

# An Integrated Architecture for Future Studies in Data Processing for Smart Cities

Cristian Chilipirea<sup>a</sup>, Andreea-Cristina Petre<sup>a</sup>, Loredana-Marsilia Groza<sup>a</sup>,  
Ciprian Dobre<sup>a</sup>, Florin Pop<sup>\*a</sup>

<sup>a</sup>Computer Science Department, Faculty of Automatic Control and Computers, University Politehnica of Bucharest, Romania

---

## Abstract

Data processing for Smart Cities become more challenging, facing with different handling steps: data collection from different heterogeneous sources, processing sometimes in real-time and then delivered to high level services or applications used in Smart Cities. Applications used for intelligent transportation systems, crowd management, water resources management, noise an air pollution management require different steps of data processing. The main subject of this paper is to proposed an architecture for data processing in Smart Cities. The architecture is oriented on the flow of data from the source to the end user. We describe seven steps of data processing: collection of data from heterogeneous sources, data normalization, data brokering, data storage, data analysis, data visualization and decision support systems. We consider two case studies on crowd management in smart cities and on Intelligent Transportation Systems (ITS) and we provide experimental highlights.

*Keywords:* arhitecture; big data; data processing; crowd sensing; crowd dynamics; intelligent transportation systems

---

## 1. Introduction

More and more applications today use, generate and handle very large volumes of data. In particular, this is true for Smart City applications, which attract a rapidly increasing interest from government, companies, citizens, developers, scientists, etc. They cover a large spectrum of needs in public safety, water and energy management, smart buildings, government and agency administration, social programs, transportation, health, education. They are fed with huge amounts of input data, in various formats, from a continuously increasing number of sources (sensors, governmental, regional, and municipal sources, cit-

---

\*Corresponding author, Tel.: +40-723-243-958; Fax: +40-318-145-309; *Email address:* florin.pop@cs.pub.ro.

10 izens, public open data sources, etc.), are describe by complex workflow and in  
11 many cases impose real-time processing capabilities, useful in decision taking.

12 The large volume of data coming from a variety of sources and in various  
13 formats, with different storage, transformation, delivery or archiving require-  
14 ments, complicates the task of context data management. At the same time,  
15 fast responses are needed for real-time applications. Despite the potential im-  
16 provements of the Smart City infrastructure, the number of concurrent appli-  
17 cations that need quick data access will remain very high. With the emergence  
18 of the recent cloud infrastructures, achieving highly scalable data management  
19 in such contexts is a critical challenge, as the overall application performance is  
20 highly dependent on the properties of the data management service.

21 Extracting valuable information from raw data is especially difficult con-  
22 sidering the velocity of growing data from year to year and the fact that 80%  
23 of data is unstructured. In addition, data sources are heterogeneous (various  
24 sensors, users with different profiles, etc.) and are located in different situa-  
25 tions or contexts. This is why the Smart City infrastructure runs reliably and  
26 permanently to provide the context as a “public utility” to different services.  
27 Context-aware applications exploit the context to adapt accordingly the timing,  
28 quality and functionality of their services. The value of these applications and  
29 their supporting infrastructure lies in the fact that end-users always operate in  
30 a context: their role, intentions, locations and working environment constantly  
31 change.

32 As the scale, complexity and dynamism of distributed systems is dramati-  
33 cally growing, their configuration and data management have started to become  
34 a limiting factor of their development. This is particularly true in the case of  
35 Cloud is used for data storage and also for data processing, where the task  
36 of managing hundreds or thousands of nodes while delivering highly reliable  
37 services entails an intrinsic complexity. Furthermore, Cloud computing intro-  
38 duces another challenge that impacts on the resource management decisions.  
39 In these contexts, self-management mechanisms have to take into account the  
40 cost-effectiveness of the adopted decisions.

41 Considering all of these aspects, the main subject of this paper is to pro-  
42 posed an architecture for Big Data processing in Smart Cities. The architecture  
43 is oriented on the flow of data from the source to the end user. We describe  
44 seven steps of data processing: collection of data from heterogeneous sources,  
45 data normalization, data brokering, data storage, data analysis, data visual-  
46 ization and decision support systems. We describe two case studies on crowds  
47 management in smart cities and on Intelligent Transportation Systems (ITS).

48 The paper are structured as follows. Section 2 presents the related work on  
49 crowd data smart cities and on ITS. The proposed architecture is presented in  
50 Section 3. Two use cases are described in Section 4. Then, the experiments  
51 obtained for these use cases are presented in Section 5. The paper ends with  
52 conclusions and future work presented in Section 6.

## 53 2. Related Work

54 Smart Cities [1] represent an important goal that can dramatically improve  
55 the life of citizens. There is a lot of research that aims to get us closer and closer  
56 to this goal. The idea of a smart city is in accordance to other movements in  
57 research such as Internet of Things [2] and Big Data [3]. New York times actually  
58 declared this period the "Age of Big Data" [4].

59 In order to enable Smart Cities technologies such as Internet of Things and  
60 even Wireless Sensor Networks [5] or Crowd Sensing [6] are the catalysts that  
61 provide all the data about our cities. The need for sensing in Smart Cities is  
62 explored in [7]. This data needs to be processed often using Big Data techniques  
63 in order to extract the information required to make decisions about the cities.  
64 This information and the decisions are then used in order to inform the citizens  
65 to take certain actions or to activate actuators that enable automatic processes.  
66 A good example where actuators can improve Smart Cities is given by the  
67 management of green spaces [8].

