



COSMOS

Cultivate resilient smart Objects for Sustainable city applicatiOns

Grant Agreement Nº 609043

D6.1.3 Reliable and Smart Network of Things: Design and Open Specification (Final)

WP6: Reliable and Smart Network of Things

Version: 1.0

Due Date: 30th April 2016

Delivery Date: 30th April 2016

Nature: Report

Dissemination Level: PUBLIC

Lead partner: 5 (UNIS)

Authors: F. Carrez (editor-UNIS) & A. Akbar (editor-UNIS), P. Bourelos (ICCS/NTUA), J. Sancho (ATOS), J. Rico (ATOS)

Internal reviewers: Achilleas Marinakis (ICCS/NTUA), Paula Ta-Shma (IBM)

www.iot-cosmos.eu



The research leading to these results has received funding from the European Community's Seventh Framework Programme under grant agreement n° 609043

Version Control:

Version	Date	Author	Author's Organization	Changes
V0.1	08/02/2016	François Carrez, Adnan Akbar	UNIS	Initial version
V0.2	26/02/2016	Panagiotis Bourellos	NTUA	Structural changes in the Experience Sharing chapter
V0.3	29/02/2016	Adnan Akbar	UNIS	Structural changes in section 2
V0.4	04/04/2016	Adnan Akbar	UNIS	New contributions in section 2,4 and 5
V0.5	22/04/2016	Achilleas Marinakis	NTUA	Final adjustments of NTUA contribution
V0.6	25/04/2016	François Carrez	UNIS	proof reading and fine tuning of Exsum / Intro/Concl
V0.7	26/04/2016	Achilleas Marinakis & Paula Ta-Shma	ICCS/NTUA & IBM	Review
V0.8	27/04/2016	Adnan Akbar	UNIS	Revised version w.r.t. review comments & Revised conclusion
V0.9	28/04/2016	François Carrez	UNIS	Final checks

Table of Contents

1	Introduction	10
2	Inference/Prediction Functional Component.....	12
2.1	Introduction	12
2.2	Background.....	12
2.2.1.	Machine Learning	12
2.2.2.	Classification Analysis	13
2.2.3.	Clustering Analysis	14
2.2.4.	Regression Analysis.....	14
2.2.5.	Statistical Inference:	15
2.2.6.	Large-Scale IoT Data	17
2.3	Functional Overview.....	19
2.4	Connection with other Components.....	19
2.5	Interfaces.....	21
2.6	Use-Case Scenario Extension from Year 1.....	22
2.6.1.	Background	22
2.6.2.	Results and explanations	22
2.7	Use-Case Scenario 2 – Pro-Active Traffic Management.....	24
2.7.1.	Introduction	24
2.7.2.	Proposed solution.....	24
2.7.3.	Adaptive Moving Window Regression.....	25
2.7.4.	Other components.....	27
2.8	Conclusion	27
3	Pre-Processing Functional Component.....	29
3.1	Introduction	29
3.1.1.	Data Cleaning.....	29
3.1.2.	Data Transformation.....	29
3.1.3.	Data Reduction	29
3.2	Functional Overview.....	30
3.3	Connection with other Components.....	30
3.4	Interfaces.....	31
3.5	Use-case scenario.....	32
3.6	Conclusion	32

4	Event (Pattern) Detection Functional Component	33
4.1	Introduction	33
4.1.1.	Background	33
4.1.2.	Complex Event Detection	34
4.1.3.	Anomaly Detection	35
4.1.4.	Complex Event Processing.....	36
4.1.5.	Probabilistic Rules for Event Processing.....	37
4.2	Functional Overview.....	37
4.3	Connection with other components	37
4.4	Interfaces.....	38
4.5	Use-Case Scenarios	39
4.5.1.	Scenario 1: Detecting Traffic State (Madrid)	39
4.5.2.	Scenario 2: Anomaly Detection (Taipei)	41
4.5.3.	Scenario 3: Detecting an Event from Twitter Data (Madrid).....	42
4.6	Conclusion	42
5	Situational Awareness Functional Component.....	43
5.1	Introduction	43
5.2	Functional Overview.....	43
5.2.1.	Data Sources	44
5.2.2.	Data Flow and Processing Model.....	45
5.2.3.	Information sharing and retrieval.....	47
5.3	Connection with other components	48
5.3.1.	Complex Event Processing Engine	48
5.3.2.	Machine Learning	48
5.3.3.	Message Bus	51
5.3.4.	Forging CEP and ML	51
5.4	Interfaces.....	51
5.5	Use-case Scenarios.....	52
5.5.1.	Scenario 1	52
5.5.2.	Scenario 2	53
5.6	Conclusion	55
6	Experience Sharing Functional Component.....	56
6.1	Introduction	56
6.2	Functional Overview.....	57



6.2.1.	Mechanism Improvements and Restating Achievements	58
6.2.2.	Architectural/Deployment Changes	58
6.2.3.	Usage of ML based VE Similarity Calculations	59
6.3	Communication with other Components.....	59
6.4	API / Interfaces for Experience Sharing.....	60
6.5	Use Cases.....	61
6.5.1.	Use Case for Year 3 Implementation	61
6.6	Conclusion	62
7	Conclusions	63
8	References	64

Table of Figures

Figure 1: COSMOS Function View (w/ WP6 FC emphasis)	11
Figure 2: Efficient Implementation of Machine Learning	13
Figure 3: Hidden Markov Model	16
Figure 4: Inference/Prediction on raw data	20
Figure 5: Inference/Prediction on VE Data	21
Figure 6: Aggregated data comparison with non-aggregated data for Machine Learning algorithms	23
Figure 7: Aggregated data comparison with non-aggregated data for HMM	23
Figure 8: Proposed solution for pro-active traffic management	25
Figure 9: Connection of pre-processing FC with other components	31
Figure 10: Pre-Processing component	32
Figure 11: Proposed solution for Event Detection	34
Figure 12: μ CEP interfaces	36
Figure 13: Connection of Event Detection FC	38
Figure 14: Flow chart for adaptive clustering	41
Figure 15: Situational Awareness data exchanges	44
Figure 16: An Example of Platform Centric Situation Awareness	49
Figure 17: Prediction of Events for SA Projection	50
Figure 18: SA in Decision-making process	52
Figure 19: Situation Awareness process	53
Figure 20: Bayesian Network Representation for Madrid Scenario	55
Figure 21: Problem Solution structure of Knowledge	56
Figure 22: New Followee-Follower Structure	57
Figure 23: Sequence of actions for Experience Sharing between two VEs	60

Table of Acronyms

Acronym	Meaning
AMWR	Adaptive Moving Window Regression
ANN	Artificial Neural Network
API	Application Programming Interface
ARIMA	AutoRegressive Integrated Moving Average
ARMA	AutoRegressive Moving Average
CB	Case Base
CBR	Case-based Reasoning
CEP	Complex Event Processing
CRUD	Create/Read/Update/Delete
DAG	Directive Acyclic Graph
EDA	Event Driven Architecture
EM	Expectation Maximization
FC	Functional Component
FFT	Fast Fourier Transform
GMM	Gaussian Mixture Model
GPS	Global Positioning System
GUI	Graphical User Interface
HMM	Hidden Markov Model
HTTP	Hyper-Text Transfer Protocol
IGR	Information Gain Ratio
IoE	Internet of Everything
IoT	Internet of Things

IP	Internet Protocol
JDK	Java Development Kit
JSON	Java Script Object Notation
KNN	K-Nearest Neighbour
LSSA	Least-Squares Spectral Analysis
MaL	Maximum Likelihood
MAP	Maximum A Posteriori
MAPE-K	Monitor/Analyse/Plan/Execute - Knowledge
MD	Model Developer
ML	Machine Learning
MLP	Multi-Layer Perceptron
MQTT	Message Queue Telemetry Transport
MTBF	Mean-Time Between Failure
NILM	Non-Intrusive Load Monitoring
OWL	Ontology Web Language
PA	Predictive Analysis
PAA	Piece-wise Aggregation Approximation
PCA	Principle Component Analysis
PMML	Predictive Model Markup Language
QoS	Quality of Service
RBF	Radial Basis Function
RDD	Resilient Distributed Datasets
RDF	Resource Description Framework
SA	Situation Awareness
SAX	Symbolic Aggregation approxImation



SPARQL	SPARQL Protocol and RDF Query Language (Recursive acronym)
SQL	Simple Query Language
SSID	Service Set IDentifier
SVM	Support Vector Machine
SVR	Support Vector Regression
UC	Use-case
URL	Unified Resource Locator
VE	Virtual Entity
XML	Extensible Markup Language
μCEP	Micro Complex Event Processing

1 Introduction

The main objective of this work package is to provide methods for inferring high-level knowledge from raw IoT data, to provide the means for situational awareness and understanding how things have behaved in comparable situations previously. Different data mining methods based on machine learning and statistical analysis techniques were explored and developed for extracting events from raw data. The focus of Year 1 was to explore state of the art methods and highlight the limitations of these methods whereas in Year 2 we addressed the limitations and extended the methods beyond the state of the art. Finally, in Year 3 we have improved and combined all the components in order to provide a global view for situation awareness. Complex Event Processing engines were deployed for knowledge inference in complex deployment situations based on near real-time analysis of complex events. We explored several Machine Learning techniques in conjunction with *Complex Event Processing* (CEP) in order to provide adaptive solutions to dynamic scenarios in order to optimize the performance of CEP for situation awareness. Finally different methods for experience sharing between virtual entities were investigated so that things can be made smarter and able to make decisions using the experience of entities, which have already faced identical situations. All of the above mentioned objectives are elaborated in this document with relevant examples.

This document is intended to provide a generic approach to meet the high-level objectives which are summarized below:

- 1) To address the limitations of existing knowledge extraction techniques and extend it to provide solutions for large-scale IoT applications;
- 2) To provide adaptive methods for inferring complex events from raw data streams;
- 3) To provide means for Situation Awareness for IoT networks;
- 4) To develop Experience sharing mechanisms between virtual entities and thus making them autonomous.

What is new in this deliverable compared to previous D6.1.2 version.

- Added a new use-case scenario (Section 2.6) which describes our prediction algorithm and its application on Madrid traffic data for proactive traffic management;
- Added new functionalities for detecting anomalies and probabilistic rules in Event Detection component (Section 4). A new use-case scenario for detecting anomalies has also been added;
- Added a new use-case scenario in Situational awareness component (Section 5) which describes the achievement of three levels of situational awareness using other components of the work package;
- All other Functional Components come in their third iteration aligned with Year 3 objectives;
- Added new functionalities in Experience Sharing FC (Section 6.2.2)
- Integrated Experience Sharing FC with Privelets FC (Sections 6.3 and 6.4)

The document is organized on the basis of description of WP6 components which are shown with thick edges within the COSMOS Functional View shown in Figure 1 below:

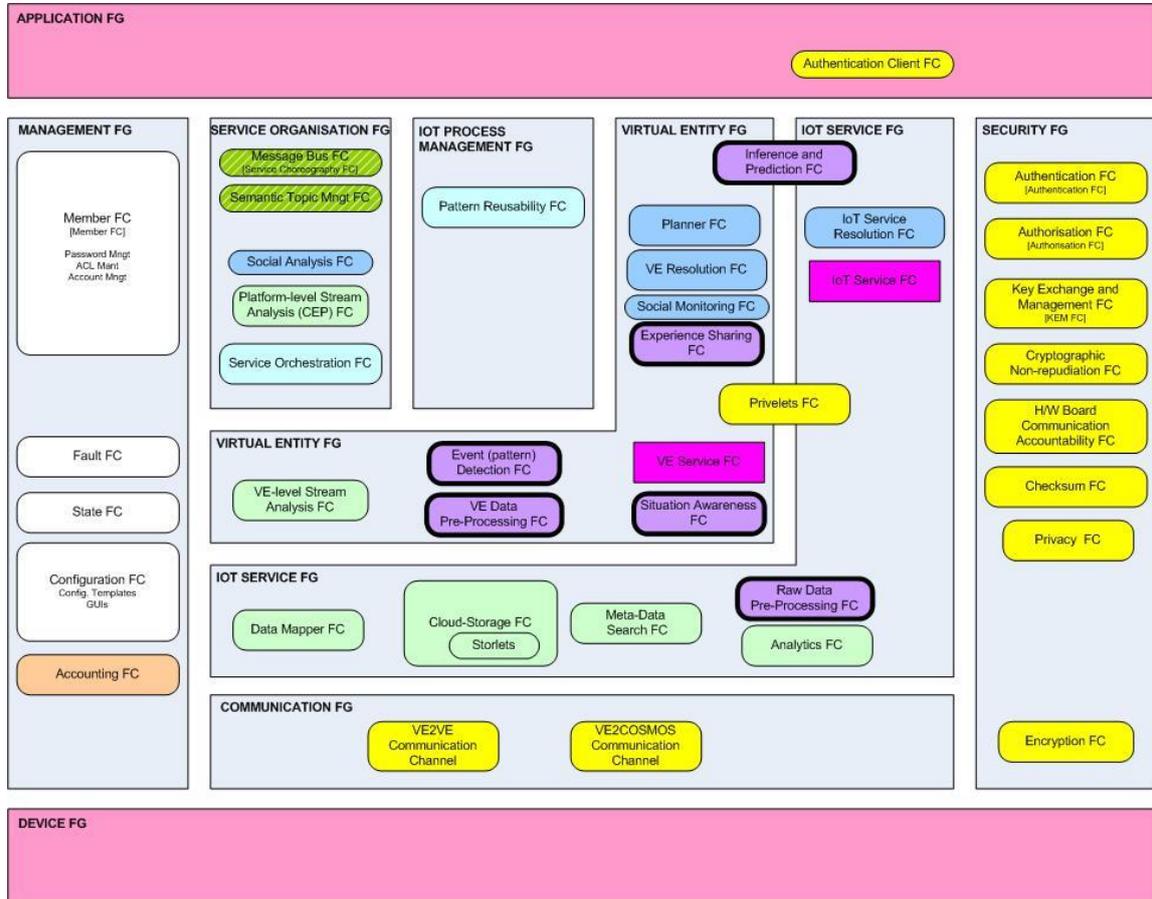


Figure 1: COSMOS Function View (w/ WP6 FC emphasis)

2 Inference/Prediction Functional Component

2.1 Introduction

In the world of *Internet of Things* (IoT), devices and sensors are deployed or used in varying conditions and different situations. Mostly they are deployed in remote places and are connected using less reliable wireless links. In order to prolong their battery life, data provided by these devices may be sporadic, less reliable and incomplete. Data itself is of no value until it is processed intelligently to extract high-level knowledge which can be used to make decisions. Data mining methods based on machine learning and statistical analysis techniques have the potential to extract knowledge from unreliable and incomplete data. Pattern recognition techniques based on data mining methods have long been used in speech and image processing applications but their use in IoT world is still in its early years.

This section explains different state of the art data mining methods which were explored for extracting high level knowledge from raw IoT data and we have demonstrated how these methods have the potential to contribute for novel applications. A high-level knowledge can be any event, some actionable knowledge or some statistical property of the data depending on the application. Data from different sensors in IoT form patterns, and different patterns represent different events. However, only certain events are interesting depending on the context of application and require further action. Pattern recognition methods based on data mining algorithms enable to extract these interesting events which have the potential to form the basis for many interesting applications. We have demonstrated in our work, how these algorithms can be applied in IoT domain for novel and innovative applications.

The performance of most of the Machine Learning models deteriorates when applied on real-time dynamic environments. One of the main factors contributing in deterioration of *Machine Learning* (ML) models is concept drift which occurs when the statistical properties of the target variable or the statistical distribution of the observation data changes over time. Secondly, the complexity of machine learning algorithms increases exponentially with the amount of data. Most of the innovative solutions found in literature based on ML and statistical analysis techniques are implemented on small data sets and are not usable for large scale IoT applications. In this regard, the aim of inference/prediction component was to focus on ML solutions which are able to cope with the dynamic nature of underlying IoT data in near real-time and also able to deal with the largeness of IoT data for smart city applications.

2.2 Background

2.2.1. Machine Learning

ML represents any algorithm or process which learns from the data. The availability of data from large number of diverse devices which represent real-time events has motivated the researchers to explore more intelligent solutions using ML algorithms. Whether it is in the context of a health department, supply chain management system, transportation or the environment, methods based on ML enables providing more intelligent and autonomous solutions by analysing and providing further insight. Different data mining algorithms have been used in very diverse contexts, detecting traffic congestion [2], predicting the energy demand of the buildings [3] and predicting the behaviour of customers in an online shop [4] are few examples of using data mining algorithms in context of IoT. There are many ML algorithms found in the literature, but in a broader context they fall into following three main categories.

- Classification Analysis;
- Clustering Analysis;
- Regression Analysis.

An efficient implementation of any ML algorithm involves different steps ranging from filtering, interpolation, and aggregation on one end to optimization methods on the other. A proper understanding of a problem is mandatory to apply the right choice of steps. The different steps involved in machine learning algorithms with possible options are shown in Figure 2. A brief introduction to few of these techniques is given in this section.

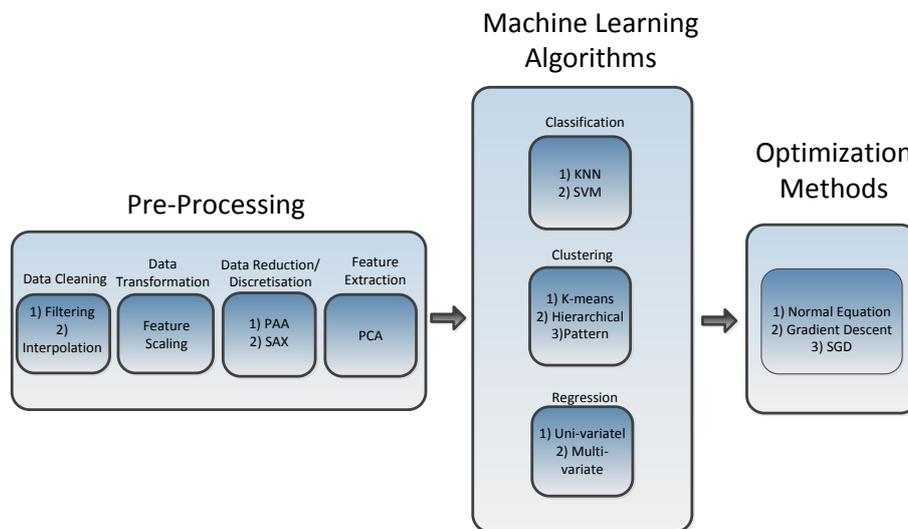


Figure 2: Efficient Implementation of Machine Learning

2.2.2. Classification Analysis

Classification is a supervised machine learning technique which requires labeled data in learning phase used widely for pattern recognition and categorization of new observations. The classification of emails as SPAM or NOT-SPAM (a.k.a. spam filters) represents a well-known example of classification analysis [5]. The server learns from the user's behavior, whenever the user marks (label) emails as spam. It looks for the key words in the marked spam email and if it detects the repetition of those keywords in new mails, it will mark them as spam. There are many variants of classification tools found in the literature. The authors in [6] included several classification algorithms in top 10 data mining algorithms including *Support Vector Machine* (SVM), *K-Nearest Neighbour* (KNN), decision trees and Naive Bayes'. Each classifier has certain advantages and disadvantages and the selection of any particular classifier depends on the data characteristics.

SVM is one of the most widely used classification algorithm. The two main advantages which give SVM an edge on others are:

1. Its ability to generate nonlinear decision boundaries using kernel methods;
2. It gives a large margin boundary classifier.

