

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

Cristian Chilipirea^a, Andreea-Cristina Petre^a, Loredana-Marsilia Groza^a,
Ciprian Dobre^a, Florin Pop^{*a}

^a*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 and air pollution management, require different data processing techniques. The main subject of this paper is to propose 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: architecture; big data; data processing; crowd sensing; crowd dynamics; intelligent transportation systems

1. Introduction

2 More and more applications today use, generate and handle very large vol-
3 umes of data. In particular, this is true for Smart City applications, which
4 attract a rapidly increasing interest from government, companies, citizens, de-
5 velopers, scientists, etc. They cover a large spectrum of needs in public safety,
6 water and energy management, smart buildings, government and agency admin-
7 istration, social programs, transportation, health, education. They are fed with
8 huge amounts of input data, in various formats, from a continuously increasing
9 number of sources (sensors, governmental, regional, and municipal sources, cit-
10 izens, public open data sources, etc.), and are described by a complex workflow

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

11 and in many cases impose real-time processing capabilities, useful in decision
12 taking.

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

22 Extracting valuable information from raw data is especially difficult con-
23 sidering the velocity of growing data from year to year and the fact that 80%
24 of data is unstructured. In addition, data sources are heterogeneous (various
25 sensors, users with different profiles, etc.) and are located in different situa-
26 tions or contexts. This is why the Smart City infrastructure runs reliably and
27 permanently to provide the context as a “public utility” to different services.
28 Context-aware applications exploit the context to adapt accordingly the timing,
29 quality and functionality of their services. The value of these applications and
30 their supporting infrastructure lies in end-users always operating in a context:
31 their role, intentions, locations and working environment constantly 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 which 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 propo-
42 pose 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 visualiza-
46 tion 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 is 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 which can dramatically improve
55 the life of citizens. There is a lot of research aiming 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, 3, 4] and Big Data [5]. New York times
58 actually declared this period the "Age of Big Data"¹.

59 In order to enable Smart Cities technologies such as Internet of Things, Wire-
60 less Sensor Networks [6] and Crowd Sensing [7] are the catalysts providing data
61 about our cities. The need for sensing in Smart Cities is explored in [8]. This
62 data needs to be processed often using Big Data techniques in order to extract
63 the information required to make decisions about the cities. This information
64 and the decisions are then used in order to inform the citizens to take certain
65 actions or to activate actuators for enabling automatic processes. A good ex-
66 ample where actuators can improve Smart Cities is given by the management
67 of green spaces [9].

68 Probably the most important issues addressed in order to build Smart Cities
69 are the ones of Crowd Dynamics [10]. In order to understand Crowd Dynamics,
70 we need data on the movements of as many people as possible. These move-
71 ments need to be recorded for both pedestrians [11] and for vehicles [12]. The
72 problem of tracking is not solved in any of the two scenarios. This is surprising,
73 considering the problem of tracking a particular individual is usually solved by
74 the use of GPS [13]. However, GPS requires user participation which is difficult
75 to obtain, in contrast WiFi [14] or cellular methods [15] can be used to gather
76 data on large crowds. These systems also do not work indoors and require
77 the cooperation of the individual being tracked in order to generate a position
78 estimate.

79 It is important to treat both indoor and outdoor cases when considering
80 human mobility. This is because modern vital facilities, such as hospitals, which
81 are part of the backbone of many cities consist of large areas with multiple
82 buildings. An example of how the dynamics inside these facilities can be used
83 in order to improve the layout of the facility is given by Ruiz et al [16]. Similarly,
84 Universities campuses, another type of large facilities at the core of cities, are
85 analyzed [17], [18] in order to better understand the dynamics inside them.

86 Crowd tracking experiments are taking place in a wider variety of places
87 like mass events [19] or festivals [20]. They are also used in order to measure
88 queues using only WiFi signals [21]. This queue can represent waiting time at
89 a counter, which directly affects customer experience or the movement through
90 security lines at an airport [22].

