrientDB utilizes could be the next step in NoSQL. Performance evaluations ( regarding reads, writes, updates...) indicate that improvements are always needed and welcomed but are at a level that is satisfying for most needs. This means that the choice of the database system that is going to be used in an application shouldn't be done from a performance perspective, but from a more general and data type oriented perspective along with evaluating the ecosystem that the main solution is a part of. More significant improvements could come in the area of data processing and analyzing. The Apache Hadoop ecosystem of open source projects is currently the most fertile ground for improvements in this area.

## ***2.5 Access Control***

Access Control is a general term that can be described as a way of securely granting, limiting or denying access to resources therefore protecting the resources from potentially malicious parties. This section a few of the most significant methodologies of enforcing access control will be describe, a general overview, comparison and evaluation will be given and some available solutions will be explored and evaluated.

### **2.5.1 Types of access control**

In this section a few of the most important and significant types of access control will be described and evaluated.

#### ***DAC and MAC***

DAC (Discretionary Access Control) is a type of access control that enforces a evaluation criteria (or policy) that mainly restricts access to the resource owner or group and all control over that resource including passing access to other subjects, and is commonly done automatically or indirectly. The users/resource owners have control over the access to the resource. An example would be in operating systems. Users that create files have ownership over them and can limit access to themselves or allow access to other users. DAC is often compared to the MAC (Mandatory Access Control).

MAC is a type of access control where a external entity like an administrator decides on giving or denying access to resources. A example would be database management systems. They allow access to resources depending on the permissions users have. The system administrator has the ownership over all resources.

#### ***IBAC***

IBAC (Identity Based Access Control) is based on giving access to a resource depending on a users identity. Implementations often become RBAC because identities are grouped and the groups can also be grouped as roles. An example of true IBAC are simple username and password based authentication systems

such as signing up and logging into accounts on websites. IBAC is simple to implement but has limited applications because of limited flexibility regarding modifying access to resources for a number of users. Most adequate implementations would be as mentioned before, password (or authentication card, fingerprint...) based authentication systems that limit access to a resource which could be a account, secure room, etc.

## ***RBAC***

RBAC (Role Based Access Control) is an access control system where users are assigned to groups or roles which have access to resources. An administration entity has power over changing access policies for resources. This solution is one of the most popular because of the simplicity of implementation, intuitive way of working, good flexibility and intuitive administration. Different roles have different authority levels or different areas of authority. Depending on the implementation, the roles can be organized in a hierarchy tree or a simple independent list of roles or groups. Almost every website system has a RBAC system implemented limiting access to information, changing and adding data, etc. Users are typically divided into groups with different levels of access depending on the type of user (buyer, seller, visitor...) and administrators are divided into different levels and areas for efficient maintaining of the system. A problem that occurs with this kind of system is that a large number of roles create a hard system to administer and maintain. An administrator needs to have deep understanding of each role and dependencies on other roles to successfully maintain the system.

## ***ABAC***

ABAC (Attribute Based Access Control) is a type of access control that evaluates attributes to decide if a user has access to a resource. Attributes are typically divided into three categories:

- subject - user attributes (examples: age, postal code, IP address...)
- object - resource attributes (examples: type, value, age...)
- environment (examples: day of the week, hour of the day...)

Policies contain condition and after comparing these attributes against policies the access to a resource is granted or denied. RBAC can be viewed as a subset of ABAC simply by evaluating everything based on one attribute that contains a users role. ABAC is therefore more complex than RBAC and IBAC and is typically more difficult to implement. Another issue is that the policies that contain the conditions is not as intuitive and RBAC and therefore more complex for an administrator to understand. The maintaining of large RBAC systems usually require a large amount of "hands on" modifications and monitoring or a very complex and well built system to change roles automatically. If some attributes of a user change or environmental conditions change, a pure RBAC system is not flexible to change automatically. Although ABAC can be more complex to build the benefit is that is much more flexible than RBAC. ABAC policies can specifically define which actions are allowed on which resources depending on the users, resources and environmental attributes. For instance, if a user becomes older than the legal age limit they can automatically have additional access granted without changing the users role. Subjects can also have a role attribute and mainly

function as a RBAC system and have policies with attributes as additional access control options that can be put in place just by adding a policy. ABAC has its place in systems that require a long term solution, complex policies and low maintenance over the policies.