68 Probably the most important issues that are addressed in order to build  
69 Smart Cities are the ones of Crowd Dynamics [9]. In order to understand Crowd  
70 Dynamics we need data on the movements of as many people as possible. These  
71 movements need to be recorded for both pedestrians [10] and for vehicles [11].  
72 The problem of tracking is not solved in any of the two scenarios. This is  
73 surprising, considering the problem of tracking a particular individual is usually  
74 solved by the use of GPS [12]. GPS systems do not yet offer the desired accuracy,  
75 even if many research projects make significant improvements in this direction  
76 [13]. These systems also do not work indoors and require the cooperation of the  
77 individual being tracked in order to generate a position estimate.

78 It is important to treat both indoor and outdoor cases when considering  
79 human mobility. This is because modern vital facilities, such as hospitals, which  
80 are part of the backbone of many cities consist of large areas with multiple  
81 buildings. An example of how dynamics inside these facilities are considered  
82 in order to improve them is available in the work of Ruiz et al [14]. Similarly  
83 Universities campuses, another type of large facilities at the core of cities, are  
84 analyzed [15], [16] in order to better understand the dynamics inside them.

85 Crowd tracking experiments are taking place in a wider variety of places  
86 like mass events [17] or festivals [18]. They are also used in order to measure  
87 queues using only WiFi signals [19]. This queue can represent waiting time at a  
88 counter, which directly affects customer experience or the movement through  
89 security lines at an airport [20].

90 Crowd Sensing can be used in order to extract all types of data for smart  
91 cities. A powerful example is given by the authors of [21] where students are  
92 asked to take pictures of plants around the campus. The pictures are then  
93 analyzed by scientists in order to better understand the status of flora. Projects  
94 like this could potentially be used at the scale of a city in order to measure a  
95 large variety of features. It is not always necessary for people to be active in their  
96 participation of data gathering. Passive systems require only their presence in  
97 the monitored location, which can even be obtained in an opportunistic manner.

98 Whenever any citizens carrying the scanner walks or drives on a specific street  
99 data about the street can be gathered. In this way maps can be enhanced  
100 with features [22] such as roundabouts or pot-holes. Diverse uses include even  
101 earthquake detection [23] and soon maybe even the detection of effects produced  
102 by this large natural disasters.

103 There are many projects and platforms targetted directly at crowd sensing:  
104 Medusa [24], Matador [25], Mosden [26] and mCrowd [27]. And these platforms  
105 already implement important features for Smart Cities such as crowd sources  
106 new reporting [28] but they do not yet combine the data sets or offer a method  
107 to analyze the data in order to extract information that is hidden inside it. This  
108 type of information represent answers to questions tht we don't yet have and  
109 they can currently only be obtained by using Big Data techniques.

110 The data gathered from all these systems is usually analyzed by experts or  
111 scientist manually. This is the case for [14], where categorization of individuals  
112 into different groups such as patients or staff is done by using rules built by  
113 experts. More information can be extracted from these data sets if they are  
114 combined and Big Data systems are used to process them.

115 Real-time processing is used to designate a category where the job outcome  
116 is needed as fast as possible, and usually the task itself is not something that  
117 will take a long time to process. These systems can be categorized as hard  
118 or soft. A Hard real-time system is an OS for a nuclear plant or a plane.  
119 Tasks must be scheduled and completed fast because otherwise a catastrophe  
120 could happen. These systems are usually governed by hard deadlines and the  
121 scheduler must make sure they get met. Soft real-time systems are the ones  
122 like hotel booking or video streaming sites. The answers must be delivered fast  
123 to the customers, but a delayed frame now and then can not lead to disastrous  
124 results. One article that explores this type of hard real-time scheduling is [29].  
125 In the paper the authors try to improve the scheduling capabilities of a system by  
126 also adding security checks to the incoming jobs. The added module can detect  
127 threats brought by snooping, alteration of spoofing and can be easily added to  
128 any real-time scheduler. Their security module name SAREC (security-aware  
129 real-time heuristic strategy for clusters) integrates with the popular Earliest  
130 Deadline First algorithm to create a security aware scheduler named SAEDF.  
131 Although the matter of securing the interactions between the users and the  
132 cluster infrastructure is important, in our case a large portion of these measures  
133 could be implemented in an intermediate cluster proxy module if needed, with  
134 little overhead to the job itself. By using a proxy to mediate all user-cluster  
135 interactions we can alleviate a large number of security risks. If a user has a  
136 malicious intent and manages to submit a job that poses a security risk, the  
137 fact that all jobs are run in virtual machines on the cluster infrastructure will  
138 limit the damages to only the users task.

139 Another example of real-time processing and scheduling [30]. The authors  
140 talk about the problem of soft real-time scheduling in rendering 3D images  
141 inside the Google Earth software. The Google Earth software allows one to  
142 navigate anywhere in the world and has multiple viewing modes from virtual  
143 3D renderings to satellite imagery. A frame is a static 2D representation that

144 is rendered on the screen at a given time. To ensure a smooth navigation  
145 experience, at least 60 of these frames must be rendered on the users screen in a  
146 second. When a scheduling deadline is not met, the previous frame is redisplayed  
147 causing the application to "stutter". In order to alleviate the problems the  
148 authors have devised a new algorithm that better estimates rendering time on  
149 multiple devices, in order to improve scheduler accuracy. We are in particular  
150 interested in their scheduling model and discovered they also abstracted some of  
151 the events into "single-active sporadic tasks" (triggered by a specific rendering  
152 phase) and "soft real-time aperiodic tasks (triggered by receiving new imagery  
153 through network)". We will use similar terms to define the submit patterns and  
154 properties of different types of jobs.