91 Crowd Sensing can be used in order to extract all types of data for smart
92 cities. A powerful example is given by the authors of [23] where students are
93 asked to take pictures of plants around the campus. The pictures are then
94 analyzed by scientists in order to better understand the status of flora. Projects

¹"The age of big data" - Steve Lohr, New York Times, 11, 2012

95 like this could potentially be used at the scale of a city in order to measure a
96 large variety of features. It is not always necessary for people to be active in their
97 participation of data gathering. Passive systems require only their presence in
98 the monitored location, which can even be obtained in an opportunistic manner.
99 Whenever any citizens carrying the scanner walks or drives on a specific street
100 data about the street can be gathered. In this way maps can be enhanced
101 with features [24] such as roundabouts or pot-holes. Diverse uses include even
102 earthquake detection [25] and soon maybe even the detection of effects produced
103 by these large natural disasters.

104 There are many projects and platforms targeted directly at crowd sensing:
105 Medusa [26], Matador [27], Mosden [28] and mCrowd [29]. And these platforms
106 already implement important features for Smart Cities such as crowd sources
107 new reporting [30] but they do not yet combine the data sets or offer a method
108 to analyze the data in order to extract information hidden inside it. This type
109 of information represents answers to questions that we don't yet have and they
110 can currently only be obtained by using Big Data techniques.

111 The data gathered from all these systems is usually analyzed by experts or
112 scientist manually. This is the case for [16], where categorization of individuals
113 into different groups such as patients or staff is done by using rules built by
114 experts. More information can be extracted from these data sets if they are
115 combined and Big Data systems are used to process them.

116 Real-time processing is used to designate a category where the job outcome
117 is needed as fast as possible, and usually the task itself is not something taking
118 a long time to process. These systems can be categorized as hard or soft. A
119 Hard real-time system is an OS for a nuclear plant or a plane. Tasks must be
120 scheduled and completed fast because otherwise a catastrophe could happen.
121 These systems are usually governed by hard deadlines and the scheduler must
122 insure they are achieved. Soft real-time systems are the ones like hotel book-
123 ing or video streaming sites [31]. The answers must be delivered fast to the
124 customers, but a delayed frame now and then cannot lead to disastrous results.
125 One article which explores this type of hard real-time scheduling is [32]. In
126 the paper the authors try to improve the scheduling capabilities of a system by
127 also adding security checks to the incoming jobs. The added module can detect
128 threats brought by snooping, alteration of spoofing and can be easily added to
129 any real-time scheduler. Their security module name SAREC (security-aware
130 real-time heuristic strategy for clusters) integrates with the popular Earliest
131 Deadline First algorithm to create a security aware scheduler named SAEDF.
132 Although the matter of securing the interactions between the users and the
133 cluster infrastructure is important, in our case a large portion of these measures
134 could be implemented in an intermediate cluster proxy module if needed, with
135 little overhead to the job itself. By using a proxy to mediate all user-cluster
136 interactions we can alleviate a large number of security risks. If a user has a
137 malicious intent and manages to submit a job that poses a security risk, runing
138 all jobs in virtual machines on the cluster infrastructure will limit the damages
139 to only the users task.

140 Another example of real-time processing and scheduling [33]. The authors

141 talk about the problem of soft real-time scheduling in rendering 3D images inside
142 the Google Earth software. The Google Earth software allows one to navigate
143 anywhere in the world and has multiple viewing modes from virtual 3D ren-
144 derings to satellite imagery. A frame is a static 2D representation, rendered on
145 the screen at a given time. To ensure a smooth navigation experience, at least
146 60 of these frames must be rendered on the users' screen in a second. When
147 a scheduling deadline is not met, the previous frame is redisplayed causing the
148 application to "stutter". In order to alleviate the problems, the authors have
149 devised a new algorithm that better estimates rendering time on multiple de-
150 vices, in order to improve scheduler accuracy. We are in particular interested in
151 their scheduling model and discovered they also abstracted some of the events
152 into "single-active sporadic tasks" (triggered by a specific rendering phase) and
153 "soft real-time aperiodic tasks (triggered by receiving new imagery through net-
154 work)". We will use similar terms to define the submit patterns and properties
155 of different types of jobs.

