



Data Descriptor

Classification of actual sensor network deployments in research studies from 2013 to 2017

Janis Judvaitis ^{1,‡}, Artis Mednis ¹, Valters Abolins ¹, Ansis Skadins ¹, Didzis Lapsa ¹, Raimonds Rava ¹, Maksims Ivanovs ¹ and Krisjanis Nesenbergs ^{1,‡,*} 

¹ Institute of Electronics and Computer Science, 14 Dzerbenes St., LV-1006, Riga, Latvia

* Correspondence: krisjanis.nesenbergs@edi.lv

‡ These authors contributed equally to this work.

Version August 20, 2020 submitted to Data

Abstract: Due to breakthroughs in embedded system development, sensing technologies, and ubiquitous connectivity in recent years, technologies such as Wireless Sensor Networks (WSN) and Internet of Things (IoT) have captured the imagination of researchers, businesses, and general public. That resulted in the emergence of an enormous, difficult-to-navigate body of work related to WSN and IoT. To highlight trends and developments in these technologies and to see whether they are actually deployed rather than subjects of theoretical research with presumed potential use cases, we gathered and codified a dataset of scientific publications from a 5-year period from 2013 to 2017 involving actual sensor network deployments. In the first iteration, 15010 potentially relevant articles were identified in SCOPUS and Web of Science databases; after two iterations, 3059 actual sensor network deployments were extracted from those articles and classified in a consistent way according to different categories such as type of nodes, field of application, communication types, etc. We publish the resulting dataset with the intent that its further analysis may identify prospective research fields and future trends in WSN and IoT.

Dataset: http://git.edi.lv/CPS_Lab/sn_deployment_mapping_review

Dataset License: CC0

Keywords: WSN; IoT; Deployment; Review;

1. Summary

As we are heading into the 21st century, digitalization trends in transportation [?] [?], in-house logistics [?], education [?] [?], agriculture [?], banking [?] [?] and other fields are providing new and engaging ways for technology to improve our daily lives. This naturally leads to the emergence of applications of Wireless Sensor Networks (WSN) and Internet of Things (IoT) in large number of different domains. The WSN and IoT popularity is growing rapidly and, according to Grand View Research, the Narrow Band IoT(NB-IoT) market size will reach more than \$6 billion by 2025 [?]. Yet the majority of researchers still use simulation tools to validate their theories [?] rather than deploy actual devices; as a consequence, it is unclear to what extent the vast majority of the available WSN/IoT devices are actually used instead of theorized as being applicable and what design choices drive the selection of devices.

The aim of this work was to gather the available information about actual WSN and IoT deployments used by the research community and present an overview about trends from the 5-year period from 2013 to 2017. The presented dataset can be further used for various statistical and contextual analysis as well as further extended to cover a broader time frame. As the complete marked

32 data set is available together with intermediate collection results, the authenticity of the data can be
33 verified.

34 Altogether 15010 data articles were identified as potential candidates, from which after two
35 iterations of screening 3059 actual sensor network deployments were extracted and codified according
36 to multiple categories as described in the next sections.

37 The data acquisition, analysis and validation took around two years for a team of 12 volunteer
38 researchers of which 8 provided significant value.

39 The rest of the document is structured as follows - Section ?? describes the data set as such, Section
40 ?? discusses methods used in acquiring the data as well as data validation and quality, and finally
41 Section ?? contains some practical data usage notes.

42 2. Data description

43 The dataset contains data files resulting from the data acquisition process as shown in Table ??
44 and described below in detail. The files are in one of three formats:

- 45 • **.bib** - BibTeX format containing entries representing published articles;
- 46 • **.json** - JSON format containing structured human readable data object entries;
- 47 • **.txt** - Text files containing TAB delimited tabular data with a header row.

Table 1. Data files in the dataset.

ID	File name	Number of data entries	Description
(A)	0_Merged_Full_Collection_15010.bib	15010	Identified candidate articles from database keyword search as a BibTeX file
(B)	1_Screened_4915.bib	4915	Candidate articles left after screening as a BibTeX file
(C)	1_Screening_statistics.txt	12	Screening statistics per person as Tab delimited text file
(D)	1_Screening_timeline.txt	17	Screening timeline per week of screening as Tab delimited text file
(E)	2_Eligibility_statistics.txt	12	Eligibility check statistics per person as Tab delimited text file
(F)	2_Eligible_3017.json	3017	Candidate articles left after eligibility check as a JSON file
(G)	2_Ineligible_1898.json	1898	Articles excluded in the eligibility check and reason for exclusion as a JSON file
(H)	2_Mistaken_as_ineligible_47.json	47	Articles mistakenly excluded in the eligibility check and re-included during validation as a JSON file
(I)	3_Eligibility_and_extraction_timeline.txt	35	Eligibility check and data extraction weekly timeline as Tab delimited text file
(J)	3_Extracted_data_3059.json	3059	Main dataset - Deployments identified and codified data extracted as a JSON file
(K)	3_Extraction_statistics.txt	12	Data extraction statistics per person as Tab delimited text file
(L)	3_Mistaken_as_eligible_15.json	15	Articles mistakenly included in the eligibility step and excluded during data extraction as a JSON file
(M)	Timeline.txt	14	Overall timeline of dataset building process as a Tab delimited text file
(N)	README.md	-	Readme file with short description of the dataset
(O)	Notebooks	-	Folder containing Jupyter notebook files for easier data visualisation with processing examples

48 In the subsections below the technical description of data entries with possible data types and
49 values are described in detail. Verbatim data values in this description will be formatted like `this`.

50 For the eager reader interested in the main resulting dataset, please refer to dataset (J) on page ??.

51 2.1. (A) - identified candidate articles

52 This file contains 15010 BibTeX entries, which have the following types: `@article` (7137), `@book`
53 (74), `@conference` (1861), `@incollection` (67) and `@inproceedings` (5871).

54 Each entry in the file starts on a new line, and can continue over multiple lines. The basic structure
55 of an entry is `@type{id,metadata}` where `type` is name of the document type e.g. `article` or `book`,
56 `id` is a unique string in the document identifying that specific entry and `metadata` is a list of comma

57 separated key/value pairs describing the entry. Not all entries contain the same metadata entries, but
58 most have the following: `abstract`, `author`, `doi`, `title` and `year`. Also depending on the entry
59 type additional metadata like `page`, `volume` or `url` could be present.

