

*Data Descriptor*

# Sigfox and LoRaWAN Datasets for Fingerprint Localization in Large Urban and Rural Areas

Michiel Aernouts<sup>1\*</sup>, Rafael Berkvens<sup>1</sup>, Koen Van Vlaenderen<sup>2</sup> and Maarten Weyn<sup>1</sup>

<sup>1</sup> University of Antwerp - imec, IDLab - Faculty of Applied Engineering, Groenenborgerlaan 171, 2020 Antwerp, Belgium

<sup>2</sup> Sensolus NV, Rijsenbergstraat 148, 9000 Ghent, Belgium

\* Correspondence: {michiel.aernouts; rafael.berkvens; maarten.weyn}@uantwerpen.be

**Abstract:** Because of the increasing relevance of the Internet of Things and location based services, researchers are evaluating wireless positioning techniques, such as fingerprinting, on LPWAN communication. In order to evaluate fingerprinting in large outdoor environments, extensive, time-consuming measurement campaigns need to be conducted to create useful datasets. This paper presents three LPWAN datasets which are collected in large-scale urban and rural areas. The goal is to provide the research community with a tool to evaluate fingerprinting algorithms in large outdoor environments. During a period of three months, numerous mobile devices periodically obtained location data via a GPS receiver which was transmitted via a Sigfox or LoRaWAN message. Together with network information, this location data is stored in the appropriate LPWAN dataset. The results of our fingerprinting implementation, which is also clarified in this paper, indicate a mean location estimation error of 214.58 m for the rural Sigfox dataset, 688.97 m for the urban Sigfox dataset and 398.40 m for the urban LoRaWAN dataset. In the future, we will enlarge our current datasets and use them to evaluate our fingerprinting methods. Also, we intend to collect additional datasets for Sigfox, LoRaWAN as well as NB-IoT.

**Keywords:** IoT; LPWAN; Sigfox; LoRaWAN; Localization; Fingerprinting.

**Data Set:** 10.5281/zenodo.1193563

**Data Set License:** CC-BY

## 1. Introduction

The Internet of Things (IoT) its growing importance creates a rapidly increasing necessity for wide area communication standards that guarantee reliable connectivity between a multitude of IoT devices. For this purpose, researchers have been developing various Low Power Wide Area Network (LPWAN) standards. IoT requires LPWAN standards to support long-range communication and high scalability of end-devices at a low cost. Also, ubiquitous indoor and outdoor connectivity as well as ultra-low power consumption are crucial aspects for reliable, transparent IoT applications that work for years on small batteries [1]. To meet these requirements, numerous measures have to be considered for LPWAN design, such as modulation techniques, network topology, hardware complexity, use of radio frequent spectrum and regulations. In general, a trade-off between these measures and data rate has to be made [2].

Context-awareness is an important aspect of IoT applications. This means that depending on the application, an IoT device can alter its behavior based on the measurements it has conducted in its environment. In order to create context-awareness for IoT applications, the location of the

32 device has to be obtained with minimal location error. Currently, Global Navigation Satellite Systems  
33 (GNSS) are the most commonly used method to do so. Although GNSS systems provide accurate  
34 location estimations, they have a few drawbacks which oppose some of the aforementioned IoT  
35 requirements. Firstly, GNSS receivers generally consume a lot of power, which limits the overall  
36 battery lifetime significantly. For instance, Global Positioning System (GPS) receivers consume 30 -  
37 50 mA while obtaining a GPS fix, which can take tens of seconds [3]. Furthermore, the GNSS location  
38 data will only be available on the device itself; forwarding the data through wireless communication  
39 involves additional power consumption. On the other hand, wireless positioning techniques can be  
40 applied to LPWAN communication messages without having to send additional messages which  
41 contain location information. Consequently, a device can be located without increasing power  
42 consumption. Secondly, GNSS systems tend to lose connectivity in indoor environments. Since  
43 many LPWAN standards operate in sub-GHz ISM bands, they can be used outdoor as well as  
44 indoor. Of course, GNSS remains a favorable solution for applications which require continuous high  
45 accuracy localization. However, many IoT use cases do not require such accurate location estimations  
46 and have more interest in long battery life-time. Therefore, wireless positioning based on LPWAN  
47 communication is an interesting alternative for long-term, low power localization.

48 Wireless positioning has been a prominent research topic for decades [4-6]. Many techniques  
49 which were developed over the years are still suitable for localization with modern wireless  
50 technologies. These techniques estimate the location of a transmitter or receiver by analyzing physical  
51 characteristics of the communication link such as Received Signal Strength (RSS), timing information,  
52 signal phase, etc. This paper presents datasets which can be used for a RSS localization method called  
53 fingerprinting. With this method, a training database of communication messages is built by storing  
54 their transmitter location as well as the Received Signal Strength Indicator (RSSI) for all receiving  
55 base stations. Afterwards, RSSI measurements of new messages are matched to the fingerprints in  
56 the training database to estimate the transmitter's location, e.g. by applying a  $k$ -Nearest-Neighbor  
57 ( $k$ NN) analysis [5]. In order to minimize the location estimation error, an extensive site survey has to  
58 be conducted to create a complete training database. Therefore, fingerprinting techniques are mostly  
59 used in constricted, indoor areas [7,8].

60 This paper presents three datasets. Firstly, we collected a Sigfox dataset in a large rural area  
61 between Antwerp and Ghent, Belgium. Secondly, an urban Sigfox dataset was built in the city center  
62 of Antwerp. Lastly, an urban LoRaWAN dataset was also built in the city center of Antwerp. In the  
63 near future, we also intend to create a rural LoRaWAN dataset in the large area between Antwerp  
64 and Ghent.

65 The remainder of the paper is structured as follows. Section 2 describes the LPWAN standards  
66 which were used to collect the dataset. Section 3 explains how the messages were collected and stored  
67 in a fingerprint database. Section 4 illustrates how the dataset can be used by the research community.  
68 In Section 5, the results of our analysis are listed. These results are then discussed in Section 6. Finally,  
69 Section 7 concludes the paper and states the intended future work.

## 70 2. Low Power Wide Area Networks

71 In recent years, numerous operators have been investing in the IoT by rolling out nation-wide  
72 LPWAN networks such as Sigfox [9], LoRaWAN [10] and NB-IoT [11]. In this section, these  
73 technologies and their constraints are discussed. Table 1 displays a summary of the aforementioned  
74 LPWAN standards.