### **2.5.2 OASIS, XACML and JSON**

XML (Extensible Markup Language) is a markup language that is used to format data in such a way that it can be readable by both humans and machines. It is defined by a set of rules for encoding documents by the W3C's XML 1.0 Specification. Although the focus of XML's design is on documents, it can be and is used for storing various kinds of structured data [36].

JSON (JavaScript Object Notation) is a lightweight data-interchange format. Same as XML it is readable by both humans and machines. Although it was built around JavaScript it is completely language independent it is ideal for mapping information stored in objects (on object oriented languages) or structures and therefore is best suited for exchanging data between languages [37].

OASIS is a non-profit consortium produces open standards in the areas of security, IoT, cloud computing, energy, content technologies, emergency management, and other areas. It has more than 5,000 participants representing over 600 organizations and individual members in more than 65 countries. The goal of OASIS is to provide standards which offer efficient and open solutions for common problems the areas mentioned before, drive development and building of standardized open source solutions that can easily be used by many, therefore accelerating innovation and development in general [38].

XACML (eXtensible Access Control Markup Language) is a declarative access control policy language implemented in XML created by OASIS. It defines a way to evaluate requests for resources according to rules defined in policies. Put simply it is a taught out and standardized solution for implementing access control in software applications. It provides common ground regarding terminology and workflow between multiple vendors building implementations of access control using XACML and interoperability between the implementations. It is primarily based on ABAC but can also implement RBAC as RBAC can logically be viewed as a subset of ABAC as explained earlier in this document [39] [40]. Along with the XML version there is also a profile utilizing JSON. That profile utilizes the latest 3.0 version of XACML and as stated in the official document of the JSON profile [41] is equivalent to the standard XML. Requests and responses can be translated in either direction (and back) without loss of information and equivalent requests result in equivalent responses.

### **2.5.3 Current Technologies and open source solutions**

TODO: describe any solutions used in the development of the solution

## **2.6 SMARTIE**

SMARTIE (Smart City) is a European project with the goal of solving security, privacy and trust issues in IoT, in a Smart City implementation. Partners include companies, universities and cities from

Germany, Serbia, Spain, Portugal and UK. As stated on the official website [42] the project officially started on September 1st 2013. and is scheduled to end on August 31st 2016. and has a total budget of 4,862,363 € with the contribution from EU in the amount of 3,286,144 €

## **2.6.1 Brief Overview**

### ***Vision [42]***

The vision of SMARTIE is to create a distributed framework to share large volumes of heterogeneous information for the use in smart-city applications, enabling end-to-end security and trust in information delivery for decision-making purposes following data owner's privacy requirements. A secure, trusted, but easy to use IoT system for a Smart City will benefit the various stakeholders of a smart city: The City Administration will have it easier to get information from their citizens while protecting their privacy. Furthermore, the services offered will be more reliable if quality and trust of the underlying information is ensured. Privacy and Trust are a key prerequisite for citizens to participate in Smart City activities. A Smart City can improve the Smart and Comfort Live of their citizens enormously. Enterprises benefit from the securely provided information. They can optimize their business processes and deal with peak demands introduced by the dynamics of the Smart City. Furthermore, they can offer more tailored solutions for their customers based on the status of the Smart City.

### ***Main goals [42]***

- Understanding requirements for data and application security and creating a policy-enabled framework supporting data sharing across applications.
- Developing new technologies that establish trust and security in the perception layer and network layer.
- Develop new technologies for trusted information creation and secure storage for the information service layer.
- Develop new technologies for information retrieval and processing guided by access control policies in the application layer.
- Demonstrate the project results in real use cases

### ***Key Challenges [42]***

The idea of the IoT brings new challenges regarding security and in consequence also for privacy, trust and reliability. The major issues are:

- Many devices are no longer protected by well-known mechanisms such as firewalls and can be attacked via the wireless channel directly. In addition devices can be stolen and analysed by attackers to reveal their key material.
- Combining data from different sources is the other major issue since there is no trust relationship between data providers and data consumers at least not from the very beginning.

- Secure exchange of data is required between IoT devices and consumers of their information

### ***Expected Impact [42]***

- The SMARTIE IoT platform will allow the virtualization of the functionalities of discovery, secure information access, processing and privacy-aware distribution between the consumer and producer of the data generated by the smart objects.
- It will demonstrate the applicability in scenarios linked to the green behaviour and sustainability of smart cities like efficient transport/mobility and energy management.
- SMARTIE will facilitate new companies for developing and providing services over its IoT infrastructure/platform.
- SMARTIE allows heterogeneous and multiple source of data to interact in reliable and secure manner providing third parties developers in Europe to enter the Smart Cities area and in that way increase the share of the IoT market.