156 Talking about the arrival patterns of the jobs, the authors from [34] build
157 a common approach to schedule static and dynamic tasks, in a system which
158 also has to deal with hard real-time deadlines. They divide their tasks in 3
159 categories, based on their arrival pattern and number of instances they require
160 for running. Aperiodic tasks need only one instance to run, and can enter the
161 system at any time. Both periodic and sporadic tasks require multiple instances
162 to run, but while the former come at a specific interval of time, the latter can be
163 submitted like the aperiodic tasks, at any time, but no sooner than a specified
164 interval. The authors have extended a previous static time-based scheduling
165 algorithm into a dynamic version which constantly changes the expected start
166 and end time of jobs while still keeping the end time in the necessary deadline.
167 They have thus provided two versions of their scheduler, one accepting aperiodic
168 tasks without affecting the existing task instances deadline, and another, with
169 the same properties, accepting periodic and sporadic tasks. Before accepting
170 any task, a formal schedulable test is run, to see if the system can handle the
171 tasks deadline. If not, it is rejected. The scheduled tasks are considered to
172 be preemptive, and a list of static tasks known beforehand is expected to be
173 provided at system startup. To account for dependencies between tasks, start
174 and end times are parameterized instead of being given a fixed value.

175 We also investigated solutions related to intelligent transport systems, since
176 this is the type of workload we are going to test our scheduler on. The [35]
177 project tries to act as a hub for such endeavors in order to help each of the
178 individual current ITS system grow and communicate through a common point
179 of contact. These systems are increasingly important since optimizing traffic
180 can also reduce CO_2 emissions along with the benefits brought to all the in-
181 habitants of a city. Current implemented solutions are mostly proprietary and
182 involve infrastructure changes. There are a number of existing solutions trying
183 to estimate the state of the traffic, ranging from sensors in the road, to GPS
184 systems on cars, to cameras interpreting images. Indifferent of the chosen so-
185 lution, all of these systems will generate a large amount of data which has to
186 be interpreted city-wide. Although our solution uses a small part of this data,

187 it could grow and adapt to provide the necessary analysis needed to drive an
 188 intelligent city of today.

189 3. Data Flow based Architecture

190 We propose an architecture which contains several steps represented by the
 191 flow of data from the source to the end user. The data represents the input of
 192 our system, it is used to create valuable good for the users, usually in the form
 193 of information or automatic actions.

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

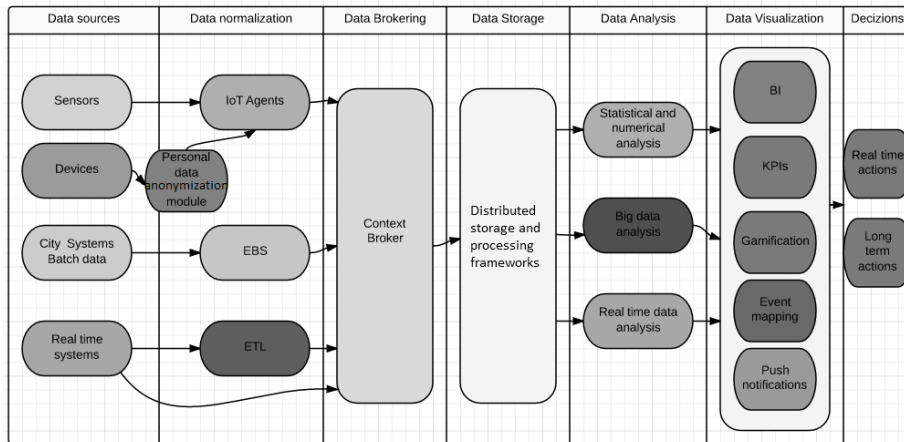


Figure 1: Proposed Architecture.