SVM requires a good knowledge and understanding about how they work for efficient implementation. SVM works by mapping the training data into a high dimensional feature space, and then separates the classes of data with a hyper plane, and maximizes the distance which is called the margin. It is possible to maximize the margin in feature space by using kernels into the method, which result into non-linear decision boundaries in the original input space. The decisions about pre-processing of the data, choice of a kernel, and setting parameters of SVM and kernel influence the performance of SVM greatly and incorrect choices may severely reduce the performance of SVM as discussed in [7]. It is also possible to use mixture of different kernel functions for optimized performance and one such example is given in [8], where authors used SVM for image classification with kernel function which is mixture of *Radial Basis Function* (RBF) and polynomial kernel. The SVM algorithm requires extensive time in training but once the model is trained, it makes prediction on new data very fast.

Another efficient and simple classification algorithm is KNN, which is one of the simplest and instance-based learning techniques used for classification. It is a non-parametric algorithm, which means it does not make any prior assumptions on the data set. It works on the principle of finding predefined number of labeled samples nearest to the new point, and predicts the class with the highest votes. KNN memorizes labeled data in the training set and then compares the new data features to them. The advantage of KNN lies in simple implementation and reduced complexity. Despite its simplicity, it works quite well in situations where decision boundary is very irregular. Its performance is also very good when different classes do not overlap in feature space [9]. KNN is also called lazy algorithm as it takes zero effort for training. However, it requires full effort for predicting for new data points [9].

2.2.3. Clustering Analysis

Clustering refers to the unsupervised machine learning technique which is used for grouping similar objects on the basis of pre-defined metric such as distance or density. As opposed to the classification analysis, it does not involve training period. Clustering is a general technique used widely for knowledge discovery in big data sets and several variants of algorithms are found in the literature. The choice of a particular clustering algorithm and parameters (such as type of metric, threshold and number of clusters) is governed by the problem scenario and the data sets involved in the problem.

Clustering generates a hypothesis about the given data, whether the data under observation has distinct classes, or overlapping classes or classes with fuzzy nature. With the advent of big data, clustering applications have even increased. Clustering serve as the basic of many data mining techniques and finds applications in almost every industry. For example market researchers apply clustering techniques to group the customers into different segments, such that customers with common interests are in a same group. In this way, they can target different groups with different focused offers which are more related to them. Social networks apply clustering to group similar people into communities in order to help in recognizing people with similar interests from a larger set of people. Another interesting application of clustering is used by insurance companies to group different areas on the basis of crime rates so that they can offer different rates for insurances to the different areas depending on the level of insecurity. For more details, please refer to state of the art deliverable D2.2.2.

2.2.4. Regression Analysis

Regression analysis is one of the most widely used ML tool; it is used to define a relation between variables so that the values can be predicted in future using the defined relation. In

general, the dependent variable is called a *target variable* and the set of independent variables which form the input are called *predicted variables*. In order to elaborate a bit more, let us consider one simple and yet practical example of regression analysis e.g. to identify a relation between energy use and weather in order to predict the energy demand according to weather changes. In this example, energy usage is a dependent variable (target variable) and weather parameters such as outside temperature, wind and humidity are independent variables (predicted variables) that act as the input. It is a general observation that if the weather is cold, the energy consumption will be higher due to usage of heating systems as compared to mild or warm weather. If the dependent variable is relying on only one factor, it is called uni-variate model whereas like in our example if it depends on more than one variable (temperature, wind and humidity) it is called multi-variate model. The energy demand can vary linearly along the weather parameters; or there can be a non-linear relation between them. Non-linear regression models are generally more accurate but they come with the increased cost of complexity. Over fitting is also a common problem in non-linear models.

If one of the dependent variables is time, then the regression analysis can also be called *time series analysis*. Traditionally, statistical methods like ARMA and ARIMA were used for time series regression but recently the trend has shifted towards more sophisticated ML models such as *Support Vector Regression (SVR)* and *Artificial Neural Networks (ANN)* because of their robustness and ability to provide more accurate solutions.

SVR is an extension of SVM which is widely used for regression analysis. The main idea is the same as in SVM, it maps the training data into higher feature space using kernel functions and find the optimum function which fits the training data using hyper-plane in higher dimension. There are many examples found in the literature where SVR has been successfully applied for prediction such as predicting stock market feeds [10], electricity demands [11] and traffic flow [12]. The prediction of travel-time is an essential aspect of intelligent transportation systems. Different algorithms from statistical theory and recently from ML domain have been applied for predicting the travel time as the random nature of traffic makes it difficult to predict. In [13], the authors applied SVR for the same problem and demonstrate that it outperforms other statistical methods performance-wise. Recently, the hybrid regression models have also been in use in IoT and one such example is given in [14] where a combination of several regression algorithms was made to predict the short-term sensor readings.

2.2.5. Statistical Inference:

Classification Analysis methods described in previous section is an example of supervised ML model where parameters of the model are estimated from training data which consist of pairs of input and annotations of the output. The performance of supervised models is often quite good but the labeling of data poses an extra task. Data is often labeled manually, and although the process of labeling may be simple but it limits the amount of examples that can be classified using manual methods. In addition, as the annotation task becomes more complicated such as for labeling data for traffic congestion, annotations become far more challenging as well. In this regards, we explored a statistical inference approach based on *Hidden Markov Models (HMM)* for predicting an event in scenarios where labeled or annotated data is not available. For such scenarios, the training of a model for predicting output without the availability of annotated output data seemed counter intuitive, but it is possible using Expectation Maximization algorithms.

2.2.5.1. Hidden Markov Model

HMM is an extension of Markov model in which the states of the system are not directly observable as contrast to simple Markov models like Markov chain where state is visible to the observer. Markov model is an example of statistical model, which assumes Markov property to define and model the system. In simple words, Markov property states that the future state of the system only depends on the present state and does not depend on the past values. Each state has some probability distribution related to the output, and therefore sequence of outputs generated gives indirectly information about the sequence of states as shown in Figure 3 where at every sampling time instant, we have an observation which is related to hidden state. HMM is widely used for temporal pattern recognition.

There are many applications in literature where Markov model and HMM have been applied to find temporal relations within data. One such example is given in [15], in which authors showed the use of HMM to find the possible transition of events (states) depending on their temporal relation. The authors demonstrated their approach on the energy data which they gathered using the test bed installed in their research centre. They first applied clustering to find the possible hidden concepts or states and name it as weekday or weekend; and then applied HMM on top of it in order to show the possible state sequence and probability of transition from one state to another. For example, if there were consecutive four weekdays registered, then there will be very high probability that the next day will be a weekend. HMM can also be used as a classification tool for inferring events. For example, the authors in [16] applied it for accident detection. It is a quite generic tool and used extensively for classification purpose in different fields as evident by their use in [17]. The main difference from the classification analysis techniques discussed in previous section is that HMM is a parametric technique which is based on the statistical properties of the underlying data with respect to time. HMM can also be used for predicting the optimal and most probable sequence of states or outcomes using *Expectation Maximization* (EM) algorithm as used by authors in [18].

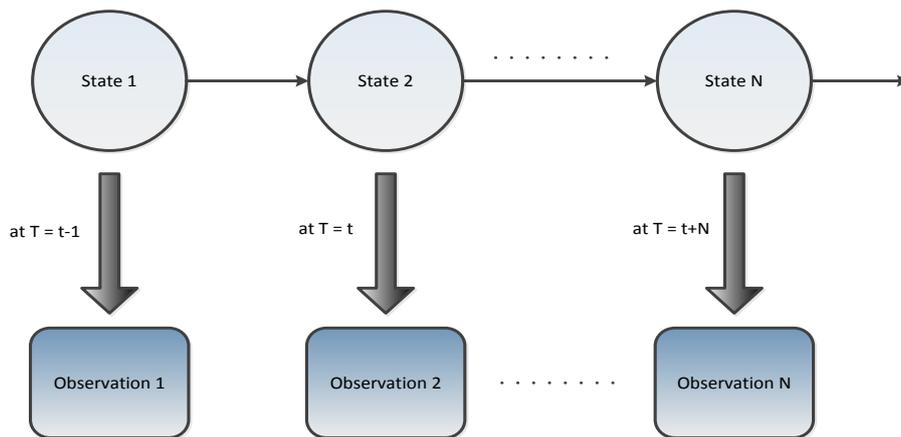


Figure 3: Hidden Markov Model

A complete description of HMM requires specification of total number of hidden states N , and the specification of three probability measures represented by $\lambda = (A, B, \pi)$, where:

- π represents the set of prior probabilities for every state;
- A represents the set of transition probabilities a_{ij} to go from state i to state j : $a_{ij} = P(q_t = S_j | q_{t-1} = S_i)$ and collected in matrix A ;

- B represents emission probabilities which is the probability of observation in a particular state: $b_{i,k} = P(x_n = v_k | q_n = S_i)$, the probability to observe v_k if the current state is $q_n = S_i$. The number $b_{i,k}$ is collected in a matrix B .

There are three basic problems of HMM in order to apply for real world applications. They are:

- **Problem 1:** Given the observation sequence $O = O_1, O_2, \dots, O_t$ and a model $\lambda = (A, B, \pi)$, how to efficiently compute the probability of observation sequence given the model $P(O|\lambda)$?
- **Problem 2:** Given the observation sequence $O = O_1, O_2, \dots, O_t$ and a model $\lambda = (A, B, \pi)$, how do we choose a corresponding state sequence $Q = q_1, q_2, \dots, q_t$ which best explains the state sequence?
- **Problem 3:** For the given observation sequence $O = O_1, O_2, \dots, O_t$, how to find a model parameters which maximize $P(O|\lambda)$?

The problem 3 is of particular interest for unsupervised learning in which only observations are given. In this case, the hidden states are estimated using Expectation Maximization (EM) algorithms based on Maximum Likelihood theory.

2.2.6. Large-Scale IoT Data

Big data is no more a fantasy and has now become a necessity for many applications. Across different applications in IoT, we have seen more data results into more effective algorithms and more meaningful abstraction of knowledge. Most of the applications in IoT are of large scale: smart buildings, intelligent transportation system and real-time management of supply chain logistics are just few examples dealing with vast repositories of data. Therefore, any practical application must fulfil the requirement of scaling up to datasets of interest.

The ever-increasing data has resulted into paradigm shift for finding new optimized techniques in contrast to traditional methods. Distributed computing is one research area which is explored to carry complex analytics in parallel on multiple machines and form the basis of MapReduce [19]. Whereas, cloud storage becomes the main candidate for providing scalable solutions for managing large data sets. In our work, we have explored the use of Apache Spark [20] for providing parallel, scalable and distributed platform for carrying analytics on large IoT data and integrated it with Openstack Swift Object storage and used both storlets [21] and metadata search for carrying out analytics close to the storage. More details regarding these aspects can be found in deliverable D4.1.3.

In the following section, we briefly describe the requirements which drive our design decisions for large scale analytics.

2.2.6.1. Requirements

The conventional approaches of dealing with large data problems are proving to be insufficient and require a distinct approach. In this section, we briefly summarized the requirements for dealing with large data problems:

- **Time latency:** Time latency is an important requirement for large data processing solutions. The proposed method should be able to finish processing data in a given time frame. The value of insights gained from the data may lose its importance if the underlying process takes longer time. The statistical properties of the data underlying may change over longer period of time analysis.

- **Scaling out:** For large data problems, scaling out (using large number of low-end servers) is the preferred choice as compared to scaling up (using small number of high-end servers). The latter approach of using machines with very high processing capabilities is not an optimum choice as the cost of such machines does not scale linearly i.e. a machine with twice as memory and processing capabilities will be of significantly more than the double of machine).
- **Reliability:** A large data programming model/framework should be able to deal with failures of the underlying hardware. In big data applications, failures of the underlying hardware are inevitable due to the large number of physical machines. Consider a simple example; a large cluster is made of 1000 highly reliable machines with *Mean-Time Between Failures* (MTBF) of 1000 days. Even with these reliable servers, the system will experience roughly one failure/day.
- **Analytics close to data:** It represents one of the fast gaining requirements for large scale data processing applications. The ability to gather and store data is developing overwhelmingly as compared to the ability to process it. Big data is no more a fantasy and has now become a necessity for many applications. Across different applications in IoT, we have seen more data results into more effective algorithms and more meaningful abstraction of knowledge. Most of the applications in IoT are of large scale; smart buildings, intelligent transportation system and real time management of supply chain logistics are just few examples dealing with vast repositories of data. Therefore, any practical application must fulfil the requirement of scaling up to datasets of interest.

2.2.6.2. ***Distributed and Scalable Machine Learning using Apache Spark***

In many ML algorithms, the optimization algorithm (such as gradient descent) involves accessing the same dataset iteratively to optimize certain parameters for finding optimum decision boundary or an appropriate function. MapReduce [22] enables to represent each iteration as a MapReduce job, but each job must reload the data from a disk which incurs a significant amount of performance penalty.

Spark is designed to address the limitations of MapReduce in order to overcome the bottleneck issues caused by disk I/O retrieval and to improve the performance for iterative algorithms. Spark provides the ability to run computations in the memory which enables it to provide much faster computation times for complex and iterative applications as compared to systems based on traditional MapReduce, such as Hadoop etc. Spark is designed to be fast and general purpose at the same time. It is designed for different data analysis applications which include iterative algorithms such as ML applications, batch applications and real-time streaming applications. By providing a generic optimized framework, it provides a generic data analysis platform which makes it suitable for a wide variety of IoT applications.

Spark [20] achieved fast distributed computing with the help of the *Resilient Distributed Dataset* (RDD), which represents a read only collection of objects which can be split into many partitions. The different partitions of an RDD can be computed on the different machines across the cluster. The programs can cache an RDD in memory which can be reused in multiple parallel operations like in MapReduce, which is the main reason for the fast performance of Spark. RDDs are quite generic and can contain objects of type Scala, Python and Java. It can run on Hadoop Yarn manager and can read data from Hadoop Distributed File System thus making it extremely portable for running on different systems.

2.2.6.3. Analytics Close to data using Storlets and Metadata Search

Although object storage can store objects, manage them, protect them in a highly scalable way, it does not itself dramatically increase the rate at which we can extract value from objects. A new research prototype developed by IBM called storlets is developed on the concept of converting the software-defined object store into a smart storage platform. It aims at greatly increasing the value of the data that can be retrieved from the storage and the speed at which we can access what we need. Storlets allow the object storage to move the computation close to the data, as opposed to the traditional approach of moving the data to the server by the system to run computation.

The impact of storlets is of considerable importance. Stored data can be processed locally, and no longer needs to be transferred over the network to a remote computer, processed and then put back onto the storage server, all of which incurs both network transfer latencies and extra costs. Storlets aim at reducing costs, providing enhanced flexibility and improving security at the same time by turning the object store into a platform and extending the functionality of the object store using software. Saving bandwidth by avoiding unnecessary data transfers is not the only advantage which storlets provide; storlets provide perfect ways to introduce new services: storlets can analyze each object and extract its metadata including size, subject, format and more.

In Year 2 we also introduced another method for moving computation closer to the storage and reducing network bandwidth. When Spark SQL accesses the object storage, metadata search on the object metadata is used to restrict the queries to only those containing data relevant to the query. This reduces the number of objects read from disk as well as sent across the network.

More details can be found in COSMOS deliverable D4.1.3.

2.3 Functional Overview

This component is responsible for providing high-level knowledge from raw IoT data using different pattern recognition techniques. In this context, we have explored several supervised ML techniques including different variants of SVM and KNN and statistical techniques such as HMM which were explained earlier. Also, it provides the capability for predicting ahead using time series historical data. In short, it provides the two following main functionalities.

- 1) If labelled historical data is available (raw data with labelled high-level knowledge), it provides the functionality to train the model and provides capability to deploy the model in order to predict the output for real-time data;
- 2) If the labelled historical data is not available (incomplete data), it exploits the temporal patterns of the data and learns using statistical properties of the data to train the model which can be used to predict the output for real-time data;
- 3) It provides time series prediction mechanisms which enable to predict ahead and form the basis for pro-active IoT applications.

All functionalities are provided in the context of large-scale IoT data.

2.4 Connection with other Components

The inference and prediction FC can take both raw data and VE data as input depending on the scenario and the usage required by application developer. Figure 4 shows the needed inter-component interactions steps for carrying inference/prediction on raw data. Raw data can be pre-processed before publishing on the Message Bus FC (WP4) under specific topic, and then stored on the Cloud Storage FC (WP4) using Data Mapper FC (WP4). Inference and Prediction FC retrieves the data from the Cloud Storage FC using the provided interface.

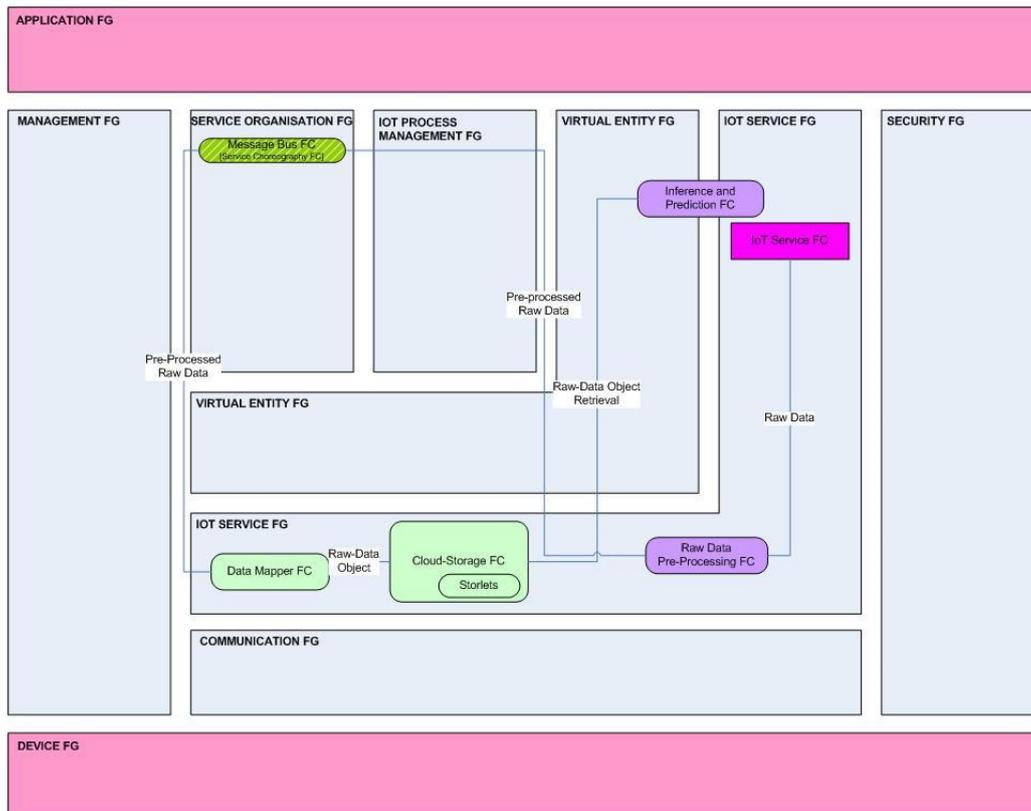


Figure 4: Inference/Prediction on raw data

Figure 5 below shows the inter-component interactions steps needed for inference and prediction on VE data. The flow of data is almost the same with the additional capability of running storlets for pre-processing on the object storage and within Spark for complex ML tasks. Pre-processing becomes an important step when running analytics on large-scale complex IoT data.

