

Cite as: Tekler, Z. D., Ono, E., Peng, Y., Zhan, S., Lasternas, B., & Chong, A. (2022, December). ROBOD, room-level occupancy and building operation dataset. In Building Simulation (Vol. 15, No. 12, pp. 2127-2137). Tsinghua University Press. doi: <https://doi.org/10.1007/s12273-022-0925-9>

ROBOD, Room-level Occupancy and Building Operation Dataset

Zeynep Duygu Tekler¹, Eikichi Ono¹, Yuzhen Peng¹, Sicheng Zhan¹, Bertrand Lasternas¹, and Adrian Chong^{1,*}

¹Department of the Built Environment, National University of Singapore, 4 Architecture Drive, 117566, Singapore
*corresponding author(s): Adrian Chong (adrian.chong@nus.edu.sg)

ABSTRACT

The availability of the building's operation data and occupancy information has been crucial to support the evaluation of existing models and development of new data-driven approaches. This paper describes a comprehensive dataset consisting of indoor environmental conditions, Wi-Fi connected devices, energy consumption of end uses (i.e., HVAC, lighting, plug loads and fans), HVAC operations, and outdoor weather conditions collected through various heterogeneous sensors together with the ground truth occupant presence and count information for five rooms located in a university environment. The five rooms include two different-sized lecture rooms, an office space for administrative staff, an office space for researchers, and a library space accessible to all students. A total of 181 days of data was collected from all five rooms at a sampling resolution of 5 minutes. This dataset can be used for benchmarking and fostering data-driven approaches in the field of occupancy prediction and occupant behaviour modelling, building simulation and control, energy forecasting and various building analytics.

Background & Summary

The building sector is currently responsible for more than one-third of the global energy consumption and approximately 40% of the total direct and indirect CO₂ emissions in the world¹. As energy demand from the building sector continues to rise due to rapid urbanisation around the globe, significant efforts have been dedicated to improving building energy efficiency while maintaining reliable building operations and high indoor environmental quality for the occupants.

To achieve this goal, researchers have relied on various modelling approaches and simulation tools to help model and quantify building energy use based on different factors, such as climatic regions², architectural design³, environmental conditions⁴, occupancy and occupant interactions with building systems^{5,6}. These models often require the systematic collection and analysis of various real-world inputs such as the buildings' operational data, energy use and occupancy information to derive meaningful insights and develop effective strategies for reducing building energy use⁷. For instance, the availability of various building data is necessary for physics-based energy models to define model assumptions and inform model calibration⁸. At the same time, data-driven or machine learning-based models require a sufficient amount of training data to produce reliable prediction results^{9,10}.

However, the collection of such real-world datasets is often challenging in reality. Firstly, it requires the installation of different sensors within each room in the building, which can incur a considerable cost depending on the number of rooms and the size of the target building. After the sensors have been deployed, another significant cost comes from the regular maintenance of these sensors to ensure they stay operational and the data storage services procured to safely store and manage the data collected. Secondly, the integration of the sensor data collected can also create additional hurdles due to the issues related to intermittent sensor failure and nonstandard sampling frequencies used by different sensor manufacturers. Lastly, the collection of occupancy data, which is often performed in person or through surveillance cameras, is labour intensive and may also encounter resistance from the building occupants due to privacy concerns¹¹. Despite these challenges, there has been a sustained effort within the building science community to encourage the release of public building datasets to facilitate collaborative and reproducible research. Some examples of these public datasets include: the Building Data Genome Project¹², which contains the energy metering data for 1,636 non-residential buildings; BLOND¹³, an energy consumption dataset for appliances in an office building; fLEECe¹⁴, an energy use and occupant behaviour dataset for residential buildings; CU-BEMS¹⁵, which contains the electrical consumption and indoor environmental sensor data for a smart office building; as well as other commercial and residential datasets containing energy consumption data, building operation data, occupancy data, indoor environmental quality data or different combinations of these data categories¹⁶⁻¹⁸. So far, no public dataset contains all the highlighted data categories alongside high-resolution occupancy data across multiple room types.

In this paper, we release ROBOD, a Room-level Occupancy and Building Operation Dataset. To the best of our knowledge,

39 this is the most comprehensive dataset that contains room-level occupant presence and count information integrated with
40 building operation data from different room types in a university environment. The dataset consists of a wide range of data
41 categories, including indoor environmental conditions, Wi-Fi connected devices, building energy end-uses (i.e., HVAC, lighting,
42 plug loads, and fans), HVAC operations, and local outdoor weather conditions collected through various heterogeneous sensors
43 together with the ground truth occupant presence and count information for five different rooms. Through the use of this dataset,
44 researchers from different fields can benefit from various applications, including but not limited to occupancy prediction and
45 occupant behaviour modelling, building simulation and control, energy forecasting, and building analytics.