203 Data comes from different sources and we need to collect it from everywhere:
 204 smart sensors, personal devices, batch data from city systems, real-time systems
 205 in order to be able to extract as much knowledge as we can from it. When we
 206 combine different data sources related to the same context then we can get more
 207 insights from it and this empowers decision makers to minimize 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 he remains anony-
 211 mous and other users cannot trace him back starting from the pieces of data
 212 provided by him. We need to make sure he cannot be identified from a group

213 of users which contribute to the system with their data by using techniques
214 and algorithms for data anonymization. Data collection needs to be done while
215 considering data privacy. This can be achieved by using a data anonymization
216 module. Anonymization needs to be done as close to the source as possible to
217 avoid any data leaks, that could identify an individual as a data provider.

218 Real-time systems data needs to be normalized by using ETL - Extract,
219 Transform and Load - representing 3 database functions joined in one tool in
220 order to get the data out from a database and introduce it into another one.

221 City systems batch data needs to be normalized by EBS - Electronic Batchload
222 Service –which is an “Online Computer Library Center” service permitting to
223 the batch load participants to send data to it over the Internet [36, 37].

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

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

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

253 When processing the data, we can perform big data analytics, statistical
254 and numerical analysis or real-time data analysis to gain valuable assets. In the
255 end, we need advanced data visualizations to enable the user to take the right
256 decisions, or to make long term actions based on historical data.

257 Real-time data processing and analytics allows decision makers the oppor-
258 tunity to take immediate action when it is required and batch data processing

259 makes the results to be more accurate due to the patterns discovered and then
260 applied in real-time to get more relevant data.

261 We need to combine the data from multiple sources to be able to predict
262 future events in order to respond in an efficient way to make a difference. It
263 is important to engage the users and to achieve this we need to empower and
264 motivate them. A way of smart user engagement and advanced visualization is
265 gamification. For example, users of a mobile application can share status about
266 how much do they recycle different things, or how much CO₂ they produced
267 based on how many km they were driving in a day, and enter in competitions
268 with others on social media.

269 4. Architecture Use Cases

270 4.1. *Intelligent Transportation Systems*

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

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

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

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

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

324 *4.2. Smart Cities and Crowds*

325 As to our knowledge, there are no complete architectures for crowd sensing or
326 crowd tracking taking into account processing of the data information extraction
327 using Big Data techniques.

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

332 Crowd tracking using WiFi scanners is based on the ubiquitousness of smart-
333 phones. These devices now have powerful processing, a large variety of sensor
334 and communication capabilities. Most importantly for our application they
335 almost always have a WiFi module. The WiFi module sends 802.11 packets
336 in order to perform communication or auxiliary functions such as searching for
337 networks. Because most of these packets contain a device identifier in the form
338 of the MAC address, this means a device can be tracked by deploying WiFi
339 scanners which record packets [19].

340 By looking at the architecture the WiFi scanners represent the sensors which
341 gather data about the movements of crowds. This data needs to first be cleaned
342 and filtered [20] as not all packets can be considered useful detections of a
343 device. This initial cleaning and filtering procedures take place both at the
344 scanners themselves in order to minimize bandwidth usage and at the central

345 server that gathers data from all the scanners. This represents the second step
346 in the architecture.

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

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

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

368 5. Experiments

369 The first use case considers the Intelligent Transportation systems. The
370 application model is as follow. The cabs are viewed as clients, which generate
371 data with a sporadic schedule in a variety of sizes. The car GPS position is
372 recorded every 15 seconds, and by default, the cluster client on the car sends
373 the last 4 known positions every minute. However, if a car experiences a loss
374 of connectivity, it may exhibit a pause in generating jobs and submit a larger
375 data task when connectivity is reestablished. These tasks are considered as
376 real-time ones, and the aggregated data is computed as soon as possible. More
377 experimental results have been presented in [38].

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

- 380 • A client sends position information to Cluster Proxy;
- 381 • The Cluster Proxy writes cab data to distributed file system;
- 382 • The Cluster Proxy encapsulates client data and puts a job in appropriate
383 scheduler queue;
- 384 • The scheduler finds available cluster resources and creates a job container
385 on a node;

- 386 • The map process reads data from the distributed file system and processes
387 it;
- 388 • The map process aggregates new data with old data from distributed
389 database or creates new DB entry if client is at its first report;
- 390 • The map process writes data back on the distributed database environ-
391 ment;
- 392 • The 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.