60 2.2. (B) - screened candidate articles

61 This file contains 4915 BibTeX entries of the same format as described in previous section. These
62 entries represent candidate articles left from the (A) dataset after first step of screening and the file
63 contains the following entry types: `@article` (2385, 33% left after screening), `@book` (12, 16% left),
64 `@conference` (569, 31% left), `@incollection` (12, 18% left) and `@inproceedings` (1937, 33% left).

65 2.3. (C) - screening statistics

66 This file is formatted as a table in a TAB delimited text file. It has 12 entries, each pertaining to
67 one of the 12 volunteer researchers involved in the screening process.

68 Each row has the following 3 headers/columns with corresponding data types:

- 69 1. `Screener` - two letter code uniquely identifying each of the researchers. Example of data in
70 column: KN;
- 71 2. `Articles_screened` - number of articles processed by the corresponding researcher in the
72 screening step. This is an integer value in range from 0 to 5473;
- 73 3. `Percentage_screened` - the percentage of the total number of articles in dataset (A) that were
74 processed by the researcher in the screening step. This number is formed as percentage value
75 rounded to 2 decimal places and has values from 0.00% to 36.46%.

76 2.4. (D) - screening timeline

77 The file is formatted as a table in a TAB delimited text file. It has 17 entries, each representing one
78 of the 17 weeks during which the screening process took place.

79 Each row has the following 3 headers/columns with corresponding data types:

- 80 1. `Week` - number of the week in screening process, represented by an integer value in range from
81 1 to 17;
- 82 2. `Screened_per_week` - number of articles processed during the specific screening week by all
83 researchers involved. This is an integer value in range from 50 to 2068;
- 84 3. `Total_screened` - cumulative number of articles processed up to and including the specific
85 screening week by all researchers involved. This is an integer value in range from 1046 to 15010.

86 2.5. (E) - eligibility statistics

87 This file is formatted as a table in a TAB delimited text file. It has 12 entries, each pertaining to
88 one of the 12 volunteer researchers involved in the eligibility checking process.

89 Each row has the following 8 headers/columns with corresponding data types:

- 90 1. `Tested_by` - two letter code uniquely identifying each of the researchers. Example of data in
91 column: KN;
- 92 2. `Marked_eligible` - number of articles processed and marked as eligible by the corresponding
93 researcher in the eligibility checking step. This is an integer value in range from 1 to 708;
- 94 3. `Marked_Ineligible` - number of articles processed and marked as ineligible by the
95 corresponding researcher in the eligibility checking step. This is an integer value in range
96 from 0 to 475;
- 97 4. `Eligibility_percentage` - the percentage of the total number of articles checked by the
98 researcher in eligibility checking step that were marked as eligible. This number is formed
99 as percentage value rounded to 2 decimal places and has values from 54.66% to 100.00%.

- 100 5. `Mistaken_as_ineligible` - number of articles mistakenly marked as ineligible by the
 101 corresponding researcher in the eligibility checking step. This is an integer value in range
 102 from 0 to 20;
- 103 6. `Error_rate` - the percentage of the total number of articles checked by the researcher in eligibility
 104 checking step that were mistakenly marked as ineligible. This number is formed as percentage
 105 value rounded to 2 decimal places and has values from 0.00% to 33.33%. Additionally one value
 106 is NaN or "not a number" representing value resulting from division by zero;
- 107 7. `Total_processed` - number of articles processed by the corresponding researcher in the eligibility
 108 checking step. This is an integer value in range from 1 to 1183;
- 109 8. `Percentage_processed` - the percentage of the total number of articles in dataset (B) that were
 110 processed by the researcher in the eligibility checking step. This number is formed as percentage
 111 value rounded to 2 decimal places and has values from 0.02% to 24.07%.

112 2.6. (F) - candidate articles marked as eligible

113 This file contains a JSON data object with 3017 entries, each representing a single article marked
 114 as eligible in the eligibility checking step.

115 The object is structured as follows: `{entry1, entry2, ..., entry3017}`. Each of the entries
 116 have the following structure: `id:{key1:value1, ..., key8:value8}` where `id` is a unique string
 117 identifier of the entry (e.g. "42") and each key/value pair represents one of 8 metadata entries from
 118 the Table ?? below. In some cases, where a specific metadata value was not available for an entry, the
 119 value can also be `null`.

Table 2. Metadata format in dataset (F).

Key	Value type	Description	Example/possible values
<code>"included_by"</code>	string	Two character unique identifier of researcher who marked this article as eligible.	<code>"KN"</code>
<code>"title"</code>	string	Text string containing full title of the eligible article	<code>"Some title"</code>
<code>"year"</code>	integer	Year, when the article was published	from 2013 to 2017
<code>"authors"</code>	string	Text string containing list of authors of the article in BibTeX format	<code>"Author, A and Author, B"</code>
<code>"DOI"</code>	string	Text string representing Digital Object Identifier (DOI) of the publication	<code>"10.1109/IE.2015.30"</code>
<code>"scihub"</code>	boolean	Whether article was publicly available (false) or had to be acquired through other channels (true)	<code>true</code> or <code>false</code>
<code>"pdf_url"</code>	string	Contains URL of the publicly available article if it exists	<code>"https://someurl.com"</code>
<code>"description"</code>	string	Text string containing a short description of the contents of the article	<code>"Self calibrating WSN"</code>

120 2.7. (G) - candidate articles marked as ineligible

121 This file contains a list of 1898 JSON data objects, each representing a single article marked as
 122 ineligible in the eligibility checking step.

123 The list is structured as follows: `[entry1, entry2, ..., entry3017]`. Each of the entries have
 124 the following structure: `{key1:value1, ..., key8:value8}` where each key/value pair represents
 125 one of 8 metadata entries from the Table ?? below. In some cases, where a specific metadata value was
 126 not available for an entry, the value can also be `null`.