155 Talking about the arrival patterns of the jobs, the authors from [31] build a  
156 common approach to schedule static and dynamic tasks, in a system that also  
157 has to deal with hard real time deadlines. They divide their tasks in 3 categories,  
158 based on their arrival pattern and number of instances they require for running.  
159 Aperiodic tasks need only one instance to run, and can enter the system at any  
160 time. Both periodic and sporadic tasks require multiple instances to run, but  
161 while the former come at a specific interval of time, the latter can be submitted  
162 like the aperiodic tasks, at any time, but no sooner than a specified interval.  
163 The authors have extended a previous static time-based scheduling algorithm  
164 into a dynamic version that constantly changes the expected start and end time  
165 of jobs while still keeping the end time in the necessary deadline. They have  
166 thus provided two versions of their scheduler, one that accepts aperiodic tasks  
167 without affecting the existing task instances deadline, and another, with the  
168 same properties, that accepts periodic and sporadic tasks. Before accepting any  
169 task, a formal schedulability test is run, to see if the system can handle the  
170 tasks deadline. If not, it is rejected. The scheduled tasks are considered to be  
171 preemptive, and a list of static tasks that are known beforehand is expected  
172 to be provided at system startup. To account for dependencies between tasks,  
173 start and end times are parameterized instead of being given a fixed value.

174 We also investigated solutions related to intelligent transport systems, since  
175 this is the type of workload we are going to test our scheduler on. The [32]  
176 project tries to act as a hub for such endeavors in order to help each of the  
177 individual current ITS system grow and communicate through a common point  
178 of contact. These systems are increasingly important since optimizing traffic can  
179 also reduce  $CO_2$  emissions along with the benefits brought to all the inhabitants  
180 of a city. Current implemented solutions are mostly proprietary and involve  
181 infrastructure changes. There are a number of existing solutions that try to  
182 estimate the state of the traffic, ranging from sensors in the road, to GPS  
183 systems on cars, to cameras that interpret images. Indifferent of the chosen  
184 solution, all of these systems will generate a large amount of data that needs to  
185 be interpreted city-wide. Although our solution uses a small part of this data,  
186 it could grow and adapt to provide the necessary analysis needed to drive an  
187 intelligent city of today.

188 **3. Data Flow based Architecture**

189 We propose an architecture that is detailed on steps that represent the flow  
 190 of data from the source to the end user. We look to the data as it is raw material  
 191 that goes into the factory to be processed and becomes a valuable good for the  
 192 users.

193 We created 7 steps architecture to accomplish our goal to make from data a  
 194 value (see Figure 1): fist we need to aggregate the data sources, then we need  
 195 to perform data normalization, but before doing that we need to anonymize the  
 196 data that comes from personal devices. The next step is to create a context for  
 197 the gathered data and after that we should send it to be stored and processed  
 198 in a parallel and distributed way. The result of the processing will provide the  
 199 starting point for data analysis that will generate the patterns and discover the  
 200 insights we need. In the end all the findings need to be visualized in an advanced  
 201 style to empower the decision makers.

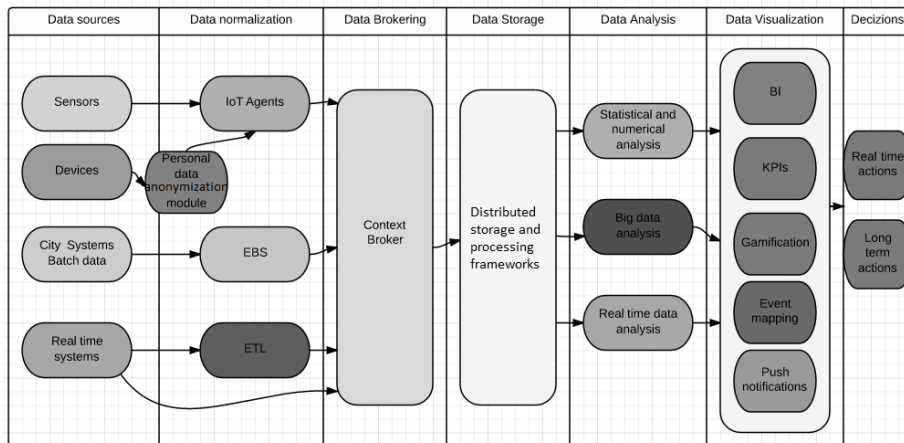


Figure 1: Proposed Architecture.

202 Data comes from different sources and we need to collect it from everywhere:  
 203 smart sensors, personal devices, batch data from city systems, real time systems  
 204 in order to be able to extract as much knowledge as we can from it. When we  
 205 combine different data sources that are related to the same context then we can  
 206 get more insights from it and this empowers the decisions makers to minimize  
 207 the risks.

208 We have plenty of law constraints and each country has its own regulations  
 209 regarding data privacy so we need to address this important issue because when  
 210 a user decided to contribute to a system he needs to be sure that he remains  
 211 anonymous and other users cannot trace him back starting from the pieces of  
 212 data that he provided. We need to make sure that he cannot be identified form a  
 213 group of users that contribute to the system with their data by using techniques  
 214 and algorithms for data anonymization. Data collection needs to deal with

215 data privacy, and in order to face that, we need a personal data anonymization  
216 module before sending to IoT agents for normalization. We need to handle this  
217 as close to the source as possible to avoid any data leaks, that could identify an  
218 individual as data provider.

219 Real time systems data needs to be normalized by using ETL - Extract,  
220 Transform and Load - that represent 3 database functions that are put together  
221 in one tool in order to get the data out from a database and introduce it into  
222 another one.

223 City systems batch data needs to be normalized by EBS - Electronic Batchload  
224 Service –which is an “Online Computer Library Center” service that permits to  
225 the batchload participants to send data to it over the Internet.