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 hyper threading at 2.4/max 3.4
405 GHz, SSD drive and 16 GB of memory. Out of the 4 virtual machines, 3 were
406 kept purely for computation and storage needs (Data nodes 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 20 GB of storage assigned, 2 CPUs, and 3 GB of memory, out of
410 which 2 GB 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. In conclusion, we can combine
418 these type of jobs, without any performance decreasing. Moreover, by inter-
419 preting the average processing time, it is clearly that the performance of the
420 cluster is best when the pattern of the input is similar to the one it was de-
421 signed 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

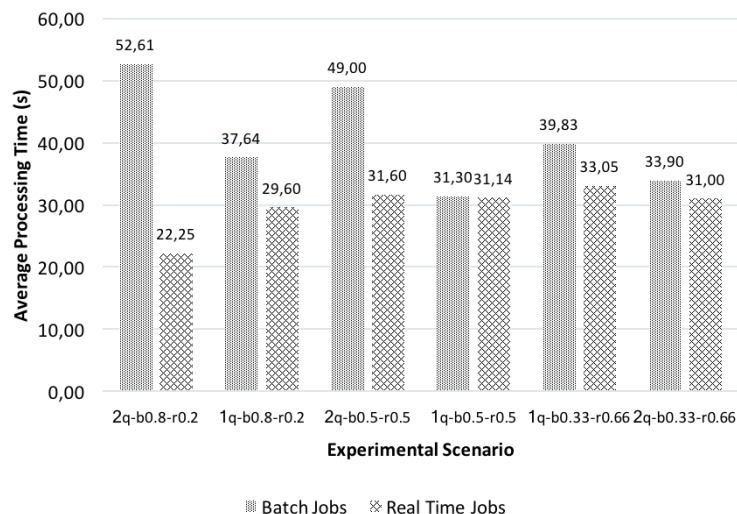


Figure 2: Comparison of average time for real-time and batch processing for different scenarios. These are total times, including data transfers time, time of data writing in HDFS and processing time, which require access to large data-sets collected from cabs.

423 be solved with greater flexibility in resource limitations, as we imposed a rather
 424 fixed margin of resource distribution in configuring the scheduler.

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 Crowdax. This data set
 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 which 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. Red marks areas with high density and yellow the ones
 441 with low density.

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

445 Using these visualizations an individual can then start the decision processes.

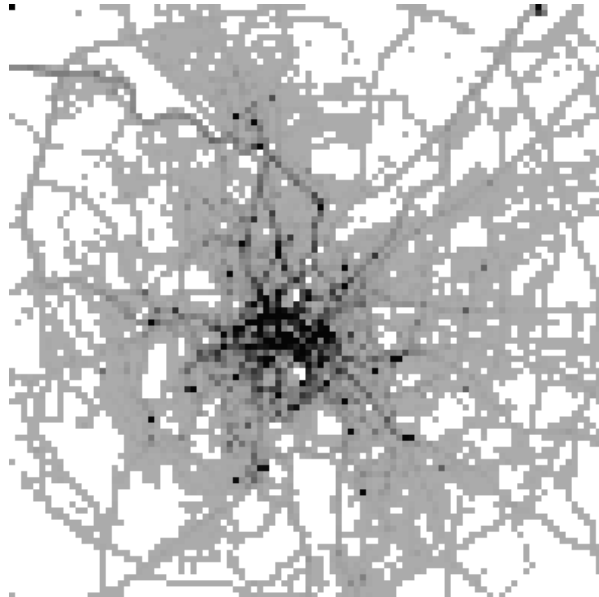


Figure 3: Rome - regions visited by taxis.