46 Methods

47 Building & Room Characteristics

48 The building considered in this dataset is the School of Design and Environment 4 (SDE4) building located at the National
49 University of Singapore. SDE4 is a 6-story academic building spanning 8,588 square meters. It is the first newly-built net-zero
50 energy building in Singapore and the first building in South Asia that obtained a Zero Energy Certification. We collected the
51 room occupancy and building operation data for five rooms located at different building levels, as visualized in Figure 1. The
52 five rooms include two different-sized lecture rooms (Room 1 and Room 2), an office space for administrative staff (Room 3),
53 an office space for researchers (Room 4), and a library space for students (Room 5). The detailed description of each room is
54 provided in Table 1.

55 Data Categories Overview and Collection

56 A building management system (BMS) is currently deployed in the building to help monitor and manage the building's
57 mechanical and electrical systems. As part of BMS, various sensors are installed throughout the study building to collect
58 information about the building's energy consumption, HVAC conditions and outdoor weather conditions. The BACnet Protocol
59 is used to retrieve these sensor measurement data to be stored in the PI Data Archive, which is a feature within the OSISoft PI
60 system. The PI Data Archive serves as an industry-standard data management system for storing time-series data and allows
61 users to perform remote data extraction using various RESTful API services. On top of that, we have installed standalone
62 indoor environmental quality (IEQ) sensors to measure the indoor environmental conditions within each room. Apart from
63 these data categories, we have also tapped into the surveillance cameras and Wi-Fi access points within the study building to
64 obtain room-level occupancy information and the number of Wi-Fi-connected devices, respectively. All data measurements
65 from different sensors were queried with a sampling frequency of 5 minutes before they are integrated to form ROBOD. A
66 5-minute sampling interval was chosen to strike a good balance between data representativeness and data collection cost.

67 The following section describes the details of each data category found within the dataset. More detailed information about
68 the data units, sensor types, sensor range and accuracy specifications from the manufacturers are provided in Table 2.

69 Indoor Environmental Quality Data

70 The indoor environmental data represent the measurements for indoor environmental quality, which include VOC (volatile
71 organic compound), sound pressure level, relative humidity, indoor air temperature, illuminance, PM2.5 (particulate matter),
72 and CO₂ concentration levels. A dedicated IEQ monitoring unit is installed in each room and its location within the room is
73 provided in Table 3.

74 Wi-Fi Data

75 The Wi-Fi data represents the number of Wi-Fi-enabled devices connected to the routers installed in each room. Some examples
76 of these devices include mobile devices (i.e., smartphones and laptops), which connect to the nearest routers depending on
77 their users' movement patterns and location, and stationary devices whose location remains fixed mainly within the room (i.e.,
78 printers and desktops). Based on this, number of Wi-Fi connected devices are higher than the number of occupants in the room
79 as the Wi-Fi data contain both mobile and stationary devices. However for further processing, any filtering logic could be
80 introduced to filter out the stationary devices within the room to identify the number of mobile devices which could be used
81 to infer the occupant count. The raw Wi-Fi dataset contains the logs for every device that connects to different access points
82 across the campus and is stored in a Hive SQL database. By querying the relevant logs through the Open Database Connectivity
83 (ODBC) API, the raw Wi-Fi logs are processed to extract the number of connected devices by counting the number of unique
84 MAC addresses recorded during a 5-minute interval for each room.

85 Energy Data

86 The energy data represents the energy consumption values of the building's electrical end uses such as HVAC, lighting, plug
87 loads, and ceiling fans. For HVAC energy consumption,

- 88 • Room 1 and Room 2 are conditioned by Fan Coil Units (FCU), with the chilled water supplied by a district chiller plant
89 and the supply airflow rate controlled by variable speed fans.

- 90 • Room 3, Room 4, and Room 5 are conditioned by Air Handling Units (AHU), which are connected to multiple rooms in
 91 the building. It should be noted that the chilled water energy and AHU fan energy consumption are same for the rooms
 92 which share the same AHU. In this dataset, Room 4 and Room 5 share the same AHU and therefore have identical energy
 93 consumption values for chilled water and AHU fan energy.