### ***Use Cases [42]***

#### **1. Frankfurt/Oder (GER)**

- Traffic management with the possibility to influence real traffic.
- Focus on authentication, trust, data security, interoperability

#### **2. Murcia (ES)**

- Monitoring energy efficiency in the campus
- Users can interact with the system to improve energy efficiency.

#### **3. Belgrade (RS)**

- Provide smart transportation using location of busses and travellers
- Focus on data security and privacy using developed access rights and policies

## **2.6.2 Conclusion**

There are a number of smart city project currently being developed but unlike SMARTIE most are focused on specific aspects and specific test scenarios and rarely grow bigger after they are complete. SMARTIE is a solution for smart cities that unlike other efforts is being developed with security concerns and standardisation in mind. The fact that it is founded by the European Union gives even more significance to the project because if successful, it will surely be utilized in many European cities. Along with having the focus on solving this problem for cities today, the project also focuses on utilizing standards, stimulating innovation and producing improvements and developments in the area of IoT.

## 3 Project description

This section will describe in detail the work that has been done along with describing the development process.

TODO: restructure this to make more sense

### 3.1 *Introduction*

TODO: rewrite to be easier to understand

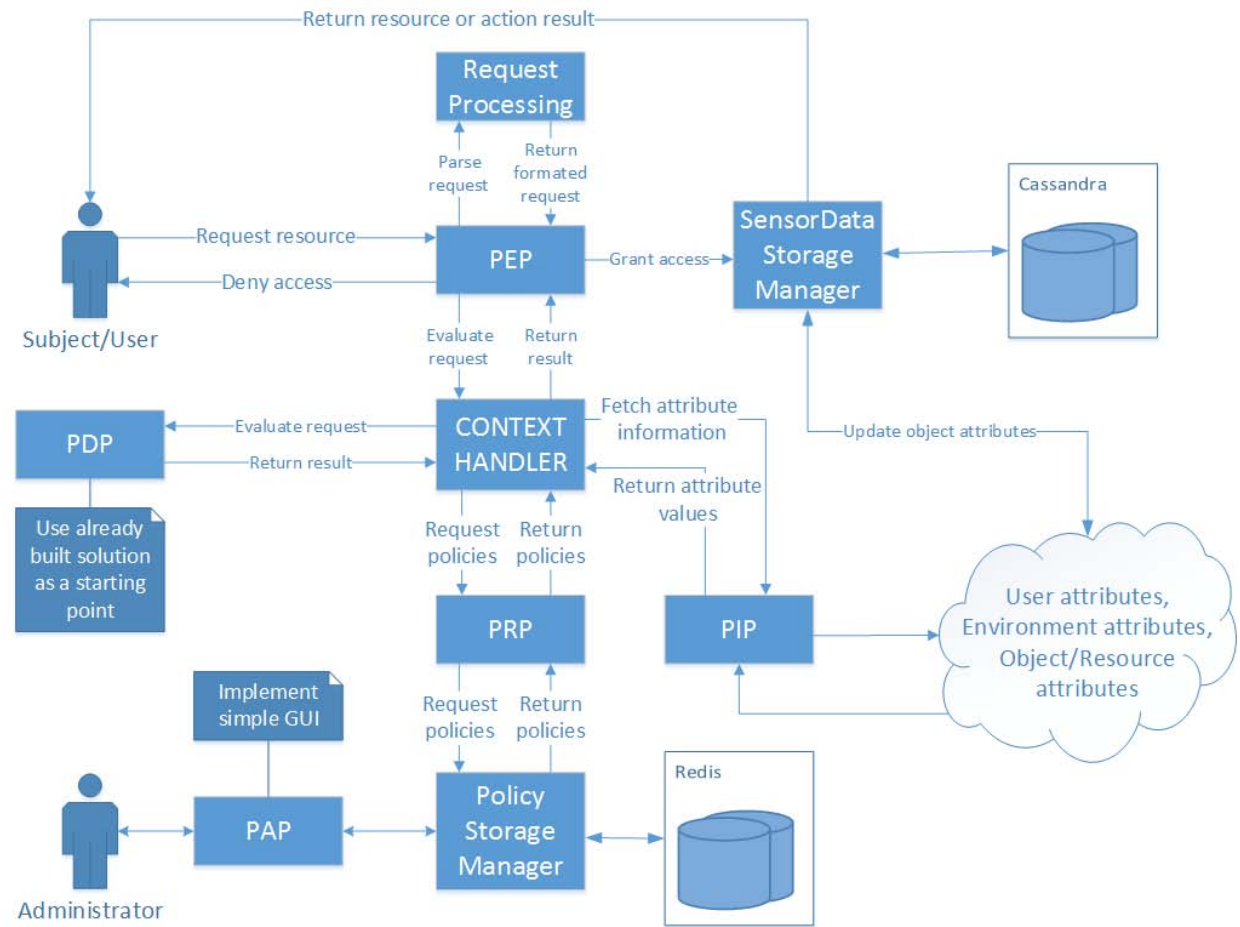
### 3.1 *General architecture and design choices*