446 Automatic systems can monitor the flow of people or cars and can decide which
447 areas are over-crowded and need assistance. They can also be used to identify
448 expected 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
453 management and intelligent transportation systems. We highlight experimental
454 results from applications developed using the model proposed in our architec-
455 ture. As further work we will analyze self-adaptive optimization methods used
456 in this architecture, focusing on data reduction and data cleaning, patter ex-
457 traction 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*:

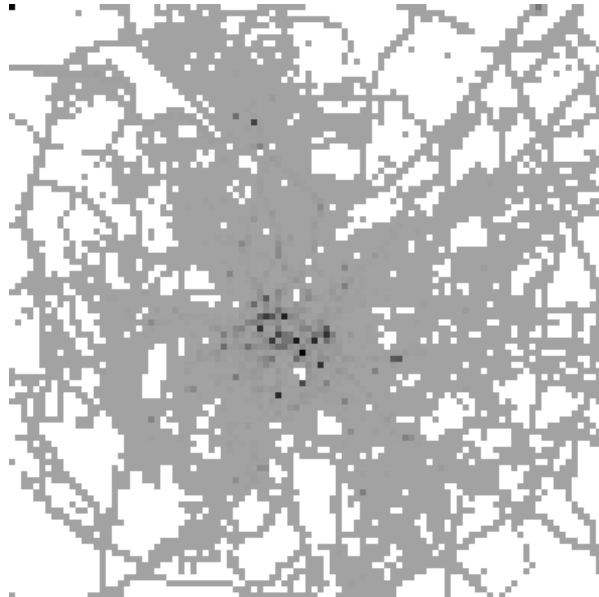


Figure 4: Rome - most important regions.

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] V. Albino, U. Berardi, R. M. Dangelico, Smart cities: Definitions, dimen-
472 sions, performance, and initiatives, *Journal of Urban Technology* 22 (1)
473 (2015) 3–21.
- 474 [2] A. Whitmore, A. Agarwal, L. Da Xu, The internet of things—a survey of
475 topics and trends, *Information Systems Frontiers* 17 (2) (2015) 261–274.
- 476 [3] J. M. Batalla, P. Krawiec, Conception of id layer performance at the net-
477 work level for internet of things, *Personal and Ubiquitous Computing* 18 (2)
478 (2014) 465–480.
- 479 [4] J. M. Batalla, M. Gajewski, W. Latoszek, P. Krawiec, C. X. Mavromous-
480 takis, G. Mastorakis, Id-based service-oriented communications for unified
481 access to iot, *Computers & Electrical Engineering* 52 (2016) 98–113.
- 482 [5] A. Gandomi, M. Haider, Beyond the hype: Big data concepts, methods,
483 and analytics, *International Journal of Information Management* 35 (2)
484 (2015) 137–144.