94 The energy consumption data of lighting, plug loads, and ceiling fans are collected through electrical meters and the number
 95 of each end use (i.e., lighting, plug loads and ceiling fans) found in each room is listed in Table 4. In this case, the number
 96 of lighting units refers to the number of luminaries in each room. Similarly, plug load units are represented as the number of
 97 inbuilt 13A double electrical sockets available in each room. It is also worth highlighting that each room may contain different
 98 types of plug loads depending on its space function. For instance, Room 1 and Room 2 contain mostly laptops and projectors,
 99 Room 3 and Room 4 contain different number of monitors, laptops, desktops, and printers; while Room 5 contains mostly
 100 laptops and printers.

101 **HVAC Operations Data**

102 The HVAC operations data represent the different parameters and settings that the building's HVAC system operates within.
 103 Some of these measurements include supply airflow, damper position, temperature setpoint, cooling coil valve position and
 104 cooling coil valve command, AHU/FCU fan speed, offcoil air temperature, offcoil temperature setpoint, supply air humidity,
 105 pressure across filter, supply air static pressure and supply air temperature. It should be noted that the building uses a dedicated
 106 outdoor air system for air supply, so the CO_2 level of incoming air is identical to the outdoor CO_2 level.

107 The temperate setpoint in all rooms is conditioned by Proportional Integral Derivative (PID) control against the thermostat
 108 temperature setpoint set by the room occupants. As Room 1 and Room 2 are conditioned by FCUs, they do not contain data
 109 measurements related to VAV. The availability of the HVAC operations is also indicated in Table 2 as a footnote.

- 110 • Room 1 and Room 2 have dedicated FCUs supplying airflow rates at 3,375 and 2,025 cubic meter per hour (CMH),
 111 respectively.
- 112 • Room 3 has a VAV airflow rate of 900 CMH and is air-conditioned by a AHU with a supply airflow rate of 1,3470
 113 CMH, serving five other rooms in the building. Room 4 has a VAV airflow rate of 3,192 CMH, while Room 5 has a VAV
 114 airflow rate of 1,944 CMH. Both rooms are air-conditioned by the same AHU with a supply airflow rate of 14,560 CMH,
 115 supplying chilled air to eleven other rooms in the building.

116 **Outdoor Weather Data**

117 The outdoor weather data is measured by a local weather station installed on the roof of the study building. Measurements
 118 include barometric pressure, dry bulb temperature, global horizontal solar radiation, wind direction and speed, outdoor CO_2 ,
 119 and relative humidity.

120 **Occupancy Data**

121 The occupancy data contains both the occupant presence and number of occupants present in each room. This information was
 122 collected by monitoring the occupants' movement through surveillance camera footage and manually counting the number of
 123 occupants. At any point in time during the data collection process, any identifiers (i.e., names and personal details) that reveal
 124 occupants' identity were not collected nor stored in this dataset. The protocols for the data collection has been approved by the
 125 host university's Institutional Review Board (NUS-IRB-2021-31).

126 **Data Pre-processing**

127 This section details the data pre-processing steps performed when merging the data categories described above to form ROBOD.
 128 These steps involve formatting the timestamp information for each data category to follow the same ISO 8601 date-time format
 129 (i.e., YYYY-MM-DD HH:MM +-HH:MM), starting with the year information, followed by the month, day, hour, minute,
 130 and time zone offset from UTC. Each data measurement follows a 5-minute sampling interval, starting with the 0th minute,
 131 followed by the 5th minute, the 10th minute, and so on till the 55th minute during each hour. After these standardisation steps
 132 are performed, the six categories are merged within the same timestep using their timestamp information as the primary key.

133 **Data Records**

134 ROBOD is currently hosted on figshare¹⁹ and consists of five comma-separated value (CSV) files. Each file contains the
 135 combined data for each room for all six data categories described in Table 2. Each data measurement also contains the
 136 timestamp information corresponding to the time when the data measurement was recorded and followed the date-time format:
 137 YYYY-MM-DD HH:MM +08:00. The last component (i.e., +08:00) indicates a UTC offset of +8 hours as the data collection

138 was conducted in the tropical island of Singapore. Given that the data measurements followed a sampling interval of 5 minutes,
139 this corresponds to 288 data points recorded per day. The data collection period spanned between September 2021 and
140 December 2021, where the sensor data collected during the weekends were excluded. Furthermore, there were also specific
141 days during the data collection period when several of the sensors were not working correctly for certain rooms, leading to the
142 data collected during these periods being dropped from the final dataset. In the end, a total of 181 days of data was collected
143 from the five rooms, where Room 1, Room 2 and Room 3 contributed 29 days of data separately while Room 4 and Room 5
144 contributed 47 days of data each. Apart from the timestamp information that is stored in the string format, the occupancy count
145 and presence information is stored as integers, while the rest of the data fields are represented as floating numbers.