**Table 1.** An overview of three main LPWAN standards

	<b>Sigfox</b>	<b>LoRaWAN</b>	<b>NB-IoT</b>
<b>Band</b>	EU: 868 MHz US: 902 MHz	EU: 433, 868 MHz US: 915 MHz	Cellular (LTE)
<b>Bandwidth</b>	100 Hz	125 - 500 kHz	180 kHz
<b>Modulation</b>	UL: UNB DBPSK DL: GFSK	LoRa	BPSK / QPSK
<b>Data rate</b>	UL: 100 bps DL: 600 bps	300 bps - 37.5 kbps	UL: 20 kbps DL: 250 bps
<b>Max. payload</b>	UL: 12 bytes DL: 8 bytes	250 bytes	125 bytes
<b>MAC</b>	Unslotted ALOHA	LoRaWAN	UL: SC-FDMA DL: OFDMA
<b>Topology</b>	Star	Star	Cellular network

### 75 2.1. Sigfox

76 Sigfox is a proprietary technology which operates on 868 MHz in Europe and on 902 MHz in  
77 the US. Because it uses an Ultra-Narrow Bandwidth (UNB) modulation technique which is called  
78 Differential Binary Phase Shift Keying (DBPSK), numerous devices can communicate over ranges up  
79 to 10 to 50 km with low power consumption, low-cost hardware setup and easy implementation [9].  
80 However, this modulation technique has a few downsides as well. Due to the narrow bandwidth, the  
81 maximum uplink throughput that can be achieved is only 100 bps. Also, regional regulations impose  
82 a limited duty cycle of 1%, i.e. 36 seconds per hour and six seconds per message [12]. Therefore, the  
83 daily limit for a Sigfox device equals 140 uplink messages (twelve bytes each) and four downlink  
84 messages (eight bytes each). Downlink messages are modulated with Gaussian Frequency Shift  
85 Keying (GFSK) and have to be requested by the end-device itself, direct communication between  
86 user and device is not possible. To improve the likelihood of successful reception, messages can be  
87 sent multiple times (three by default) on random frequency channels. Due to the Sigfox network  
88 architecture, messages from a single end-device are received by multiple base stations.

89 Since Sigfox uses UNB modulation, timing methods for localization are not a suitable option [13].  
90 Therefore, several other localization methods have been researched. Firstly, Sigfox Geolocation is  
91 built on probabilistic distance calculation using RSSI measurements. If a message is received by three  
92 or more base stations, the location of the transmitter can be estimated via triangulation. For a location  
93 accuracy below 500 m, Sigfox advises to use Sigfox Geolocation as a complementary method to other  
94 positioning systems such as WiFi localization or GPS. Secondly, Sallouha et al demonstrated RSSI  
95 fingerprinting as a feasible approach [14]. By providing their Sigfox nodes with GPS modules, the  
96 researchers were able to apply fingerprinting to subdivide a large region into classes. Within these  
97 classes, location estimation is improved by distance estimation between end-devices and GPS nodes  
98 within the same class. However, the radius of a class was limited to 200 m to minimize the location  
99 error, as errors of over 60 m were measured for larger radii. This approach would demand many GPS  
100 nodes, which impedes scalability. Lastly, Janssen et al implemented a WiFi fingerprinting method in  
101 a large urban area [15]. A mobile device was carried around in the city of Antwerp, Belgium while  
102 sniffing for WiFi BSSIDs. Every 10 minutes, the device simulated a Sigfox message containing the two  
103 BSSIDs with the highest RSSI, which were compared to open source WiFi BSSID databases to obtain  
104 a location estimate. This approach led to estimation errors of 23 m to 45 m. However, the Sigfox  
105 transmissions had to be dedicated for localization purposes, whereas other methods benefit from the  
106 physical characteristics of the Sigfox communication link.

## 107 2.2. LoRaWAN

108 Another widely used proprietary long range technology is LoRaWAN. Contrary to Sigfox,  
109 LoRaWAN applies a proprietary Chirp Spread Spectrum (CSS) modulation technique called LoRa,  
110 which stands for 'Long Range' [16]. In Europe, LoRa operates in the license free 433 and 868 MHz  
111 bands; in the US, 915 MHz is used. Symbols are encoded using a number of chirps, which causes  
112 the signal to be spread over a wider channel bandwidth. This technique reduces interference with  
113 other signals, protects the signal against jamming and cancels out multi-path and fading effects. The  
114 number of chirps is determined by the spreading factor, which can range from 7 to 12. By altering  
115 this parameter, a balance between data rate and range can be mediated. With a high spreading factor,  
116 long ranges can be achieved at the expense of a low 300 bit/s data rate. On the other hand, the data  
117 rate can be raised up to 37.5 kbit/s for shorter distances by using a low spreading factor. Because  
118 LoRaWAN supports multiple channels simultaneously, end-devices are able to exchange data for a  
119 longer period of time without violating duty cycle regulations. Subsequently, payloads of up to 250 B  
120 can be sent.

121 Compared to Sigfox, LoRa works with a higher bandwidth which enables localization with Time  
122 Difference Of Arrival (TDoA). However, very accurate synchronization between the receiving base  
123 stations is required for such a method. Distances between the transmitter and the base stations are  
124 estimated based on the time of flight of a signal, a location estimation is obtained via triangulation.  
125 Semtech has implemented a proprietary geolocation feature in LoRaWAN, which is based on TDoA.  
126 The LoRa Alliance claims that this feature achieves an estimation error of 20 to 200 m, depending on  
127 the conditions [17]. Fargas et al evaluated a TDoA method, and concluded that it is possible to obtain  
128 a location accuracy of around 100 m for static transmitters [3]. Further research has to be conducted  
129 to improve the location accuracy for real-time tracking of dynamic targets.