After extensive research and regarding the technologies currently available, a number of design choices were made and a solution architecture is proposed.

The NoSQL databases that will be used are Cassandra for storing sensor data and Redis for storing policies. Cassandra was chosen because of its scaling characteristics and Redis for its speed. A DataManager module will be built for each of these allowing easier use of these for other parts of the solution. The access control framework will have the OASIS architecture containing the: PEP, PDP, PRP, PIP and PAP but will be implemented with a subset of simpler functionalities. The limitations that will be taken into account come directly from the current state of the SMARTIE project which means that all or most use-case scenarios will be fulfilled but the overall solution can be simplified making the delivery more realistic (taking in account the time constraints).

In its current stage, SmartData doesn't need support for complicated policies. Policies should mainly give access according to the users id, group and role and the resources id and type. There are no conditions regarding environment conditions, range of sensor values or other more complex conditions. This greatly simplifies the PAP and PDP needed in the solution.





**Figure 3. Architecture of the proposed solution**

In Figure 3. a more detailed architecture of the proposed solution can be seen. This solution is a slight modification to the one proposed in OASIS XACML 3.0.

TODO: Describe the diagram and implementation, benefits...

### 3.1 Development plans and overview

TODO: rewrite to be easier to understand

The NoSQL database used for storing sensor data in the first stages will be setup for simple and generic sensor readings. The database will contain the following data: sensorID, sensorType, timestamp, valueOfReading. The structure will most likely mirror that list of values. In later stages the database will either be expanded to support the real sensor data currently used in the SMARTIE project or the policy list will be expanded and a connection to the current database used in SMARTIE (which is also a Cassandra DB).

The policies that are going to be generated and used are going to be simple and access to data will be denied by default. Options for policies in the first stages will include allowing access according to sensorID, sensorType, userGroup, userType (role). In later stages they will support basic environment

conditions like day of the week and time of the day and other resource and user attributes (not yet decided which ones they will be).

The main focus most of the time spent in the development will be for the access control part. For the PDP which is one of the more complicated parts, an open source solution (Balana or AT&T XACML open source GitHub project) will be used as a starting point. It supports XACML and should simplify the implementation of a PDP significantly. For the XACML/JSON policy building, a framework (<https://github.com/att/XACML>) with a REST service will be utilized to help with the integration. The PRP and PIP will connect to the policy database and will retrieve the attributes of the users, resource and environment. Implementation of these should be simple and straightforward. In the first stages User and Resource attributes will be stored in a separate database (not decided which one yet) and a number of test users and test resources and data will be generated for testing and development purposes. Later stages will include integration with SmartData regarding fetching user attributes and generating and storing resource attributes. The PAP will be the last component made in the Access Control framework and will support the creation of simple policies described near the beginning of this section.

TODO: Database schema diagrams

This solution doesn't rely on other parts of the SMARTIE project and SmartData platform. This is good because independence will allow for good testing and easier development. However it is going to be built knowing that it has to be integrated or that it should be able to be integrated without much effort and/or modification.

### ***3.2 Potential improvements/achievements***

In the current state, SmartData has access control implemented by encrypting data and allowing decryption to users that have access to it. The policies that contain the information on which users can access the data are stored together with the data itself ("sticky policies"). This means that the all data needs to be fetched in order to evaluate the policy on it. This has a drawback in case of requests that result in denying access. If a user or set of users send a large number of requests for a lot of data, for which they don't have access to, an opening for a DoS type of attack could happen. This also means that in case of changing access policies for past data means fetching all the data that we want the change to affect and updating of the policy attached to it. This solution therefore isn't flexible if policies that are meant to be updated are to be integrated.