146 Technical Validation

147 This section presents the technical validity of our dataset starting with a preliminary analysis of missing data and various
148 visualisations involving occupant count, outdoor environmental condition, room air temperature, room temperature setpoint,
149 and energy consumption based on the raw dataset.

150 **Missing data** A preliminary analysis of the dataset highlighted a small number of missing data points for each room in
151 ROBOD due to issues related to intermittent sensor failure. Table 5 presents a detailed breakdown of the amount of missing
152 data found in each column and for each room. The temporal relationship of the missing data is also presented in Figure 2. It
153 should be reiterated that the datasets for Room 1 and 2 do not contain columns related to VAV (i.e., Supply Air Flow, Damper
154 Position, Cooling Coil Valve Position and Command, Offcoil Temperature Setpoint, Offcoil Air Temperature, Pressure across
155 Filter, and Supply Air Humidity) as they are conditioned by FCUs, therefore they are not included in the dataset.

156 **Occupant count** Figure 3 presents the average occupant count for each room on an average weekday. Based on the
157 occupancy fluctuations, it can be observed that the occupant count patterns differ slightly among different rooms. More
158 specifically, the occupant count for Room 1 and Room 2 experience heavy fluctuations throughout the day compared to other
159 rooms. In particular, we observed three distinct peaks in Room 2 that occur at 11 am, 1 pm, and 3 pm, which can be explained
160 by the block lectures that are regularly scheduled during these periods. Room 3 presents a regular office schedule with the
161 office workers arriving at work between 8 to 10 am and leaving the office at the end of the workday between 6 to 8 pm. The
162 occupants in Room 4 are observed to follow a flexible work schedule where the last departure times for some occupants can
163 stretch late into the night after midnight. Lastly, we can observe a sharp increase in occupancy levels in Room 5 from zero at 9
164 am, followed by a sharp drop back to zero at 9 pm every day, corresponding with the operational hours of the library space.

165 **Outdoor environmental condition** Figure 4 shows the monitored outdoor conditions of dry-bulb temperature, global
166 horizontal solar radiation, relative humidity, and CO₂. As the data was collected from the study building located in the tropic,
167 the outdoor dry-bulb temperature ranges from 22.6°C to 35.5°C, where temperatures tend to rise to higher levels in the afternoon
168 (i.e., 12 pm to 4 pm). The global horizontal solar radiation can reach over 1000 W/m² between 11 pm and 3 pm. At the same
169 time, the relative humidity ranges from 40% to 100%, of which over 98% accounts for the primary ratio (25%). The cooling
170 systems process dry-bulb temperature and relative humidity to deliver the required supply air flow to cool down the internal
171 thermal zones within the building, while removing thermal energy generated by the solar radiation. The outdoor CO₂ levels
172 span between 439 ppm to 510 ppm, which is used as the basis of maintaining the indoor CO₂ levels at a standard or comfortable
173 range.

174 **Room temperature setpoint** Figure 5 depicts the distributions of temperature setpoints for each of the five rooms. As the
175 occupants' thermal sensation is subjective, the temperature setpoints may differ among the rooms and during different periods
176 of the day. Room 1 and Room 2 show a wide range of temperature setpoints, ranging from 22°C to 27.2°C, and from 22°C to
177 27.7°C, respectively. Room 4 shifted the setpoints to the range of 25.3°C to 28°C. Unlike the other rooms, Room 3 and Room 5
178 kept the temperature setpoints consistently at 25°C and 26°C, respectively.

179 **Room air temperature** Figure 6 shows a heatmap of the average indoor air temperature or thermal distribution at different
180 time periods during the day for each room. The vertical axis indicates each of the five rooms, and the horizontal axis shows the
181 time of the day. For example, it can be observed that the air temperatures in Room 1, Room 2, and Room 3 tends to be cooler
182 than Room 4 and Room 5. Moreover, air temperatures in the afternoon are also warmer than in the morning for all five rooms.

183 **Energy consumption** Figure 7 summarizes energy consumption of space cooling, plug load, and lighting in each room.
184 The cooling energy consumption that combines the energy consumed by chilled water and AHU/FCU fans. Since Room 4
185 and Room 5 are air-conditioned by the same AHU, their cooling energy consumption is not separated for this analysis. The
186 difference in the cooling demand among the rooms can be explained by the differences in room functions and room area.
187 For instance, Room 1 and Room 2 are used as lecture spaces with similar indoor areas resulting in identical cooling energy
188 consumption and schedules. Similarly, the cooling energy and schedules for Room 3 and Rooms 4+5 are similar in terms of
189 its pattern as all three rooms function as multi-occupant offices (i.e., Room 3 and Room 4). Devices that are connected to
190 electrical sockets can be classified into two groups: non-mobile devices located in the rooms and portable devices. The former
191 contributes 24-hour plug load consumption, including the small energy consumption when the devices enter into idle modes.