127 Because only articles describing actual physical deployment of sensor network devices (more
 128 than one and networked) were included, several groups of articles were excluded as illustrated by the
 129 `"reason"` metadata field, which can take one of the following values (number of matching entries in
 130 the dataset in brackets):

- 131 • `"Article not available"` (438 entries) - we were not able to access full text of the article;
- 132 • `"Theoretical"` (160 entries) - the article described theoretical aspects not practical deployment;
- 133 • `"Not deployed"` (293 entries) - no deployment was described even though device might be
 134 developed;

- 135 ● "Article not English" (88 entries) - article not available in English language;
 136 ● "Simulation" (485 entries) - experiment was simulated thus not using actual deployment;
 137 ● "No network" (183 entries) - non-networked devices (usually data loggers) or a single device
 138 deployed;
 139 ● "Review" (23 entries) - a review article of other deployment articles, excluded to avoid
 140 duplication;
 141 ● "Other" (166 entries) - some other reason for exclusion - usually not related to sensor networks
 142 at all.

Table 3. Metadata format in dataset (G).

Key	Value type	Description	Example/possible values
"excluded_by"	string	Two character unique identifier of researcher who marked this article as ineligible.	"KN"
"reason"	string	Text string containing a short reason for exclusion of the contents of the article	See explanation above
"title"	string	Text string containing full title of the eligible article	"Some title"
"year"	integer	Year, when the article was published	from 2013 to 2017
"authors"	string	Text string containing list of authors of the article in BibTeX format	"Author, A and Author, B"
"DOI"	string	Text string representing Digital Object Identifier (DOI) of the publication	"10.1109/IE.2015.30"
"scihub"	boolean	Wether article was publicly available (false) or had to be acquired through other channels (true)	true or false
"pdf_url"	string	Contains URL of the publicly available article if it exists	"https://someurl.com"

143 2.8. (H) - candidate articles mistaken as ineligible

144 This file contains a list of 47 JSON data objects, each representing a single article marked as
 145 ineligible by mistake, even though it was actually eligible, during the eligibility checking step.

146 The list is structured as follows: [entry1, entry2, ..., entry47]. Each of the entries have
 147 the following structure: {key1:value1, key2:value2, key3:value3} where each key/value pair
 148 represents one of 3 metadata entries from the Table ?? below.

Table 4. Metadata format in dataset (H).

Key	Value type	Description	Example/possible values
"article_id"	integer	An integer identifier of the article, matching the id field in dataset (F)	from 2971 to 3017
"mistake_by"	string	Two character unique identifier of researcher who made the mistake in marking the article as ineligible	"KN"
"comment"	string	Text string containing short explanation why the article should be included	"Some comment"

149 2.9. (I) - timeline of eligibility check and data extraction phase

150 This file is formatted as a table in a TAB delimited text file. It has 35 entries, each representing one
 151 of the 35 weeks during which the eligibility checking and data extraction phase took place.

152 Each row has the following 5 headers/columns with corresponding data types:

- 153 1. **Week** - the number of week for which statistics is given. This is an integer value in range from 1
 154 to 35;
- 155 2. **Processed_per_week** - the number of articles processed per week in the eligibility checking and
 156 data extraction phase. This is an integer value in range from 56 to 302;
- 157 3. **Total_processed** - the cumulative number of articles processed up to and including that week.
 158 This is an integer value in range from 73 to 4915;
- 159 4. **Included_and_extracted_per_week** - the number of articles included and actually used for data
 160 extraction per week. This is an integer value in range from 33 to 186;

161 5. `Total_included_and_extracted` - the cumulative number of articles included and actually used
 162 for data extraction up to and including that week. This is an integer value in range from 48 to
 163 2970.

164 2.10. (J) - extracted codified data

165 This file contains a JSON data object with 3059 entries, each representing a single deployment from
 166 the previously identified articles and containing extracted codified data related to this deployment.

167 The object is structured as follows: `{entry1, entry2, ..., entry3059}`. Each of the entries
 168 have the following structure: `id:{key1:value1, ..., key12:value12}` where `id` is a unique string
 169 identifier of the entry (e.g. "42") and each key/value pair represents one of 12 metadata entries from
 170 the Table ?? below. In some cases, where a specific metadata value was not available for an entry, the
 171 value can also be `null`.

172 In the context of this data a device is considered to be a "rich" device instead of ordinary sensor
 173 network device, if it is an interactive computer like system with some multimedia capabilities, e.g.
 174 smartphone, personal computer, Raspberry PI etc.

Table 5. Metadata format in dataset (J).

Key	Value type	Description	Example/possible values
<code>"related_article_id"</code>	integer	Identifier number of the article from dataset (F) in which the specific deployment was described	from 1 to 3017
<code>"extracted_by"</code>	string	Two character unique identifier of the researcher who extracted data about this deployment	"KN"
<code>"year"</code>	integer	Year, when the article containing the deployment was published	from 2013 to 2017
<code>"has_goal_network"</code>	boolean	Whether deployment is made with goal application in mind (true, 1825 entries) instead of just technology focused (false, 1234 entries)	true or false
<code>"goal_network"</code>	data object	Contains data object describing the goal network if it exists	See below
<code>"deployment_notes"</code>	string	Text string containing optional comments by extracting researcher on the deployment	"Vague description..."
<code>"node_connection"</code>	string	Type of connection between sensor nodes - wireless (2860 entries), wired (94 entries), hybrid using both types (89 entries), or not defined (16 entries)	"Wireless", "Wired", "Hybrid", or null
<code>"node_mobility"</code>	string	Mobility type of sensor nodes - static only (2286 entries), mobile only (580 entries), mixed - some static and some mobile (140 entries), or not defined (53 entries)	"Static", "Mobile", "Mixed", or null
<code>"rich_nodes"</code>	string	Which sensor nodes are "rich" devices - none of the nodes are rich (2374 entries), only base station nodes are rich (416 entries), all nodes are rich (226 entries), mixed - some simple and some rich nodes (27 entries), or not defined (16 entries)	"None", "Base_stations", "All", "Mixed", or null
<code>"deployed_as_tool_or_subject"</code>	string	Whether the deployment in the article is used as a tool in the research described (1618 entries), or is the subject of the research itself (1441 entries)	"Tool" or "Subject"
<code>"testbed"</code>	string	Whether a sensor network testbed is used for the described sensor network deployment (478 entries), isn't used (2516 entries) or whether the sensor network itself is part of a testbed (65 entries)	"Used", "No" or "Part of"
<code>"deployment_trl"</code>	string	The Technology Readiness Level of the deployment: 3-bench tested concept (103 entries), 4-validated in laboratory (682 entries), 5-tested in artificial environment/testbed (888 entries), 6-demonstrated on close-to-real environment (479 entries), 7-demonstrated in real environment or 8-final system in real environment (81 entry)	"3-Bench", "4-Lab", "5-Test", "6-Demo", "7-Target" or "8-Final"

