



1

2 **ARTICLE INFORMATION**

3 **Article title**

4 *DigitalExposome: A Dataset for Wellbeing Classification using Environmental Air Quality*
5 *and Human Physiological data.*

6 **Authors**

7 Thomas Johnson

8

9 **Affiliations**

10 Department of Computer Science, Nottingham Trent University, UK.

11 **Corresponding author's email address and Twitter handle**

12 *Thomas Johnson: thomas.johnson@ntu.ac.uk - Twitter: @tomwjohnson*

13 **Keywords**

14 Urban environment; air pollution; sensors; physiological wellbeing assessment; wearable
15 devices; participatory research; air quality monitoring

16 **Abstract**

17 Urban environments play a critical role in shaping mental wellbeing, yet their impact remains
18 understudied, particularly in relation to environmental air quality and human physiology.
19 Despite this growing awareness of the importance of mental health in urban planning,
20 challenges in integrating diverse datasets, spanning environmental, physiological, and self-
21 reported mental wellbeing data limit the scope of research in this area. The DigitalExposome
22 dataset addresses this gap by providing a comprehensive resource for understanding the
23 relationship between these factors. The resulting data was collected from October 2021 to
24 September 2022 in Nottingham, UK with the dataset including over 42, 437 samples from 40
25 participants aged between 18-50. Participants conducted a walk through diverse urban
26 environments including polluted and green spaces, while carrying a custom-built
27 environmental monitoring system (Enviro-IoT), wearing an Empatica E4 wearable, and using
28 a smartphone mobile application to self-label mental wellbeing via emojis. Environmental
29 variables (e.g., a range of particulates and gases including particulate matter and nitrogen
30 dioxide), physiological metrics (e.g., HR, HRV, EDA, BVP), and mental wellbeing labels
31 were recorded. Data was processed following collection through resampling and
32 interpolation, and normalization for analysis. This novel dataset lays the groundwork for
33 exploring the relationships between air quality, physiological changes, and mental wellbeing,
34 offering valuable insights for urban planning and public health.

35

36

37

38

39

40

41

42

43

44

45 **SPECIFICATIONS TABLE**
46

Subject	Computer Science, Artificial Intelligence, Pollution
Specific subject area	Environmental air quality, human physiological monitoring, and mental wellbeing analysis in urban and green space environments
Type of data	12 columns (.csv format)
Data collection	Data was collected from October 2021 to September 2022 in Nottingham, United Kingdom, using a custom-built air quality monitoring station (Enviro-IoT), an Empatica E4 wearable, and a self-report smartphone mobile application. The Enviro-IoT tracked air pollutants (e.g., PM1.0, PM2.5, PM10, CO, NO2, NH3) and noise, while the E4 Empatica measured physiological metrics (such as HR, HRV, EDA, BVP). Participants (n=40) aged between 18-50 years old completed a 40-minute pre-defined urban route. Data was resampled following collection to 1Hz and normalised for analysis, excluding variables with logging issues.
Data source location	Institution: Department of Computer Science, Nottingham Trent University, Clifton Lane, Nottingham, United Kingdom.
Data accessibility	Repository name: Mendeley Data Data identification number: 10.17632/mbwxy48223.2 Direct URL to data: https://data.mendeley.com/datasets/mbwxy48223/2
Related research article	Johnson, T., Kanjo, E. & Woodward, K. DigitalExposome: quantifying impact of urban environment on wellbeing using sensor fusion and deep learning. <i>Comput. Urban Sci.</i> 3, 14 (2023). https://doi.org/10.1007/s43762-023-00088-9

47
48

49 **VALUE OF THE DATA**

- 50 • The data comprises 12 columns of 42, 437 samples involving self-labelled mental
51 wellbeing captured with emojis, environmental air quality and on-body human
52 physiological data. It offers greater understanding of the link between the environment,
53 wellbeing and emotions by collecting data which is obtained at the point-of-exposure
54 within an urban environment.
- 55 • DigitalExposome dataset can be used to train, validate and test various deep learning
56 models, such as CNNs, aiding in the development of tools and enhancing the robustness
57 and accuracy of mental wellbeing classification.
- 58 • The dataset will be used by researchers, scientists and engineers for investigating the
59 relationship and consequences of air pollution factors towards wellbeing, creating novel
60 pollution control techniques, and evaluating the efficacy of treatments.