This Access Control framework should solve this issues offering a safer and more flexible solution and could potentially be integrated as access control "before" (regarding the point in communication) the current access control point.

### ***3.3 Conclusion***

Using the OASIS ABAC implementation proposal architecture will allow the solution to be easily be built upon and allow the adding of other essential functionalities. This solution will therefore retain the structure of the OASIS ABAC implementation but will not have the full set of features/functionalities implemented. Just a subset determined to be enough taking in account the limitations set. The integration in SMARTIE or SmartData is not yet certain but will helpfully be done and prove to have significant benefits.

..... TODO: elaborate a bit more

## **4 (Main) Development implementation**

### ***4.1 Final structure description in detail***

#### **4.1.1 Global architecture**

#### **4.1.2 Policy database**

#### **4.1.3 Sensor data database**

### ***4.2 Data storage modules***

### ***4.3 Access control module***

#### **4.3.1 PAP**

#### **4.3.2 PIP & PRP**

#### **4.3.3 PEP**

#### **4.3.4 PDP**

## **5 Testing**

### ***5.1 Test scenatio1***

### ***5.2 Test scenario2***

## **6 Integration in SMARTIE**

### ***6.1 Proposed implementation***

### ***6.2 Challenges***

### ***6.3 Solution***

### ***6.4 Testing results***

## 7 Conclusion

Companies such as Google, Facebook, Amazon... etc. are investing a lot of time, effort and money in solving the challenges regarding IoT, Big Data, Access Control, Security in all of these areas and more because future steps in the commercial technological aspect in general point to the integration of microcontrollers or computers in general into almost every device and everyday item or can be integrated into infrastructure, factories, farms, production facilities ect. This leads to an IoT which leads to Big Data problems which are being addressed by NoSQL data management systems. Along with the fundamental problems like infrastructure and having solutions that can cope and handle everything, security and privacy of all of that data is an issue that needs to be addressed for "real world" applications to be acceptable, safe and successful. ABAC using the XACML/JSON standard from OASIS is one of many possible solution/routes that can be taken regarding the solving of this problem. Compared to other solutions it offers exceptional security if taken into account the minimum amount of human effort needed for managing and updating policies. The implementation of a standardized solution means that modules can be taken and ported to other applications easily and used in other areas such as web and mobile applications and effectively addressing half of the issues from the OWASP current top ten list without relying on individuals to take in consideration security issues one by one.

TODO: finish work and elaborate

## 8 References

- [1] A. I. G. M. Luigi Atzori, "The Internet of Things: A survey," *Computer Networks* 54, p. 2787–2805, 14 June 2010.
- [2] N. B. A. C. V. M. Z. Andrea Zanella, "Internet of Things for Smart Cities," 2014 April 2014. [Online]. Available: <http://eprints.networks.imdea.org/id/eprint/740>. [Accessed 12 February 2015].
- [3] A. V. V. J. W. J. L. D. Q. Qi Jing, "Security of the Internet of Things: perspectives and challenges," *Wireless Netw*, p. 20:2481–2501, 17 June 2014).
- [4] J. W. C. Z. J. L. Hui Suo, "Security in the Internet of Things: A Review," in *International Conference on Computer Science and Electronics Engineering*, 2012.
- [5] A. R. L. G. A. C.-P. S. Sicari, "Security, privacy and trust in Internet of Things: The road ahead," *Computer Networks* 76, p. 146–164, 15 July 2014.
- [6] Y. S. J. W. Prasant Misra, "Towards a Practical Architecture for the Next Generation Internet of Things," 3 February 2015. [Online]. Available: <http://arxiv.org/pdf/1502.00797v1.pdf>. [Accessed 10 February 2015].
- [7] "Wikipedia M2M," Wikimedia Foundation, Inc., 13 February 2015. [Online]. Available: [http://en.wikipedia.org/wiki/Machine\\_to\\_machine](http://en.wikipedia.org/wiki/Machine_to_machine). [Accessed 16 February 2015].
- [8] S. M. I. J. W. M. I. S. G. M. I. X. L. M. I. a. V. C. L. F. I. Min Chen, "A Survey of Recent Developments in Home M2M Networks," *IEEE COMMUNICATIONS SURVEYS & TUTORIALS*, VOL. 16, NO. 1, pp. 98-114, 2014.
- [9] X. Z. Xiao Nie, "M2M Security Threat and Security Mechanism Research," in *3rd International Conference on Computer Science and Network Technology*, 2013.
- [10] M. A. P. O. S. N. M. L. t. H. David S. Watson, "Machine to Machine (M2M) Technology in Demand Responsive Commercial Buildings," in *2004 ACEEE Summer Study on Energy Efficiency in Buildings*, Pacific Grove, CA, 2004.
- [11] S. P. S. V. K. Roshni Bajpayee, "Big Data: A Brief investigation on NoSQL Databases," *International Journal of Innovations & Advancement in Computer Science, IJIACS*, Volume 4, Issue 1, p. 2347 – 8616, January 2015.
- [12] P. P. V. K. Rajendra Kumar Shukla, "Big Data Frameworks: At a Glance," *International Journal of Innovations & Advancement in Computer Science, IJIACS*, Volume 4, Issue 1, p. 2347 – 8616, January 2015.