- 485 [6] C. S. Raghavendra, K. M. Sivalingam, T. Znati, *Wireless sensor networks*,
486 Springer, 2006.
- 487 [7] H. Ma, D. Zhao, P. Yuan, Opportunities in mobile crowd sensing, *Commu-*
488 *nications Magazine*, IEEE 52 (8) (2014) 29–35.
- 489 [8] G. P. Hancke, G. P. Hancke Jr, et al., The role of advanced sensing in smart
490 cities, *Sensors* 13 (1) (2012) 393–425.
- 491 [9] K. Su, J. Li, H. Fu, Smart city and the applications, in: *Electronics,*
492 *Communications and Control (ICECC)*, 2011 International Conference on,
493 IEEE, 2011, pp. 1028–1031.
- 494 [10] G. K. Still, *Crowd dynamics*, Ph.D. thesis, University of Warwick (2000).
- 495 [11] H. Du, Z. Yu, F. Yi, Z. Wang, Q. Han, B. Guo, Group mobility classification
496 and structure recognition using mobile devices, 2016 IEEE International
497 Conference on Pervasive Computing and Communications, PerCom 2016.
- 498 [12] D. Kumar, H. Wu, Y. Lu, S. Krishnaswamy, M. Palaniswami, Understanding
499 Urban Mobility via Taxi Trip Clustering, 2016 17th IEEE International
500 Conference on Mobile Data Management (MDM) (2016) 318–324.
- 501 [13] M. S. Grewal, L. R. Weill, A. P. Andrews, *Global positioning systems,*
502 *inertial navigation, and integration*, John Wiley & Sons, 2007.
- 503 [14] Y. Chon, S. Kim, S. Lee, D. Kim, Y. Kim, H. Cha, Sensing WiFi packets
504 in the air, Proceedings of the 2014 ACM International Joint Conference
505 on Pervasive and Ubiquitous Computing - UbiComp '14 Adjunct (2014)
506 189–200.
- 507 [15] M. Dash, K. K. Koo, S. P. Krishnaswamy, Y. Jin, A. Shi-Nash, Visualize
508 people’s mobility - Both individually and collectively - Using mobile phone
509 cellular data, Proceedings - IEEE International Conference on Mobile Data
510 Management 2016-July (2016) 341–344.
- 511 [16] A. J. Ruiz-Ruiz, H. Blunck, T. S. Prentow, A. Stisen, M. B. Kjaergaard,
512 Analysis methods for extracting knowledge from large-scale wifi monitoring
513 to inform building facility planning, in: *Pervasive Computing and Commu-*
514 *nications (PerCom)*, 2014 IEEE International Conference on, IEEE, 2014,
515 pp. 130–138.
- 516 [17] L. Vu, K. Nahrstedt, S. Retika, I. Gupta, Joint bluetooth/wifi scanning
517 framework for characterizing and leveraging people movement in univer-
518 sity campus, in: *Proceedings of the 13th ACM international conference on*
519 *Modeling, analysis, and simulation of wireless and mobile systems*, ACM,
520 2010, pp. 257–265.
- 521 [18] M. Zhou, Z. Tian, K. Xu, X. Yu, X. Hong, H. Wu, Scanme: location
522 tracking system in large-scale campus wi-fi environment using unlabeled
523 mobility map, *Expert systems with applications* 41 (7) (2014) 3429–3443.

- 524 [19] B. Bonne, A. Barzan, P. Quax, W. Lamotte, Wifipi: Involuntary tracking
525 of visitors at mass events, in: World of Wireless, Mobile and Multime-
526 dia Networks (WoWMoM), 2013 IEEE 14th International Symposium and
527 Workshops on a, IEEE, 2013, pp. 1–6.
- 528 [20] C. D. C. Chilipirea, A.-C. Petre, M. v. Steen, Filters for wi-fi generated
529 crowd movement data, in: 10th International Conference on P2P, Parallel,
530 Grid, Cloud and Internet Computing, IEEE, 2015, pp. 285–290.
- 531 [21] Y. Wang, J. Yang, H. Liu, Y. Chen, M. Gruteser, R. P. Martin, Measuring
532 human queues using wifi signals, in: Proceedings of the 19th annual inter-
533 national conference on Mobile computing & networking, ACM, 2013, pp.
534 235–238.
- 535 [22] L. Schauer, M. Werner, P. Marcus, Estimating crowd densities and pedest-
536 rian flows using wi-fi and bluetooth, in: Proceedings of the 11th In-
537 ternational Conference on Mobile and Ubiquitous Systems: Computing,
538 Networking and Services, ICST (Institute for Computer Sciences, Social-
539 Informatics and Telecommunications Engineering), 2014, pp. 171–177.
- 540 [23] K. Han, E. Graham, D. Vassallo, D. Estrin, et al., Enhancing motivation
541 in a mobile participatory sensing project through gaming, in: Privacy, Se-
542 curity, Risk and Trust (PASSAT) and 2011 IEEE Third International Con-
543 ference on Social Computing (SocialCom), 2011 IEEE Third International
544 Conference on, IEEE, 2011, pp. 1443–1448.
- 545 [24] H. Aly, A. Basalamah, M. Youssef, Map++: A crowd-sensing system for
546 automatic map semantics identification, in: Sensing, Communication, and
547 Networking (SECON), 2014 Eleventh Annual IEEE International Confer-
548 ence on, IEEE, 2014, pp. 546–554.
- 549 [25] M. Faulkner, M. Olson, R. Chandy, J. Krause, K. M. Chandy, A. Krause,
550 The next big one: Detecting earthquakes and other rare events from
551 community-based sensors, in: Information Processing in Sensor Networks
552 (IPSN), 2011 10th International Conference on, IEEE, 2011, pp. 13–24.
- 553 [26] M.-R. Ra, B. Liu, T. F. La Porta, R. Govindan, Medusa: A programming
554 framework for crowd-sensing applications, in: Proceedings of the 10th in-
555 ternational conference on Mobile systems, applications, and services, ACM,
556 2012, pp. 337–350.
- 557 [27] I. Carreras, D. Miorandi, A. Tamin, E. R. Ssebagala, N. Conci, Matador:
558 Mobile task detector for context-aware crowd-sensing campaigns, in: Perva-
559 sive Computing and Communications Workshops (PERCOM Workshops),
560 2013 IEEE International Conference on, IEEE, 2013, pp. 212–217.
- 561 [28] P. P. Jayaraman, C. Perera, D. Georgakopoulos, A. Zaslavsky, Efficient op-
562 portunistic sensing using mobile collaborative platform mosden, in: Collab-
563 orative Computing: Networking, Applications and Worksharing (Collabo-