175 In addition to the overall description of the sensor network deployment itself, such as type of
 176 connection of sensor nodes and technology readiness level of the deployment as described in the
 177 article, an additional group of metadata was extracted related to the potential future goal network
 178 that the reserach is building towards. Although a major part of the deployments are driven by
 179 technology development (1234 entries) not application and don't have such a goal network, for those
 180 deployments which have some practical application in mind (1825 entries), the following metadata

181 object is stored under the key "goal_network": {key1:value1, . . . , key4:value4}. In this object
182 for each deployment four keys with these possible values and number of entries in dataset are provided:

- 183 • "field" - The target field of application with one of the following values:
- 184 1. "Health & wellbeing" including patient, frail and elderly monitoring systems, sports
185 performance, and general health and wellbeing of both body and mind (349 entries);
 - 186 2. "Education" including systems meant for educational purposes and serious games (10
187 entries);
 - 188 3. "Entertainment" including computer games, AR/VR systems, broadcasting, sporting and
189 public events, gambling and other entertainment (17 entries);
 - 190 4. "Safety" including anti-theft, security, privacy enhancing, reliability improving, emergency
191 response and military applications and tracking people and objects for these applications
192 (163 entries);
 - 193 5. "Agriculture" including systems related to farming, crop growing, farm and domesticated
194 animal monitoring, precision agriculture (229 entries);
 - 195 6. "Environment" monitoring of environment both in wild life and city, including weather,
196 pollution, wild life, forest fires, aquatic life, volcanic activity, flooding, earthquakes etc. (297
197 entries);
 - 198 7. "Communications" general communications like power lines, water and gas pipes, energy
199 consumption monitoring, internet, telephony, radio etc. (51 entries);
 - 200 8. "Transport" including intelligent transport systems (ITS) smart mobility, logistics and
201 goods tracking, smart road infrastructure etc. (123 entries);
 - 202 9. "Infrastructure" general infrastructure, such as tunnels, bridges, dams, ports, smart
203 homes and buildings etc. (413 entries);
 - 204 10. "Industry" anything related to industrial processes, production and business in general
205 like coal mine monitoring, production automation, quality control, process monitoring etc.
206 (143 entries);
 - 207 11. "Research" not related to other fields, but to support future research - better research tools
208 and protocols, testbeds etc. (20 entries);
 - 209 12. "Multiple" the deployed network will serve multiple of the previously described fields (10
210 entries).
- 211 • "scale" - The target deployment scale of the sensor network with one of the following values
212 (from smallest to largest):
- 213 1. "Single actor" including such single entities as a person (e.g. body area network), animal,
214 vehicle or robot (345 entries);
 - 215 2. "Room" include such relatively small territories as rooms, garages, small yards (131 entries);
 - 216 3. "Building" include larger areas with separate zones, like houses, private gardens, shops,
217 hospitals (530 entries);
 - 218 4. "Property" include even larger zones capable of containing multiple buildings, like city
219 blocks, farms, small private forest or orchard (447 entries);
 - 220 5. "Region" include areas of city or self-government scale like a rural area, forest, lake, river,
221 city or suburbs (317 entries);
 - 222 6. "Country" include objects of scale relative to countries, like national road grid, large
223 agricultural or forest areas, smaller seas (27 entries);
 - 224 7. "Global" include networks of scale not limited to a single country, such as oceans, jungle or
225 space (24 entries);
 - 226 8. null - no scale information of target deployment provided or it is not clearly defined (4
227 entries).
- 228 • "subject" - The main target subject meant to be monitored by the goal network with one of the
229 following values:
- 230 1. "Environment" includes all types of environmental phenomena, like weather, forests, bodies
231 of water, habitats etc. (728 entries);

- 232 2. "Equipment" includes all types of inanimate objects, including industrial equipment,
 233 buildings, vehicles or robots as systems not actors in environment, dams, walls etc. (498
 234 entries);
- 235 3. "Opposing actor" include all types of actors in environment, which do not want to be
 236 monitored, thus including security and spying applications, tracking and monitoring of
 237 perpetrators or military opponents, pest control etc. (126 entries);
- 238 4. "Friendly actor" includes actors that do not mind to be tracked or monitored for some
 239 purpose, like domestic or wild animals (tagging), elderly or frail, people in general if
 240 compliant (456 entries);
- 241 5. "SELF" includes cases where the sensor network monitors itself - location of nodes,
 242 communication quality etc. (1 entry);
- 243 6. "Mixed" - this includes target deployments with multiple subjects from the previouslu
 244 stated values (16 entries).
- 245 • "interactivity" - The interactivity of the goal sensor network with the following values:
- 246 1. "Passive" includes passive monitoring nodes and data gathering for decision making
 247 outside the system or for general statistics purposes (1448 entries);
- 248 2. "Interactive" includes sensor networks providing some kind of feedback, control or
 249 interactivity within the loop or confines of the system, like automated irrigation systems,
 250 real time alarms etc. (375 entries);
- 251 3. null no specific interactivity of target deployment is provided or clearly defined in the
 252 article (2 entries).

253 2.11. (K) - statistics of extraction process

254 This file is formatted as a table in a TAB delimited text file. It has 12 entries, each pertaining to
 255 one of the 12 volunteer researchers involved in the data extraction process.