- 61 • The dataset includes environmental air quality measurements, human physiology data and
62 self-reported mental wellbeing, and, to our knowledge, is the largest publicly accessible
63 dataset of its kind collected in a real-world setting.
64
65

66 BACKGROUND

67 The air we breathe is a familiar environmental hazard not only to our health [3] but in recent
68 years has led to new studies focused on the impact towards our behaviour [4], mental health
69 [5], wellbeing [6] and emotions [7]. As the world continues to grow in significant numbers
70 these issues remain one of the key factors to reduce in order to improve the quality of life, not
71 only in the UK but worldwide. Approximately, 99% of the world's society in 2019 were
72 living in areas where the air quality guidelines are below the recommended levels and
73 because of non-clear fuels and household emissions are causing over 4.2 million deaths each
74 year [8]. Also, individuals who live within urban environments in the UK are as a result more
75 than likely to develop an increased level of blood pressure, asthma, allergy related illnesses
76 and behavioural issues [9]. The current datasets available predominately focus on the issue of
77 poor air quality towards health making it difficult to draw conclusions on the impact towards
78 mental wellbeing and emotions. So, there arises a need to solve the problem for which a real-
79 world, real-time extensive dataset is required which collects data at the point-of-exposure.
80

81 DigitalExposome Dataset [1] was created to provide an open, accessible and high-quality
82 resource to explore the relationship between the urban environment, human behaviour and
83 on-body physiology and mental wellbeing. Measuring a range of environmental factors such
84 as particulate matter (1.0, 2.5 & 10), noise, carbon monoxide, ammonia, and nitrogen
85 dioxide, the dataset provides a more thorough view and understanding into the entirety of the
86 environment. Precisely labelled datasets are essential for building effective and practical deep
87 learning applications [13]. The dataset was collected in the real-world and so seeks to provide
88 high volume and variety of results which are more diverse than synthetic datasets.
89

90 Alternative works such as Datasets by Air Quality and Health Impact Assessment [10] and
91 CitiS-Health: Air Pollution and Mental Health include aspects of air quality and mental
92 health [11] include aspects of air quality and mental health. However, they are generally
93 unrelated to other domains and have limited relevance to wellbeing classification. The Air
94 Quality and Health Impact Assessment dataset provides relevant data on the urban
95 environment, however the focus is much more on the health impact such as towards
96 respiratory, cardiovascular and hospital admissions, while CitiS-Health: Air Pollution and
97 Mental Health dataset is limited to only part of an individual's mental health attributes such
98 as physical activity and diet habits. In contrast, the DigitalExposome dataset focuses on real-
99 time data collection, at the point of exposure to enhance wellbeing classification by
100 combining environmental air quality factors and on-body human physiological data.
101

102 DATA DESCRIPTION

103 This article describes the dataset which includes environmental air quality, human on-body
104 physiological and self-report mental wellbeing collected from the October of 2021 to

105 September 2022. The analysis of the collected data in this article, include the descriptive
 106 statistics of both the collected environmental and physiological variables obtained during the
 107 study; including (min, max, mean, median, quartiles). Additionally, the shape and distribution
 108 of the data, including skewness and kurtosis are also described. This is evident in Table 1
 109 below.

110

111 Table 1: Summary of descriptive statistics of the collected environmental air quality and
 112 physiological factors.
 113

Variables	Mean	Median	Min	1st Qu.	2nd Qu.	3rd Qu.	Max	Skewness	Kurtosis
BVP (μV)	-1.5	0	-1050	-36.0	0	34.7	1070.0	-0.02	11.0
EDA (μS)	0.3	0.2	0	0.1	0.2	0.3	4.5	3.90	15.8
HR (bpm)	100.0	101.0	0.7	91.2	100.6	109.0	174	0.13	1.7
HRV (s)	0.5	0.6	0	0.2	0.6	0.6	1.3	-0.52	-0.5
NH ₃ (ppm)	879.0	686.0	15.0	509	686.0	1060.0	3800.0	1.30	1.4
Noise (dB)	97.4	96.4	47.2	94.5	96.4	100.0	140.0	-1.70	20.0
Nitrogen Dioxide ($\mu\text{g}/\text{m}^3$)	38.0	38.0	2.0	30.0	38.0	42.3	88.0	0.08	0.2
PM 1.0 ($\mu\text{g}/\text{m}^3$)	4.4	3.0	0	0	3.0	7.0	65.0	3.20	19.0
PM 2.5 ($\mu\text{g}/\text{m}^3$)	5.8	0	0	3.0	3.0	9.0	65.0	2.00	7.1
PM 10 ($\mu\text{g}/\text{m}^3$)	7.3	3.0	0	0	4.0	12.0	65.0	19.00	4.4
Carbon Monoxide (ppm)	453.0	509.0	47.0	341.0	509.0	548.0	1201.0	-1.40	0.5

114

115

116 The data was obtained at a sub-urban area of Clifton, UK which involved using land located
 117 at Nottingham Trent University (Clifton Campus). The route selected took participants on a
 118 journey into different urban environments next to a dual carriageway and several green
 119 spaces.