## 130 2.3. NB-IoT

131 In 2016, the 3rd Generation Partnership Project (3GPP) has introduced the Narrowband-IoT  
132 standard (NB-IoT). A main advantage of this new standard is its compatibility with traditional  
133 cellular networks. A 200 kHz GSM carrier can be re-purposed as an NB-IoT carrier. Also, NB-IoT can  
134 be implemented in the guardband of the LTE carrier or inside an LTE carrier by assigning a 180 kHz  
135 Physical Resource Block (PRB). Numerous properties such as MAC protocols are inherited from LTE.  
136 For instance, Orthogonal Frequency Division Multiple Access (OFDMA) is used for downlink and  
137 Single Carrier Frequency Division Multiple Access (SC-FDMA) is implemented for uplink. As a  
138 result, data rates are limited to 250 kbps for downlink, and 20 kbps for uplink [11]. Raza et al  
139 summarize a few flaws of NB-IoT [2]. Firstly, not all messages are acknowledged due to limited  
140 downlink capacity. Secondly, implementation of packet aggregation introduces additional latency  
141 which is undesirable for time-critical communication. Thirdly, the performance of NB-IoT may  
142 decline when the network is being heavily used for voice- and data traffic. Lastly, NB-IoT is a very  
143 recent LPWAN technology. Thus, the amount of commercial applications has been limited up to now.  
144 More research has to be conducted in order to obtain more knowledge about real world performance  
145 and battery life of and NB-IoT end-device.

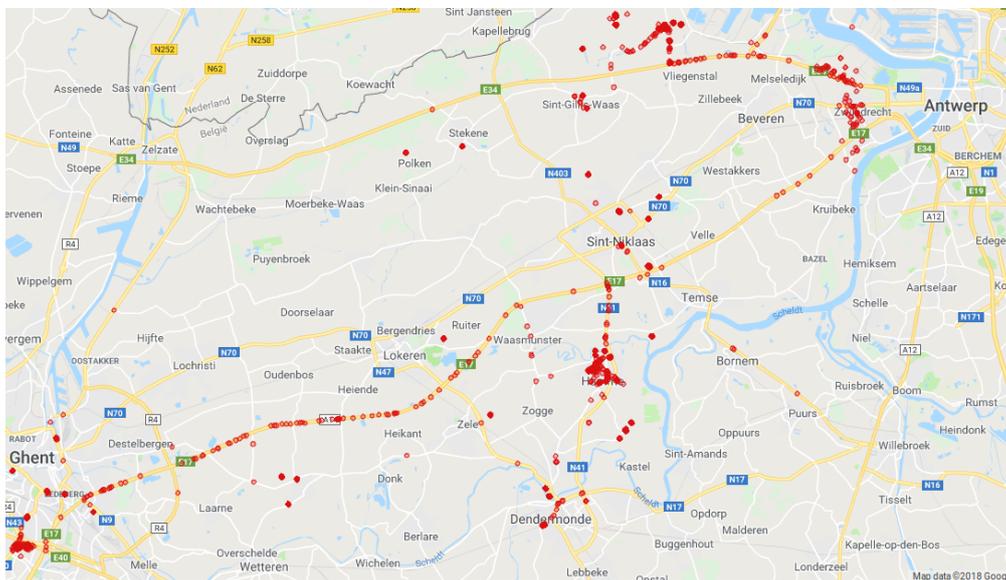
146 In order to locate an NB-IoT end device, Observed Time Difference of Arrival (OTDoA)  
147 localization could be used. A number of synchronized base stations transmit a Positioning Reference  
148 Signal (PRS), which is received by the end device. The end device forwards the TOA per transmitting  
149 base station to a localization server, where the difference between these TOA's and the PRS is used to  
150 calculate the estimated location of the end device [18].

151 In the remainder of this paper, we describe how our Sigfox and LoRaWAN datasets were  
152 collected and analyzed. Future work includes a similar measurement campaign and analysis for  
153 NB-IoT messages.

### 154 3. Collection Methodology

155 From 16 November 2017 until 5 February 2018, three LPWAN datasets were collected in two  
 156 different environments. Firstly, a Sigfox dataset was recorded in a large area between Antwerp and  
 157 Ghent ( $\pm 1068 \text{ km}^2$ ), which was delineated because of its overall rural characteristics. Secondly,  
 158 another Sigfox dataset was recorded in an urban area in and around the city center of Antwerp  
 159 ( $52.97 \text{ km}^2$ ). Lastly, a LoRaWAN dataset was collected in the same urban area. The messages in these  
 160 datasets contain GPS coordinates, which allows the user to correlate accurate location information  
 161 with RSSI measurements on that location. For the Sigfox datasets, a proprietary, nation-wide Sigfox  
 162 network which was deployed by Engiem2M was used. The LoRaWAN dataset was collected over a  
 163 proprietary, nation-wide network that was rolled out by Proximus. In this section, the collection and  
 164 structure of each dataset is discussed.

#### 165 3.1. Sigfox Rural

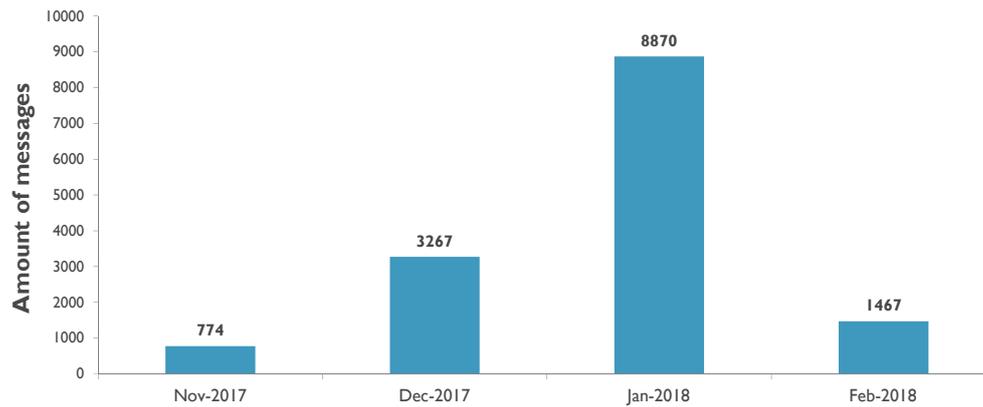


**Figure 1.** All messages in the rural Sigfox dataset were obtained in a large area between the cities of Antwerp and Ghent. Transmission locations are indicated by red dots on the map.