256 Each row has the following 11 headers/columns with corresponding data types:

- 257 1. **Extractor** - two letter code uniquely identifying each of the researchers. Example of data in
 258 column: KN;
- 259 2. **Total_articles_processed** - number of articles processed by the corresponding researcher in
 260 the data extraction step. This is an integer value in range from 1 to 708;
- 261 3. **Total_deployments_extracted** - number of actual sensor network deployments extracted from
 262 these articles by the corresponding researcher in the data extraction step. This is an integer value
 263 in range from 1 to 708;
- 264 4. **Not_sure_goal_deployment** - number of articles in which the extractor was not sure about the
 265 goal deployment of the sensor network and required peer input to get the final value. This is an
 266 integer value in range from 0 to 12;
- 267 5. **Error_goal_deployment** - number of articles in which the extractor mistakenly marked a wrong
 268 goal deployment value, which was later corrected in validation stage. This is an integer value in
 269 range from 0 to 79;
- 270 6. **Total_goal_deployment_mistakes** - Sum of two previous values representing the total amount
 271 of errors related to the goal deployment - made by the specific extractor. This is an integer value in
 272 range from 0 to 87;
- 273 7. **Goal_deployment_error_rate** - the percentage of the total number of deployments processed by
 274 the researcher in data extraction stage that contained some sort of error related to goal deployment
 275 data extraction. This number is formed as percentage value rounded to 2 decimal places and has
 276 values from 0.00% to 30.00%;
- 277 8. **Not_sure_other** - number of articles in which the extractor was not sure about the some other
 278 metadata value not related to goal deployment and required peer input to get the final value. This
 279 is an integer value in range from 0 to 60;
- 280 9. **Error_other** - number of articles in which the extractor mistakenly marked a wrong metadata
 281 value not related to goal deployment, which was later corrected in validation stage. This is an
 282 integer value in range from 0 to 10;

- 283 10. `Total_other_mistakes` - Sum of two previous values representing the total amount of errors
 284 not related to the goal deployment made by the specific extractor. This is an integer value in range
 285 from 0 to 66 ;
- 286 11. `Other_error_rate` - the percentage of the total number of deployments processed by the
 287 researcher in data extraction stage that contained some sort of error not related to goal deployment
 288 data extraction. This number is formed as percentage value rounded to 2 decimal places and has
 289 values from 0.00% to 50.00% .

290 2.12. (L) - candidate articles mistaken as eligible

291 This file contains JSON data object with 15 entries, each representing a single article marked as
 292 eligible by mistake, even though it was actually ineligible, discovered during the data extraction step.

293 The object is structured as follows: {entry1, entry2, ..., entry15}. Each of the entries have
 294 the following structure: `id:{key1:value1, key2:value2, key3:value3}` where `id` is a unique
 295 string identifier of the article and each key/value pair represents one of 3 metadata entries from the
 296 Table ?? below.

Table 6. Metadata format in dataset (L).

Key	Value type	Description	Example/possible values
<code>"related_article_id"</code>	integer	An integer identifier of the article, matching the <code>id</code> field in dataset (F)	from 279 to 2676
<code>"included_by"</code>	string	Two character unique identifier of researcher who made the mistake in marking the article as eligible	"KN"
<code>"error_type"</code>	string	Text string describing the reason for exclusion of the mistakenly included deployment, including simulation (3 entries), use of previously existing data (5 entries), no network between devices (4 entries) and sensor network not actually deployed (3 entries)	"simulation", "existing data", "no network" or "no deployment"

297 2.13. (M) - overall timeline of dataset creation

298 This file is formatted as a table in a TAB delimited text file. It has 14 entries, each pertaining to a
 299 milestone date in the progress of dataset creation and has no column headers.

300 Each row has the following 2 columns with corresponding data types:

- 301 1. Date in format of `yyyy-mm-dd` with values in the range from 2018-06-12 to 2020-07-02 ;
- 302 2. Milestone event description in the form of a text string.

303 2.14. (N) - readme file

304 This file contains a short human readable description of the data in this dataset in the form of a
 305 Markdown document.

306 2.15. (O) - notebook folder

307 In this folder several Jupyter notebook files are stored for easy loading of and access to the data
 308 files. These contain example Python 3 code for opening the files and extracting the data within.

309 3. Methods

310 To acquire this dataset, first the scope of the problem was defined as follows: to gather and codify
 311 all scientific peer reviewed publications describing original practical sensor network deployments
 312 from a 5-year period from 2013 to 2017. The scope was narrowed for practical purposes as follows:

- 313 • Only publications in English language were considered;
- 314 • Only publications that could be accessed by the research team without use of additional funds
 315 were considered;

- 316 ● To be considered a network, the deployment had to have at least two actually deployed sensor
 317 devices;
- 318 ● Devices didn't have to be wireless, to be considered sensor network - also wired, acoustic or other
 319 networks were considered;
- 320 ● Only research doing the deployment themselves was considered - no use of ready datasets from
 321 other deployments was included;
- 322 ● No simulated experiments were included;
- 323 ● the timeframe was selected as 2013-2017, because the data acquisition was started in the middle
 324 of 2018, and only full years were chosen for comparability;
- 325 ● To avoid duplicates only original deployments were included instead of review articles.

326 Based on this scope a systematic literature review methodology was devised and followed
 327 consisting of the following steps (note that in dataset, the files related to these steps are enumerated
 328 starting from 0 not 1):

- 329 1. Candidate article acquisition
 330 2. Screening (exclusion)
 331 3. Screening (inclusion/eligibility)
 332 4. Codification and data extraction
 333 5. Verification

334 3.1. Candidate article acquisition

335 Due to their popularity and wide access in the institutions represented by the authors two main
 336 indexing databases were selected for querying articles: SCOPUS and Web of Science.

337 For each of these databases a query with the same information based on the scope defined above
 338 was prepared:

- 339 ● **SCOPUS:** KEY(sensor network OR sensor networks) AND TITLE-ABS-KEY(test* OR experiment*
 340 OR deploy*) AND NOT TITLE-ABS-KEY(review) AND NOT TITLE-ABS-KEY(simulat*) AND (
 341 LIMIT-TO (PUBYEAR,2017) OR LIMIT-TO (PUBYEAR,2016) OR LIMIT-TO (PUBYEAR,2015)
 342 OR LIMIT-TO (PUBYEAR,2014) OR LIMIT-TO (PUBYEAR,2013))
- 343 ● **Web of Science:** TS = ("sensor network" OR "sensor networks") AND TS = (test* OR
 344 experiment* OR deploy*) NOT TI="review" NOT TS=simulat* **with additional parameters:**
 345 Indexes=SCI-EXPANDED, SSCI, A&HCI, CPCI-S, CPCI-SSH, BKCI-S, BKCI-SSH, ESCI,
 346 CCR-EXPANDED, IC Timespan=2013-2017

347 The querying was done on **2018-06-12** and yielded the following results:

- 348 ● **SCOPUS:** 11536 total articles identified of which 4814 were not found in Web of Science database;
 349 ● **Web of Science:** 10204 total articles identified of which 3636 were not found in SCOPUS database;
 350 ● After checking for duplicates **15010 unique candidate articles** were identified of which 6560
 351 articles were found in both databases. Duplicates were checked both automatically using features
 352 provided by Mendeley software and manually by title/author/year combination.