192 The latter only need the electricity from electrical sockets when their owners occupy the rooms. For instance, the plug load
193 consumption in five rooms are nearly constant before 7 am. Furthermore, this consumption in Room 1, Room 2, and Room
194 5 increased simultaneously from 9 am. Similar to the plug load consumption, the energy consumption of lighting is closely
195 related to occupants' room usages. Therefore, the lighting demand increases from 9 am in most of the rooms.

196 **Usage Notes**

197 The dataset provided in this paper is in the CSV format for all rooms and has a total file size of 20 MB. The CSV data format
198 allows the files to be easily imported by most spreadsheet programs and databases. It is also easy to work with due to its
199 human-readable format and can be readily processed and analysed by most popular programming languages such as Python,
200 Java, Javascript, and R.

201 Due to the presence of missing data in the dataset, we have also included several data post-processing steps as a reference
202 for researchers who would like to use the existing dataset. These steps involves imputing the dataset's missing or erroneous
203 sensor data by using the *missingpy* imputation library. While different imputation algorithms have been utilised in past studies²⁰,
204 a Random Forest-based imputation algorithm (i.e., MissForest²¹) is adopted in this case by performing column-wise imputation
205 in an iterative fashion. The algorithm begins by imputing the column with the least number of missing values (i.e., candidate
206 column) and filling the missing values in the remaining columns with an initial guess, such as the column's mean. Following
207 this, a Random Forest (RF) model is trained by setting the candidate column as the output variable and the remaining columns
208 as the model's input for those rows that do not contain missing values in the candidate column. After the RF model has been
209 trained, it is used to impute the missing values in the candidate column before moving on to the next candidate column with the
210 second smallest number of missing values. This process is repeated for each column containing missing values over multiple
211 iterations until the difference between the dataset imputed in the previous round and the newly imputed dataset increases for the
212 first time.

213 **Code availability**

214 All data post-processing steps and visualisations performed in this manuscript are implemented using Python 3.6 and public
215 libraries including Numpy and Pandas for data manipulation, Matplotlib, Seaborn, and Missingno for data visualisation, and
216 Missingpy for data imputation. A step-by-step guide has been compiled within a single Jupyter notebook and made available
217 on Github (<https://github.com/ideas-lab-nus/robod.git>).

218 **References**

- 219 1. Buildings, I. T. Iea: Paris. <https://www.iea.org/reports/tracking-buildings-2020>, Accessed:2022-01-21 (2020).
- 220 2. Orouji, P. *et al.* Atlas of heating: Identifying regional climate-dependent heat demands in residential buildings of iran. In *Building Simulation*, vol. 14, 857–869 (Springer, 2021).
- 221 3. Ataman, C. & Dino, İ. G. Performative design processes in architectural practices in turkey: architects' perception. *Archit. Eng. Des. Manag.* 1–15, <https://doi.org/10.1080/17452007.2021.1995315> (2021).
- 222 4. Aydin, E. E. & Jakubiec, J. A. Sensitivity analysis of sustainable urban design parameters : Thermal comfort , urban heat
223 island , energy , daylight , and ventilation in singapore (2018).
- 224 5. Tekler, Z., Low, R. & Blessing, L. Using smart technologies to identify occupancy and plug-in appliance interaction
225 patterns in an office environment. In *IOP Conference Series: Materials Science and Engineering*, vol. 609, 062010,
226 <https://doi.org/10.1088/1757-899X/609/6/062010> (IOP Publishing, 2019).
- 227 6. Peng, Y., Nagy, Z. & Schlüter, A. Temperature-preference learning with neural networks for occupant-centric building
228 indoor climate controls. *Build. Environ.* **154**, 296–308, [10.1016/j.buildenv.2019.01.036](https://doi.org/10.1016/j.buildenv.2019.01.036) (2019).
- 229 7. Ding, Y. *et al.* Review on occupancy detection and prediction in building simulation. In *Building Simulation*, 1–24
230 (Springer, 2021).
- 231 8. Chong, A., Gu, Y. & Jia, H. Calibrating building energy simulation models: A review of the basics to guide future work.
232 *Energy Build.* **253**, 111533, <https://doi.org/10.1016/j.enbuild.2021.111533> (2021).
- 233 9. Zhan, S. & Chong, A. Data requirements and performance evaluation of model predictive control in buildings: A modeling
234 perspective. *Renew. Sustain. Energy Rev.* 110835, <https://doi.org/10.1016/j.rser.2021.110835> (2021).
- 235 10. Peng, Y., Rysanek, A., Nagy, Z. & Schlüter, A. Using machine learning techniques for occupancy-prediction-based cooling
236 control in office buildings. *Appl. Energy* **211**, 1343–1358, [10.1016/j.apenergy.2017.12.002](https://doi.org/10.1016/j.apenergy.2017.12.002) (2018).