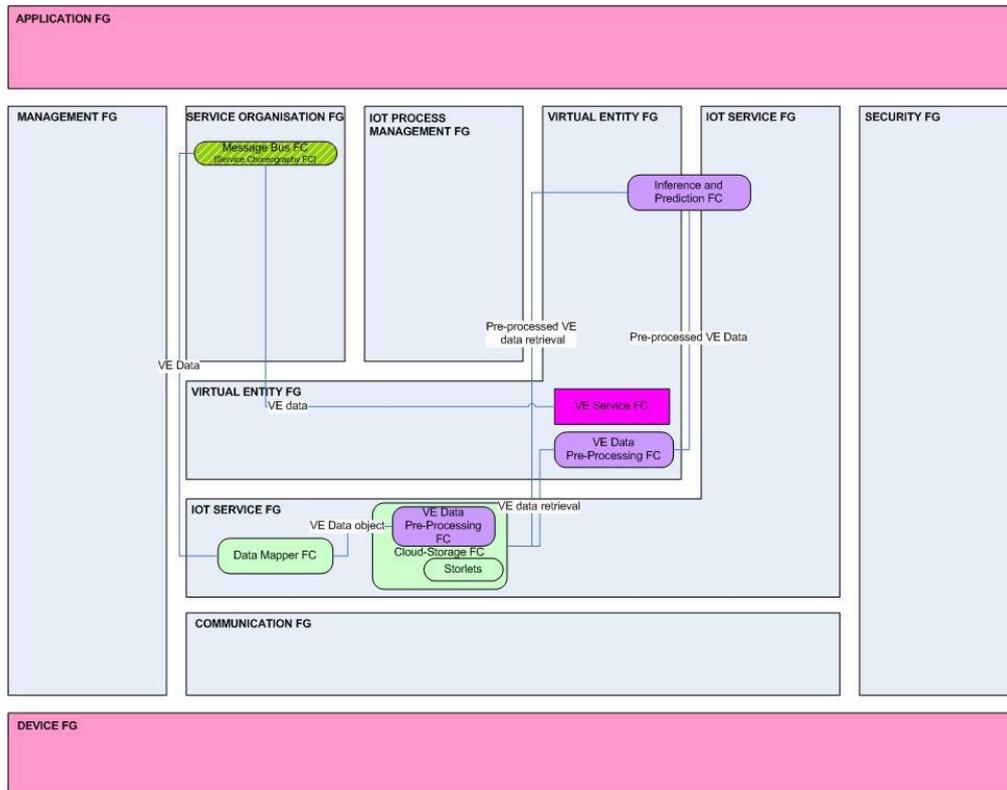


Figure 5: Inference/Prediction on VE Data

2.5 Interfaces

The application developer can use the models provided by COSMOS for inference of high-level knowledge and for predicting events. In order to use the existing models, the application developer will have to define the input features and output entity in which application developer is interested in. In order to achieve this, the application developer will use the COSMOS storage services in order to access historical data for the construction of off-line model. Then when a model is available on-line update of the model can be done under certain circumstances. Once enough data is collected, and therefore a model available, the Prediction component/functionality will be instructed for making a prediction based on the available model. The resulting model will be persisted and semantically annotated in order to make the model retrievable for later use. Using this approach, the semantic description of VEs and IoT services or data bus topics, could also include a reference to a prediction model (if available).

Since COSMOS is intended to be used by different actors and forge cooperation and reuse, prediction models can be built by different parties (for instance in the case of public data) provided that they are semantically described and linked to the data sources. Once stored and annotated, other actors will be able to query the semantic store for prediction models and use them according to their needs. The interfaces are summarized below:

Application Developer-Inference/prediction Block API:

An API will be provided to connect the application developer to Inference/Prediction block for the following purposes:

- 1) In order to select the particular prediction model and to define the input and output for the model;
- 2) To select the specific type of pre-processing required for the application

Client/VE-Inference/Prediction Block API:

An interface will be provided between the client/VE and inference/prediction block which will serve the following two purposes:

- 3) A Client or a VE will send a request using the API to register its interest in particular topic stating the interested characteristics/services (Occupancy state of a room, Traffic conditions on a road);
- 4) A Client or VE can also use API to get required prediction value at particular instant. An example can be a client sending a query to prediction block to find the traffic state at particular location.

Storage-Modeling Block Interface:

All the historical data is stored in the form of objects in object storage. Machine Learning models require an access to historical data for training purpose. In this regards, modeling block should be connected to the object storage.

2.6 Use-Case Scenario Extension from Year 1

This use-case demonstrates the use of pre-processing techniques and knowledge extraction techniques at the same time. We have extended the use-case scenario from year 1 for detecting occupancy state from electricity consumption data with the use of storlets for pre-processing and providing the means for extracting knowledge from large-scale IoT data in a distributed way using apache Spark. The work done in this section is based on the idea of moving computations to the storage in order to avoid the bottleneck caused by the bandwidth limitations of the network. We have proposed a novel solution based on our initial findings which show encouraging results and are summarized in this section.

2.6.1. Background

The amount of data in IoT is increasing exponentially and the trends indicate further increase in data size. The ability to store data is developing overwhelmingly as compared to the ability to process it. The fact that the improvements in the storage have outmatched the improvements in processing data to the extent where even the ability to read could not match what we store is itself distressing. The tendency of data storage has increased from tens of megabytes to magnitude of terabytes in last couple of decades whereas the increase in latency and bandwidth performance has merely improved to the factor of tens or perhaps hundreds.

It is a common practice for most of the existing data analysis platforms to have separate processing nodes and storage nodes such as amazon platform which has EC2 for computing and S3 for storage. Many data-intensive workloads are not very much processing demanding and the transfer of data from storage nodes to the processing nodes proves to be a bottleneck in such scenarios. In our work, we have explored the possibility of moving computations to the storage in order to reduce the amount of data transferred over the network. The initial results have shown significant potential towards reduced complexity and providing scalable solutions.

2.6.2. Results and explanations

We have extended the same scenario of occupancy detection discussed in the previous section. The sampling rate of the data gathered for occupancy detection scenario was 10 seconds. The sampling rate of 10 seconds might seem redundant for occupancy detection as it is highly unlikely that the state of the user will change in the period of 10 seconds. We explored the possibility of applying aggregation techniques and analyse their effect on the accuracy and time complexity of different data mining algorithms which were discussed in 2.1.

Figure 6 shows the comparison of accuracy and time complexity of the following classification algorithms.

- 1) K Nearest Neighbor (KNN)
- 2) Support Vector Machine (SVM) with Radial Basis Function (RBF) kernel
- 3) Support Vector Machine (SVM) with linear kernel

Whereas, Figure 7 shows the results for the *Hidden Markov Model* (HMM). Although the accuracy of algorithms has dropped a little but the gain in reducing time complexity is huge.

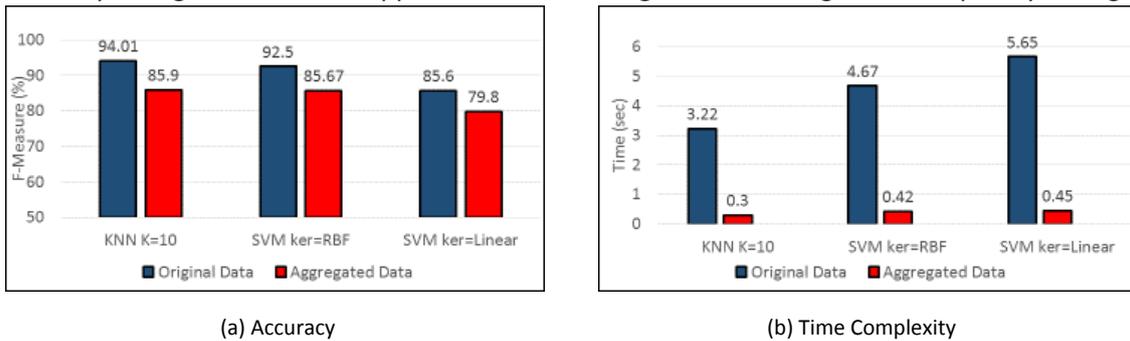


Figure 6: Aggregated data comparison with non-aggregated data for Machine Learning algorithms

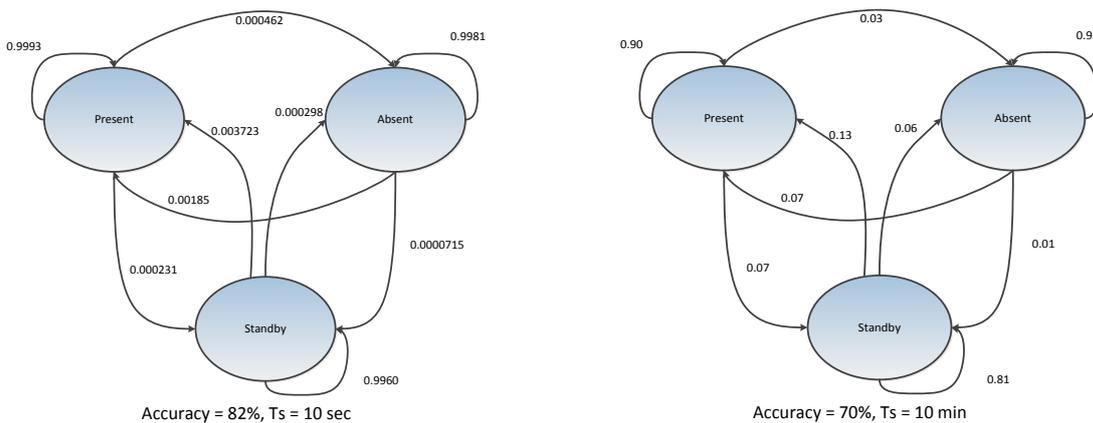


Figure 7: Aggregated data comparison with non-aggregated data for HMM

One might argue that applying pre-processing tasks before the data is stored using middleware platform is a better option. In our case this is not an optimum choice as occupancy detection is one application which can be inferred from raw data at low sampling rate. There are other applications such as *Non-Intrusive Load monitoring* (NILM) [23], for detecting which devices are being used and requires very high sampling rate. In such a scenario, if we apply pre-processing at middleware platform, it will limit the use of data to specific applications and valuable information would be lost. In this regard, we proposed to store all raw data and apply pre-processing techniques at the object storage side. The two main advantages of our approach are:

- 1) The amount of data traveling over the network will be reduced to the factor of aggregation. In the example above, we applied *Piece-wise Aggregation Approximation* (PAA) with the fixed window size of 6 samples which directly resulted into the six times less data transfer over network.
- 2) The storlets can also be used extracting basic form of knowledge directly, thus reducing the time and data over network and increasing the web services provided by

the system. For example, if we are interested in the readings at particular day and time, a simple storlet can be written which takes day and time as input and will return only the required data. Whereas, in a conventional method, system will have to load the file in memory and scan it to get the required information. For large data files, it results in latency and large amount of data transferred over network.

2.7 Use-Case Scenario 2 – Pro-Active Traffic Management

2.7.1. Introduction

Although CEP provides solutions for dealing with data streams in real-time, it lacks the predictive power provided by both *Machine Learning* (ML) and statistical data analysis methods. Most of the CEP applications are intended to provide reactive solutions by correlating data streams using predefined rules as the events happen but do not exploit historical data due to limited memory. However in many applications, early prediction of an event would be more useful than detecting it when it has already eventually occurred. For example, it will be more useful to predict the traffic congestion as compared to detecting it.

On the other hand there are several methods from ML and statistical analysis domain which have the ability to provide innovative and predictive solutions; however they are unfortunately not suitable for analysing data in real-time. ML methods exploit historical data and apply diverse disciplines such as probability and artificial intelligence in order to train the models for making predictions about future states. They do have the potential to provide the basis for proactive solutions for IoT applications but they lack the power of scalability and processing multiple data streams which is provided by CEP.

In our work, we exploit both approaches and propose an architecture based on CEP and ML in order to provide a proactive solution for IoT applications that leverage best of both approaches. The promise behind our work is that if the input to the CEP is predicted data, then the complex event detected by CEP using causal and temporal pattern recognition techniques will be a predicted complex event. In contrast to the current prediction methods which are based on static model parameters, we propose an adaptive prediction algorithm called *Adaptive Moving Window Regression* (AMWR) for dynamic IoT environments, which utilizes moving window for training the model and updates the model as new data arrives. It tracks down errors and prevents their propagation. The size of the training window can be found automatically which optimizes the performance for specific dataset and the size of prediction window is also adaptive in nature in order to ensure certain accuracy in the prediction.

2.7.2. Proposed solution

The proposed architecture illustrating our approach is shown in Figure 8. It consists of several components from different functional groups of COSMOS architecture and is showed here in order to demonstrate the overall picture. Prediction capabilities using AMWR were developed and evaluated in this work package. Node-Red provides the front-end for our architecture where data from different sources such as MQTT or REST are accessed and after performing pre-processing such as converting to Json and filtering redundant data is published under a specific topic on the Kafka message broker. More details about the Kafka and Node-Red can be found in the next section. AMWR block accesses the real-time data from the Kafka topic and publishes the predicted data under the specific topic on Kafka. The event collector module of CEP is listening to events and as soon as new events are available, it collects the events, performs pattern matching using CEP engine and produces the resulting complex events.

Node-Red provides us a visual tool for performing pre-processing such as filtering or data conversion and converting the external data sources in a data format defined for our architecture. In this way, we can add any external data sources without changing the configuration of the internal components. Kafka is an open source message broker where messages can be transmitted.

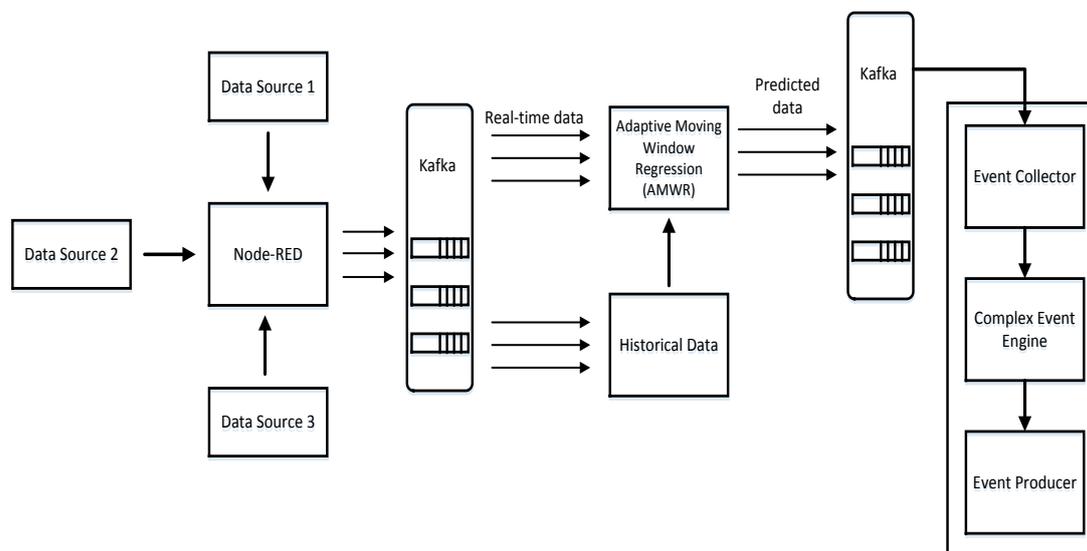


Figure 8: Proposed solution for pro-active traffic management

2.7.3. Adaptive Moving Window Regression

We proposed and developed an adaptive prediction algorithm called AMWR for dynamic IoT data. In general, prediction models are trained using large historical data and once the model is trained it is not possible to further update the model. In real-time dynamic environments, the performance of the model may deteriorate over time due to change in statistical properties of underlying data. The context of the application may change resulting in the degradation of prediction model performance. For such scenarios, we propose a prediction model which utilises moving window of data for training the model and once new data arrives, it calculates an error and re-trains the model accordingly. The optimum size for training window is found graphically and is specific to the underlying data stream. Our proposed approach is adaptive in nature as it tracks down errors and prevents it from propagating by retraining the model periodically. The size of prediction window or forecast horizon is also adaptive which is derived by the performance of the model in order to ensure a proper reliability in the prediction. There are three main steps involved in the implementation of AMWR which are described below.

- 1) Selection of regression algorithm;
- 2) Optimum Training window size;
- 3) Size of prediction horizon.

2.7.3.1. Selection of regression algorithm

Traditionally, statistical methods like ARMA and ARIMA were used for time series regression but recently the trend is shifted towards more sophisticated ML models such as different

variants of *Support Vector Regression (SVR)* and *Artificial Neural Networks (ANN)* because of their robustness and ability to provide more accurate solutions.

We have implemented our approach using SVR due to its ability to model non-linear data using kernel functions. SVR is an extension of SVM which is widely used for regression analysis. The main idea is the same as in SVM, it maps the training data into higher feature space using kernel functions and find the optimum function which fits the training data using hyper-plane in higher dimension.

2.7.3.2. Optimum training window size

The choice of optimum training window size for machine learning models is an open research issue. In general, the accuracy of prediction model increases as the size of training data increases which reflects to have large historical data for training prediction models so that it covers all possible patterns spanning time series. Although this approach generates generic and accurate model for prediction, there is one major drawback associated with it: if the behavior or statistics of the underlying data changes, the trained model is unable to track the changes and results into erroneous readings; the error will then start accumulating in future predictions.

In contrast to this approach, researchers have proposed to use the moving window for training the models in which most recent data is feed into the models [37]. The size of the optimum window is an open research problem with no generic solution. A large window size can have more accurate results but it increases the complexity of the model making it unsuitable for real-time applications whereas a small window size can result into an increased error and hence effecting the reliability of the system.

In order to overcome this issue, we proposed a novel and generic method based on time series analysis -called Lomb Scargle method- to find the optimum window size and validate our results on a real-world data. In our method, we exploited the inherent periodic nature of most of the time series data of real-world. Time series data consists of three main components; trend, seasonal and the random component. The use of the moving window minimizes the error caused by trend component for prediction as the most recent data is used for training the model. The error induced by the random component is Gaussian by nature which is averaged out over a set of predictions. In order to minimize the seasonal component, we proposed the window size equivalent to the seasonal component of the time series and proposed to find it by spectral analysis of the data stream. *Fast Fourier Transform (FFT)* algorithm is the most commonly used method for finding spectral components by searching for the sharp peaks in the periodogram calculated by Fourier transform of the time series. FFT requires the time series to be evenly spaced which is not always possible for most of the IoT data streams. Missing values is a common phenomenon in IoT and the inability of FFT to deal with it makes it unsuitable for our system.

For such systems, another method called *Least-Squares Spectral Analysis (LSSA)* or more commonly known as Lomb Scargle can be used to find the periodicities in a time series data. Lomb first proposed the method while studying variable stars in astronomy which is based on the least squares fit of sinusoids to data samples.

2.7.3.3. Prediction horizon

Prediction horizon is the number of steps ahead for which the prediction is made. In our work, we propose to have an adaptive size for prediction horizon in order to always ensure a proper level of accuracy. The intuition behind our approach is to increase the size of prediction window if the accuracy of model is high and to decrease it if the performance of prediction model decreases. The performance of the model is evaluated by comparing the predicted data with the actual data when it arrives.

2.7.4. Other components

2.7.4.1. Kafka

In our architecture, we have used apache Kafka [38] as the message broker. It is also an open source tool for real-time publishing and subscribing of messages or data. It provides a scalable architecture for high throughput data feeds with very low latency. It was developed by LinkedIn and was open sourced later in 2011. Like other publish-subscribe messaging systems, Kafka maintains feeds of messages in topics. Producers write data to topics and consumers read from topics. Since Kafka is a distributed system, topics are partitioned and replicated across multiple nodes. A single topic can have one or more consumers. Messages are simply byte arrays and the developers can use them to store any object in any format – with String, Json, and Avro the most common. In our architecture, all the messages are published in Json format. What makes Kafka unique w.r.t. other available systems is its persistent nature to hold the messages for a set amount of time in the form of a logs (ordered set of messages). For more information see deliverable D4.1.3.

2.7.4.2. Node-Red

Node-Red serves as the front-end interface for our architecture. Node-RED is an open source visual tool which is used extensively for wiring the Internet of Things. It provides APIs for connecting different components, and with the help of user provided Java code, it can be used to filter the data change the format as well.