353 The resulting dataset was saved as BibTeX file (see dataset (A)) and imported in Mendeley
 354 software for collaborative screening for exclusion.

355 3.2. First screening iteration - exclusion

356 During this stage, research team was instructed to exclude articles conservatively - only exclude
 357 those which definitely match the exclusion criteria and leave all others for more thorough examination
 358 in the next stage.

359 The exclusion criteria was defined as follows: *The article does not feature a real life deployment of a
 360 sensor network or is not in English language.*

361 The team of volunteers participating in the screening process were all provided access to a shared
 362 Mendeley group with the dataset of 15010 articles and given the following instructions:

- 363 1. Spend no more than 10 minutes on a single article;
 364 2. Only look at article title and abstract for exclusion;
 365 3. When processing an article mark it as "read" (a gray/green circle mark in Mendelay);
 366 4. If the article matches exclusion criteria move onto the next article;
 367 5. Otherwise, include it for the next stage by marking it with favourite/star icon in Mendeley;
 368 6. Regularly synchronize progress and follow randomized article slots based on alphabetic order of
 369 article titles to avoid collisions of multiple reviewers;
 370 7. In case of doubt, articles could be tagged for second opinion by another reviewer.

371 To ensure consistent understanding of the exclusion process, one researcher took lead of it and
 372 re-evaluated first 100 articles processed by all other researchers, and discussed any differences or
 373 problems. Weekly discussion on progress, problematic articles etc. was held.



Figure 1. Weekly progress of first screening phase.

374 The screening for exclusion took place from 2018-06-13 till 2018-10-08 with the weekly progress
 375 shown in Figure ???. Then validation stage started, during which randomized sample or articles was
 376 double checked by other researchers and 142 articles identified as requiring second opinion were
 377 discussed and marked appropriately. The validation of screening/exclusion phase ended on 2018-11-03
 378 with 4915 articles left for the next stage (thus 10095 articles were excluded in this phase).

379 3.3. Second screening iteration - inclusion/eligibility

380 After the first stage of screening the second screening iteration phase started. This phase required
 381 opening and reading the full text of the articles, thus, for time conservation it was done in parallel with
 382 the next phase - codification and data extraction (See next section).

383 First, from 2018-11-04 till 2019-01-13 an instruction for full text eligibility validation was developed
 384 together with data codification and data extraction methodology. An online spreadsheet was developed
 385 with the 4915 articles from the previous screening phase, with columns for the required data as
 386 dropboxes.



Figure 2. Weekly progress of second screening phase.

387 Then, from 2019-01-14 till 2019-09-15 data inclusion/eligibility and codified data extraction stage
 388 took place - the weekly progress can be seen in Figure ???. The main steps in this stage for all researchers
 389 involved were as follows:

- 390 1. Mark the row corresponding to the selected article with unique identifier of the researcher, so that
 391 noone else accidentally takes the same article for analysis;
- 392 2. Locate the full text of the article - if it is not available in English language from any source
 393 (indexing pages, preprint publishing pages, author pages, Researchgate, Sci-hub, general google
 394 search etc.) then exclude the article from data extraction, otherwise move to the next step;
- 395 3. Read the article to identify any sensor network deployments in it. If there are no deployments,
 396 the article must be excluded. If there are several deployments, insert new lines in the table, thus
 397 describing each deployment separately;
- 398 4. Do not include any articles that should have been excluded in the previous stage (review articles,
 399 articles without actual deployments or using old data from previous deployments, or even
 400 deployments with single sensor device or multiple devices, which have no sensors or network
 401 between them);
- 402 5. For each included row leave a comment on which/how many actual sensor network deployments
 403 are there - these deployment rows in the spreadsheet table are then filled by the same researcher
 404 as part of the next phase (see next Section).

405 During the second screening stage 2970 articles were first included and codified. Then on
 406 2019-09-17 a verification phase of excluded articles was begun, and involved both randomized reviews,
 407 as well as multiple reviews of any article marked as uncertain by the original researcher. After this
 408 phase ended on 2019-12-02 additional 47 articles were found in the mistakenly excluded article list and
 409 included thus leading to 3017 total articles eligible for extraction.

410 From the excluded $4915 - 3017 = 1898$ articles, the reasons for exclusion from most frequent to least
 411 frequent were: (1) article describes simulation not actual deployment - 485 articles; (2) Article full text
 412 not available - 438 articles; (3) Sensor network only described, but not actually deployed - 293 articles;
 413 (4) only separate sensor devices with no network/local data logging - 183 articles; (5) Theoretical
 414 article with no practical experiments - 160 articles; (6) Article not available in English - 88 articles; (7)
 415 Article uses existing data gathered from a previous deployment or public data set - 62 articles and (8)
 416 Article is a review article of other deployments - 23 articles. Additionally 166 articles were excluded
 417 due to other reasons, that didn't correspond to one of the above mentioned categories (e.g. nothing to
 418 do with sensor devices or disqualified due to multiple categories).

419 Till 2020-01-05 all deployments in these articles were identified and codified and a thorough
420 validation phase of codified data was carried out during the process in which 15 articles were identified
421 as mistakenly included for codification leaving only 3002 articles.

422 The total number of identified sensor network deployments in these articles was 3059.

423 3.4. Data codification and extraction

424 For all of the 3059 deployments the researchers involved had to extract two codified groups of
425 data:

- 426 1. Details on the actual sensor network deployment described in the article;
- 427 2. If exists - the goal deployment towards which this research is aimed in the future.

428 The specific codification values are described in detail in the data description of dataset (J) as
429 shown in section ???. In addition to these values all researchers were allowed to provide `null` value if
430 the article did not mention or describe the specific value of interest and `OTHER` value if the researcher
431 did not think the value could fit in any previously defined category. Additionally on every field the
432 researchers could leave comments asking for second opinion or leaving discussion points about the
433 codification system.