120

121 A total of 40 participants made up from 25 male and 15 female took part in the study, all aged
 122 between 18 and 50 years old. All individuals who took part were screened prior to the study
 123 to ensure they were fit and healthy with a questionnaire. The questionnaire is available via the
 124 Mendeley data repository to view.

125

126 The dataset folder ‘DigitalExposome Dataset (40 users combined)’ in the Mendeley Data
 127 repository comprises a total of 12 columns and 42, 437 samples, with an average of 1,061
 128 samples collected per participant presented in a Microsoft Excel format (.csv) comprising of
 129 1 worksheet. The DigitalExposome dataset comprises data from the 40 participants, with all
 130 data integrated as one which has been normalised for analysis. To normalise the data, each
 131 value collected was scaled using min-max normalization approach. This technique is

132 presented at Equation 1. Specifically, the minimum and maximum values for each variable
 133 were identified, and each data point was transformed by subtracting the minimum value and
 134 dividing by the range (maximum – minimum). This process ensures that all values are within
 135 a uniform scale of 0 to 1.

$$z_i = (x_i - \min(x)) / (\max(x) - \min(x))$$

138 Equation 1. Min-max normalisation formula for the scaling data to a range between 0 and 1.

139
140 Where:

- 141 - Z_i is the i th normalised value in the dataset,
- 142 - X_i is the i th value in the dataset,
- 143 - $\min(x)$ is the minimum value in the dataset,
- 144 - $\max(x)$ is the maximum value in the dataset.

145
146 The ‘DigitalExposome’ data collected at each sample point, includes the following aspects
 147 (1) the environmental air quality factors (including Nitrogen Dioxide, Carbon Monoxide,
 148 Ammonia, PM1.0, PM2.5 and PM10), (2) physiological on-body factors (including HRV,
 149 HR, EDA and BVP) and (3) labelled mental wellbeing data (label). In total there are 12
 150 column names with data under each.

153 EXPERIMENTAL DESIGN, MATERIALS AND METHODS

154 Data for the study was collected using three devices as depicted in Fig. 1 which involved a
 155 custom-built environmental air quality monitoring station (Enviro-IoT), industry standard
 156 physiological measurement wrist wearable (Empatica E4 model) and a self-report custom-
 157 developed mobile application.



159
160 Figure 1. Integrated System for Monitoring Air Quality, Physiology and Wellbeing: The
 161 Enviro-IoT system combines air quality sensors, the E4 Empatica wearable, and a



162 smartphone app to track environmental factors, physiological data and self-reported user
163 wellbeing.
164

165
166
167
168
169

170 **Data Collection Preparation**

171 Firstly, the Enviro-IoT is an air quality monitoring system that employs low-cost sensors to
172 capture pollutants and gases; housed within a rucksack which has been validated in line with
173 industry standard equipment as detailed in our other work [14]. The collected variables
174 included the following variables: Particulate Matter (PM1), (PM2.5), (PM10), Nitrogen
175 Dioxide (NO₂), Carbon Monoxide (CO), Ammonia (NH₃) and Noise (dB)).
176

177
178
179
180
181

177 Secondly, an E4 Empatica to measure and record on-body physiological changes involving
178 ElectroDermal Activity (EDA), Heart-Rate (HR), Heart-Rate Variability (HRV), Body
179 Temperature, Blood Volume Pulse (BVP) and movement). Finally, an Android smartphone
180 with a pre-loaded mobile application was used to record wellbeing changes in-situ within the
181 environment.

182
183
184
185
186
187
188
189
190
191

184 The labelling process in our study makes use of the five-point Likert SAM scale which has
185 been adapted using emojis from the work carried out from the ‘Personal Wellbeing Index for
186 Adults’ [12]. This approach involves asking users how they are feeling with their life as a
187 whole. The process behind this step is that during the data collection participants will be
188 constantly prompted by the researcher to ascertain how they are feeling. In the pre-installed
189 mobile application this previous work is applied in the way that participants are met with five
190 well-known emojis, displayed on buttons which equate to a score of 1 = very negative/ very
191 low to 5 = very positive/ very high.