IoT has provided the researchers with a global view enabling access to truly heterogeneous data sources for the very first time. These data sources can be RestFul web service, MQTT data feed or any other external data source. Data format is not limited to any specific format in IoT. XML and Json are two most commonly used formats which are used extensively for transmitting IoT data. Also, different data feeds from different sources may contain data which might be redundant for a specific application which needs to be filtered out. Node-Red provides pre-processing functionalities such as filtering and conversion of data into specific format using graphical interface. It accesses the heterogeneous data source, filters the relevant data and converts it into Json format and publishes it under the specific topic on the internal Kafka message bus.

2.8 Conclusion

In this component, we have explored several Machine Learning methods for extracting high-level knowledge and providing basis for innovative IoT applications. We have adopted an incremental approach where the functionalities of the component were gradually improved every year. We started with the use of supervised machine learning methods and highlighted their drawbacks in Year 1. In Year 2, we addressed those drawbacks using unsupervised methods and demonstrated their application for innovative applications. In Year 3, we



explored time series prediction mechanism and demonstrated how they can be optimized for IoT data streams and provide the basis for proactive IoT applications.

3 Pre-Processing Functional Component

3.1 Introduction

Pre-processing is an important step for applying ML algorithms in IoT due to many reasons. Most of the devices are connected with wireless links in a dynamic environment and resource constraint nature of these devices affects the communication link and their performance. The deployment of cheap and less reliable devices is a common practice in IoT for bringing the overall cost of a system down with the unfortunate resulting flaw of missing values, out of range values or impossible data contributions. The sentence “*garbage in garbage out*” fits perfectly to many ML algorithms. The amount of data is increasing exponentially in IoT and the processing of such large data with minimum time latency is an important factor which can be optimized by the use of proper pre-processing methods. Several aggregation techniques are commonly used in IoT for reducing the total amount of data travelling through the network. In this context, *Piecewise Aggregation Approximation* (PAA) and *Symbolic Aggregation approximation* (SAX) are the most common techniques. Data collected in IoT can be a combination of many features. The number of features plays an important role in the complexity and the performance of a ML algorithm. In few scenarios, these features are correlated with each other and it is possible to reduce the features by representing the data using new uncorrelated features using statistical analysis such as *Principal Component Analysis* (PCA). A brief introduction about the most commonly used pre-processing techniques is given below.

3.1.1. Data Cleaning

In real time dynamic environment, a faulty or a missing sensor reading may occur due to bad communication channel or loss of service. The missing values can result in irregular time series or incompatible vector size as compared to other devices connected. A simple data cleaning method involves then filtering out-of-range values and filling out missing values. Missing values can be filled by the mean value of the sensor over some time window, by last recorded value or by simple interpolation methods using historical data.

3.1.2. Data Transformation

Data transformation involves transforming the data into the form that is optimum for ML process. Feature scaling is an important example which is used extensively as a pre-processing step for ML algorithms [24]. The range of values of different features is on different scales. If one of the features has considerably wider range as compared to other features, the optimization function for the ML algorithm will be governed by that particular feature and will affect the performance of algorithm. Also, it will take much longer time for the optimization objective to converge in such cases. Feature scaling is therefore a method that brings all the features on the same scale making sure they contribute equally to classification algorithm. Standardization is most commonly method used for feature scaling which involves having each feature represent as a Gaussian like curve with zero mean and unit variance.

3.1.3. Data Reduction

Data reduction is perhaps the most important pre-processing step when dealing with very large data sets. Several variants of aggregation techniques are used in order to reduce the data size without any loss of information. PAA and extended version of PAA called SAX are the most commonly used aggregation techniques in IoT for data reduction.

PAA is a simple dimensionality reduction method for time series analysis. In order to reduce time series from n dimensions to m dimensions, the data is divided into m equal sized frames and the mean value is calculated for the data lying in that frame. The vector of mean values calculated will represent new series with m dimensions.

Feature extraction is another technique used widely for data reduction where the number of features of data is large and mostly correlated to each other. Feature extraction enables to extract most relevant and uncorrelated features in order to perform optimum analysis. PCA [25] is one of the most famous feature extraction technique which form new uncorrelated features based on the statistical properties in order to represent the same data with less dimensions. PCA is widely used in image processing and pattern recognition systems.

3.2 Functional Overview

This component provides several functions at VE and platform level which are summarized below:

- 1) It provides the capability for basic pre-processing on raw data streams including filtering, selecting or aggregation on real-time data using light weight CEP running on VE;
- 2) It provides the capability to run pre-processing on the object storage using storlets. Storlets can be both generic and domain specific as well. For domain specific storlets, an application developer can write his own storlets which can be run using restful web interface;
- 3) It provides the functionality to run different pre-processing techniques including aggregation, interpolation, and data cleaning and feature scaling using Apache Spark on the platform level.

3.3 Connection with other Components

Pre-processing FC provides the capability to pre-process raw data using light weight CEP running at VE level and to pre-process VE data using storlets and Spark. Both connections are shown in the Figure 9. Pre-processing on raw data provides simple functionalities like filtering or aggregating the data whereas VE data pre-processing is more complex and provides a fundamental step for ML algorithms. Also, application developers can write domain specific pre-processing using storlets.

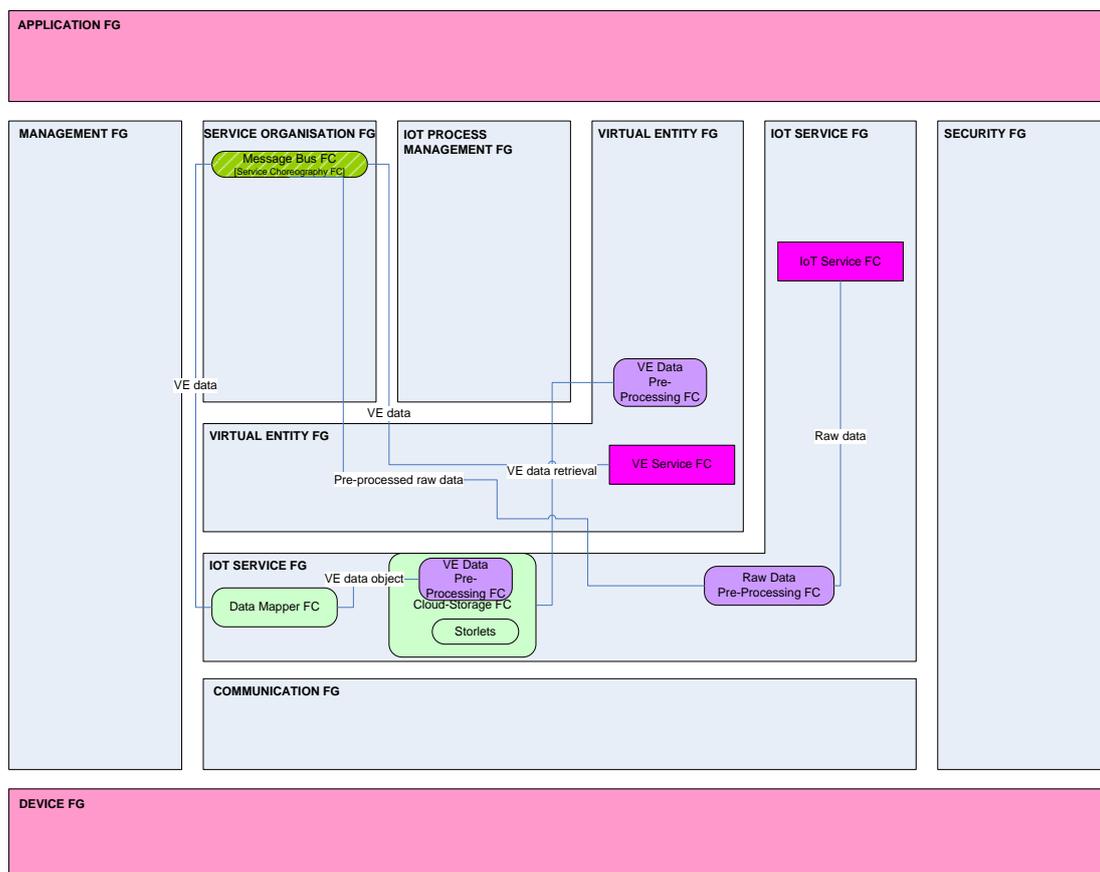


Figure 9: Connection of pre-processing FC with other components

3.4 Interfaces

As shown in the Figure 10, Pre-processing component is present at three different levels and the interfaces for every level are described below:

VE Level using CEP

An interface will be provided for application developer to define the pre-processing using DOLCE language for the CEP. The output will be published to under the specific topic on message bus.

Object Storage using Storlets

An application developer can code their own storlets in java or can use generic storlets provided by COSMOS for aggregation. Aggregation storlet will take the aggregation factor and the input features as input.

Platform level using Apache Spark

Different pre-processing methods are provided which can be specified when choosing a particular data mining method.

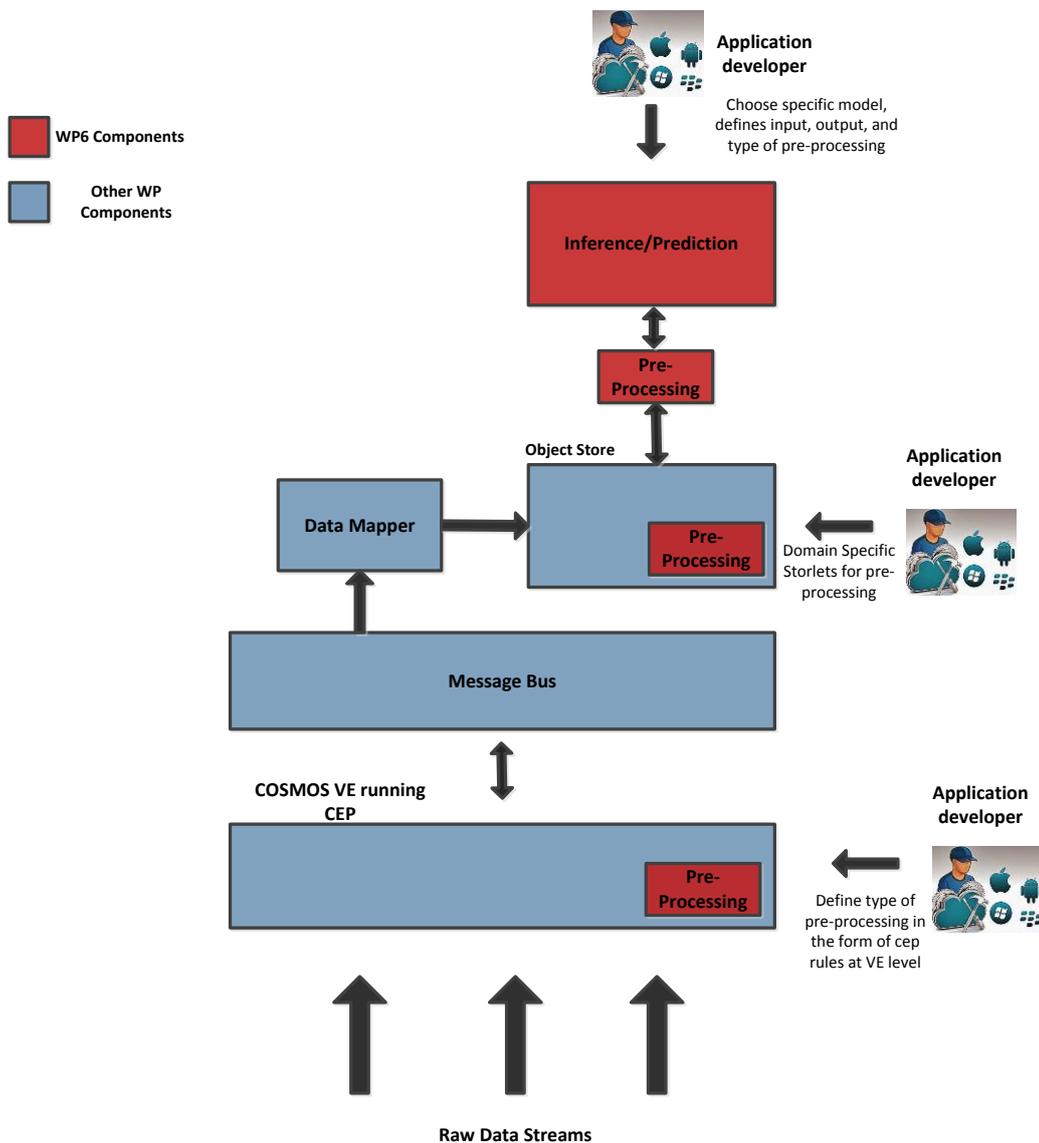


Figure 10: Pre-Processing component

3.5 Use-case scenario

A use-case scenario for the pre-processing has already been described in section 2.6 where we discussed using aggregation as a pre-processing step using storlets for reducing the amount of data travelling over the network and hence reducing the time complexity for machine learning algorithms.

3.6 Conclusion

Pre-processing is an important step for running analytics on large IoT data sets. It is possible to carry pre-processing at different levels within the platform depending on the application. Pre-processing on raw data can filter the data and only store the relevant data, while pre-processing after storage can be useful for running different types of complex analytics.

4 Event (Pattern) Detection Functional Component

4.1 Introduction

This component is intended to provide context-aware stream processing for extracting complex events from raw IoT data streams. In this regard, we proposed and developed a novel method for combining machine learning with CEP for detecting events automatically in an optimized manner. The combination of heterogeneous data sources in IoT form complex patterns where different patterns represent different hidden events. The objective of this component is to extract those hidden events by learning and matching these patterns in near real-time and to contribute towards more context-aware system. In this regard, this component is aimed to detect following two types of events.

- 1) Complex Events by correlating different data streams
- 2) Unusual events by exploring Anomaly detection methods

4.1.1. Background

There are many applications in IoT which require real-time processing of data such as intelligent transportation systems and smart buildings. As an example, consider an intelligent transportation system: analytics on the historical data can enable the researchers to have an idea about the average traffic flow which can assist when taking important decisions like managing traffic signals and the location of important buildings such as hospitals and fire stations (as they should be located around the places where less dense traffic provides easier access). But real-time traffic flow is a random and dynamic process and any incident or external factor such as bad weather or rain can affect the traffic flow and result into traffic congestion. Historical data analysis does not provide the mean to detect traffic congestion in real-time. Real-time monitoring of traffic system is a very important aspect of smart city applications. Traffic has to be managed in an optimized way by analysing parameters like traffic flow and traffic density in real-time. There are thousands of sensors deployed in smart cities generating gigabytes of data per day. The process of analysing, inferring and correlating these data streams in real-time requires therefore alternate solutions.

The Research area of CEP includes processing, analyzing and correlating event streams from different data sources to infer more complex events in real-time. The main objective of CEP was to provide the processing capability of big data engines which enables it to analyze the data and extract patterns on the run in real-time with distributed architecture. The main difference from the traditional big data engines was that CEP can handle multiple events which are seemingly unrelated and can correlate them in order to provide a desired and meaningful output. In the early years of data processing, data was at rest but now data processing is taking place as the information is in motion. Sensors data, stock market feeds and GPS data from vehicles produce data in the form of real-time data streams.

CEP engines require rules or patterns to detect an event from data streams which have to be given manually by the administrators of the system. Based on this, there is an assumption that administrators have the required background knowledge which sometimes is neither available nor so precise. The manual setting of rules and patterns limits the use of CEP only for expert's domain and poses a weak point and even though with prior expertise and knowledge, experts are prone to make errors in choosing optimized parameters for dynamic systems. Systems based on CEP deploy static rules and there is no means to update the rules automatically. In

real-time dynamic scenarios, parameters of a system may change and performance of CEP may deteriorate in such dynamic scenarios.

4.1.2. Complex Event Detection

In our work, we have addressed the above mentioned limitations of CEP using our novel approach based on adaptive clustering in order to exploit historical data and find optimized parameters for CEP rules automatically for extracting complex events from real-time data streams. The parameters for CEP rules depend on the threshold values in order to differentiate between different events. We explored adaptive clustering approach for finding threshold values and update the rules accordingly. Our approach is able to follow the changes in the distribution of data and adapt to it in run time. Our proposed architecture is able to infer complex events from raw data streams in a distributed manner and is able to provide adaptive solutions.

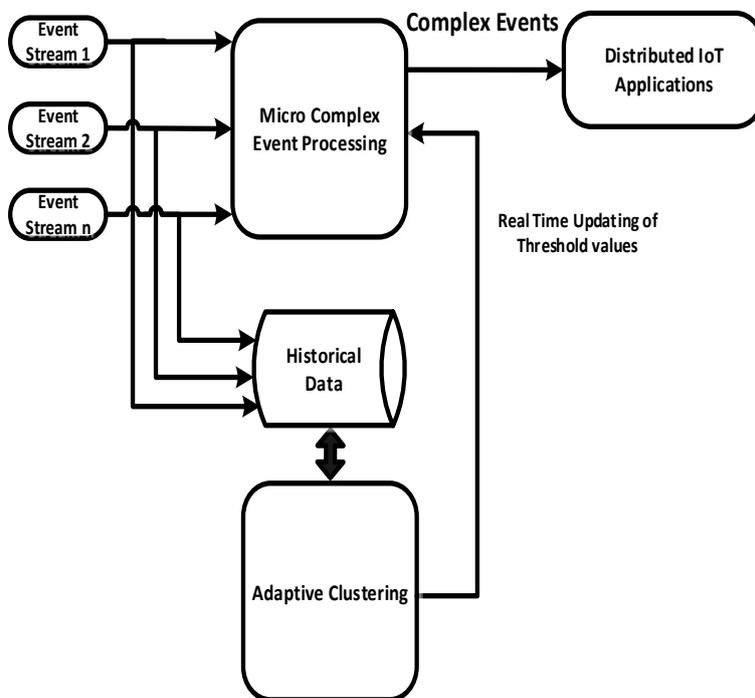


Figure 11: Proposed solution for Event Detection

Clustering refers to the unsupervised ML technique which is used for grouping similar objects on the basis of pre-defined metrics such as distance or density. Cluster analysis is a general technique used widely for knowledge discovery in big data sets and several variants of algorithms are found in the literature. The choice of a particular clustering algorithm and parameters (such as type of metric, threshold and number of clusters) is governed by the problem scenario and the data sets involved in the problem.

Cluster analysis generates a hypothesis about the given data, whether the data under observation has distinct classes, or overlapping classes or classes with fuzzy nature. With the advent of big data, clustering applications have even increased. Clustering serves as the basic of many data mining techniques and finds applications in almost every industry. For example

market researchers apply clustering techniques to group the customers into different segments, such that customers with common interests are in a same group. In this way, they can target different groups with different focused offers which are more related to them. Social networks apply clustering to group similar people into communities in order to help in recognizing people with similar interests from a larger set of people. Another interesting application of clustering is used by insurance companies to group different areas on the basis of crime rates so that they can offer different rates for insurances to the different areas depending on the level of insecurity.

4.1.2.1. K-means

K-means [26] is an iterative process which forms the clusters by finding centroids such that the sum of the squares of distance from centroids to data is minimized. For a data set X with n number of samples, it divides them into k number of disjoint clusters, each described by the centroid value. In the first initialization step, it assigns initial centroids, most commonly by selecting k number of samples randomly from the data set X . In a second step it assigns each sample to its nearest centroids and forms k clusters. In the third and final step, it creates new centroids by calculating the mean value of the samples assigned to each previous centroid and calculates the difference between the old and new centroids. It keeps on iterating the process until the difference is less than a threshold value. K-means algorithm finds the local optimum solution which depends on the initial positions of centroids it has chosen. Hence usually it is run multiple times with different random initializations and chooses the best result from multiple runs. The number of clusters k have to be specified in advance and is a major drawback in the approach as the distribution of data may be unknown.

4.1.2.2. GMM

Gaussian Mixture Models (GMM) [27] is a probabilistic clustering technique. GMM is a probabilistic model which assumes that all the observation points are generated from a mixture of finite number of Gaussian distributions and the parameters of Gaussian distribution are not known. It can also be taken as the generalization of k-means clustering which also incorporate the information about the covariance structure of the data. K-means is an example of hard clustering whereas GMM is an example of soft clustering where assignments of underlying data points to particular Gaussian distribution is done in terms of probability.