- 239 11. Tekler, Z. D., Low, R., Gunay, B., Andersen, R. K. & Blessing, L. A scalable bluetooth low energy approach to identify
240 occupancy patterns and profiles in office spaces. *Build. Environ.* **171**, 106681, <https://doi.org/10.1016/j.buildenv.2020.106681> (2020).
- 242 12. Miller, C. *et al.* the building data genome project 2, energy meter data from the ashrae great energy predictor iii competition.
243 *Sci. data* **7**, 1–13, <https://doi.org/10.1038/s41597-020-00712-x> (2020).
- 244 13. Kriechbaumer, T. & Jacobsen, H.-A. Blond, a building-level office environment dataset of typical electrical appliances. *Sci.*
245 *data* **5**, 1–14, <https://doi.org/10.1038/sdata.2018.48> (2018).
- 246 14. Paige, F., Agee, P. & Jazizadeh, F. fleece, an energy use and occupant behavior dataset for net-zero energy affordable
247 senior residential buildings. *Sci. data* **6**, 1–9, <https://doi.org/10.1038/s41597-019-0275-3> (2019).
- 248 15. Pipattanasompon, M. *et al.* Cu-bems, smart building electricity consumption and indoor environmental sensor datasets.
249 *Sci. Data* **7**, 1–14, <https://doi.org/10.1038/s41597-020-00582-3> (2020).
- 250 16. Tekler, Z. D. *et al.* Near-real-time plug load identification using low-frequency power data in office spaces: Experiments
251 and applications. *Appl. Energy* **275**, 115391, <https://doi.org/10.1016/j.apenergy.2020.115391> (2020).
- 252 17. Schwei, J. H. *et al.* Room-level occupant counts and environmental quality from heterogeneous sensing modalities in a
253 smart building. *Sci. data* **6**, 1–11, <https://doi.org/10.1038/s41597-019-0274-4> (2019).
- 254 18. Li, H., Wang, Z. & Hong, T. A synthetic building operation dataset. *Sci. data* **8**, 1–13, <https://doi.org/10.1038/s41597-021-00989-6> (2021).
- 256 19. Tekler, Z. D. *et al.* Robod, room-level occupancy and building operation dataset. *figshare* <https://doi.org/10.6084/m9.figshare.19234530.v5> (2022).
- 258 20. Low, R., Tekler, Z. D. & Cheah, L. Predicting commercial vehicle parking duration using generative adversarial multiple
259 imputation networks. *Transp. Res.* **2674**, 820–831, <https://doi.org/10.1177/0361198120932166> (2020).
- 260 21. Stekhoven, D. J. & Bühlmann, P. Missforest—non-parametric missing value imputation for mixed-type data. *Bioinformatics*
261 **28**, 112–118, <https://doi.org/10.1093/bioinformatics/btr597> (2012).

262 Acknowledgements

263 This research project is supported by the National Research Foundation, Singapore, and Ministry of National Development,
264 Singapore under its Cities of Tomorrow R&D Programme (CoT Award COT-V4-2020-5). Any opinions, findings and
265 conclusions or recommendations expressed in this material are those of the author(s) and do not reflect the views of National
266 Research Foundation, Singapore and Ministry of National Development, Singapore.

267 Author contributions

268 **Zeynep Duygu Tekler:** Conceptualisation, Methodology, Software, Validation, Data Curation, Writing - Original Draft.
269 **Eikichi Ono:** Investigation, Writing - Review & Editing. **Yuzhen Peng:** Visualisation, Writing - Original Draft. **Sicheng**
270 **Zhan:** Visualisation, Writing - Review & Editing. **Bertrand Lasternas:** Resources. **Adrian Chong:** Writing - Review &
271 Editing, Supervision, Funding Acquisition.