192

193 **DigitalExposome Experiment**

194 All participants completed a declaration of health assessment to ensure they were fit and
195 healthy as well as an informed consent form before taking part. At the start of each session
196 participants were given the Enviro-IoT rucksack, smartphone with pre-loaded application and
197 asked to wear the E4 wristband on their non-dominant hand. Prior to the study, participants
198 were reminded to constantly select an emoji option displayed to them on the smartphone
199 whilst walking around. Additionally, participants were shown a map outline of the pre-
200 specified route to familiarise themselves with the outline route as depicted in Figure 2. The
201 pre-specified route took participants through a variety of different scenarios within the urban
202 environment including busy, polluted, and green space. In all, walking all the way around
203 took each participant around 40 minutes.

204
205
206
207

205 Data was collected from the three devices on a regular basis, specifically after each
206 participant’s experimental session had ended. This approach ensured that the data could be
207 promptly reviewed, cleaned and prepared ready for analysis. The regular retrieval process

208 also helped to identify any potential technical issues or inconsistencies early on, allowing for
 209 corrective measurements if necessary.
 210



211
 212 Figure 2. The pre-defined route of DigitalExposome taking participants through an array of
 213 different urban environments.
 214

215 **Cleaning and Pre-processing collected data.**

216 Between the experimental devices there is a varying sampling rate of each device which was
 217 taken into account. Physiological data collected by the E4 Empatica varies at HR = 1Hz,
 218 EDA = 4Hz, BVP = 64Hz. HRV is provided as a sequence of time intervals corresponding to
 219 detected heartbeats, rather than at a fixed sampling rate. The environmental air quality data is
 220 sampled at 0.2Hz.

221
 222 To ensure a consistent sampling rate, the physiological data produced by the Empatica was
 223 down sampled to a rate of 1Hz to match the sample rate of collected HR. The environmental
 224 data had to be up sampled to match the sampled rate of the physiological data at 1Hz.
 225 Furthermore, the self-labelled mental wellbeing data collected from the smartphone was
 226 extracted and up sampled to the same rate as the environmental and physiological data to 1Hz
 227 to remain consistent with the other data. Linear interpolation was used as a mechanism to
 228 sample the data. If the two known points are given by the coordinates (x_1, y_1) and (x_2, y_2) ,
 229 The linear interpolant is the straight line between these points. For a value x in the interval
 230 (x_2, x_1) , the value y along the straight line is given from Equation 2 of slopes as shown
 231 below:

232

$$y = y_1 + \frac{(x - x_1)(y_2 - y_1)}{x_2 - x_1}$$

233

234 Equation 2. Linear interpolation formula used to estimate the value of (y) at a given point (x)
235 based on two known data points (x1, y1) and (x2 and y2).

236

237 Normalisation was used following this on all variables to bring them within the same range
238 for both data analysis and machine learning. Outliers in the data were identified using a visual
239 inspection of the resulting participant dataset. Any extreme or unusual values that fell outside
240 the predefined range were flagged for further review. In this case outliers were carefully
241 assessed to determine whether they were due to experimental variability, participant
242 noncompliance or technical errors before deciding whether to retrain or exclude them from
243 the final analysis. It is important to note that during this process, the variables; Carbon
244 Dioxide and Volatile Organic Compound were removed from the sample as there were issues
245 with logging the data. The data recorded from these variables at data collection read
246 'ERROR' within the downloaded CSV file and hence was removed and replaced with a zero.
247 Furthermore, the values for PM2.5, PM10 and NO2 recorded as a zero indicate that there was
248 no detectable pollution for these parameters at the time of measurement. As such, these
249 values represent actual readings. Therefore, in the process of normalisation, these zero values
250 were used as they were, reflecting on the absence of pollution during those specific time
251 points.

252

253

254 LIMITATIONS

255 The DigitalExposome study has few limitations in that HRV was recorded for all but three
256 participants due to issues around the Empatica E4 wearable. For these users the HRV data
257 was removed from the dataset and replaced with a zero. Although considering the
258 practicalities of a real-world experiment being used in this study the number of participants
259 who took part could involve more to create an even more diverse set of users. Finally,
260 although Volatile Organic Compound and Carbon Dioxide sensors were used in the
261 equipment setup, these variables had to be discounted prior to normalisation due to a sensors
262 malfunction.

263

264

265 ETHICS STATEMENT

266 The ethical approval for this study was granted by Nottingham Trent University Invasive
267 Ethics Committee (Document No. 068/2020). All procedures performed in this study
268 involving the human participants were carried out by the ethical standards and guidelines of
269 Nottingham Trent University. Informed consent was obtained from all participants who took
270 part in this study. All data was handled with strict confidence and pseudonymised in line with
271 ethical agreements from NTU and to protect participants.

272



273 **CRedit AUTHOR STATEMENT**

274 *Thomas Johnson: Conceptualisation, software, validation, data curation, Writing – original*
275 *draft, Writing – Review & Editing.*

277 **ACKNOWLEDGEMENTS**

278 This research did not receive any specific grant from funding agencies in the public,
279 commercial, or not-for-profit sectors.