The source of generation of different data points is unknown which makes it harder to learn GMM for the data observations as one does not know which points belong to which source. This problem can be solved using Expectation Maximization algorithm which is well known statistical algorithm and can solve this problem by iteratively optimizing the observation data set over different distributions.

It works in two steps. In the first step, it assumes Gaussian distributions centered around random points in the data set and then it computes the probability for each point being generated by every component of the model. It then maximizes the likelihood of data for those parameters. It keeps repeating it until the likelihood reaches the local optimum.

4.1.3. Anomaly Detection

Anomalies are a special type of events which can be described as anything unusual and which does not confine to the normal behaviour. For example, a traffic accident can lead to a congestion at unusual time which is an example of anomaly, or malfunctioning of any

electronic device can result into excessive amount of power dissipation which needs to be detected and reported as soon as possible. Similarly, a temperature of 50°C reported by a sensor during winters is an anomaly which might be due to either a faulty sensor or a fire.

There are many methods from machine learning domain which are available for anomaly detection which can be used for monitoring applications. But they are not suitable for large-scale real-time applications. In contrast to these methods, event processing methods based on *Event Driven Architecture* (EDA) which are optimized for real-time analysis can be used. It involves analysing the data stream on the run using rule-based inference. For monitoring applications, it involves setting up rules using threshold values in order to make decisions. For example, if we want to monitor power measurements, a rule can be set up as if the current measurement is greater than the threshold power then generate the event. For a single house, there might be many appliances and each appliance may have a different behaviour. The normal working range of every appliance is specific and it is almost impossible to use the human knowledge to set optimized threshold values. Also, the behaviour of devices is expected to change with respect to time. The use of microwave in daytime is not an anomaly but if we detect the use of microwave at night, it is an example of anomaly. Hence, the threshold values should evolve with the time. This was just one example of monitoring; another application may involve the monitoring of sensor readings such as temperature or humidity in order to detect malfunctioning of a sensor or any unusual event such as fire.

We propose to explore the devices historical data in order to model their normal behaviour and find the threshold values automatically. We propose to have threshold values with respect to the different contexts such as morning or evening, weekday or weekend, summers or winters etc. Threshold values will be adaptive as the behaviour of devices evolves and the rules for CEP can be updated on the run in order to provide a generic solution.

4.1.4. Complex Event Processing

A number of COSMOS components are using the μ CEP engine developed by the *Internet of Everything* (IoE) Lab of ATOS Research & Innovation division. Due to the fact that this component is served as an enabler for the real-time analysis of streaming data, it is already described in deliverable D4.1.2 - Information and Data Lifecycle Management [28]. The following picture represents the architectural view of this component.

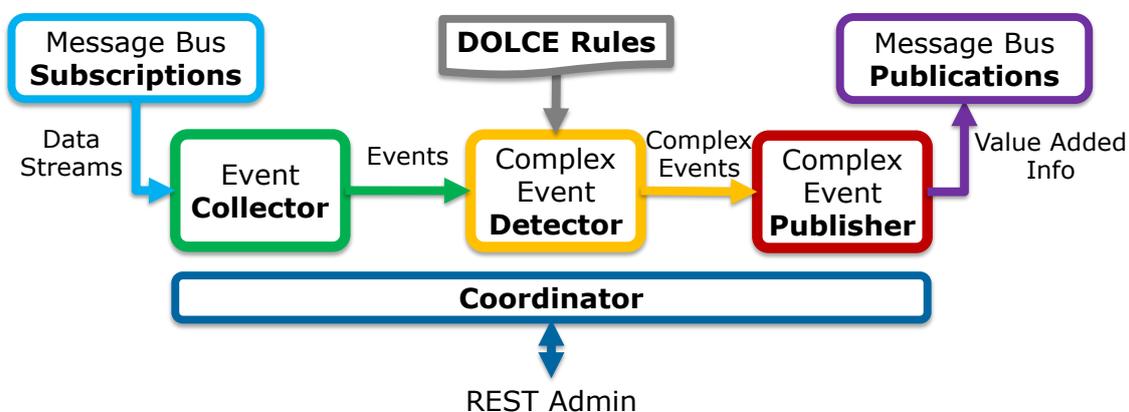


Figure 12: μ CEP interfaces

4.1.5. Probabilistic Rules for Event Processing

In CEP, rules are defined in an absolute manner, i.e. either they are true or false which means even if one of the conditions becomes false, it excludes completely the possibility of complex event happening which might not be accurate for many real-world problems. Even if one of the conditions becomes false, complex event can still happen although with reduced probability. For example, consider a complex event which is dependent on three events: Event A, Event B and Event C and a rule is defined as:

If (Event A happens) AND (Event B happens) AND (Event C happens) → Generate a Complex Event

In a general CEP, if any of the Event A, B or C is missing, it will make the whole condition 'false' and hence it will not generate any complex event. Instead, we proposed to attach weights with the conditions so even if one condition becomes 'false', still Complex event is generated but with less probability depending on the weight attached to the condition.

In this regard, we explored probabilistic graphical models using Bayesian networks to map different events to a complex event. Bayesian Networks is an example of directed acyclic graphical model that represents a set of random variables and their conditional dependencies via a *Directed Acyclic Graph* (DAG). Different events detected by CEP such as Event A, B and C in the above example are the binary random events with two possible values as 'yes' or 'no'. And their conditional dependence on a complex event can be calculated using Bayesian chain rule.

4.2 Functional Overview

This component provides the following functionalities:

- 1) It enables to correlate the data from different data sources in real time to infer complex event;
- 2) It monitors the individual data stream and detect anomalies or any unusual event in near real-time;
- 3) It exploits the historical data to automatically calculate the parameters for CEP Rule so the administrators do not require domain knowledge.

4.3 Connection with other components

Event Detection FC is using VE-level stream analysis capability provided by COSMOS. As it can be seen from the Figure 13 below, raw data is generated by IoT Services and Event Detection block correlates the data using DOLCE rules provided by application developer in real-time to infer a complex event. Also, at the same time raw data is published on the message bus and stored in the cloud storage. Inference/Prediction FC access historical data and run machine learning algorithms to estimate optimized parameters for CEP rules which will be updated automatically through interface provided between Inference and Prediction FC with Event Detection FC.

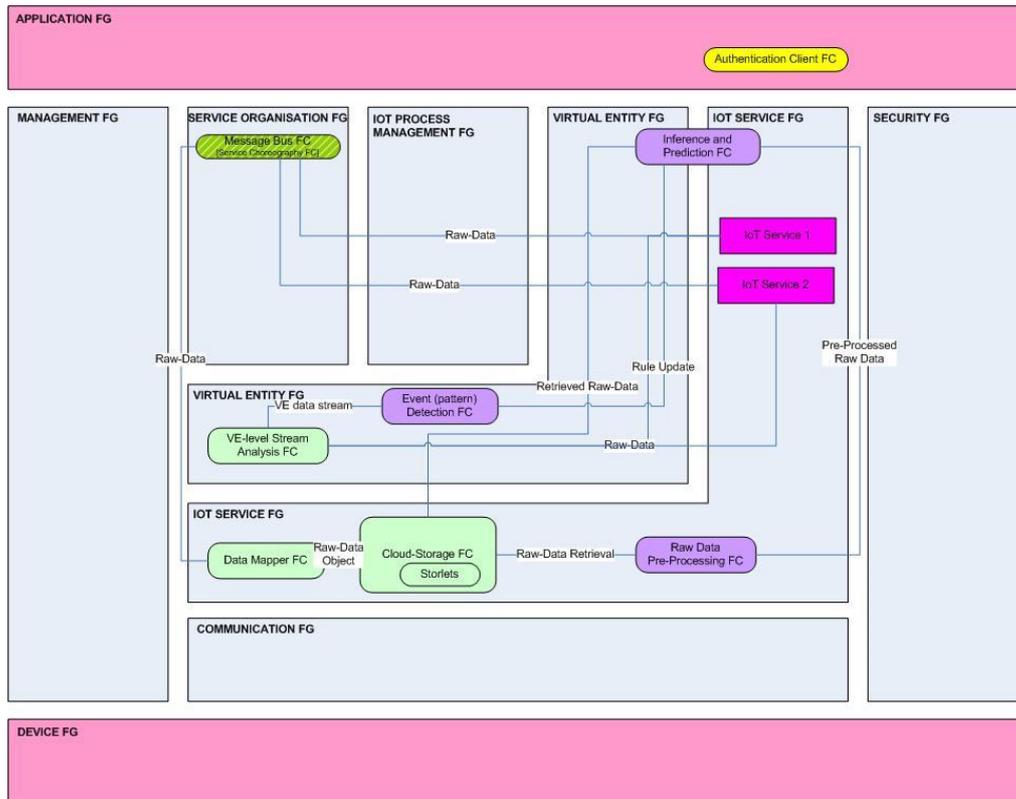


Figure 13: Connection of Event Detection FC

4.4 Interfaces

VE - μ CEP

The μ CEP engine that is being used in COSMOS utilizes the DOLCE domain language for the definition of rules. A *DOLCE Rule* file must have at least two clauses: *Events* and *Complex Events*. The former represents the data streams received by the engine, while the later represents the data that will be outputted in case a rule is triggered –i.e., when a *Detect clause* is evaluated to *True*.

The injection of data streams into the μ CEP is done using the Message Bus, and the same applies for outputting the complex events as *Value Added Info*. The *Event Detector* module will be subscribed to various *Topics*, and the *Complex Event Publisher* will be publishing to certain *Topics*.

Application Developer - μ CEP

Apart from the data streams that are willing to be analysed, there is an additional entrance point to the engine, the *DOLCE Rules* file. Application Developers are able to modify it in two different ways:

1. Editing a **.dolce* file on their own, using their preferred Text Editor. This approach provides full control over the definition of the rules
2. Using a specific Web Graphical User Interface (GUI) developed by COSMOS project, what facilitates the creation of a *DOLCE Rules* file in an easier manner, following a guided wizard

In either case, the file has to be deployed into the μ CEP running instance, what is possible to be done using a specific API function over the *REST Administration* interface. In this sense,

Event Detection-Inference/Prediction

An interface will be developed which connects the Event detection component with the Inference/Prediction component for updating of threshold values for the CEP rules automatically using ML methods.

4.5 Use-Case Scenarios

4.5.1. Scenario 1: Detecting Traffic State (Madrid)

In the city of Madrid, thousands of heterogeneous traffic sensors have been deployed on different locations across the city. These sensors provide real-time information about the traffic flow in the city such as average traffic speed, average traffic intensity, type of traffic or type of road etc.

CEP has the potential to provide a distributed solution for analyzing, correlating and inferring high-level knowledge from this large amount of data in near real-time. The core of CEP is a rule-based engine which requires rules for extracting complex patterns. These rules are based on different threshold values. For example, traffic speed can be analyzed using a simple rule as "if current speed is less than a threshold speed; generate slow speed event". The setting of these rules requires system administrators to have prior knowledge about the system which is not always available and poses a weakness. Secondly, every road segment has a different response. There might be a road segment with speed restrictions as compared to a free flow road segment so the threshold values will be different for both road segments. In this way, at city level there will be hundreds of road segments and it is almost impossible for the administrators to understand the behavior of each individual road segment.

In current systems based on CEP, rules set by system administrators are static and there exist no means to update them automatically. Static set-up of rules can effectively degrade the performance of applications as IoT mainly consists of dynamic environments where the context of application is always changing. For example in the event of bad weather or rain, traffic will move slowly and the rules which are set up for normal conditions can generate a false alarm of congestion; threshold values would have to be different in such conditions. In addition, the response of the road is also changing with respect to time. A road might have a different behavior in morning rush hours as compared to quite night hours.

In short, the following drawbacks have been identified in current technologies based on CEP:

- 1) Threshold values have to be set manually and there is no automatic way to find the optimal threshold values;
- 2) Threshold values set are static and once set, CEP system is unable to update it during run-time;
- 3) Current solutions are not context-aware and adaptive by nature.

4.5.1.1. Adaptive Clustering

Figure 14 shows the flowchart of adaptive clustering. It acquires historical data and extracts the data for specific time period from it. In real-time scenarios, the definition of context is changing with time. In morning hours, traffic intensity is usually higher with low average speed and hence the definition of bad traffic is different from the night hours when traffic is lighter. The response of traffic is different at different time periods and hence we proposed to have different threshold values for different time periods. We slice the historical data along the time line in terms of morning, afternoon, evening and night traffic hours. More detail is given in the next section.

After extracting the data, we apply feature scaling in order to optimize the clustering algorithm. Feature scaling is a method that brings all the features on the same scale so they contribute equally to the clustering algorithm. The range of values of different features of traffic data is on different scales. For example, the value of traffic speed ranges from 0 to 120 km/h whereas the range of traffic intensity is in hundreds. If one of the features has considerably wider range as compared to other features, the optimization function for clustering will be governed by that particular feature and will impact the boundaries.

We have implemented k-means clustering algorithm which is an iterative algorithm which forms the clusters by finding centroids such that the sum of the squares of distance from centroids to data is minimized. For a data set X with n number of samples, it divides them into k number of disjoint clusters, each described by the centroid value. In the first initialization step, it assigns initial centroids, most commonly by selecting k number of samples randomly from the data set X . In a second step it assigns each sample to its nearest centroids and forms k clusters. In the third and final step, it creates new centroids by calculating the mean value of the samples assigned to each previous centroid and calculates the difference between the old and new centroids. It keeps on iterating the process until the algorithm converges. We applied k-means clustering with the $k=2$, hence it resulted into two clusters. The midpoint between the two centroids divides the data into different events. We propose to use these midpoints as threshold values for CEP rules as it defines the boundary between different events.

In general, threshold values for CEP rules are static and are a major drawback when deployed in real-world dynamic environments. Statistical properties of the underlying data may change over time resulting in inaccurate threshold values. Therefore, we propose to keep track the changing in data distribution by accessing the quality of cluster as new data arrives. And as the quality of clusters deteriorate, k-means model is retrained and new threshold values are found using latest data.

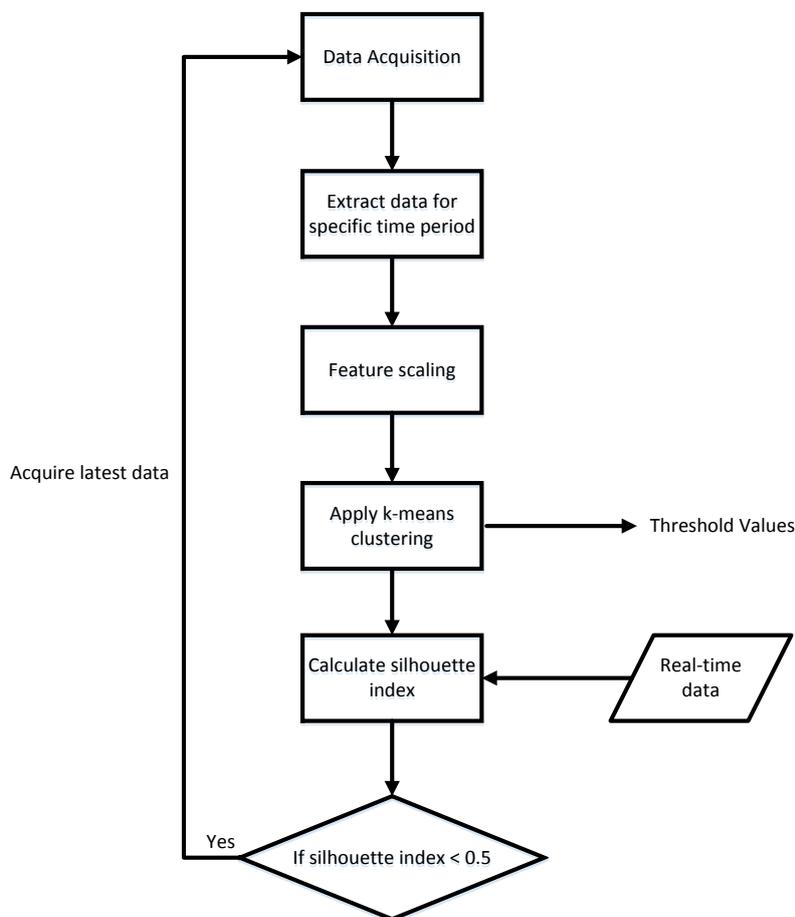


Figure 14: Flow chart for adaptive clustering

4.5.2. Scenario 2: Anomaly Detection (Taipei)

Taipei Test Scenario consist of fifty (50) volunteer households where III's In-Snergy suite (incl. home gateway, smart plugs, and smart strips) were installed during Year 3 of COSMOS which enables user to monitor energy usage and control/set up/use appliances in a smarter way. It can also provide the safety notifications by WIFI on iOS/Android app.

Different appliances are connected to smart gateway with the help of smart plugs. Smart plugs monitor real-time electricity consumption data which is provided as a web service from III. COSMOS access this data in near real-time using credentials provided by III and published on the specific topic in apache Kafka. Real-time data is being monitored with the help of μ CEP in COSMOS to detect anomalies. Historical data is analysed in apache spark using different statistical and ML methods in order to calculate threshold values for μ CEP rules.

Anomalies for a specific appliance are highly dependent on time as well. For example, switching on the TV in an evening is a normal behaviour whereas switching it on during midnight can be an anomaly. Analytics on historical data enables to model the normal behaviour of specific devices and understand the users in a better way.

4.5.3. Scenario 3: Detecting an Event from Twitter Data (Madrid)

In order to fully utilize the potential of IoT, we also analyze the twitter data to extract unusual events using our event detection component. Real-time data from twitter for the city of Madrid is extracted using Node-RED API for twitter. We stored the total number of tweets per unit time for different areas near measuring points in the Madrid city. The intuition behind our approach is that in case of any unusual event such as football match, concert or any incident, the total number of tweets increases as people tend to tweet about those events. We learn normal pattern of total number of tweets from historical data and as we detect a pattern change indicating high number of tweets using our component, we generate a complex event.

4.6 Conclusion

In this component, we have proposed a context-aware method to analyse and extract complex events from data streams in near real-time using CEP and ML. Our propose method is adaptive and is able to cope with dynamic environments as opposed to current state of the art methods. We have demonstrated the feasibility and usage of our proposed architecture with the help of different real-world use-case scenarios such as *Intelligent Transportation System* (ITS) and Smart energy management.

5 Situational Awareness Functional Component

5.1 Introduction

As it was stated in the updated version of State of the Art document [1], it is not an easy task to provide a closed definition of what *Situation Awareness (SA)* is. In spite of that, we are focusing on the definition that defines the entire SA process as “*the perception of the elements in the environment within a volume of time and space, comprehension of their meaning and the projection of their status in the near future*” [29]. Following this approach, worth going into details of each step:

- **Level 1 - Perception of the elements in the environment:** The first step in achieving SA involves perceiving the status, attributes, and dynamics of relevant elements in the environment. This will be commonly referenced as “acquiring context”, “receiving data streams”, “exploiting information from data sources”, and the like;
- **Level 2 - Comprehension of the current situation:** Based upon knowledge of Level 1 elements, particularly when put together to form patterns with other elements, a holistic picture of the environment will be formed, including a comprehension of the significance of information and events;
- **Level 3 - SA Projection of future status:** It is the ability to project the future actions of the elements in the environment, at least in the near term, that forms the third and highest level of Situation Awareness. This is achieved through knowledge of the status and dynamics of the elements and a comprehension of the situation (both Level 1 and Level 2 SA).