- 564 ratecom), 2013 9th International Conference Conference on, IEEE, 2013,
565 pp. 77–86.
- 566 [29] T. Yan, M. Marzilli, R. Holmes, D. Ganesan, M. Corner, mcrowd: a plat-
567 form for mobile crowdsourcing, in: Proceedings of the 7th ACM Conference
568 on Embedded Networked Sensor Systems, ACM, 2009, pp. 347–348.
- 569 [30] H. Väättäjä, T. Vainio, E. Sirkkunen, K. Salo, Crowdsourced news report-
570 ing: supporting news content creation with mobile phones, in: Proceedings
571 of the 13th International Conference on Human Computer Interaction with
572 Mobile Devices and Services, ACM, 2011, pp. 435–444.
- 573 [31] J. M. Batalla, Advanced multimedia service provisioning based on efficient
574 interoperability of adaptive streaming protocol and high efficient video cod-
575 ing, *Journal of Real-Time Image Processing* 12 (2) (2016) 443–454.
- 576 [32] T. Xie, X. Qin, Scheduling security-critical real-time applications on clus-
577 ters, *Computers, IEEE Transactions on* 55 (7) (2006) 864–879.
- 578 [33] J. P. Erickson, G. Coombe, J. H. Anderson, Soft real-time scheduling in
579 google earth, in: *Real-Time and Embedded Technology and Applications*
580 *Symposium (RTAS)*, 2012 IEEE 18th, IEEE, 2012, pp. 141–150.
- 581 [34] B. Sprunt, L. Sha, J. Lehoczky, Aperiodic task scheduling for hard-real-
582 time systems, *Real-Time Systems* 1 (1) (1989) 27–60.
- 583 [35] C. Dobre, G. Suci, C. Chilipirea, C. Gosman, Mobility beyond individ-
584 ualism: an integrated platform for intelligent transportation systems of
585 tomorrow, in: *ITS Romania Congress*, 2014, pp. 31–35.
- 586 [36] Y. Kryftis, G. Mastorakis, C. X. Mavromoustakis, J. M. Batalla, E. Pal-
587 lis, G. Kormentzas, Efficient entertainment services provision over a novel
588 network architecture, *IEEE Wireless Communications* 23 (1) (2016) 14–21.
- 589 [37] J. M. Batalla, M. Kantor, C. X. Mavromoustakis, G. Skourletopoulos,
590 G. Mastorakis, A novel methodology for efficient throughput evaluation in
591 virtualized routers, in: *Communications (ICC)*, 2015 IEEE International
592 Conference on, IEEE, 2015, pp. 6899–6905.
- 593 [38] C. Barbieru, F. Pop, Soft real-time hadoop scheduler for big data processing
594 in smart cities, in: *Advanced Information Networking and Applications*
595 *(AINA)*, 2016 IEEE 30th International Conference on, IEEE, 2016, pp.
596 863–870.