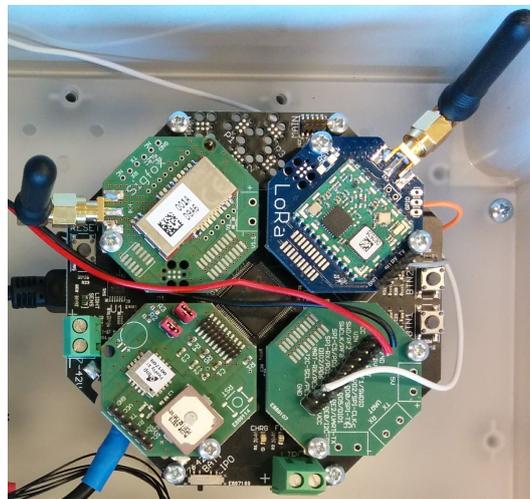
166 As shown in Figure 1, the first dataset was collected in a large rural area between the cities  
 167 of Antwerp and Ghent. This dataset contains 25 638 Sigfox messages which were received between  
 168 16 November 2017 and 5 February 2018, as indicated by the temporal spread graph in Figure 2. In  
 169 order to collect these messages, ten *Sensolus Stickntrack* devices [19], like the one in Figure 3, were  
 170 carried by people who commute by car between the cities of Antwerp and Ghent on a daily basis.  
 171 These devices contain a Ublox CAM-M8C GPS receiver, which obtains a GPS fix every ten minutes.  
 172 A Sigfox message with GPS coordinates is sent directly after obtaining this GPS fix. Due to the  
 173 fact that the total transmission time of a Sigfox message can take up to six seconds, the received  
 174 coordinates of the transmitting device could differ from the actual device coordinates at receiving  
 175 time. Therefore, the correlation between RSSI measurements and the received GPS coordinates could  
 176 hold an additional GPS location estimation error of tens of meters, depending on the ground speed of  
 177 the transmitting device. Additional information such as the ground speed of the device, or GPS signal  
 178 quality measurements such as the Horizontal Dilution Of Precision (HDOP) could not be included in  
 179 the message because of the limited Sigfox payload size.







**Figure 5.** The urban Sigfox dataset contains 14 378 messages which were obtained from 16 November 2017 to 5 February 2018. This figure illustrates the temporal spread of this dataset.



**Figure 6.** OCTA-Connect hardware modules are used to send GPS data to the local data server [21].  
Bottom left: Firefly X1 GPS receiver. Top left: TD1207R RF module for Sigfox communication.  
Top right: IM880B-L RF module for LoRaWAN communication

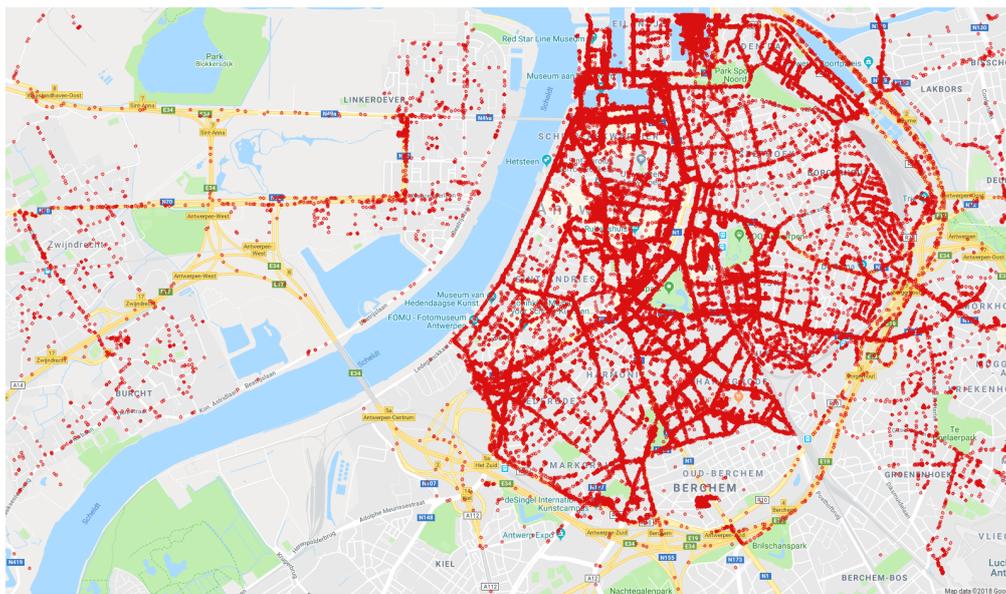
213 From this point forward, the Sigfox messages are processed and stored with the same  
 214 methodology that is used for the rural Sigfox dataset. The Sigfox backend forwards JSON strings  
 215 for all duplicates of a message to a local data server. These strings, which contain GPS coordinates,  
 216 base station ID, base station RSSI, receiving time and a message ID are then grouped to form a new  
 217 database of unique Sigfox messages. By querying this database, a list of all messages within Antwerp  
 218 is extracted.

219 As shown in Table 3, the urban Sigfox dataset is structured in the same way as the rural Sigfox  
 220 dataset. All rows represent one of the 14 378 messages in the urban dataset, the columns indicate  
 221 which of the 84 base stations in the urban area have received the message. If a base station has not  
 222 received the message, an RSSI of  $-200$  dB is inserted in the cell. The receiving time, latitude and  
 223 longitude can be seen in the last three columns of the dataset.

**Table 3.** The structure of the urban Sigfox dataset. Every row represents a Sigfox message and indicates its receiving base stations (BS), the receiving time of the message (RX Time) and the GPS coordinates at transmission time.

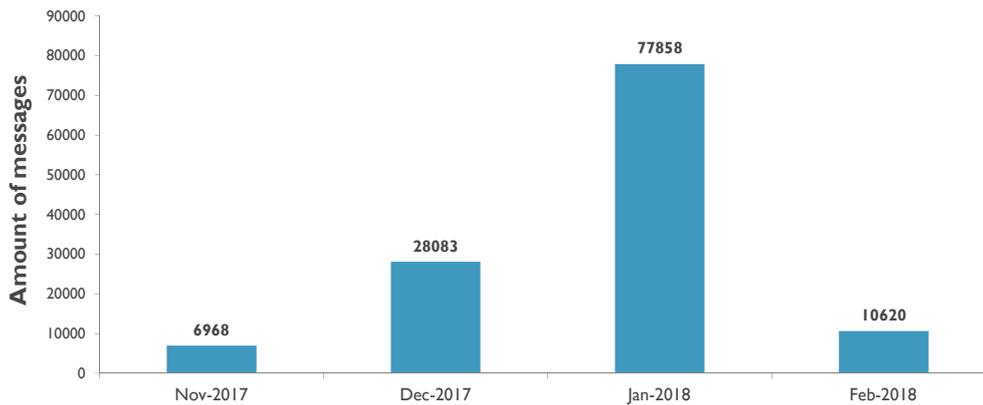
BS 1	BS 2	...	BS 84	RX Time	Latitude	Longitude
-122	-105	...	-200	2018-01-24T09:20:09+00:00	51.2111778259	4.4175133705
-114	-119	...	-200	2018-01-09T13:11:01+00:00	51.2206916809	4.3969035149
...	...	...	...	...	...	...