434 As with the exclusion stage, data extraction stage also contained coordination between the
435 researchers involved - first 10 codification efforts by each of the researchers were double checked by
436 one researcher so that everyone had common understanding. All questions and unclear values were
437 discussed weekly for clarifications.

438 Finally the codified data was verified - all comments were manually processed, outlier values,
439 `null` values and `OTHER` values were double checked by other researchers, to verify that something
440 was not missed by the original reader of the article. Additionally random validation of codified entries
441 occurred.

442 The errors during validation were labeled and counted for each of the researchers involved (as
443 can be seen in datasets (E), (H), (K) and (L). The deployments that were checked by researchers who
444 were outliers (with low amount of articles processed or high amount of specific errors) were re-checked
445 by other researchers.

446 Finally, on 2020-05-29 the dataset was completed and preparation started for publishing the data
447 set. Data set was cleaned up, formatted and submitted to an open access Git repository on 2020-07-02.
448 Afterwards, the text of this publication was prepared together with Jupyter Notebook examples on use
449 of this data.

450 3.5. Data quality

451 In addition to random validation and checking for errors as described in the previous steps,
452 additional checks on the data set were done to ensure quality of the data.

453 First, the number of article candidates from each year were compared to see if there is a bias for
454 specific years (e.g. older articles). Each year the number of articles was around the mean 3002, with
455 deviation of less than 4.5%.

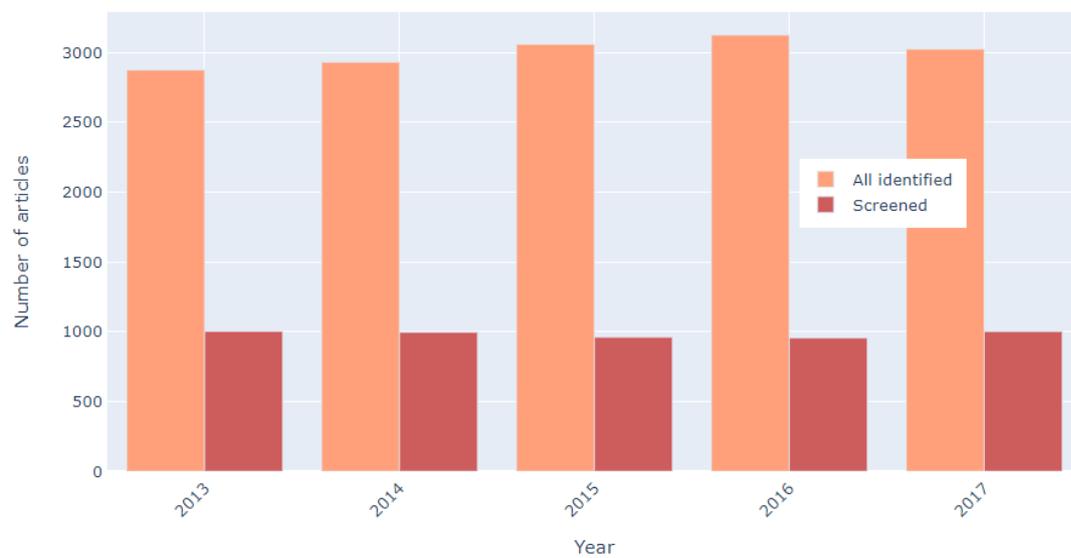


Figure 3. Number of articles initially identified per year and corresponding number of articles included in first screening phase.

456 Then, the screening phase results were analyzed to test for bias related to year. In all years 30%
 457 to 35% of articles survived the first screening/exclusion phase, with no observable bias towards any
 458 particular year (see Figure ??).

459 The approximate 1/3 inclusion rate also held true for the three most represented categories
 460 of articles: @article with 33.42%, @conference with 30.57% and @inproceedings with 32.99%
 461 inclusion rates. The two less represented groups @book and @incollection each had less than 75
 462 instances in the first dataset and thus even though their inclusion rate differed from the expected
 463 (16.22% and 17.91% respectively) this is most likely due to the small number of articles in these
 464 categories not an inherent bias towards them in the screening process.

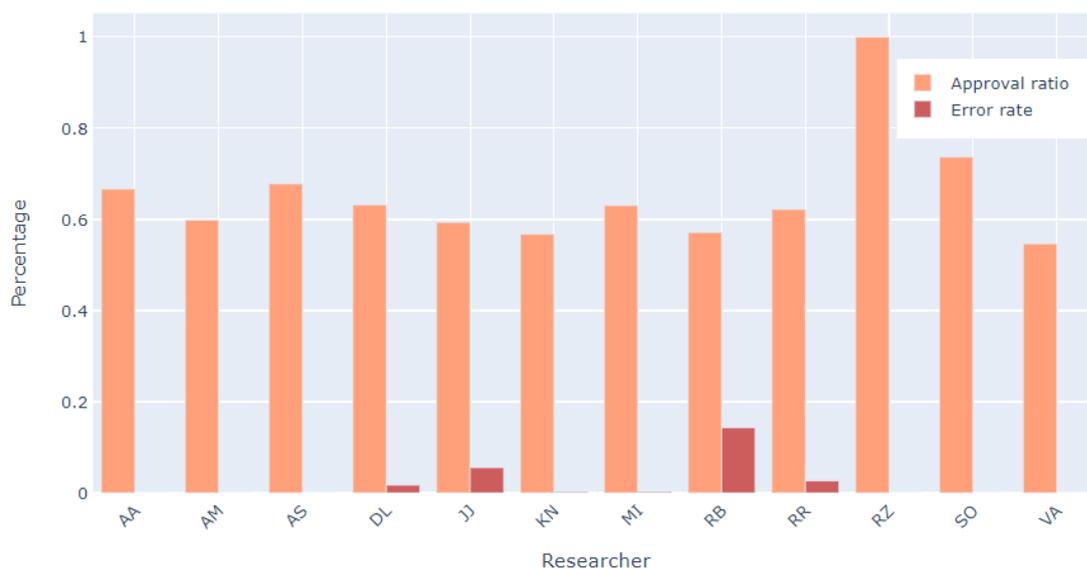


Figure 4. For each researcher - the ratio of articles marked as eligible and their detected error percentage in eligibility screening phase.