282 **DECLARATION OF COMPETING INTERESTS**

283 The author declares that they have no known competing financial interests or personal
284 relationships that could have appeared to influence the work reported in this paper.

286 **REFERENCES**

- 287
- 288 1) T, Johnson. (2025), “DigitalExposome: A Dataset for Wellbeing Classification using
289 Environmental Air Quality and Human Physiological data”, Mendeley Data, V2, doi:
290 10.17632/mbwxy48223.2
 - 291 2) T, Johnson., E, Kanjo. & K, Woodward. DigitalExposome: quantifying impact of
292 urban environment on wellbeing using sensor fusion and deep
293 learning. *Comput.Urban Sci.* **3**, 14 (2023). [https://doi.org/10.1007/s43762-023-](https://doi.org/10.1007/s43762-023-00088-9)
294 [00088-9](https://doi.org/10.1007/s43762-023-00088-9)
 - 295 3) I, Manisalidis, E, Stavropoulou, A, Stavropoulos, E, Bezirtzoglou. Environmental and
296 Health Impacts of Air Pollution: A Review. *Front Public Health.* 2020 Feb 20;8:14. doi:
297 10.3389/fpubh.2020.00014. PMID: 32154200; PMCID: PMC7044178.
 - 298 4) H, Haddad, A, Nazelle. The role of personal air pollution sensors and smartphone
299 technology in changing travel behaviour. *Journal of Transport & Health.* Volume 11.
300 (2018) Pages 230-243. <https://doi.org/10.1016/j.jth.2018.08.001>.
 - 301 5) J, King. Air pollution, mental health, and implications for urban design: a review.
302 (2018). *Journal of Urban Design and Mental Health, March*, 4:6.
 - 303 6) L, Sarmiento, N, Wagner, A, Zaklan. *The air quality and well-being effects of low*
304 *emission zones. Journal of Public Economics. Volume 227 (2023)*
305 <https://doi.org/10.1016/j.jpubeco.2023.105014>.
 - 306 7) T. Johnson, K. Woodward and E. Kanjo, "Emotion on the Edge: Air Quality Sensors
307 Decoded as a Real-World Emotion Indicator," *2024 IEEE International Conference on*
308 *Pervasive Computing and Communications Workshops and other Affiliated Events (PerCom*
309 *Workshops)*, Biarritz, France, 2024, pp. 267-272, doi:
310 10.1109/PerComWorkshops59983.2024.10502563.
 - 311 8) *World Health Organisation, 2024. Ambient Air Quality and Health. Available at:*
312 [https://www.who.int/en/news-room/fact-sheets/detail/ambient-\(outdoor\)-air-](https://www.who.int/en/news-room/fact-sheets/detail/ambient-(outdoor)-air-quality-and-health)
313 [quality-and-health](https://www.who.int/en/news-room/fact-sheets/detail/ambient-(outdoor)-air-quality-and-health) . (accessed 7th January 2025)

- 314 9) Department for Environment, Farms and Rural Affairs (2019). *Clean Air Strategy*.
315 Available at: <https://www.gov.uk/government/publications/clean-air-strategy-2019>.
316 (accessed 7th January 2025)
- 317 10) R. Kharoua. Air Quality and Health Impact Dataset. Kaggle Datasets. 2024.
318 10.34740/kaggle/dsv/8675842
- 319 11) *CitiesHealth*. *Air Pollution and Mental Health: identifying Short-Term Human Impacts*
320 *of Air Pollution*. 2023. Kaggle Datasets.
321 <https://www.kaggle.com/datasets/thedevastator/air-pollution-and-mental-health>
- 322 12) International Wellbeing Group (2013). *Personal Wellbeing Index: 5th Edition*.
323 Melbourne: Australian Centre on Quality of Life, Deakin University.
324 (<http://www.deakin.edu.au/research/acqol/instruments/wellbeing-index/index.php>)
- 325 13) Woodward, K., Kanjo, E., Oikonomou, A., & Chamberlain, A. (2020). LabelSens:
326 enabling real-time sensor data labelling at the point of collection using an artificial
327 intelligence-based approach. *Personal and Ubiquitous Computing*, 24, 709–722.
328 <https://doi.org/10.1007/s00779-020-01427-x/Published>
- 329 14) Johnson, T., & Woodward, K. (2025). *Enviro-IoT: Calibrating Low-Cost*
330 *Environmental Sensors in Urban Settings*. <http://arxiv.org/abs/2502.07596>