### 224 3.3. LoRaWAN Antwerp



**Figure 7.** The urban LoRaWAN dataset was collected in the city center of Antwerp.

225 Lastly, a large dataset of LoRaWAN messages was obtained in the city center of Antwerp. It  
 226 holds 123 529 messages which are collected from 17 November 2017 until 5 February 2018, the spatial  
 227 and temporal spread of the dataset can be found in Figures 7 and 8 respectively. The City of Things  
 228 hardware that is used to create the urban Sigfox dataset (Figure 6) is also used to collect the urban  
 229 LoRaWAN dataset: twenty cars of Antwerp's postal service drove around in the city center while  
 230 carrying this hardware; the Firefly X1 GPS receiver continuously acquired the current latitude and  
 231 longitude of the car, as well as the HDOP of the GPS signal. This location information is sent in a  
 232 LoRaWAN message via the IM880B-L radio module. Compared to Sigfox, LoRaWAN implements  
 233 a much wider bandwidth and higher data rate. Therefore, the radio modules were able to send a  
 234 new message every minute. Also, additional information such as the HDOP could be included in the  
 235 message because of the larger maximum payload. Concisely, more data could be sent more frequently,  
 236 which is why the urban LoRaWAN dataset is significantly larger than the Sigfox datasets. Including  
 237 the HDOP in the message allows the user of this dataset to omit messages with poor GPS signal  
 238 quality. However, the GPS coordinates at receiving time could still differ from the GPS coordinates at  
 239 transmission time, especially if the transmitter is moving at high speeds.



**Figure 8.** The urban LoRaWAN dataset holds 123 529 messages that were collected from 17 November 2017 to 5 February 2018. This figure illustrates the temporal spread of this dataset.

240 On the LoRaWAN backend server, a callback function is configured to forward the payload of  
 241 each message, with additional network information attached, to the local data server. Table 4 shows  
 242 how this information was stored in the urban LoRaWAN dataset. Similar to the Sigfox datasets, the  
 243 first columns indicate the receiving base stations and their respective RSSI. In the urban area that was  
 244 explored for this dataset, 64 LoRaWAN base stations are detected. In the next columns, the receiving  
 245 time, spreading factor, HDOP, latitude and longitude are stored.

**Table 4.** The structure of the urban LoRaWAN dataset. Every row represents a LoRaWAN message and indicates its receiving base stations (BS), the receiving time of the message (RX Time), the LoRa spreading factor, and the HDOP, latitude and longitude of the transmitter at transmission time.

BS 1	BS 2	...	BS 68	RX Time	SF	HDOP	Latitude	Longitude
-101	-95	...	-200	2018-01-09T23:42:19+00:00	9	0.60	51.19404602	4.41862487
-200	-111	...	-200	2018-01-31T10:01:27+00:00	12	1.08	51.20004272	4.4116702
...	...	...	...	...	...	...	...	...

#### 246 4. Analysis

247 This section describes how the aforementioned LPWAN datasets can be used for basic  
 248 fingerprint-based localization. Before localization, a dataset should be split up in 3 subsets: a training  
 249 set, an evaluation set and a test set. In order to ensure an unbiased spatial spread across the subsets,  
 250 the main dataset should be divided randomly. The training set, which should be the largest of the  
 251 three subsets, serves as a reference database for the fingerprinting algorithm. Afterwards, these  
 252 subsets can be used for the localization algorithm as follows. Firstly, a distance matrix between the  
 253 training and evaluation sets has to be computed by calculating the Euclidean distance  $d$  between  
 254 every fingerprint in the evaluation set and every fingerprint in the training set, as illustrated in  
 255 Equation 1. In this equation, the RSSI measurements for the evaluation fingerprints and the training  
 256 fingerprints are represented by  $RSSI_{eval}$  and  $RSSI_{training}$  respectively.

$$d = \sqrt{\sum (RSSI_{eval} - RSSI_{training})^2} \quad (1)$$

257 Every distance  $d$  is stored in a matrix with  $m$  rows and  $n$  columns, with  $m$  matching the amount  
 258 of samples in the evaluation set and  $n$  the amount of samples in the training set. Secondly, this  
 259 distance matrix is used to find the  $k$  Nearest Neighbors ( $k$ NN) of an evaluation fingerprint: for every  
 260 row in the matrix, the column indexes of the  $k$  smallest distances  $d$  are extracted and compared to the  
 261 real world locations of the corresponding row indexes of the training set. These rows in the training  
 262 set are the training fingerprints which correspond best to the evaluation sample. The centroid of

263 the locations of these  $k$  training fingerprints can be used as the location estimate for the evaluation  
 264 sample. By comparing this estimate to the actual GPS coordinates of the evaluation sample, the  
 265 location estimation error can be quantified. This location estimation process is repeated for every row  
 266 in the distance matrix with different values for  $k$ , to calculate the mean estimation error for a given  $k$   
 267 and subsequently, to find an optimal value for  $k$ . Lastly, the test subset is used to validate this optimal  
 268  $k$  with fingerprints that do not occur in the training or evaluation subset. A new Euclidean distance  
 269 matrix between the training subset and the test subset is created and used to calculate the mean  
 270 location estimation error for the best  $k$ . In the list below, the main steps of this basic fingerprinting  
 271 algorithm are summarized, whereas Figure 9 displays a concise visual representation.