- [13] "Google Images," [Online]. Available: <http://1.bp.blogspot.com/-4sM6wRgPBUA/T9FiKnMhNFI/AAAAAAAAACc/bOO8PGQ0rDs/s1600/BigData+-+V5+Lens.JPG>. [Accessed 17 February 2015].
- [14] G. Lafuente, "The big data security challenge," *Network Security*, January 2015.
- [15] "Alternatives to relational database: Comparison of NoSQL and XML approaches for clinical data storage," *computer methods and programs in biomedicine 110*, p. 99–109, 2013.
- [16] S. A. H. A B M Moniruzzaman, "NoSQL Database: New Era of Databases for Big data Analytics - Classification, Characteristics and Comparison," *International Journal of Database Theory and Application*, 2013.
- [17] "MongoDB HomePage," MongoDB, Inc., [Online]. Available: <http://www.mongodb.org/>. [Accessed 7 February 2015].
- [18] "Apache Hadoop," The Apache Software Foundation, 12 12 2014. [Online]. Available: <http://hadoop.apache.org/>. [Accessed 7 February 2015].
- [19] "Apache Hbase," The Apache Software Foundation, 11 Februaray 2015. [Online]. Available: <http://hbase.apache.org/book.html>. [Accessed 15 February 2015].
- [20] "Apache Hive," The Apache Software Foundation, 12 January 2015. [Online]. Available: <https://cwiki.apache.org/confluence/display/Hive/Home>. [Accessed 7 February 2015].
- [21] "Apache Spark," The Apache Software Foundation, [Online]. Available: <http://spark.apache.org/>. [Accessed 7 February 2015].
- [22] "Apache Spark summary," 2015 Black Duck Software, Inc., January 2015. [Online]. Available: <https://www.openhub.net/p/apache-spark>. [Accessed 16 February 2015].
- [23] "Apache Cassandra," The Apache Software Foundation, [Online]. Available: <http://cassandra.apache.org/>. [Accessed 16 February 2015].
- [24] "Wikipedia MongoDB," Wikimedia Foundation, Inc, 9 March 2015. [Online]. Available: <http://en.wikipedia.org/wiki/MongoDB>. [Accessed 12 March 2015].
- [25] "Tutorialspoint Neo4j," [Online]. Available: [http://www.tutorialspoint.com/neo4j/neo4j\\_features\\_advantages.htm](http://www.tutorialspoint.com/neo4j/neo4j_features_advantages.htm). [Accessed 18 February 2015].
- [26] "Neo4j official," Neo Technology Inc., 2015. [Online]. Available: <http://neo4j.com/>. [Accessed 18 February 2015].
- [27] "OrientDB official site," Orient Technologies, 2015. [Online]. Available: <http://www.orienttechnologies.com/orientdb/>. [Accessed 18 February 2015].