226 In the next step the data reaches the context broker which takes individual  
227 pieces of data and puts them into a relevant context. A context broker is  
228 represented by a service that needs to gather context data from different types,  
229 sources and velocity, then it needs to create the conditions, integrate the data,  
230 create the rules to be able to provide prepared context data. A certain piece of  
231 data is meaningful only in appropriate relation to other pieces of data, which  
232 happens only in a given context.

233 After we created the links between the data and each context then we are  
234 ready to send the data to the distributed storage and processing frameworks. In  
235 order to create a powerful platform for big data processing we need to combine  
236 the patterns extracted from the batch data processing with the speed from  
237 real time data processing. The main idea is to bring together real time data  
238 processing and batch processing when dealing with large data sets.

239 We proposed two well known frameworks to be used in this step of the data  
240 flow: Hadoop, which is focused on batch data processing and Storm, which  
241 handles real time data processing. Hadoop, architected around batch processing,  
242 is the most popular open-source software framework for distributed storage and  
243 distributed processing of big data on clusters. The main advantages are given  
244 by the fact that was designed to be fault-tolerant, it is highly scalable and  
245 cost effective. The main components of Hadoop are the storage called Hadoop  
246 Distributed File System (HDFS), and the processing part called MapReduce.  
247 Real time data processing involves a continuous input, process and output of  
248 data. Data must be processed in a small time period (or near real time) so we  
249 recommend to use Storm because it is a free, open source, distributed real-time  
250 system that can compute over a million tuples per second on each node. Other  
251 big advantages are given by the scalability, fault-tolerance and the fact that it  
252 guarantees that the data will be processed. Also it is simple to set up, utilize,  
253 and integrate with other queueing and database technologies, which is a big plus  
254 especially when you need to create a big data platform for smart cities.

255 Now that we processed the data, we can perform big data analytics, statisti-  
256 cal and numerical analysis or real time data analysis to transform the data into  
257 valuable assets. In the end, we need advanced data visualizations that could  
258 offers us the opportunity to take the right decisions at the right time , or to  
259 take long term actions based on historical data.

260 Real time data processing and analytics allows decision makers the oppor-

261 tunity to take immediate action when it is required and batch data processing  
262 makes the results to be more accurate due to the patterns that are discovered  
263 and then applied in real time to get more relevant data.

264 We need to combine the data from multiple sources to be able to predict  
265 future events in order to respond in an efficient way that can make the difference.  
266 It is important to engage the users, in our case the citizens, and to do that we  
267 need to empower them and motivate them. A way of smart user engagement and  
268 advanced visualization is gamification. For example users of a mobile application  
269 can share status about how much do they recycle different things, or how much  
270 CO<sub>2</sub> they produced based on how many km they were driving in a day, and  
271 enter in competitions with others on social media.

## 272 4. Architecture Use Cases

### 273 4.1. *Intelligent Transportation Systems*

274 Large cities present many problems with their systems, but only transporta-  
275 tion system entertains the dynamics of this environment. Currently, it cannot  
276 cope anymore with the enormous number of cars driven on its streets using  
277 classical traffic systems. Any problems like congestion, accidents, high fuel con-  
278 sumption, pollution, etc. which affect us daily in a city can have as root causes  
279 the bad usage of current infrastructure or not enough streets for current traffic  
280 flow.

281 Trying to solve the second cause is a temporary solution due to the con-  
282 tinuously increasing cars' number, because any new street added will move the  
283 problems from one street to another or in short time if the city area which  
284 presents these issues will bring more traffic to it once new streets will be added.  
285 Also, adding new streets for vehicular traffic is very hard to be done in cities'  
286 centers. The majority of the problems encountered by citizens of a metropolis  
287 in traffic are especially determined by the bad traffic planning or by the lack  
288 of traffic control systems. Before to try to extend current streets infrastructure  
289 of a city, it has to be checked if the largest part of its roads are used at their  
290 maximum traffic flow as much as possible and then to try another expensive  
291 solutions.

292 The congestion type presents its particularities for each city not having a  
293 predefined pattern, but using different information for city infrastructure layout,  
294 drivers' behavior and habits etc. together with proper traffic prediction systems,  
295 it can be realized a generic traffic system which diminishes the congestion. The  
296 majority of navigators guides the user during its ride based on the decisions  
297 taken locally not having the global perspective about traffic from the areas that  
298 are crossed by the vehicle or about the other participants' decisions. All routing  
299 applications from cars see the same traffic events in the above scenario and  
300 all from the same area choose locally the same optimum alternative road. For  
301 instance, if there is a congestion event in same area of a city, all cars being  
302 around see it and compute their routes to the same alternative roads moving  
303 indirectly the congestion to the routes' roads.



304 Intelligent Traffic Control Systems (ITCS) are designed to reduce the global  
305 level of congestion in a city, by sensing the city environment through streets in-  
306 frastructure and the traffic participants counting on Inter-Vehicle-Communication  
307 (IVC) and Road-to-Vehicle Communication (RVC) in order to exchange data  
308 about roads congestion level, cars' speed, cars' direction, etc. ITCS is able to  
309 collect complete information about traffic in a large city, because it exchanges  
310 data with various entities from road infrastructure and traffic participants. In  
311 order to perform traffic optimization, this system is realized to support three  
312 phases (traffic monitoring and data collecting; driving conditions perspective  
313 built using the traffic model; traffic controlling by offering to participants feed-  
314 back/new routes and controlling the WTLs to improve the traffic flow.