In order to complement the above, the following summarizes the categorization of *Context Information* that best represents the Use Cases that are being developed in our work:

- **System Context:** Corresponds, for a given application, to the context information of the system (software and hardware) on which it runs as well as contextual information of the used communication system (for example, the wireless network type).
- **User Context:** It is any information that can characterize a user. This may be the age, location, medical history, biometric information, emotions, etc. It may also be user activities, social relationships, family, friends, colleagues, ...
- **Environment Context:** Represents information describing the physical environment that is not covered by the system and user contexts. Particularly the context information from external sources such as temperature, climate, lighting, ...
- **Temporal Context:** Defines all information related to time, especially: hour, day, month, year...

5.2 Functional Overview

The large volume of data which is made available in an IoT environment does not necessarily mean that applications can take effective decisions directly from the data or even make correct interpretations about the data. For example, a single vehicle reporting a low speed is not always an indication of a traffic jam. On the other hand, a large number of vehicles reporting slow speeds on a particular highway section can be interpreted as an indication of such a traffic condition.

COSMOS intends to support the transition from raw data to value added information by providing mechanisms which facilitate SA at two distinct levels: the *Virtual Entity (VE) centric SA* and the *platform level SA*. Both of them will be described in the following sections.

Despite their differences both levels are targeting the extraction of value added information from raw data in order to place VEs into the relevant context when required. Also targeted is the information sharing among VEs.

The sharing of information among VEs is supported through the use of semantic descriptors for the resulting SA information. This allows VEs to query for SA information either from the platform or from other VEs.

5.2.1. Data Sources

While analysing the prerequisites for the generation of SA information we have identified four distinct data sources which VEs should be able to access.

Figure 15 below depicts the data sources as well as the flows involved in the SA information generation and sharing.

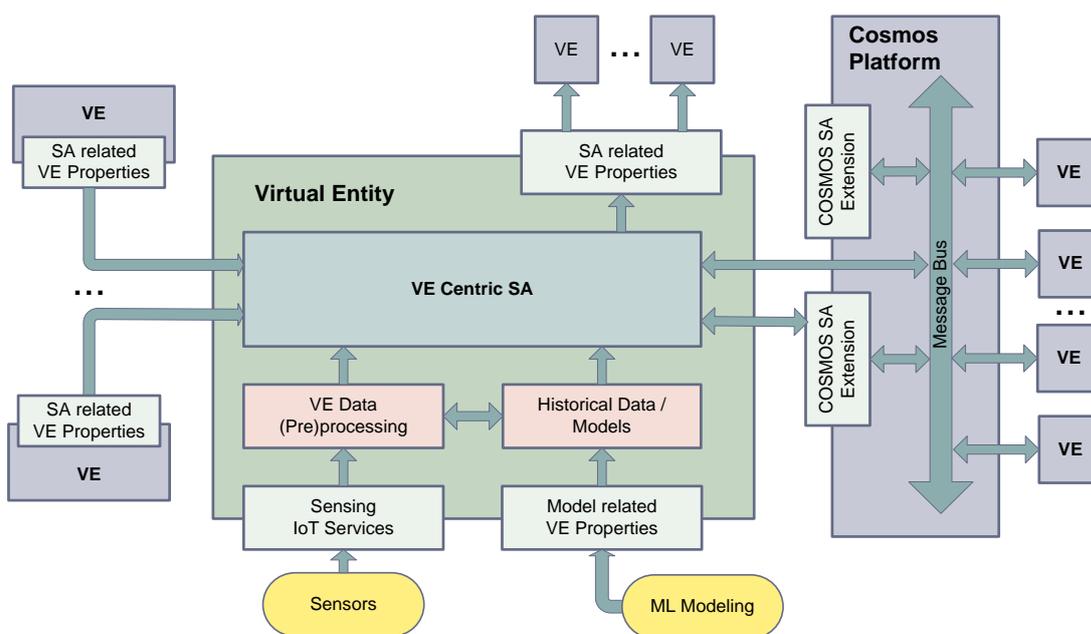


Figure 15: Situational Awareness data exchanges

Internal sources – raw data

The most basic and direct way VEs are able to obtain data for the generation SA information is by using direct sensing through the sensors which are associated to each VE. In other words, the code implementing the VE centric SA is accessing its own sensors, and applies to the cases where for the SA assessment there is no need for further data processing.

Internal sources – processed data

In this case the SA context is generated by mechanisms which could combine data from multiple sensors of the VE, an internal data processing mechanism, and even VE level historical data.

A simple example would be the computation of the expected battery lifetime for a device used in telemetry, for instance a remote weather station which could read the battery voltage and drawn current, the measured charge up to the moment, the device and the ambient temperature and also use historical data about the battery charging /discharging cycles.

External sources – COSMOS platform SA extensions

This is mainly the case where the VEs are not able to extract directly information about their environment through their own sensing capabilities but can use publicly accessible services. These might be any service, provided that the VE is able to access it and has a compatible interface. Apart from external data sources, COSMOS itself will support the development of platform extensions. These are able to process raw data submitted by VEs and to generate reusable and shared SA information which any VE can use.

Such an extension could provide traffic related information, as presented in the example described at the beginning of the section, based on the raw location data submitted by vehicle associated VEs.

External sources – other VEs

This is the case where VEs require SA information about conditions which are beyond their own sensing range and which can be provided by other VEs.

This is another sharing mechanism COSMOS will provide and facilitates VE to VE direct communication without the data being relayed through the platform.

For example, in a logistics scenario, lorry VEs could be able to find other lorry VEs operating on the same route and ask about the current restrictions in place or about queues at border crossing points. The way these data sources are built and used is described in the following sections.

5.2.2. Data Flow and Processing Model

The key concept of COSMOS is the VE which adds a new level of abstraction and functionality on top of IoT services. As mentioned before, the SA approach in COSMOS relies not only on the VEs' capability to sense the operating environment but also on its social characteristics and the ability to share information either in raw or processed state.

The ability to produce and consume data using such sharing mechanisms has to be considered at design time and COSMOS facilitates this task by providing the necessary infrastructure and design patterns. The following paragraphs describe this from the perspective of the relevant COSMOS actors.

VE developers

The VE developer is the one which builds VEs around existing or new IoT services and adds new functionality with the support of COSMOS. While using only internal sources, SA computation requires no or limited interaction with the platform or other VEs. Still, as mentioned above, determining a VE's SA context can be achieved by also using information from other VEs or from the platform SA extensions. This is only possible if the conditions listed below are met.

The first condition is that VEs share completely or partially their own raw data or even their own SA information. This data can be made accessible directly, by the VEs using their own interfaces, or indirectly, by publishing data to a COSMOS platform instance. The task of the VE developer is to implement these COSMOS compatible interfaces based on the provided design patterns and to properly describe the VE's capabilities when publishing it to the VE registry.

The second condition is that VE developers design and develop VEs considering the VEs' ability to access information which is beyond their sensing capabilities by relying on external data sources. This means that once a VE instance is deployed, a binding between such external sources and its own state variables has to be created based on a set of requirements, as later described.

Both conditions are easy to be met since COSMOS supports these requirements by providing an uniform way of accessing data (by using interfaces and data formats compatible between VEs) and by providing the means to search and address this data through the use of the semantic descriptions published and stored into the VE registry.

Application developers

Application developers are using VEs and the platform in order to create VE enabled applications. VEs are deployed either by VE developers or by application developers depending on the scenario. With regards to the SA assessment based on external data sources, at the moment of the deployment, the mapping of the requirements for the bindings mentioned above have to be configured so that the VE is able to retrieve the relevant data sources.

These requirements specify the semantic type (e.g. a VE needs outside temperature information from a specific location) as well the data type (e.g. which kind of schema is required). This is needed in order to provide semantic data compatibility between the producer and the consumer, but also to build the VE registry query parameters for finding the appropriate data sources.

Thanks to COSMOS, the task of the developer is limited to the proper description of these requirements and to the use of the VE registry query mechanism.

Platform owners

Each COSMOS platform instance is operated by a platform owner (e.g. a municipality, a telecom company, an internet service provider, etc.). A platform owner can attract VE developers and application developers by providing extensions of the COSMOS instance they own. These extensions provide additional functionality such as publicly available SA data sources.

Platform extension developers

The above mentioned platform extensions are domain specific implementations which can be plugged into a COSMOS instance in order to provide additional functionality. This extension mechanism is required since the COSMOS platform core functionality is meant to be domain independent.

A COSMOS platform extension is nothing else than an application which is able to consume data published by VEs to the platform or data from other data sources, process it so that it adds value to that data and then share the results by publishing them back to the platform. These extensions can run as stand-alone components, independently of the platform but can also use, if required, platform level functionality such as the event processing or the ML models and prediction mechanisms exposed by the platform.

An SA extension follows the same pattern and is meant to provide value added SA information to VEs by processing high volumes of raw data published by other VEs.

SA in COSMOS is available, as mentioned before at two distinct levels: VE centric and platform level respectively. The main difference between the two lies in the data processing model.

At VE centric level the SA assessment is performed mainly based on the VE's internal data sources. Also, the data processing takes place internally, as part of the VE implementation. The VE developer bears the responsibility for the design, development and testing of this functionality as well as for the proper annotation of the results (to support further reuse based on information sharing) once the VE is published into the VE registry.

Platform level SA involves the processing of large quantities of mainly raw data submitted by VEs. It uses techniques which might require more computing power and are not suitable for VE centric deployment. Some of these processing components are also provided by the COSMOS platform such as those for complex event processing, prediction using machine learning models, advanced storage and processing using servlets.

In contrast with the VE centric approach, the platform level SA does not require any VE developer effort, except for the construction of the queries used to retrieve the information exposed by the platform extensions. This is because in this case, VEs are using SA information which is made publicly available by the platform. SA extensions are built for the platform owner so that it can extend the core functionality of the platform it operates. Depending on the processing requirements, these extensions can be deployed together with the platform or as stand-alone applications. It is the task of the platform extension developer to properly design, develop, test and annotate these services which facilitate information reuse.

The backbone for the platform level information sharing is a message bus which follows a publish/subscribe pattern. Choosing such a design pattern was natural given the requirements mentioned for information sharing. VEs act as publishers (when submitting their own data) as well as subscribers (when receiving reusable information such as SA assessments). The same applies for different components of the platform as well as for the platform extensions. A description of chosen message bus solution is found in section 5.3.3.

5.2.3. Information sharing and retrieval

As mentioned above, information sharing is a key concept in COSMOS. COSMOS provides the components which facilitate data processing and storage at both VE as well as platform level. Nevertheless, the ability to produce and store raw data and value added information does not suffice when it comes to information sharing, since such scenarios involve different data sources and sinks employed by different owners.

Such data needs to be retrievable and properly addressable. To support these requirements, COSMOS employs the use of semantic descriptions of the VEs and of the platform level data sources and sinks (implemented as message bus topics).

Based on the light-weight, domain independent COSMOS ontology, and a dedicated API and tools, VE or platform extension developers are able to describe each data source and sink. This description seeks not only interface compatibility with regards to the data types but also from the semantic type so that the matching between the producer and consumer is guaranteed. As a result, a VE will not only be able to request a data source returning a Double data type but a Double representing the temperature in °C from specific room of a office building.

The ability to describe the data from both the data type as well as the semantic type perspective, combined with the querying mechanism which includes such constraints, extends the ability to share and reuse data among VEs and among VEs and the platform.

These mechanisms are described in detail in the Revised Architecture Deliverable D2.3.2 [30].

5.3 Connection with other components

Apart from the semantic technologies used for the description of the COSMOS entities, such as the ones listed above, the platform makes use other technologies to extend its functionality and facilitate the use of IoT services.

In order to achieve the process of real-time situational awareness, the combined used of this technologies should comply with the following facts:

- Be very efficient to deal with huge amounts of events;
- Predict the occurrence of 'interesting' situations;
- Be tolerant to various types of noise;
- Deal with dynamically changing situations.

5.3.1. Complex Event Processing Engine

A CEP engine is designed as a component to be fed by an amount of data sources, whose ultimate goal is to generate value added information in the form of complex events, that are the output generated after processing many small, independent incoming input events, which can be understood as a given collection of parameters at a certain temporal point.

In relation to the three levels of Situation Awareness, μ CEP engine can provide up to Level 1 and Level 2 SA working on its own. To do so, the μ CEP instance features one or more Event Collectors modules to acquire information from heterogeneous Data Sources, as it was presented in section 5.2.1 Internal Source providing raw data (Level 1 SA). Together with a specific set of DOLCE Rules, this Context Information (system, user, environment or temporal) is processed at the time of arriving to the engine, where historical data is maintained by means of Time Sliding Windows and Tuple Windows. When a certain rule is triggered, relevant Complex Events are generated by aggregating input data streams along with post-processed information, becoming an Internal Source delivering processed data (Level 2 SA). Finally, the Complex Event Publisher module transfers the value added information to the VE or component which will take advantage of it. Furthermore, when the μ CEP is used in combination with other techniques, such as ML, it is possible to get Level 3 SA, as will be explained later on this document.

5.3.2. Machine Learning

ML represents –as said in previous chapter already- any algorithm which learns from the data. Methods based on ML have the potential to extract high-level knowledge from raw IoT data and to contribute for novel applications. We have described several ML methods in section 2. In this part, we will explore how the high-level knowledge can be used for Situation Awareness at both VE centric level and platform level.

At platform level, different type of data sources -as discussed previously- are generating huge amounts of data which is heterogeneous, mobile and ever changing along time. Data from different data sources form complex patterns and analysing, inferring and correlating these complex patterns in order to access the current situation is by no means a trivial task. Consider a simple example of Situation Awareness at platform level as shown in Figure 16, where an instance of different events happening at the same time is shown. The task of SA component is to combine this information from the external data sources to extract high-level knowledge that represents the current situation. Pattern recognition techniques such as classification and clustering can be explored in order to detect particular situation.

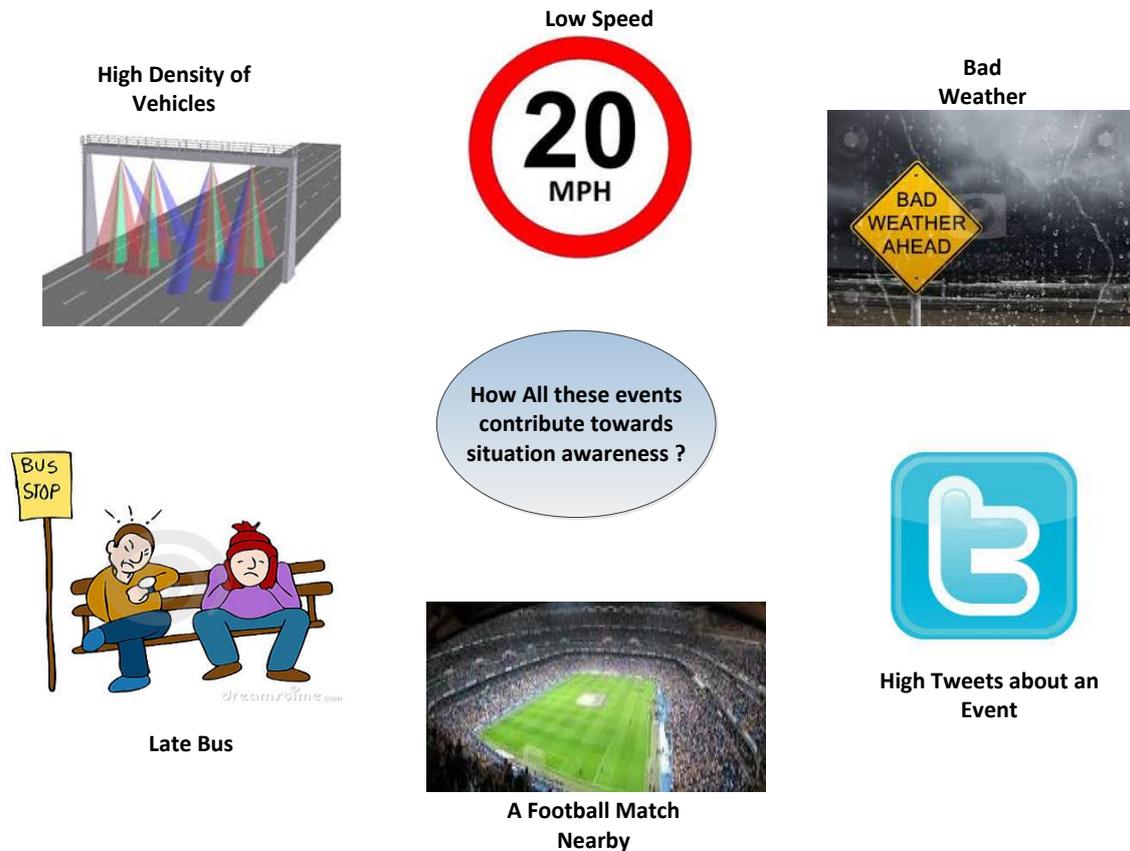


Figure 16: An Example of Platform Centric Situation Awareness

At VE-centric level, the amount of data generated will be less, but at the same time the resource constrained nature of VEs limits the use of complex ML tasks. One proposed solution is to exploit the historical data of VEs in COSMOS platform in order to train ML models and then extract high-level knowledge. The work done in Year 1, in which occupancy state of a user is inferred from his electricity usage pattern is an example of VE-centric SA.

The discussion until now falls into perception and comprehension of current situation but the vision of SA in COSMOS is beyond it. We intend to explore the projection of Situation Awareness for potential future events. The value of predicting future events is of immense importance as it enables the administrators of the platform to take pro-active approach to potential situations. We intend to explore prediction techniques from ML domain in order to provide projection of SA. One possible approach for the projection of events is shown in Figure 17 below.

5.3.2.1. Projection of SA using ML

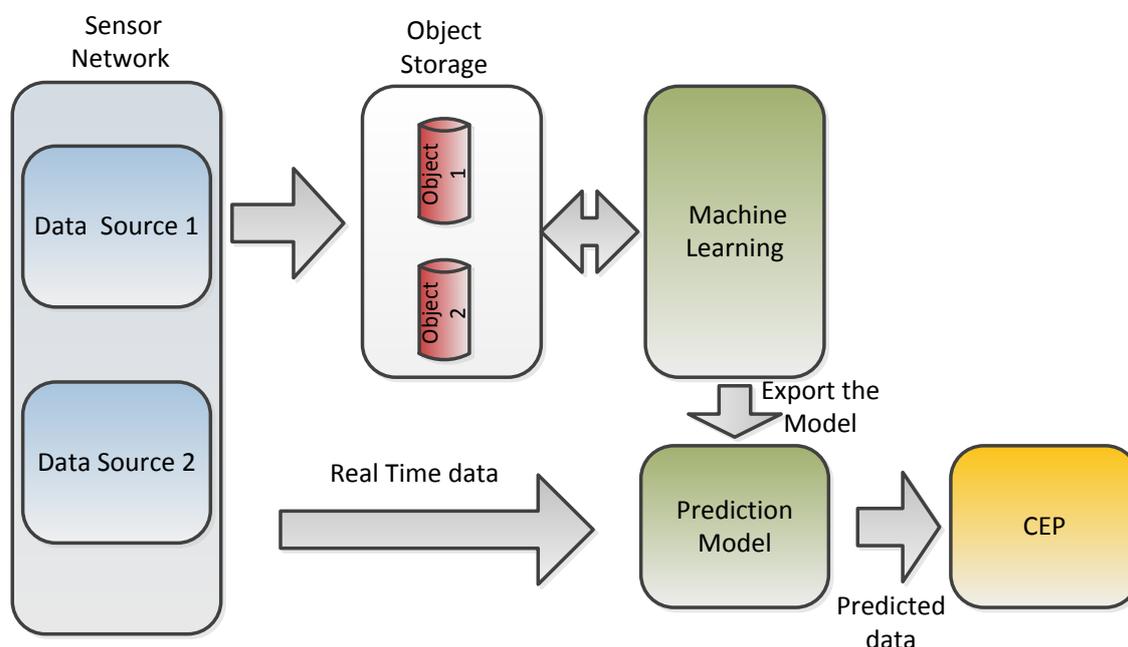


Figure 17: Prediction of Events for SA Projection

We propose to explore ML techniques for the projection of SA for potential future events. It enables the system administrators to adopt a pro-active approach and take measures before the actual event happens.