465 Another potential source of bias is the differences in researchers doing the screening, so all
 466 involved researchers were analyzed. Most had similar approval ratio of articles (articles marked
 467 eligible over all articles processed) and similar low error percentage from total articles processed. As

468 seen in Figure ?? there are three main outliers: RZ who has 100% approval ratio, which is due to the
469 fact that this researcher only processed one article in this stage, S0 whose approval ratio is closer to
470 70% instead of 60% like others, which is also due to the low number of articles processed (less than
471 40), and finally RB who had around 15% error rate in comparison to other researchers who had error
472 rate below 5%. This is also due to the low number of processed articles (7). All other researchers in
473 this phase processed several hundreds of articles and their statistics and error rates were very similar
474 showing that the efforts to reduce bias introduced by individuals were successful.

475 Overall, wherever a potential source for bias was detected due to a low number of articles
476 processed by a researcher, their work was re-validated by at least one other researcher to guarantee
477 high data quality.

478 4. User notes

479 The data set was primarily meant for easy processing using programming tools, such as
480 Python/Jupyter Notebooks, thus it is machine readable first.

481 The data is made freely accessible to everybody, although we would appreciate credit if at all
482 possible. To the best of our knowledge this is the only data set of its kind and currently only covers
483 years 2013 to 2017.

484 To access the data, use Git to clone the repository:

```
485 git clone http://git.edi.lv/CPS_Lab/sn_deployment_mapping_review
```

486 In this way you will get all of the files described in Chapter ??.

487 For examples on loading and processing this data using Python, You can access the folder
488 `Notebooks` where Jupyter notebook files with examples on data exploration are stored.

489 **Author Contributions:** For research articles with several authors, a short paragraph specifying their individual
490 contributions must be provided. The following statements should be used "conceptualization, Janis Judvaitis
491 and Krisjanis Nesenbergs; methodology, Krisjanis Nesenbergs; software, Krisjanis Nesenbergs; validation,
492 Janis Judvaitis, Artis Mednis, Valters Abolins and Krisjanis Nesenbergs; formal analysis, Krisjanis Nesenbergs;
493 investigation, Janis Judvaitis, Artis Mednis, Valters Abolins, Ansis Skadins, Didzis Lapsa, Raimonds Raba,
494 Maksims Ivanovs and Krisjanis Nesenbergs; data curation, Krisjanis Nesenbergs; writing–original draft
495 preparation, Janis Judvaitis; writing–review and editing, Krisjanis Nesenbergs, Maksims Ivanovs and Raimonds
496 Rava; visualization, Krisjanis Nesenbergs; supervision, Krisjanis Nesenbergs.

497 **Funding:** This research is co-financed by ERDF funds under the project No. 1.1.1.1/18/A/183 "iTrEMP: Intelligent
498 transport and emergency management platform"

499 **Acknowledgments:** In addition to the authors of this article also these people provided their work to the
500 acquisition and processing of this data: Rihards Balass, Reinholds Zviedris, Armands Ancans and Sandra Ose.

501 **Conflicts of Interest:** The authors declare no conflict of interest. The funders had no role in the design of the
502 study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to
503 publish the results.

504 Abbreviations

505 The following abbreviations are used in this manuscript:

506 WSN Wireless Sensor Network
IoT Internet of Things
507 TAB Tabulator character
JSON JavaScript Object Notation

508

509 . Noussan, M.; Hafner, M.; Tagliapietra, S. Digitalization Trends. In *The Future of Transport Between*
510 *Digitalization and Decarbonization*; Springer, 2020; pp. 51–70.

511 . Noussan, M.; Tagliapietra, S. The effect of digitalization in the energy consumption of passenger transport:
512 An analysis of future scenarios for Europe. *Journal of Cleaner Production* **2020**, p. 120926.

- 513 . Winkler, H.; Zinsmeister, L. Trends in digitalization of intralogistics and the critical success factors of its
514 implementation. *Brazilian Journal of Operations & Production Management* **2019**, *16*, 537–549.
- 515 . Dorofeeva, A.A.; Nyurenberger, L.B. Trends in digitalization of education and training for industry 4.0 in
516 the Russian Federation. IOP Conference Series: Materials Science and Engineering. IOP Publishing, 2019,
517 Vol. 537, p. 042070.
- 518 . Li, F.; Yang, J.; Wang, J.; Li, S.; Zheng, L. Integration of digitization trends in learning factories. *Procedia*
519 *Manufacturing* **2019**, *31*, 343–348.
- 520 . Kosareva, O.A.; Eliseev, M.N.; Cheglov, V.P.; Stolyarova, A.N.; Aleksina, S.B. Global trends of digitalization
521 of agriculture as the basis of innovative development of the agro-industrial complex of Russia. *Eurasian*
522 *Journal of Biosciences* **2019**, *13*, 1675–1681.
- 523 . Rodin, B.; Ganiev, R.; Orazov, S. «Fintech» in digitalization of banking services. International Scientific
524 and Practical Conference on Digital Economy (ISCDE 2019). Atlantis Press, 2019.
- 525 . Evdokimova, Y.; Shinkareva, O.; Bondarenko, A. Digital banks: development trends. 2nd International
526 Scientific conference on New Industrialization: Global, national, regional dimension (SICNI 2018). Atlantis
527 Press, 2019.
- 528 . NB-IoT Market Size Worth \$6.02 Billion by 2025. [https://www.bloomberg.com/press-releases/2019-07-
529 23/nb-iot-market-size-worth-6-02-billion-by-2025-cagr-34-9-grand-view-research-inc](https://www.bloomberg.com/press-releases/2019-07-23/nb-iot-market-size-worth-6-02-billion-by-2025-cagr-34-9-grand-view-research-inc). [online] [viewed
530 16.06.2020].
- 531 . Lima, L.E.; Kimura, B.Y.L.; Rosset, V. Experimental environments for the internet of things: A review. *IEEE*
532 *Sensors Journal* **2019**, *19*, 3203–3211.

533 © 2020 by the authors. Submitted to *Data* for possible open access publication under the terms and conditions of
534 the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).