315 The ITCS' key entities involved directly in the traffic are cars which are the  
316 only one component from the traffic flow having their own will according to the  
317 driver's decisions. Their main target is to collect data from the environment  
318 and then to exchange it with the other traffic participants and infrastructure.  
319 They can collect data using the sensors from incorporated navigators or using  
320 smartphones (e.g. GPS, accelerometer, barometer, etc.). Offering data to the  
321 system, they obtain feedback about traffic in real time and also new routes  
322 suggestions. The local decision capability is used only when they do not have  
323 possibility to communicate to the other system entities in order to receive a new  
324 route in exchange of the provided data, instead the global routing decisions are  
325 coordinated by servers.

#### 326 *4.2. Smart Cities and Crowds*

327 As to our knowledge there are no complete architectures for crowd sensing  
328 or crowd tracking that take into account the processing of the data and the  
329 extraction of information using Big Data processing techniques.

330 The architecture we presented in the previous section is well suited for crowd  
331 applications. In order to show this we detail each of the major parts of the  
332 architecture and show how they can be mapped for a simple crowd tracking  
333 system using WiFi scanners.

334 Crowd tracking using WiFi scanners is based on the ubiquitousness of smart-  
335 phones. These devices now have powerful processing, a large variety of sensor  
336 and communication capabilities. Most importantly for our application they  
337 almost always have a WiFi module. The WiFi module sends 802.11 packets  
338 in order to perform communication or auxiliary functions such as searching for  
339 networks. Because most of these packets contain a device identifier in the form  
340 of the MAC address, this means a device can be tracked by deploying WiFi  
341 scanners which record packets [17].

342 By looking at the architecture the WiFi scanners represent the sensors that  
343 gather data about the movements of crowds. This data needs to first be cleaned  
344 and filtered [18] as not all packets can be considered useful detections of a  
345 device. This initial cleaning and filtering procedures take place both at the  
346 scanners themselves in order to minimize bandwidth usage and at the central  
347 server that gathers data from all the scanners. This represents the second step  
348 in the architecture.

349 After the data from the WiFi scanners has a clean, normalized and standard-  
350 ized form it can be directly correlated with context data. There are numerous  
351 sources of context data freely available on the Internet. The simplest examples  
352 of context data sources are schedules or news posts. Both schedules and news  
353 posts offer a clear reasoning behind certain movements, for instance they can  
354 explain why a shop area has a lot of movement during work days and almost  
355 none in the weekend or during an important event.

356 Having multiple data sources and a continuous flow of information that can  
357 be correlated with historical events imposes the need for long term storage. Both  
358 context and sensor data is stored as well as any correlations between them. This  
359 data can be then analysed in real time or at a set time. The storage and data  
360 analysis steps match the next steps in our architecture.

361 Finally after the data is analyzed visualization tools need to be used in order  
362 to create an accessible way of making sense of the data for the individuals that  
363 need it. In the case of crowd tracking data visualization can take many forms.  
364 Usually it takes the form of a map where the density of people is shown by  
365 varying color or intensity. More information can be displayed in the form of a  
366 city map such as flows of people or events that happen at particular locations.  
367 Some decisions can skip the visualization step and directly announce the user.  
368 For instance if a traffic jam is detected people can be automatically informed in  
369 order for them to avoid the affected area.

## 370 5. Experiments

371 The first use case consider the Intelligent Transportation systems. The ap-  
372 plication model is as follow. The cabs are viewed as clients, which generate  
373 data with a sporadic schedule in a variety of sizes. The car GPS position is  
374 recorded every 15 seconds, and by default, the cluster client on the car sends  
375 the last 4 known positions every minute. However, if a car experiences a loss  
376 of connectivity, it may exhibit a pause in generating jobs and submit a larger  
377 data task when connectivity is reestablished. These tasks are considered as real  
378 time ones, and the aggregated data is computed as soon as possible.

379 A step by step workflow of the implemented application respect the model  
380 presented in Figure 1, as follows:

- 381 • Client sends position information to Cluster Proxy;
- 382 • Cluster Proxy writes cab data to distributed file system;
- 383 • Cluster Proxy encapsulates client data and puts job in appropriate sched-  
384 uler queue;
- 385 • Scheduler finds available cluster resources and creates job container on  
386 node;
- 387 • Map process on node reads data from the distributed file system and  
388 processes it;

- 389 • Map process on node aggregates new data with old data from distributed  
390 database or creates new DB entry if client is at its first report;
- 391 • Map process writes data back on the distributed database environment;
- 392 • Job is finished and resources are freed.

393 The flow of a request from inception until the end of its processing, from the  
394 technological point of view is as follows:

- 395 • Client thread reads positions from file and, depending on the profile it is  
396 assigned at startup, starts acting like a normal, mixed, or batch client;
- 397 • Proxy receives JSON through Camel;
- 398 • Proxy writes the data onto HDFS;
- 399 • Proxy triggers a new Hadoop job and submits it to the appropriate queue;
- 400 • Map process reads input HDFS data;
- 401 • Map reads existing data from HBase and aggregates it with the new data;
- 402 • Map writes end result back to HBase.