A high-level architecture for the proposed approach is shown in the Figure 17. The intuition behind our proposed approach is that if the input data streams to CEP are predicted data (in future), then CEP can be used for predicting the potential future event. The model for predicting the data can be trained and validated offline using historical data sets. The accuracy and complexity of prediction models will be the main research challenges in this work. Prediction Models trained on historical data may become inaccurate with the passage of time due to phenomena of concept drift. The statistical properties of the underlying data may change with the time. Prediction Models should be able to cope with the changes and adapt accordingly.

5.3.2.2. Automatic generation of Rules for CEP using ML

As discussed earlier, manual setting of rules and patterns limits the use of CEP only for expert's domain and poses a weak point. In this regards, we intend to apply predictive analytics principles to improve decision-making and performance of existing CEP solutions by exploring adaptive and automatic solutions. Both fields have many overlapping concepts but yet there exist only few solutions where PA concepts have been applied to CEP. One possible research area to explore is automatic generation of rules for CEP in order to cope with above-mentioned limitations. CEP language uses number of operations for describing pattern of events. In [31], authors applied the ML techniques for generating automatically the following five most commonly used operations:

1. To determine the relevant time frame to consider, i.e. the window size;
2. To identify the relevant event types and the relevant attributes;

3. To identify the predicates that select, among the event types identified above, only those notifications that are really relevant, i.e., determine the predicates for the selection operator;
4. To determine if and how relevant events are ordered within the time window, i.e., the sequences;
5. To identify which event should not appear for the composite event to happen; the negated event notification

The authors applied *Information Gain Ratio* (IGR) principle and implemented each operation in different module. This is just one example where ML can be used in conjunction with CEP. We intend to explore more predictive analysis methods in conjunction with CEP to provide more adaptive, automatic and optimized solutions that can lead to more accurate, faster and consistent performance.

5.3.3. Message Bus

The Message Bus FC entails a key integration point for the cooperation of several entities in a common scenario. Based on a publish/subscribe messaging mechanism, it allows for sharing information from heterogeneous data providers by abstracting the inherent complexity that implies dealing with different protocols and data formats. To do so, each COSMOS component provides the necessary modules to adapt to this horizontal communication bus. Furthermore, the new Message Bus implementation provided by COSMOS during Year 2 improves the way data queues are handled, thus providing new mechanisms for fault tolerance and data persistence.

5.3.4. Forging CEP and ML

CEP and ML have many overlapping concepts but yet belong to two different research fields. CEP acts on real-time data from multiple data sources in order to analyse, correlate and infer a more complex event whereas ML methods are based on exploiting historical data to learn models which can be deployed for prediction. CEP engines require rules or patterns to detect an event from data streams which have to be given manually by the administrators of the system. Based on this, there is an assumption that administrators have the required background knowledge which sometimes is not available or not precise enough. So manual setting of rules and patterns are in a sense weak point of CEP.

We proposed to explore ML methods by exploiting historical data for finding rules for CEP. As we have discussed previously, that VE-centric SA is performed mainly locally as a part of VE implementation. A light version of CEP provides a good solution for combining data from different local sensors of VE to extract high-level knowledge. The resource constrained nature and limited memory of VE are the limiting factors of exploiting historical data and using complex machine learning models. Our proposed solution fixes this problem as part of complex computations for finding optimized parameters for CEP rules can be done in the central platform and the rules for the running CEP at VE level can be updated periodically.

5.4 Interfaces

For both the VE-centric approach as well as for the platform level SA, the interfaces provided for accessing SA information follow the same pattern as those of the “standard” VE properties.

In fact, there is no difference between accessing an IoT-related VE property and an SA related one.

SA information access over both REST as well as message bus endpoints is supported at VE as well as platform level. In addition to the interface type, we need to emphasize the fact that in either case the endpoints are also semantically described as any other endpoint of the VE allowing SA information retrieval and sharing.

5.5 Use-case Scenarios

5.5.1. Scenario 1

Existing Madrid EMT transportation infrastructure provides possibility to build contextual models based on data streams related to various sources, such as the context of a particular VE Bus –position, velocity, schedule, CO₂ emission–, as well as the status of the surrounding environment of the VE –weather condition, scheduled city events, local traffic congestion–. Taking advantage of these data sources, combining them with other external related ones, and even making up relevant historical data, it is possible to foresee certain scenario evolutions like traffic jams or bus delays. Figure 18 highlights the suitability of a SA mechanism in the *Decision Making* process that is being followed.

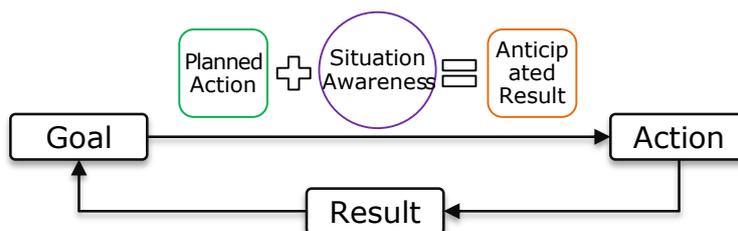


Figure 18: SA in Decision-making process

Attending to this diagram, it is crucial to have a clear identification and description of the goal to achieve. In the urban mobility scenario we are going to monitor the journey of a *person with special needs* on his/her way home, while the correspondent *caregiver* is sent notifications if any abnormal situation may occur. In this sense, the planned action starts taking the right bus and ends up stopping at the correct bus stop, what defines the temporal boundaries of the SA process. Next step consists in selecting the data sources that will feed the system components involved in the use case, such as the location of the person with special needs, the route of those buses arriving nearby user’s home and the traffic state condition in the city, to name only the most relevant ones. By acquiring, processing and maintaining this information during the required amount of time we are able to predict if the problem is going to be solved appropriately. The following Figure 19 depicts this process.

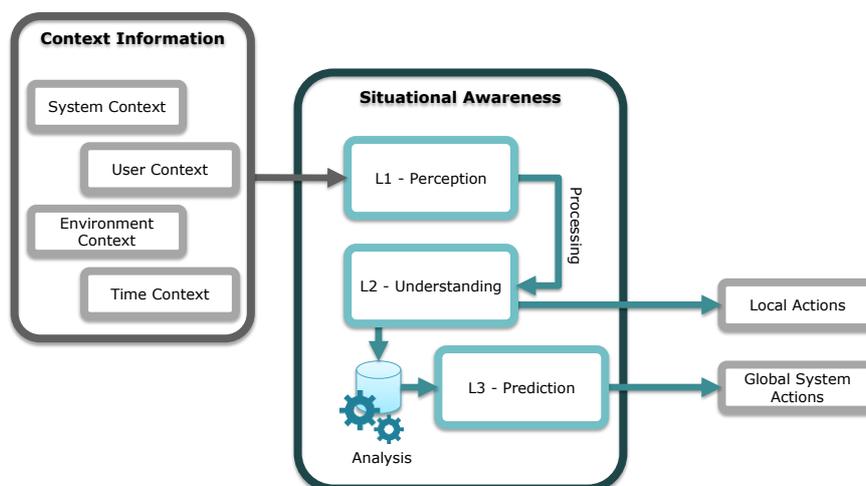


Figure 19: Situation Awareness process

Ingesting context information into the Situation Awareness functional component means that data sources are connected to the message bus publishing to topics, the ones that the μ CEP components are subscribed to. Each time a rule is triggered, the generated information may flow to one or more components depending on what level of SA is needed in each particular case. As an example, in the urban mobility scenario the user's location is continuously compared with the location of the bus so to understand if the user keeps being inside the bus or not. Moreover, an additional source of data may provide valuable information into the equation, such as the WiFi networks the user mobile is detecting and the SSID the bus is broadcasting.

5.5.2. Scenario 2

In this scenario, we demonstrated how we can achieve all the three levels of situational awareness as described in 5.1 for Madrid traffic scenario using other components of the WP.

5.5.2.1. Level 1: Perception

We consider three different types of data streams for analyzing and extracting level 1 form of situational awareness. First data stream consists of traffic parameters including average traffic speed and average traffic intensity as described in section 4.5.1 which gives us the high level knowledge about the traffic state. Second data stream consists of weather data where first level of perception is defined as sunny, cloudy or rainy weather. Third data stream consists of twitter data which can be analysed to detect whether an event is taking place in nearby region or not. A summary of level 1 perception utilizing all the other components in the work package is given below:

Input Data Stream	Perception
Traffic data (speed, intensity, occupancy etc.)	Traffic State (good/bad)
Weather data	Sunny/cloudy/rainy weather
Social media data (Twitter)	Event/no Event

5.5.2.2. Level 2: Comprehension

In the second level for situational awareness we combined all level 1 perceptions in order to provide a more global perspective. For example, if we detect low moving traffic, rainy weather and an event happening in the city, it is an indication of bad traffic or congestion due to rain and an event happening in the region. It provides the holistic picture enabling the traffic administrators to manage traffic in a better way.

5.5.2.3. Level 3: Projection

Finally, we move to the level 3 of situational awareness which is called projection. This is the highest level of situational awareness where projections are made in future using first two levels of situational awareness. We utilise the prediction component as described in Section 2.7 to predict future events. Perceptions and hence comprehension made will be in future which indicates to potential future events.

5.5.2.4. How it all fits together

We combined the functionalities provided by inference/prediction component and event detection component in order to perceive the probability of congestion. In this regard, we also explored probabilistic CEP rules as described in section 4.1.5. CEP works in a hierarchical manner where first layer of CEP detects events like good or bad traffic, rainy or sunny weather or an event from twitter data feed. In the second layer, these events are combined in probabilistic manner using Bayesian networks in order to predict the probability of congestion. Figure 20 shows the probabilistic Bayesian network representation of Madrid scenario. Traffic is time-of-the-day dependent as traffic is busier during morning and evening rush hours as compared to afternoon or night hours. Similarly, weather also affects the traffic as bad weather or rain can slow down the vehicles and increases the probability of bad traffic. Another factor is any event happening nearby in the region such as football. We exploited twitter data feed and applied event detection component functionality to detect an event. Finally, particular day has also effect on the traffic state, as weekdays are busier around office hours as compared to weekend. In order to simply the scenario, we only considered two possible states for traffic i.e. either good or bad traffic state. A bad traffic state increases the probability of congestion but still it is not a certainty. Similarly, if the traffic state is good, it does not exclude the probability that the congestion can happen in future although the probability of having congestion with good traffic state will be low. In this scenario, our goal is to find the probability of congestion given all other events as shown in the Figure.

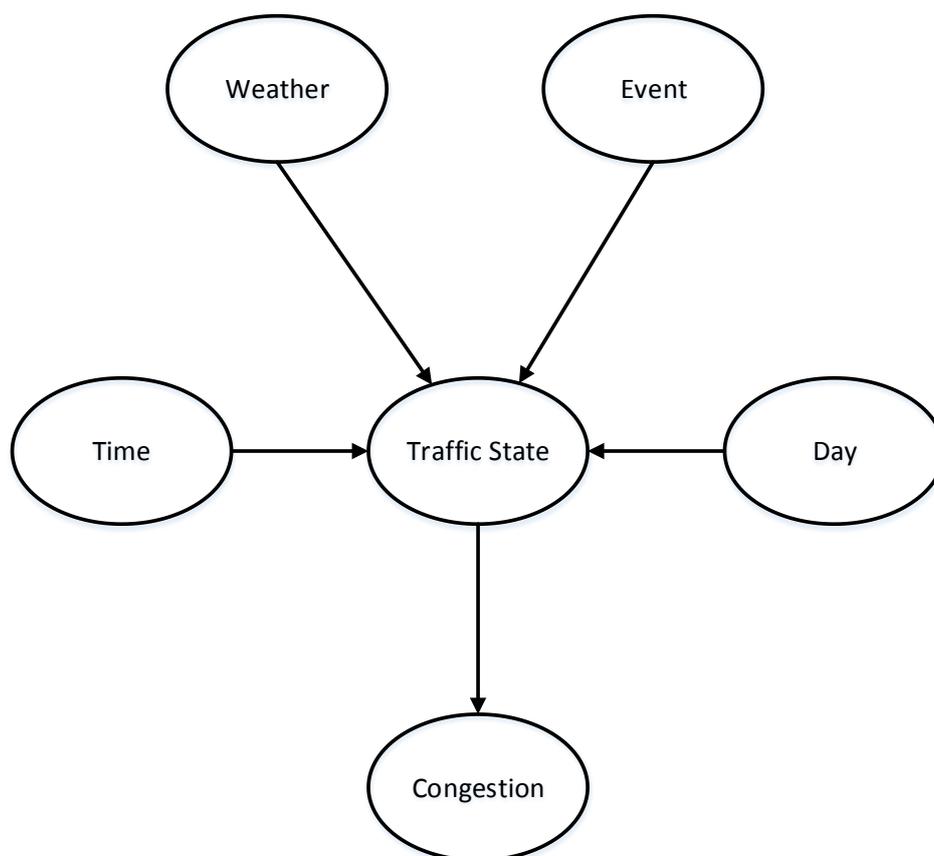


Figure 20: Bayesian Network Representation for Madrid Scenario

5.6 Conclusion

The nature of the IoT based applications (highly dynamic environment involving many times the mobility of user or even resources, the large number of data sources, their volatility, the number of involved actors etc.) brings new challenges for application developers but also opens up new opportunities. One of these opportunities is the ability to integrate situation awareness into the business logic of the application and access data about things beyond the coverage range of the things addressed by the application.

COSMOS is aware of this opportunity and therefore supports the SA through a series of data processing and sharing mechanisms which are exposed both as VE level as well at platform level. Depending on the type of the data sources and expected SA information, SA processing tasks can be executed as part of the VE functionality or in a centralized and high capacity manner when run by the platform.

These tasks can be light-weight processing jobs (especially in the case of the VEs with limited processing resources), can include advanced processing techniques relying on complex event processing or ML-based prediction models. In both cases, the output of the processing tasks (which is the SA information) is semantically described using the endpoint description model provided by the COSMOS ontology.

6 Experience Sharing Functional Component

6.1 Introduction

In the course of the COSMOS project, Task 6.3 has endeavoured to promote communication between remote VEs. This communication is engaged, by properly structuring information in an Application specific content, but always following a very strict pattern. The structure of the following section will focus on work accomplished and the final improvements which are being developed.

During the development of this concept, we structured our efforts around the codification of Knowledge in the form of Experience. This term refers to the CBR concepts of Problem-Solution as described in Deliverables of WP5 and demonstrated in Figure 21. Therefore we focused our information packaging efforts into conforming to a standard way of “objectification”, in the sense of data objects to be transmitted.

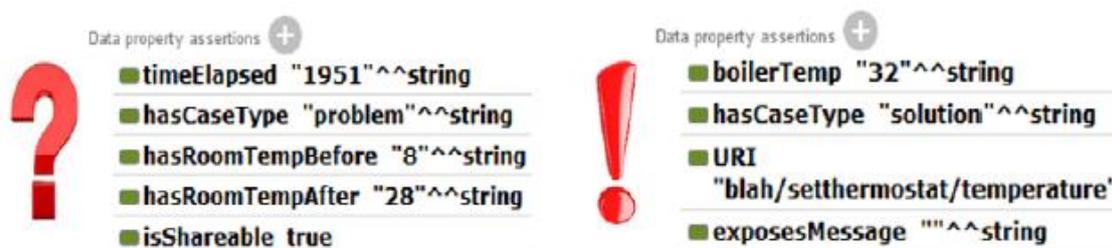


Figure 21: Problem Solution structure of Knowledge

Afterwards, the objects themselves make use of existing protocols for remote unicasting, as the Architectural deployment of the component is performed as a RESTful service. Prior to the actual communication of Experience, as a request-response pair, we focused on structuring the relationships between VEs as shown in Figure 22. This new structure incorporates new advancements enabled by the close cooperation between the service-masking oriented parts of the Privelets FC. Additionally, it also aims to make use of the Social Paradigm proposed from the Social Monitoring FC, given the successful completion of dependant work in the Project.

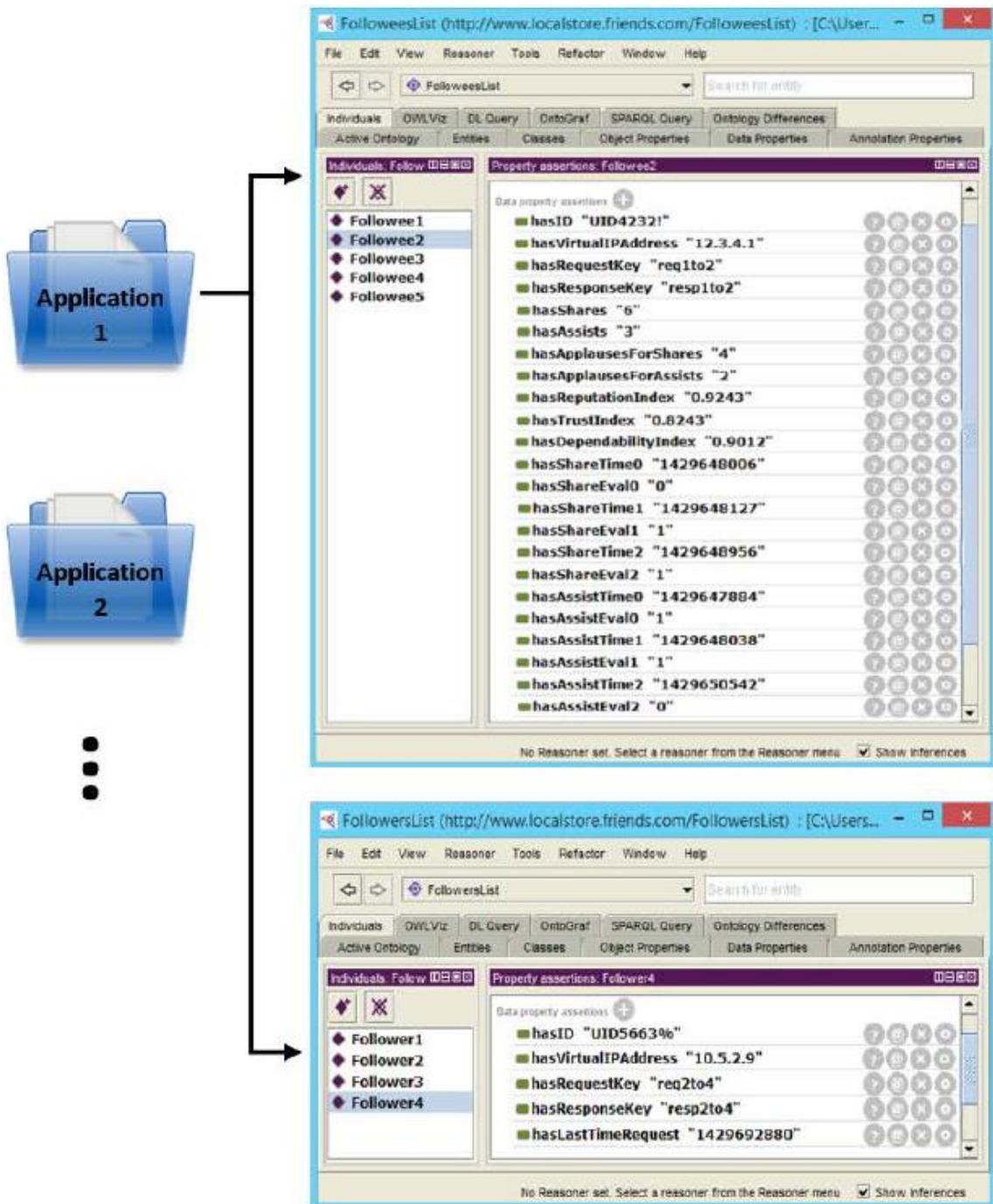


Figure 22: New Followee-Follower Structure

In the final iteration of the Experience Sharing FC we will be describing in-depth the Architectural Design as well as specific implementation details of the FC along with connecting points with other FCs.