272 Competing interests

273 The authors declare no competing interests.

274 Figures & Tables

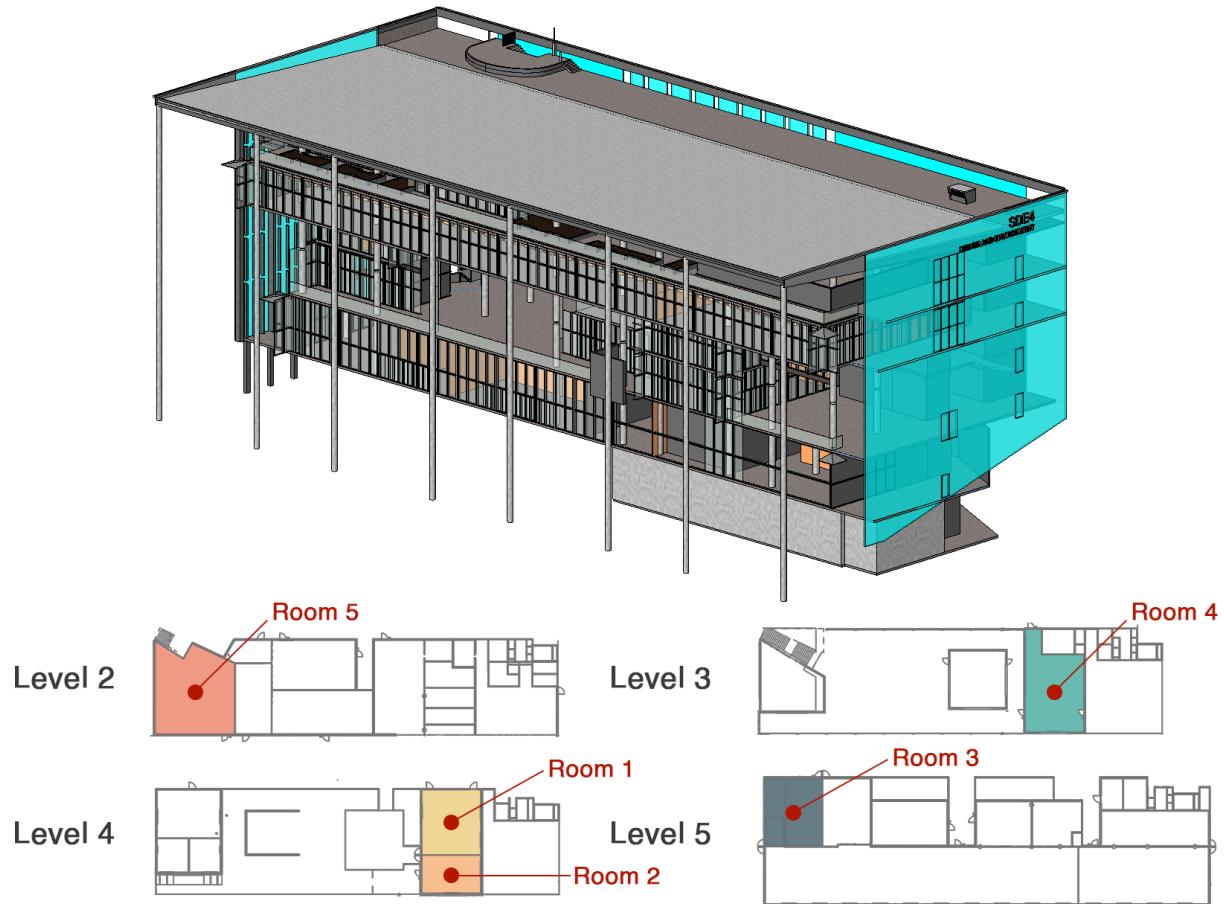


Figure 1. Study building (top) and room layouts corresponding to the building levels (bottom)

Table 1. Room descriptions.

Room	Space Function	Occupant Type	Level	Floor Area [m ²]	Floor to Ceiling Height [m]	Room Volume [m ³]	Seating Capacity [person]	Maximum Occupancy Density [m ² /person]
Room 1	Lecture room	Students	4	118.6	4.1	486.2	40	3.0
Room 2	Lecture room	Students	4	53.7	4.1	220.2	40	1.3
Room 3	Office space	Administrative staff	5	98.4	4.2	413.2	15	6.6
Room 4	Office space	Researchers	3	141.9	4.1	581.7	25	5.6
Room 5	Library space	Students	2	182.8	7.5	1363.3	36	5.0

Table 2. Data categories and sensor specifications.