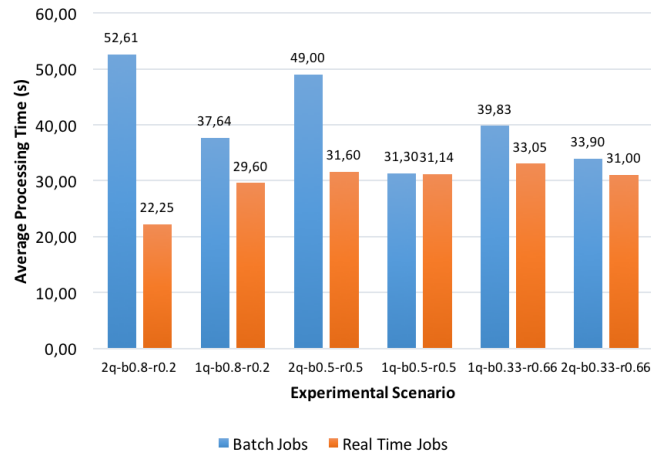


Figure 2: Comparison of average time for real-time and batch processing for different scenarios. These are total times, including data transfers time, time of data writing in HDFS and processing time, which require access to large data-sets collected from cabs.

403 Our experimental setup consists of 4 Virtual Box machines on top of a single  
404 physical host. The host has a 4-core CPU with hyperthreading at 2.4/max 3.4  
405 GHz, SSD drive and 16GB of memory. Out of the 4 virtual machines, 3 were  
406 kept purely for computation and storage needs (Datanodes in the case of HDFS)

407 and one was considered a master machine, which ran all the master nodes in  
408 the Hadoop architecture and also ran the Cluster Proxy module. The virtual  
409 machines had 20GB of storage assigned, 2 CPUs, and 3GB of memory, out of  
410 which 2GB were assigned for yarn containers in the case of the slaves.

411 The Hadoop Scheduler was configured with the following capacity parameters:  
412 The batch queue gets 30% of the capacity and may dynamically grow to no more  
413 than 60%, and the real-time queue has 70% of the capacity, but no more than  
414 90% if the batch queue is underutilized.

415 The experimental results are presented in Figure 2. The experiments were  
416 run with 1 and 2 queues. We can see that the processing time for batch jobs  
417 became comparable with time for real-time jobs. So, the conclusion is that we  
418 can combine these type of jobs, without any performance decreasing. More, by  
419 interpreting the average processing time, it is clearly that the performance of  
420 the cluster is best when the pattern of the input is similar to the one it was  
421 designed for. Although its resource limitations are flexible, they do not cater for  
422 extreme situations when the load is clearly not balanced. This problem could  
423 be solved with greater flexibility in resource limitations, as we imposed a rather  
424 fixed margin of resource distribution in configuring the scheduler.

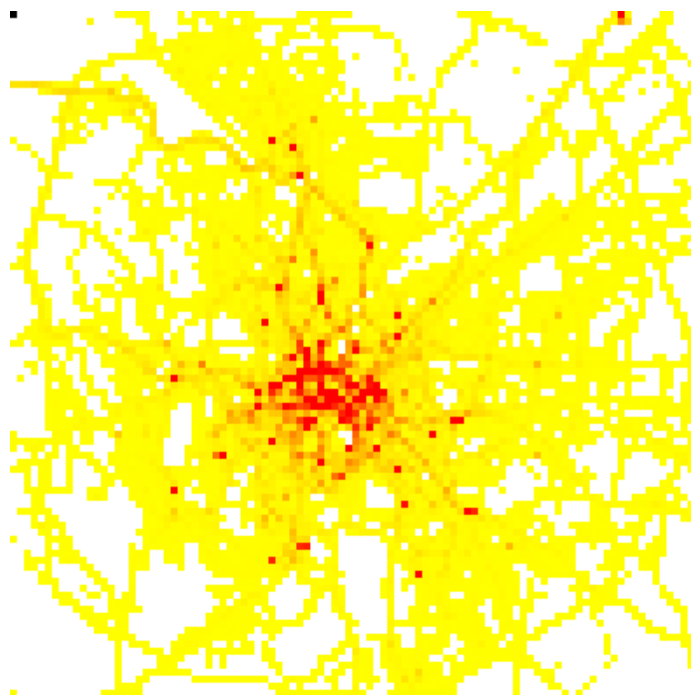


Figure 3: Rome - regions visited by taxis

425 Secondly we looked at crowd sensing data. We were interested to see what  
426 information the architecture can provide given an extensive data set. The data  
427 set we used was the roma-taxi data set available on CrowdAD. This data set

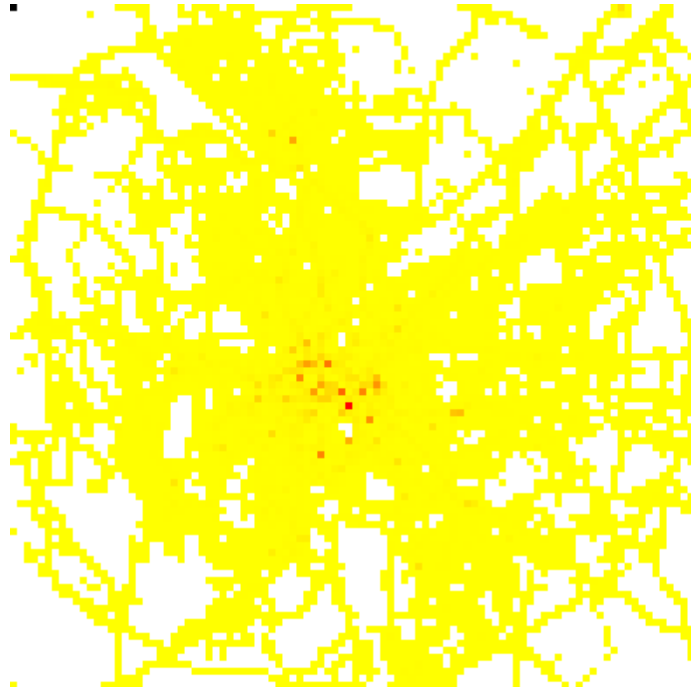


Figure 4: Rome - most important regions

428 consists of timestamped GPS data from multiple taxis that travel around the  
 429 city of Rome following their normal routines.