6.2 Functional Overview

6.2.1. Mechanism Improvements and Restating Achievements

During Year 1 the Experience Sharing component was implemented as a VE level service offered through a REST interface. Programmatically this was achieved using implementation of an embedded server and the accompanying servlet. The characteristic that promoted this approach is that since the Experience Sharing process does not involve modifying any local system resources (i.e. Case Base content) and simply accesses the persistent local store for the process of information retrieval, VEs can use the multithreaded parallel executing nature of servlets. Therefore in a decentralized environment, in which VEs demand swift access to information, a target VE can accommodate a multitude of requests. Also we made a strong connection of the Experience Sharing FC with the Planner FC, as during this time, we concluded methods of intercallability depending on the context of execution flow.

In Year 2 in continuation of this approach we focused in reducing the load overbearing, such a solution causes, in a two-fold way. At first we successfully configured the embedded Jetty Server instance in the VE code which supports deployed Services and in a second approach we performed load detection in the incoming part of the one on one communication pair of the Experience Sharing process. Using native Java capabilities introduced since version 1.5, we succeeded in detecting accurately the volume of requests, and thus eventually handle them in more efficient ways. Experimentally, we concluded that a hard limit on concurrently served requests, is around 500. Further advancements during Year 2, included the integration of Privelets code, in a way which runs in advance of any request handling, therefor saving time in the weeding out of unauthorized communications. Such controls included the knowledge of the partner VEs identity (through Virtual IP identification) initiating the request, as well as the authentication, via pre-assigned communication keys.

Focus on Year 3, is the improvement of the mechanism in order to provide a more in line approach to the modular aaS (as a Service) nature of the VE code in total, which is being presented as a redesigning principle of the VE sided FCs. The changes performed will be more thoroughly described in Subsection 6.2.2. Along with the changes in the nature of the Service, we will also include the proposed final draft of the FC's architecture in connection to the Planner FC and the Privelets FC, as well as any relevant new mechanisms/FCs influenced through integratory actions during Year 3.

6.2.2. Architectural/Deployment Changes

Driven by the changes in the structure of the VE code the Experience Sharing FC code was divided in incoming and outgoing software subcomponents. As such the code in the incoming request handling end, remains relatively unchanged barring changes in the implementation of the Privelets FC embedded code and any changes deemed necessary by the new VE code Architecture deployed aaS.

The major change was performed in the case of outgoing requests, as now the Planner FC, given a direct access by an Application, will use a Wrapper provided for the calling of the Experience Sharing functionalities. The purpose of the Wrapper is modifying the input object for use with the POST function and querying the Followees of the VE appropriately, in structure.

However, internally the Wrapper changes the way Followees are contacted as it now implements a multithreaded way of initiating a search for a candidate solution to a Case. The blocking nature of the mechanism still remains, as it is still implemented synchronously, but in this case, we achieve in cutting down on total execution times, given a populated Followees

List. This is the main effect of the parallelisation nature of the Runnable implementation in a Java environment.

6.2.3. Usage of ML based VE Similarity Calculations

In Year 3, we proposed to find similar VEs for experience sharing using clustering mechanisms from the machine learning domain. The behaviour and actions of VEs will be more similar to the VEs which have the same characteristics and hence experience sharing will be more optimized for such cases. As an example, consider a 2 bedrooms flat located on the top floor of the building and facing West. The heating plan schedule for such a flat cannot be applied efficiently to another flat which has a single bed room on the ground floor and facing north (no sun). In order to overcome this issue, we proposed to group the VEs with similar characteristics using clustering mechanisms. This is an example of grouping the VEs on the basis of their static properties. Another option is to group the VEs on the basis of the data which they are measuring. In such scenario, the flats with similar electricity and heating consumption pattern will be grouped together. Flats with high consumption of electricity and heating energy data will be in a same group whereas flats with low energy usage will be in a same group. The two possibilities are;

- 1) Similar VEs on the basis of static characteristics
- 2) Similar VEs on the basis of dynamic data

Static characteristics:

There are several options for grouping the VEs on the basis of static characteristics. The size of a flat, number of people in the flat, location such as floor number and whether facing sun or not, and the age group of residents are few of the examples which can be used for grouping the similar VEs together.

Dynamic Data:

Grouping the VEs on the basis of the data they are measuring is a totally different way of grouping. For the same example of flats, in this method the flats will be grouped on the electricity consumption and energy data patterns. Experience sharing component will utilize this information in order to share experience in an optimized manner.

6.3 Communication with other Components

The Experience Sharing Functional Component maintains in this iteration as well, its connection to the Planner FC. In this specific case, given an input from an Application in the form of a Case, the Planner will use the Wrapper to execute an Outgoing instance of the Experience Sharing mechanism. The communication is done through Json objects which properly formed, span the length of the programming flow of the VE Case retrieval process.

In the event of an incoming Experience Sharing interface use (through the forwarding of an outgoing request to the remote partner), there is a connection to the Privelets functionalities authorizing and authenticating the incoming VE's credentials and after acceptance, the Service contacts the Planner for further actions.

Pending changes to the connection with the Social Monitoring will be described in full inside the prototype Deliverable of WP6, given a renewal of the deployment model for Social Evaluations and retrievals.

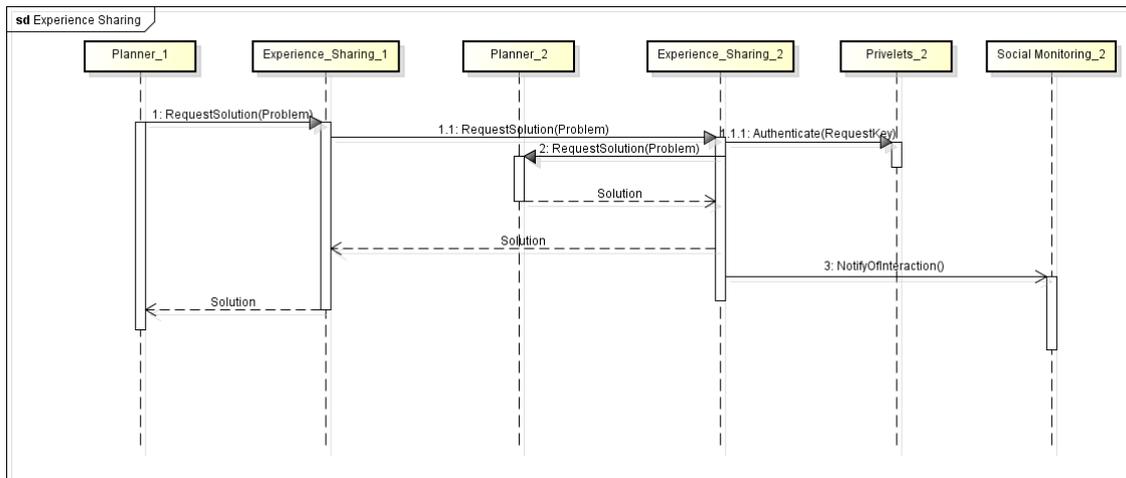


Figure 23: Sequence of actions for Experience Sharing between two VEs.

6.4 API / Interfaces for Experience Sharing

The APIs used by the Experience Sharing FC are called either implicitly or explicitly, as shown in the functionalities that must be performed below:

- Storing Experience:** Through the use of the semantic store API, VEs can modify the contents of their semantic description. Regarding Experiences (Cases), it is very important that they can establish a persistent storage of their personal Case base, which contains all their Experience. This happens by creating new instances of Problems and/or Solutions or by adding connections to already existing ones; VEs are exposed as web services and are semantically described into the semantic store of COSMOS. The semantic description of the VE, contains not only the description of the underlying IoT-Services, of the VE capabilities and constraints but includes also references to any prediction models or experiences which it uses or has created;
- Finding Experience:** The next important step in the implementation of Experience Sharing is the actual access of the semantic storage in order to find and return a certain solution that satisfies a problem. Such forays into the semantic store are made through the use of SPARQL and the API's Query Functions. The API of the semantic store, dedicated to the Experience retrieval, will include specialized methods providing as input parameters the Experience search criteria. A VE could easily create a query string, which would also be fully configurable based on information stored locally in VE variables. Therefore, any access of the VE in the semantic store can be accomplished and the actual Experience Sharing is made easier, due to the fact that Experience instances can be accessed and most importantly assessed by any VE. Jetty server is used to enable VE2VE communication through HTTP, including finding and sharing Experience, by creating a servlet for each exposed IoT-Service that is accessible by a URL;
- Choosing Experience:** Having already described the process by which an Experience can be accessed, it is also important to mention that any retrieved information (including Experience), can also be assessed, ranked, etc., giving VEs the capability to

choose one experience of those offered by other VEs. By defining a variety of social properties such as Trust & Reputation index, we enable VEs to check whether a Solution answers not only their queries but also satisfies some quality criteria. This process is executed by the Planner of the VE based on how many times a VE has shared its cases or how many times its IoT-services have been used;

- **Privelets inclusion:** The modifications during Year 3 established new object structure by adding fields for communication key fields inside the transmitted object. These keys are utilized by the VEs in order to authenticate and secure their transmissions and have been demonstrated during the Year 2 review process. Privelets functions are specifically called to access the Followers list of a VE and check keys, as well as maintain an image on the call frequency from a remote VE, therefore blocking flooding (frequently reappearing) requests.

6.5 Use Cases

6.5.1. Use Case for Year 3 Implementation

The Experience Sharing FC will exclusively focus on the Camden use-case scenario. Specifically we will make use of the more decentralized, social aspects provided by the nature of the UC namely the divergent houses and flats. In this Year's efforts for maximizing the effects of the COSMOS platform's capabilities in relation to the needs of the Camden Housing authority, we will focus on real time data on a greater scale. Provisions in the UC scenario for the eventual release of 15 flats with interactive residents are being made, therefore we anticipate to make full use of capabilities delivered in both Year 2 and Year 3. While the burden of connectivity will not be extravagant, the actual use of multiple communication nodes in the form of interacting VEs will demonstrate the validity of our approach which was tested in extensive simulations and made public in ["Heating schedule management approach through decentralized knowledge diffusion in the context of social internet of things", Panagiotis Bourellos, George Kousiouris, Orfefs Voutyras, Theodora Varvarigou].

Of specific importance in this case is not the actual handling of overbearing network traffic, rather, the speed at which incoming requests can be forwarded and eventually resolved back to the initiator with either a successful or unsuccessful conclusion. Additionally given control over the finer points of configuration we can realize observations on how a greater network of VEs handles cyclic request alerts, authorization difficulties, as well as the assignment and maintenance of the concept of Virtual IPs, introduced by work done on the integration of the FC with the Inter-VE Privelets mechanisms.

Taking major decisions on where the actual code deployment will occur is also a point of consideration as a more decentralized approach of installation and execution is preferable to showcase the portability and lightweight capabilities of the VE code. However ongoing discussions will indicate whether the actual interference with resident everyday lives outweighs the benefits of such a deployment. Other options, like deploying the code in a more Virtually Decentralized environment, are also being taken into consideration. Such a focus will provide other benefits as it proves that management authorities may have full control over the deployed software in terms of easier access and maintenance, while also preserving the concepts of autonomy and automaticity.

Finally by introducing a more Thread-independent version of the FC, we can also measure any tangible benefits and through testing, extract the optimal settings for seamless communication

between VEs. A non-blocking version of the mechanism will be of benefit to any automation system, that desires to have simultaneous access to multiple Self or Remote monitoring Applications or Services.

6.6 Conclusion

Within the context of COSMOS project, Experience has been defined as Case which consists of a problem and its corresponding solution (see COSMOS Deliverable D5.1.3). We have developed Experience Sharing Functional Component which is responsible for knowledge (Case) sharing between friends VEs, while communicating with each other using HTTP RESTful services. The component collaborates with Planner FC (WP5) for the solution retrieval based on the occurred problem, as well as with Privelets FC (WP3) to ensure authentication on top of VE2VE communication. Experience Sharing FC has been applied in Camden smart heating application scenario.

7 Conclusions

As the Work Package name suggests, the aim of our work is to make “things” more reliable and smarter. The reliability of the system is induced by utilizing historical data in order to make prediction models which can be used for decision making in real world dynamic scenarios with incomplete data. Different Machine Learning and statistical methods were explored for inferring high-level knowledge from raw IoT in order to contribute towards more situation aware and autonomous systems.

Over the years, as the technology evolves, the amount of data at our disposal has also increased rapidly. And the availability of such diverse type of data enables us to induce intelligence in the different real world applications. The use case scenarios of Madrid, Taipei and London provided us with large amount of Data which had to be processed, correlated and synthesized in order to extract high level information in near real-time, in order to help in decision making. In this regard, we explored state of the art methods in Year 1 and leveraging on Year 1 results, we proposed new novel methods in Year 2 and demonstrated it using use case scenarios. In Year 3, we improved our components (including for instance Complex Event prediction) and integrated them together in order to contribute towards more situational awareness systems.

Finally, we have explored how experience of Virtual Entities can be used to make things more autonomous and to take decisions in new situations. Additionally, an initial mechanism of experience sharing is introduced. This version of the deliverable goes further in all topics extending the state of the art and features next steps in individual conclusion sections of each Functional Component described in the previous section.

8 References

- [1] C. COSMOS, "COSMOS Project D2.2.2 State of the Art Analysis and Requirements Definition," .
- [2] G. Marfia, M. Rocchetti and A. Amoroso. A new traffic congestion prediction model for advanced traveler information and management systems. *Wireless Communications and Mobile Computing* 13(3), pp. 266-276. 2013.
- [3] T. Salsbury, P. Mhaskar and S. J. Qin. Predictive control methods to improve energy efficiency and reduce demand in buildings. *Comput. Chem. Eng.* 51pp. 77-85. 2013.
- [4] M. Limayem and C. M. Cheung. Understanding information systems continuance: The case of internet-based learning technologies. *Information & Management* 45(4), pp. 227-232. 2008.
- [5] J. Provost. Naive-bayes vs. rule-learning in classification of email. *University of Texas at Austin* 1999.
- [6] X. Wu, V. Kumar, J. R. Quinlan, J. Ghosh, Q. Yang, H. Motoda, G. J. McLachlan, A. Ng, B. Liu and S. Y. Philip. Top 10 algorithms in data mining. *Knowledge and Information Systems* 14(1), pp. 1-37. 2008.
- [7] A. Ben-Hur and J. Weston. "A user's guide to support vector machines," in *Data Mining Techniques for the Life Sciences* Anonymous 2010, .
- [8] D. Tian, X. Zhao and Z. Shi. "Support vector machine with mixture of kernels for image classification," in *Intelligent Information Processing VIA* Anonymous 2012, .
- [9] G. Shakhnarovich, P. Indyk and T. Darrell. *Nearest-Neighbor Methods in Learning and Vision: Theory and Practice* 2006.
- [10] H. Yang, L. Chan and I. King. "Support vector machine regression for volatile stock market prediction," in *Intelligent Data Engineering and Automated Learning—IDEAL 2002* Anonymous 2002, .
- [11] D. C. Sansom, T. Downs and T. K. Saha. Evaluation of support vector machine based forecasting tool in electricity price forecasting for australian national electricity market participants. *J. Electr. Electron. Eng. Aust.* 22(3), pp. 227-234. 2003.
- [12] A. Ding, X. Zhao and L. Jiao. Traffic flow time series prediction based on statistics learning theory. Presented at Intelligent Transportation Systems, 2002. Proceedings. the IEEE 5th International Conference On. 2002, .
- [13] C. Wu, J. Ho and D. Lee. Travel-time prediction with support vector regression. *Intelligent Transportation Systems, IEEE Transactions On* 5(4), pp. 276-281. 2004.

- [14] P. Ni, C. Zhang and Y. Ji. A hybrid method for short-term sensor data forecasting in internet of things. Presented at Fuzzy Systems and Knowledge Discovery (FSKD), 2014 11th International Conference On. 2014, .
- [15] F. Ganz, P. Barnaghi and F. Carrez. Automated semantic knowledge acquisition from sensor data.
- [16] S. Kamijo, Y. Matsushita, K. Ikeuchi and M. Sakauchi. Traffic monitoring and accident detection at intersections. *Intelligent Transportation Systems, IEEE Transactions On 1(2)*, pp. 108-118. 2000.
- [17] Y. Kallberg, U. Oppermann and B. Persson. Classification of the short-chain dehydrogenase/reductase superfamily using hidden markov models. *FEBS Journal 277(10)*, pp. 2375-2386. 2010.
- [18] G. Yu, J. Hu, C. Zhang, L. Zhuang and J. Song. Short-term traffic flow forecasting based on markov chain model. Presented at Intelligent Vehicles Symposium, 2003. Proceedings. IEEE. 2003, .
- [19] J. Dean and S. Ghemawat. MapReduce: Simplified data processing on large clusters. *Commun ACM 51(1)*, pp. 107-113. 2008.
- [20] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker and I. Stoica. Spark: Cluster computing with working sets. Presented at Proceedings of the 2nd USENIX Conference on Hot Topics in Cloud Computing. 2010, .
- [21] S. Rabinovici-Cohen, E. Henis, J. Marberg and K. Nagin. *Storlet Engine: Performing Computations in Cloud Storage* .
- [22] J. Dean and S. Ghemawat. MapReduce: Simplified data processing on large clusters. *Commun ACM 51(1)*, pp. 107-113. 2008.
- [23] A. Zoha, A. Gluhak, M. A. Imran and S. Rajasegarar. Non-intrusive load monitoring approaches for disaggregated energy sensing: A survey. *Sensors 12(12)*, pp. 16838-16866. 2012.
- [24] D. M. Tax and R. P. Duin. *Feature Scaling in Support Vector Data Descriptions* 2000.
- [25] I. Jolliffe. *Principal Component Analysis* 2005.
- [26] J. A. Hartigan and M. A. Wong. Algorithm AS 136: A k-means clustering algorithm. *Applied Statistics* pp. 100-108. 1979.
- [27] N. Shental, A. Bar-Hillel, T. Hertz and D. Weinshall. Computing gaussian mixture models with EM using equivalence constraints. *Advances in Neural Information Processing Systems 16(8)*, pp. 465-472. 2004.
- [28] C. COSMOS, "D4.1.2 - Information and Data Lifecycle Management ," .

- [29] M. R. Endsley, T. C. Farley, W. M. Jones, A. H. Midkiff and R. J. Hansman. *Situation Awareness Information Requirements for Commercial Airline Pilots* 1998.
- [30] C. COSMOS, "COSMOS D 2.3.2: Conceptual Model and Reference Architecture," .
- [31] A. Margarra, G. Cugola and G. Tamburrelli, "Learning from the past: Automated rule generation for complex event processing ," in *DEBS`2014*, Mumbai, India, 2014, .
- [32] C. COSMOS, "COSMOS Project D 6.1.1: Reliable and Smart Network of Things," .
- [33] <http://www.ibm.com/developerworks/library/j-jtp11234/>
- [34] Fehr, Ernst; Simon Gächter (Summer 2000). "Fairness and Retaliation: The Economics of Reciprocity". *Journal of Economic Perspectives* 14 (3): 159–181. doi:10.1257/jep.14.3.159. ISSN 0895-3309.
- [35] <http://jmeter.apache.org/>
- [36] C. COSMOS, "COSMOS Project D7.1.2: Use Cases Scenarios Definition and Design," .
- [37] Páez, Antonio, Fei Long, and Steven Farber. "Moving window approaches for hedonic price estimation: an empirical comparison of modelling techniques." *Urban Studies* 45.8 (2008): 1565-1581.
- [38] <http://kafka.apache.org/documentation.html>