- 272 **Step 1:** Split the dataset in subsets. E.g. 70% training, 15% evaluation and 15% test samples.  
 273 **Step 2:** Calculate a Euclidean distance matrix between the training set and the evaluation set.  
 274 **Step 3:** In the distance matrix, find the  $k$  nearest neighbors of every evaluation sample. Use the  
 275 centroid of the  $k$  nearest training fingerprints as the location estimate. Repeat this step for  
 276 all  $k$ 's you want to evaluate.  
 277 **Step 4:** Calculate the mean estimation error for every  $k$  of step 3. We consider the  $k$  with the smallest  
 278 mean estimation error to be the optimal parameter.  
 279 **Step 5:** Calculate a Euclidean distance matrix between the training set and the test set.  
 280 **Step 6:** In the new distance matrix, find the  $k$  nearest neighbors of every training sample. For  $k$ , use  
 281 the optimal parameter that was obtained in step 4. Use the centroid of the  $k$  nearest training  
 282 fingerprints as the location estimate.  
 283 **Step 7:** Validate the optimal value for  $k$  by calculating the estimation errors.

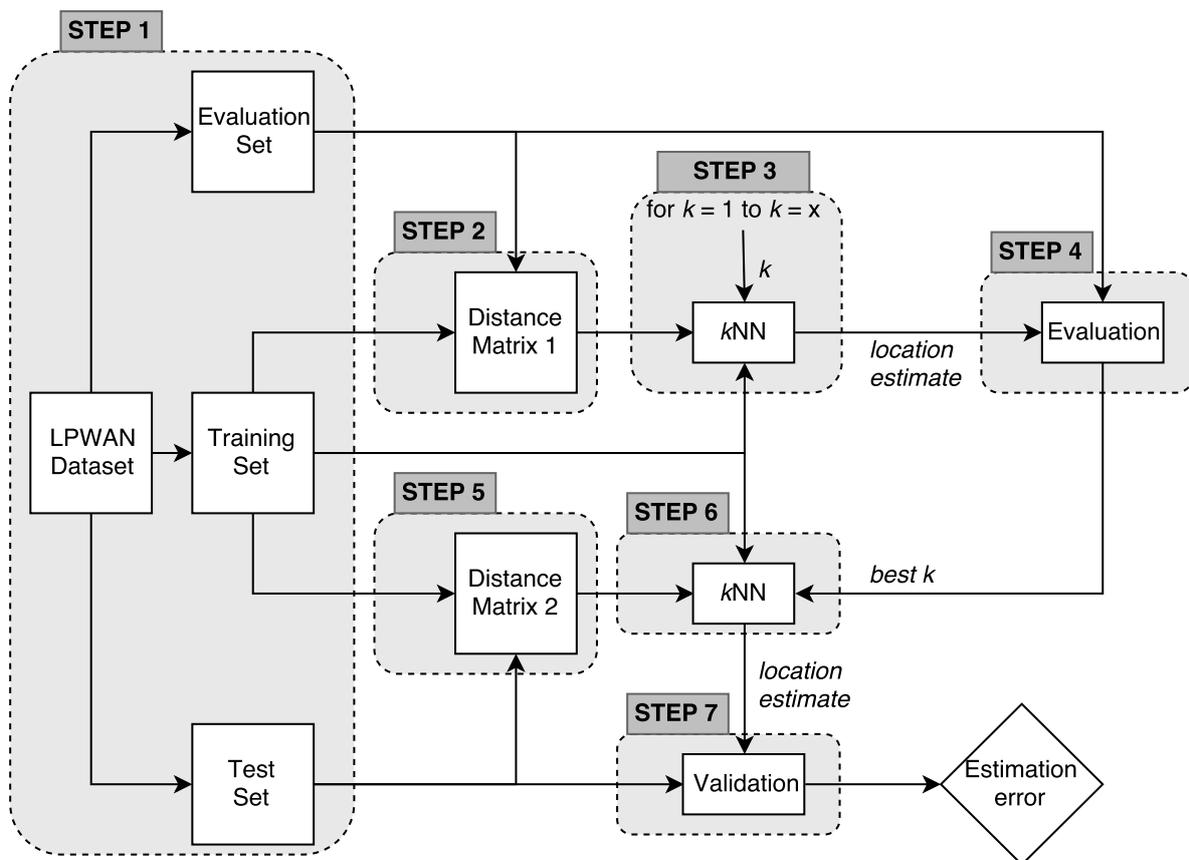
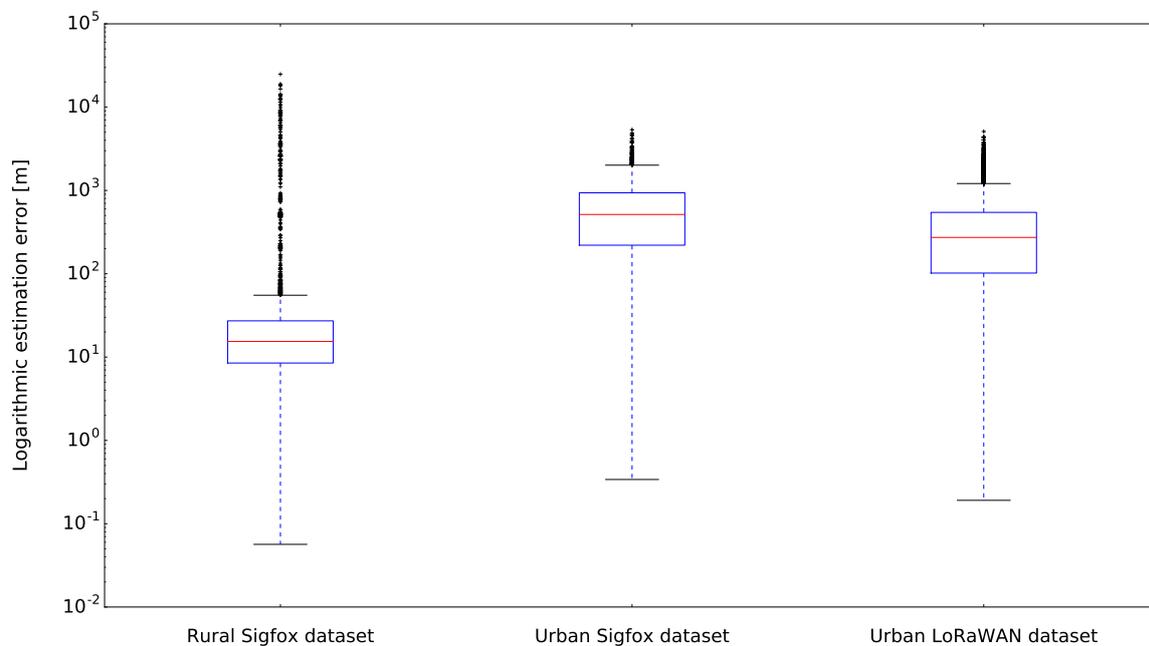


Figure 9. Visual representation of a basic fingerprinting algorithm.