430 We imagine a future on which any car and in this case any taxi is equipped  
 431 not only with the necessities of every day transport but with sensors that are  
 432 able to provide all types of data. In order to understand how the data spans  
 433 across the city we measured how popular each individual part of the city is for  
 434 taxis. The data source for our architecture is given by the taxi GPS sensors.  
 435 This data is cleaned and normalized in the second part of the architecture. For  
 436 example, all positions outside the city limits are removed. After the data is  
 437 stored we move to the data analysis. We split the city in a grid of 100x100  
 438 and count the number of items with GPS coordinates in each of the grids. In  
 439 Figure 3 we displayed the results. This is equivalent to the data visualization  
 440 part of our network. With red we mark areas with high density and with yellow,  
 441 the ones with lower density.

442 Another visualization is available in Figure 4. Here we visualize the same  
 443 data but we set the maximum values as the maximal ones in the data set. This  
 444 permits us to accurately identify the centre of the city, the most popular area.

445 Using these visualization decision processes can be started. Automatic systems  
 446 can monitor the flow of people or cars and can decide which areas are  
 447 over-crowded and need assistance. They can also be used to identify expected  
 448 behavior when a large event such as a concert takes place in town.

## 449 6. Conclusions and Future Work

450 In this paper we proposed a generic architecture for data flow handling spe-  
451 cific for Smart Cities. We describe the functions and components for each step  
452 and identify specific technologies. Then we provide two use cases on crowd man-  
453 agement and intelligent transportation systems. We highlights experimental re-  
454 sults from applications developed using the model proposed in our architecture.  
455 As further work we will analyze self-adaptive optimisation methods used in this  
456 architecture, focusing on data reduction and data cleaning, patter extraction  
457 and data aggregation.

## 458 Acknowledgment

459 The research presented in this paper is supported by projects: *DataWay*:  
460 Real-time Data Processing Platform for Smart Cities: Making sense of Big Data  
461 - PN-II-RU-TE-2014-4-2731; *MobiWay*: Mobility Beyond Individualism: an In-  
462 tegrated Platform for Intelligent Transportation Systems of Tomorrow - PN-II-  
463 PT-PCCA-2013-4-0321; *CyberWater* grant of the Romanian National Authority  
464 for Scientific Research, CNDI-UEFISCDI, project number 47/2012; *clueFarm*:  
465 Information system based on cloud services accessible through mobile devices,  
466 to increase product quality and business development farms - PN-II-PT-PCCA-  
467 2013-4-0870.

468 We would like to thank the reviewers for their time and expertise, construc-  
469 tive comments and valuable insight.

## 470 References

- 471 [1] A. Caragliu, C. Del Bo, P. Nijkamp, Smart cities in europe, Journal of  
472 urban technology 18 (2) (2011) 65–82.
- 473 [2] R. H. Weber, R. Weber, Internet of Things, Springer, 2010.
- 474 [3] J. Manyika, M. Chui, B. Brown, J. Bughin, R. Dobbs, C. Roxburgh, A. H.  
475 Byers, Big data: The next frontier for innovation, competition, and pro-  
476 ductivity.
- 477 [4] S. Lohr, The age of big data, New York Times 11.
- 478 [5] A. Rev, Wireless sensor networks.
- 479 [6] H. Ma, D. Zhao, P. Yuan, Opportunities in mobile crowd sensing, Commu-  
480 nications Magazine, IEEE 52 (8) (2014) 29–35.
- 481 [7] G. P. Hancke, G. P. Hancke Jr, et al., The role of advanced sensing in smart  
482 cities, Sensors 13 (1) (2012) 393–425.
- 483 [8] K. Su, J. Li, H. Fu, Smart city and the applications, in: Electronics,  
484 Communications and Control (ICECC), 2011 International Conference on,  
485 IEEE, 2011, pp. 1028–1031.