- [28] “Redis website FAQ,” Salvatore Sanfilippo and Pieter Noordhuis, [Online]. Available: <http://redis.io/topics/faq>. [Accessed 20 February 2015].
- [29] “Redis official website,” Redis, [Online]. Available: <http://redis.io/>. [Accessed 20 February 2015].
- [30] R. M. M. A. S. Anam Zahid, “Security of Sharded NoSQL Databases;,” in *Conference on Information Assurance and Cyber Security (CIACS)*, 2014.
- [31] C. P. N. Priya P. Sharma, “Securing Big Data Hadoop: A Review of Security, Issues, Threats and Solution,” *International Journal of Computer Science and Information Technologies*, Vol. 5 (2), pp. 2126-2131, 2014.
- [32] N. G.-O. Y. G. E. G. J. A. Lior Okman, “Security Issues in NoSQL Databases,” *International Joint Conference of IEEE TrustCom-11/IEEE ICSS-11/FCST-11*, pp. 541-547, 2011.
- [33] “A comparison between several NoSQL databases with comments and notes”.
- [34] M. G. M. I. Enrico Barbierato, “Performance evaluation of NoSQL big-data applications using multi-formalism models,” *Future Generation Computer Systems* 37, p. 345–353, 2014.
- [35] D. CORPORATION, “Benchmarking Top NoSQL Databases, A Performance Comparison for Architects and IT Managers,” 2013.
- [36] “XML wikipedia,” Wikimedia Foundation, Inc, March March 2015. [Online]. Available: <http://en.wikipedia.org/wiki/XML>. [Accessed 20 March 2015].
- [37] “JSON org,” [Online]. Available: <http://json.org/>. [Accessed 20 March 2015].
- [38] “OASIS official wabsite,” OASIS, 2015. [Online]. Available: <https://www.oasis-open.org/org>. [Accessed 20 March 2015].
- [39] “eXtensible Access Control Markup Language (XACML) Version 3.0,” 22 January 2013. [Online]. Available: <http://docs.oasis-open.org/xacml/3.0/xacml-3.0-core-spec-os-en.pdf>. [Accessed 4 November 2014].
- [40] “XACML wikipedia,” Wikimedia Foundation, Inc., 16 January 2015. [Online]. Available: <http://en.wikipedia.org/wiki/XACML>. [Accessed 20 March 2015].
- [41] “JSON Profile of XACML 3.0 Version 1.0,” 15 May 2014. [Online]. Available: <http://docs.oasis-open.org/xacml/xacml-json-http/v1.0/csprd03/xacml-json-http-v1.0-csprd03.pdf>. [Accessed 11 February 2015].
- [42] “SMARTIE Homepage,” IHP GmbH, 22 July 2014. [Online]. Available: <http://www.smartie-project.eu/project.html>. [Accessed 20 February 2015].

- [43] R. B. S. M. M. P. Jayavardhana Gubbi, "Internet of Things (IoT): A vision, architectural elements, and future directions," *Future Generation Computer Systems* 29, p. 1645–1660, 24 February 2013..
- [44] "Apache Hadoop YARN," Apache Software Foundation, 11 13 2014. [Online]. Available: <http://hadoop.apache.org/docs/current/hadoop-yarn/hadoop-yarn-site/YARN.html>. [Accessed 7 February 2015].
- [45] "Apache Hadoop HDFS," Apache Software Foundation, 13 11 2014. [Online]. Available: <http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsDesign.html>. [Accessed 6 February 2015].
- [46] "Apache Hadoop MapReduce," Apache Software Foundation, 13 11 2014. [Online]. Available: <http://hadoop.apache.org/docs/current/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html>. [Accessed 7 February 2015].
- [47] Y. H. A. C. A. G. G. H. E. N. Hanson, "Major Technical Advancements in Apache Hive," in *SIGMOD '14 Proceedings of the 2014 ACM SIGMOD international conference on Management of data*, New York, NY, USA, 2014 .
- [48] M. C. M. J. F. S. S. I. S. Matei Zaharia, "Spark: Cluster Computing with Working Sets," 2010.
- [49] M. C. F. X. D. L. a. K. Z. Jiafu Wan, "From Machine-to-Machine Communications towards Cyber-Physical Systems," *Computer Science and Information Systems Vol. 10, No. 3*, p. 1105–1128, June 2013.
- [50] C. S. DU Jiang, "A Study of Information Security for M2M of IOT," in *3rd International Conference on Advanced Computer Theory and Engineering (ICACTE)*, 2010.
- [51] "Wikipedia," Wikimedia Foundation, Inc., 13 February 2015. [Online]. Available: [http://en.wikipedia.org/wiki/Machine\\_to\\_machine](http://en.wikipedia.org/wiki/Machine_to_machine). [Accessed 16 February 2015].
- [52] "XML website," O'Reilly Media, Inc., [Online]. Available: <http://www.xml.com/>. [Accessed 20 March 2015].