## 284 5. Results

285 To evaluate the usability of our LPWAN datasets, we have executed the fingerprinting algorithm  
 286 that was described in Section 4 on all three datasets. Every dataset was randomly split up in a training  
 287 subset, an evaluation subset and a test subset; the sizes of these subsets are 70%, 15% and 15% of  
 288 the total size of the complete set. The evaluation subsets were used to find the optimal parameters  
 289 for each LPWAN dataset. A parameter sweep was conducted, varying  $k$  from 1 to 15. We decided  
 290 upon an optimal value for  $k$  per dataset by calculating which  $k$  resulted in the lowest mean location  
 291 estimation error. This optimal  $k$  was then validated with the test subset, as described in steps 6 and 7  
 292 of algorithm 9. The results of this validation are shown in the logarithmic box plots of Figure 10,  
 293 as well as in Table 5.

294 The rural Sigfox dataset has a mean error of 214.58 m and a median error of 15.4 m when looking  
 295 for the single nearest neighbor, whereas the urban Sigfox dataset has a mean error of 688.97 m and a  
 296 median error of 514.83 m when  $k$  equals 10. Lastly, the urban LoRaWAN dataset has a mean error of  
 297 398.4 m and a median error of 273.03 m when using the 11 nearest neighbors to get a location estimate.  
 298 We will discuss these results in the next section.



**Figure 10.** These box plots display the logarithmic estimation error of our first fingerprinting implementation on the LPWAN datasets.

**Table 5.** Fingerprinting results for all LPWAN datasets. The best value for  $k$  was determined by executing a parameter sweep during the evaluation phase.

	Best $k$	Mean error [m]	Median error [m]
<b>Sigfox Rural</b>	1	214.58	15.4
<b>Sigfox Antwerp</b>	10	688.97	514.83
<b>LoRaWAN Antwerp</b>	11	398.40	273.03

## 299 6. Discussion

300 In Figure 10 of the previous section, the estimation errors for the three LPWAN datasets are  
 301 shown. Table 5, lists the optimal value for  $k$  as well as the corresponding mean and median location  
 302 estimation error. When observing Figure 10, it becomes clear that the fingerprinting method has

303 the best results on the rural Sigfox dataset. This is mainly due to the spatial spread of this dataset,  
304 which can be seen in Figure 1. Measurements were conducted by people who commute by car in the  
305 rural area between Antwerp and Ghent, but a major part of these measurements were taken while  
306 the cars were parked, usually even on the same parking spot. Therefore, the rural Sigfox dataset  
307 mainly consists of dense message clusters on several small parking lots which are far apart from each  
308 other. This explains why we obtain the best results when we estimate a location based on the single  
309 nearest neighbor: there is a high probability that the nearest neighbor of a test sample is located in the  
310 same message cluster. If a test sample is not located in such a cluster, the location estimation based  
311 on a single nearest neighbor can have an error of almost 25 km, as shown in the outliers of the box  
312 plot. Concisely, the fingerprinting results for the rural Sigfox dataset are only highly accurate if a test  
313 sample is located in one of the message clusters.

314 The second box plot depicts the estimation error of the urban Sigfox dataset. As shown  
315 in Figure 4, this dataset has a more equal spatial density compared to the rural Sigfox dataset.  
316 Consequently, location estimations will be based on a higher number of nearest neighbors: we  
317 evaluated that the optimal  $k$  for this dataset equals 10.

318 Lastly, the estimation error of the urban LoRaWAN dataset can be seen in the third box plot of  
319 Figure 10. For this set, we found that location estimates based on 11 nearest neighbors resulted in  
320 the smallest estimation error. The mean and median estimation error of the LoRaWAN dataset are  
321 significantly lower than those of the urban Sigfox dataset, which is mainly a consequence of the large  
322 size of the LoRaWAN dataset (123529 messages). In the LoRaWAN set, it is more likely that the 11  
323 nearest neighbors of a test sample are close to the actual location of the test sample, as the dataset  
324 holds more messages in the same urban area. This can also be empirically evaluated by comparing  
325 Figures 4 and 7, which display the spatial spread of the urban Sigfox dataset and the LoRaWAN  
326 dataset respectively. Hence, we expect to decrease the location estimation error of the urban Sigfox  
327 dataset by adding more messages to this set.

## 328 7. Conclusion

329 This paper has described the collection methodology of three large LPWAN datasets, as well  
330 as their suitability for fingerprint-based localization. With these datasets, we intend to provide the  
331 global research community with a benchmark tool to evaluate fingerprinting algorithms for LPWAN  
332 standards. In the next months, we will keep collecting LPWAN messages in the same areas, by  
333 implementing the methodology that was described in Section 3. Apart from increasing the size of  
334 the datasets, this allows us to analyze short-term and long-term fluctuations in RSS measurements,  
335 and research how we can cope with them. We also aim to improve the spatial spread of the rural  
336 Sigfox dataset by deploying more Sigfox hardware in the rural area. Furthermore, three additional  
337 datasets will be built: a rural LoRaWAN dataset, a rural NB-IoT dataset and an urban NB-IoT dataset.

338 In Section 5, we have demonstrated the results of a basic  $k$ NN fingerprinting method which  
339 was explained in Section 4. Similar to previous research on indoor WiFi fingerprinting [22], we will  
340 analyze the effect of other distance functions, data representations and thresholding techniques using  
341 our LPWAN datasets as input.

342 **Acknowledgments:** Part of this work was funded by the MuSCLe-IoT (Multimodal Sub-Gigahertz  
343 Communication and Localization for Low-power IoT applications) project, co-funded by imec, a research  
344 institute founded by the Flemish Government, with project support from VLAIO (contract number  
345 HBC.2016.0660). The MuSCLe-IoT industry partners are Flash, Sensolus, Engie M2M, and Aertssen. Part of  
346 this research was funded by the Flemish FWO SBO S004017N IDEAL-IoT (Intelligent DENSE And Long range  
347 IoT networks) project. Part of the equipment used in this work was funded by the Flemish Hercules program.