- 486 [9] G. K. Still, Crowd dynamics, Ph.D. thesis, University of Warwick (2000).
- 487 [10] H. Zhao, R. Shibasaki, A novel system for tracking pedestrians using mul-  
488 tiple single-row laser-range scanners, *Systems, Man and Cybernetics, Part*  
489 *A: Systems and Humans*, IEEE Transactions on 35 (2) (2005) 283–291.
- 490 [11] M. Betke, E. Haritaoglu, L. S. Davis, Multiple vehicle detection and track-  
491 ing in hard real-time, in: *Intelligent Vehicles Symposium, 1996.*, Proceed-  
492 ings of the 1996 IEEE, IEEE, 1996, pp. 351–356.
- 493 [12] M. S. Grewal, L. R. Weill, A. P. Andrews, *Global positioning systems,*  
494 *inertial navigation, and integration*, John Wiley & Sons, 2007.
- 495 [13] D. Niculescu, B. Nath, Ad hoc positioning system (aps), in: *Global*  
496 *Telecommunications Conference, 2001. GLOBECOM'01. IEEE, Vol. 5,*  
497 *IEEE, 2001*, pp. 2926–2931.
- 498 [14] A. J. Ruiz-Ruiz, H. Blunck, T. S. Prentow, A. Stisen, M. B. Kjaergaard,  
499 Analysis methods for extracting knowledge from large-scale wifi monitoring  
500 to inform building facility planning, in: *Pervasive Computing and Commu-*  
501 *nications (PerCom), 2014 IEEE International Conference on, IEEE, 2014,*  
502 *pp. 130–138.*
- 503 [15] L. Vu, K. Nahrstedt, S. Retika, I. Gupta, Joint bluetooth/wifi scanning  
504 framework for characterizing and leveraging people movement in univer-  
505 sity campus, in: *Proceedings of the 13th ACM international conference on*  
506 *Modeling, analysis, and simulation of wireless and mobile systems, ACM,*  
507 *2010*, pp. 257–265.
- 508 [16] M. Zhou, Z. Tian, K. Xu, X. Yu, X. Hong, H. Wu, Scanme: location  
509 tracking system in large-scale campus wi-fi environment using unlabeled  
510 mobility map, *Expert systems with applications* 41 (7) (2014) 3429–3443.
- 511 [17] B. Bonne, A. Barzan, P. Quax, W. Lamotte, Wifipi: Involuntary tracking  
512 of visitors at mass events, in: *World of Wireless, Mobile and Multime-*  
513 *dia Networks (WoWMoM), 2013 IEEE 14th International Symposium and*  
514 *Workshops on a, IEEE, 2013*, pp. 1–6.
- 515 [18] C. D. C. Chilipirea, A.-C. Petre, M. v. Steen, Filters for wi-fi generated  
516 crowd movement data, in: *10th International Conference on P2P, Parallel,*  
517 *Grid, Cloud and Internet Computing, IEEE, 2015*, pp. 285–290.
- 518 [19] Y. Wang, J. Yang, H. Liu, Y. Chen, M. Gruteser, R. P. Martin, Measuring  
519 human queues using wifi signals, in: *Proceedings of the 19th annual inter-*  
520 *national conference on Mobile computing & networking, ACM, 2013*, pp.  
521 235–238.
- 522 [20] L. Schauer, M. Werner, P. Marcus, Estimating crowd densities and pedes-  
523 trian flows using wi-fi and bluetooth, in: *Proceedings of the 11th In-*  
524 *ternational Conference on Mobile and Ubiquitous Systems: Computing,*

- 525 Networking and Services, ICST (Institute for Computer Sciences, Social-  
526 Informatics and Telecommunications Engineering), 2014, pp. 171–177.
- 527 [21] K. Han, E. Graham, D. Vassallo, D. Estrin, et al., Enhancing motivation  
528 in a mobile participatory sensing project through gaming, in: Privacy, Se-  
529 curity, Risk and Trust (PASSAT) and 2011 IEEE Third International Con-  
530 ference on Social Computing (SocialCom), 2011 IEEE Third International  
531 Conference on, IEEE, 2011, pp. 1443–1448.
- 532 [22] H. Aly, A. Basalamah, M. Youssef, Map++: A crowd-sensing system for  
533 automatic map semantics identification, in: Sensing, Communication, and  
534 Networking (SECON), 2014 Eleventh Annual IEEE International Confer-  
535 ence on, IEEE, 2014, pp. 546–554.
- 536 [23] M. Faulkner, M. Olson, R. Chandy, J. Krause, K. M. Chandy, A. Krause,  
537 The next big one: Detecting earthquakes and other rare events from  
538 community-based sensors, in: Information Processing in Sensor Networks  
539 (IPSN), 2011 10th International Conference on, IEEE, 2011, pp. 13–24.
- 540 [24] M.-R. Ra, B. Liu, T. F. La Porta, R. Govindan, Medusa: A programming  
541 framework for crowd-sensing applications, in: Proceedings of the 10th in-  
542 ternational conference on Mobile systems, applications, and services, ACM,  
543 2012, pp. 337–350.
- 544 [25] I. Carreras, D. Miorandi, A. Tamin, E. R. Ssebagala, N. Conci, Matador:  
545 Mobile task detector for context-aware crowd-sensing campaigns, in: Perva-  
546 sive Computing and Communications Workshops (PERCOM Workshops),  
547 2013 IEEE International Conference on, IEEE, 2013, pp. 212–217.
- 548 [26] P. P. Jayaraman, C. Perera, D. Georgakopoulos, A. Zaslavsky, Efficient op-  
549 portunistic sensing using mobile collaborative platform mosden, in: Collab-  
550 orative Computing: Networking, Applications and Worksharing (Collabo-  
551 ratecom), 2013 9th International Conference Conference on, IEEE, 2013,  
552 pp. 77–86.
- 553 [27] T. Yan, M. Marzilli, R. Holmes, D. Ganesan, M. Corner, mcrowd: a plat-  
554 form for mobile crowdsourcing, in: Proceedings of the 7th ACM Conference  
555 on Embedded Networked Sensor Systems, ACM, 2009, pp. 347–348.
- 556 [28] H. Vääätäjä, T. Vainio, E. Sirkkunen, K. Salo, Crowdsourced news report-  
557 ing: supporting news content creation with mobile phones, in: Proceedings  
558 of the 13th International Conference on Human Computer Interaction with  
559 Mobile Devices and Services, ACM, 2011, pp. 435–444.
- 560 [29] T. Xie, X. Qin, Scheduling security-critical real-time applications on clus-  
561 ters, *Computers*, IEEE Transactions on 55 (7) (2006) 864–879.
- 562 [30] J. P. Erickson, G. Coombe, J. H. Anderson, Soft real-time scheduling in  
563 google earth, in: Real-Time and Embedded Technology and Applications  
564 Symposium (RTAS), 2012 IEEE 18th, IEEE, 2012, pp. 141–150.



- 565 [31] S. Choi, A. K. Agrawala, Scheduling aperiodic and sporadic tasks in hard  
566 real-time systems.
- 567 [32] C. Dobre, G. Suciu, C. Chilipirea, C. Gosman, Mobility beyond individ-  
568 ualism: an integrated platform for intelligent transportation systems of  
569 tomorrow.