348 **Author Contributions:** M.W., R.B. and M.A. conceived and designed the experiments; M.A. performed the  
349 experiments; M.A. and R.B. analyzed the data; K.V.V. contributed hardware to transmit the LPWAN messages;  
350 M.A. and R.B. wrote the paper.

351 **Conflicts of Interest:** The authors declare no conflict of interest.

352 **References**

- 353 1. Margelis, G.; Piechocki, R.; Kaleshi, D.; Thomas, P. Low Throughput Networks for the IoT: Lessons  
354 learned from industrial implementations. *IEEE World Forum on Internet of Things, WF-IoT 2015 -*  
355 *Proceedings*, 2015, pp. 181–186.
- 356 2. Raza, U.; Kulkarni, P.; Sooriyabandara, M. Low Power Wide Area Networks: An Overview. *IEEE*  
357 *Communications Surveys and Tutorials* **2017**, *19*, 855–873.
- 358 3. Fargas, B.C.; Petersen, M.N. GPS-free geolocation using LoRa in low-power WANs. 2017 Global Internet  
359 of Things Summit (GIoTS). IEEE, 2017, pp. 1–6.
- 360 4. Bensky, A. *Wireless positioning technologies and applications*, 2 ed.; Artech House, 2008; p. 305.
- 361 5. Liu, H.; Darabi, H.; Banerjee, P.; Liu, J. Survey of wireless indoor positioning techniques and systems.  
362 *IEEE Transactions on Systems, Man and Cybernetics Part C: Applications and Reviews* **2007**, *37*, 1067–1080.
- 363 6. Mao, G.; Fidan, B.; Anderson, B.D. Wireless sensor network localization techniques. *Computer Networks*  
364 **2007**, *51*, 2529–2553.
- 365 7. He, S.; Chan, S.H. Wi-Fi fingerprint-based indoor positioning: Recent advances and comparisons. *IEEE*  
366 *Communications Surveys and Tutorials* **2016**, *18*, 466–490.
- 367 8. Torres-Sospedra, J.; Montoliu, R.; Martinez-Uso, A.; Avariento, J.P.; Arnau, T.J.; Benedito-Bordonau, M.;  
368 Huerta, J. UJIIndoorLoc: A new multi-building and multi-floor database for WLAN fingerprint-based  
369 indoor localization problems. IPIN 2014 - 2014 International Conference on Indoor Positioning and Indoor  
370 Navigation, 2014, pp. 261–270.
- 371 9. Sigfox - The Global Communications Service Provider for the Internet of Things (IoT). Available online:  
372 <https://www.sigfox.com/en> (accessed on 7 March 2018).
- 373 10. Augustin, A.; Yi, J.; Clausen, T.; Townsley, W. A Study of LoRa: Long Range and Low Power Networks  
374 for the Internet of Things. *Sensors* **2016**, *16*, 1466.
- 375 11. Wang, Y.P.; Lin, X.; Adhikary, A.; Grövlén, A.; Sui, Y.; Blankenship, Y.; Bergman, J.; Razaghi, H.S. A Primer  
376 on 3GPP Narrowband Internet of Things. *IEEE Communications Magazine* **2017**, *55*, 117–123.
- 377 12. Vejlgård, B.; Lauridsen, M.; Nguyen, H.; Kovacs, I.Z.; Mogensen, P.; Sorensen, M. Coverage and Capacity  
378 Analysis of Sigfox, LoRa, GPRS, and NB-IoT. 2017 IEEE 85th Vehicular Technology Conference (VTC  
379 Spring). IEEE, 2017, pp. 1–5.
- 380 13. Gezici, S.; Zhi Tian.; Giannakis, G.; Kobayashi, H.; Molisch, A.; Poor, H.; Sahinoglu, Z. Localization via  
381 ultra-wideband radios: a look at positioning aspects for future sensor networks. *IEEE Signal Processing*  
382 *Magazine* **2005**, *22*, 70–84.
- 383 14. Sallouha, H.; Chiumento, A.; Pollin, S. Localization in long-range ultra narrow band IoT networks using  
384 RSSI. 2017 IEEE International Conference on Communications (ICC). IEEE, 2017, pp. 1–6.
- 385 15. Janssen, T.; Weyn, M.; Berkvens, R. Localization in Low Power Wide Area Networks Using Wi-Fi  
386 Fingerprints. *Applied Sciences* **2017**, *7*, 936.
- 387 16. LoRa Modulation Basics. Available online: [https://www.semtech.com/uploads/documents/an1200.22.](https://www.semtech.com/uploads/documents/an1200.22.pdf)  
388 [pdf](https://www.semtech.com/uploads/documents/an1200.22.pdf) (accessed on 7 March 2018).
- 389 17. LoRaWAN Geolocation Whitepaper. Available online: [https://docs.wixstatic.com/ugd/eccc1a\\_](https://docs.wixstatic.com/ugd/eccc1a_d43b3b29dfff4ec2b00f349ced4225c4.pdf)  
390 [d43b3b29dfff4ec2b00f349ced4225c4.pdf](https://docs.wixstatic.com/ugd/eccc1a_d43b3b29dfff4ec2b00f349ced4225c4.pdf) (accessed on 7 March 2018).
- 391 18. Hu, S.; Berg, A.; Li, X.; Rusek, F. Improving the Performance of OTDOA Based Positioning in NB-IoT  
392 Systems. GLOBECOM 2017 - 2017 IEEE Global Communications Conference. IEEE, 2017, pp. 1–7.
- 393 19. Sensolus Stickntrack. Available online: <http://www.sensolus.com/stickntrack-gps/> (accessed on 7  
394 March 2018).
- 395 20. Latre, S.; Leroux, P.; Coenen, T.; Braem, B.; Ballon, P.; Demeester, P. City of things: An integrated and  
396 multi-technology testbed for IoT smart city experiments. 2016 IEEE International Smart Cities Conference  
397 (ISC2). IEEE, 2016, pp. 1–8.
- 398 21. OCTA-Connect. Available online: <http://www.octa-connect.com/>.
- 399 22. Torres-Sospedra, J.; Montoliu, R.; Trilles, S.; Belmonte, O.; Huerta, J. Comprehensive analysis of distance  
400 and similarity measures for Wi-Fi fingerprinting indoor positioning systems. *Expert Systems with*  
401 *Applications* **2015**, *42*, 9263–9278.