Data Category	Measured Variable	Data Unit	Sensor Type (Brand)	Sensor Range	Sensor Accuracy
Indoor environmental quality	VOC	ppb	IAQ monitoring unit (Awair Omni)	0 – 60000	±10%
	Sound pressure level	dB(A)		Not specified	Not specified
	Relative humidity	%RH		0 – 100	±2%RH
	Air Temperature	°C		-40 – 125	±0.2 °C
	Illuminance	lux		0 – 64000	Not specified
	PM2.5	µg/m³		0 – 1000	±15µg/m³ or ±15%
	CO ₂	ppm		400 – 5000	±75 ppm or ±10%
Wi-Fi	Wi-Fi connected devices	Number	Wi-Fi Router (Cisco)	NA ^a	NA ^a
Energy	Ceiling fan energy	kWh	Energy meter (Schneider Electric Acti9 iEM3000)	0 – 999999	±1%
	Lighting energy		BTU meter (Integra Metering CALEC ST II)		±2%
	Plug load energy		Energy meter (Schneider Electric PM5300)		±0.5%
	Chilled water energy				
	AHU/FCU fan energy				
HVAC operations	Supply air flow ^b	CMH	VAV box (Johnson Controls)	0 – 3375	±15%
	Damper position ^b	%		0 – 100	NA ^a
	Temperature setpoint	°C	NA ^a	NA ^a	NA ^a
	Cooling coil valve position ^b	%	Valve (Johnson Controls)	0 – 100	NA ^a
	Cooling coil valve command ^b				
	AHU/FCU fan speed	Hz	Variable speed drive (ABB)	0 – 50	±0.2%
	Offcoil air temperature ^b	°C	NTC thermistor (GreyStone TSDC series)	-40 – 60	±0.2 °C
	Offcoil temperature setpoint ^b	°C	NA ^a	NA ^a	NA ^a
	Supply air humidity ^b	%RH	Capacitive (GreyStone HSDT series)	0 – 100	±2%RH
	Pressure across filter ^b	Pa	Capacitive (Setra Model 264)	Not specified	±1%
	Supply air static pressure				
	Supply air temperature	°C	NTC thermistor (GreyStone TSAP series)	-40 – 60	±0.2 °C
Outdoor weather ^c	Barometric pressure	hPa	Piezoresistive	600 – 1100	±0.5 hPa @20 °C
	Dry bulb temperature	°C	Pt100	-40 – 60	±0.15°C ±0.1%
	Global horizontal solar radiation	W/m²	Thermopile	0 – 2000	2nd class pyranometer
	Wind direction	° (Degree)	Ultrasonic	0 – 360	±2° RMSE
	Wind speed	m/s	Ultrasonic	0 – 60	±0.2 m/s or ±2%
	CO ₂	ppm	Non-dispersive infrared	0 – 2000	±(5 ppm + 2%)
	Relative humidity	%RH	Capacitive	0 – 100	±1.5%RH
Occupancy	Occupant presence	Binary (1/0)	Surveillance camera (XeronVision 2M HD IP Vari-Focal Lens Dome)	NA ^a	NA ^a
	Occupant count	Number			

^a NA refers to "Not Applicable".

^b Indicated measurements are not applicable for Room 1 & Room 2.

^c All the outdoor weather data were collected by a weather station (Delta OHM HD52.3D).

Table 3. Locations of retrofitted sensors (i.e., surveillance camera and IEQ units).

Room	Surveillance Camera	IEQ Units
Room 1	Two surveillance cameras outside of two doors	Mounted to an east side column vertically on wall
Room 2	One surveillance camera outside of the door	Mounted to an east side column vertically on wall
Room 3	On the top corner of the rooms near the entrance doors	Mounted to a west side column vertically on wall
Room 4	On the top corner of the rooms near the entrance doors	Mounted to a west side column vertically on wall
Room 5	Three surveillance cameras inside the room	Mounted to a east side column vertically on wall

Table 4. Number of end uses in each room.

Room	No. of Ceiling Fans	No. of Luminaries	No. of 13A Double Sockets
Room 1	6	20	26
Room 2	4	12	14
Room 3	4	14	9
Room 4	6	32	20
Room 5	6	11	12

Table 5. A detailed breakdown of the amount of missing data in the relevant columns of each room.

Room	Total	Missing Data	Column Name
Room 1	8352	9 10 14	supply_air_pressure and ahu_fan_speed chilled_water_energy and ahu_fan_energy voc, sound_pressure_level, indoor_relative_humidity, illuminance, pm2.5, indoor_co2
Room 2	8352	30	voc, sound_pressure_level, indoor_relative_humidity, illuminance, pm2.5, indoor_co2
Room 3	8352	13	voc, sound_pressure_level, indoor_relative_humidity, illuminance, pm2.5, indoor_co2
Room 4	13536	13	voc, sound_pressure_level, indoor_relative_humidity, illuminance, pm2.5, indoor_co2
Room 5	13536	15 2580	voc, sound_pressure_level, indoor_relative_humidity, illuminance, pm2.5, indoor_co2 supply_air_flow and damper_position

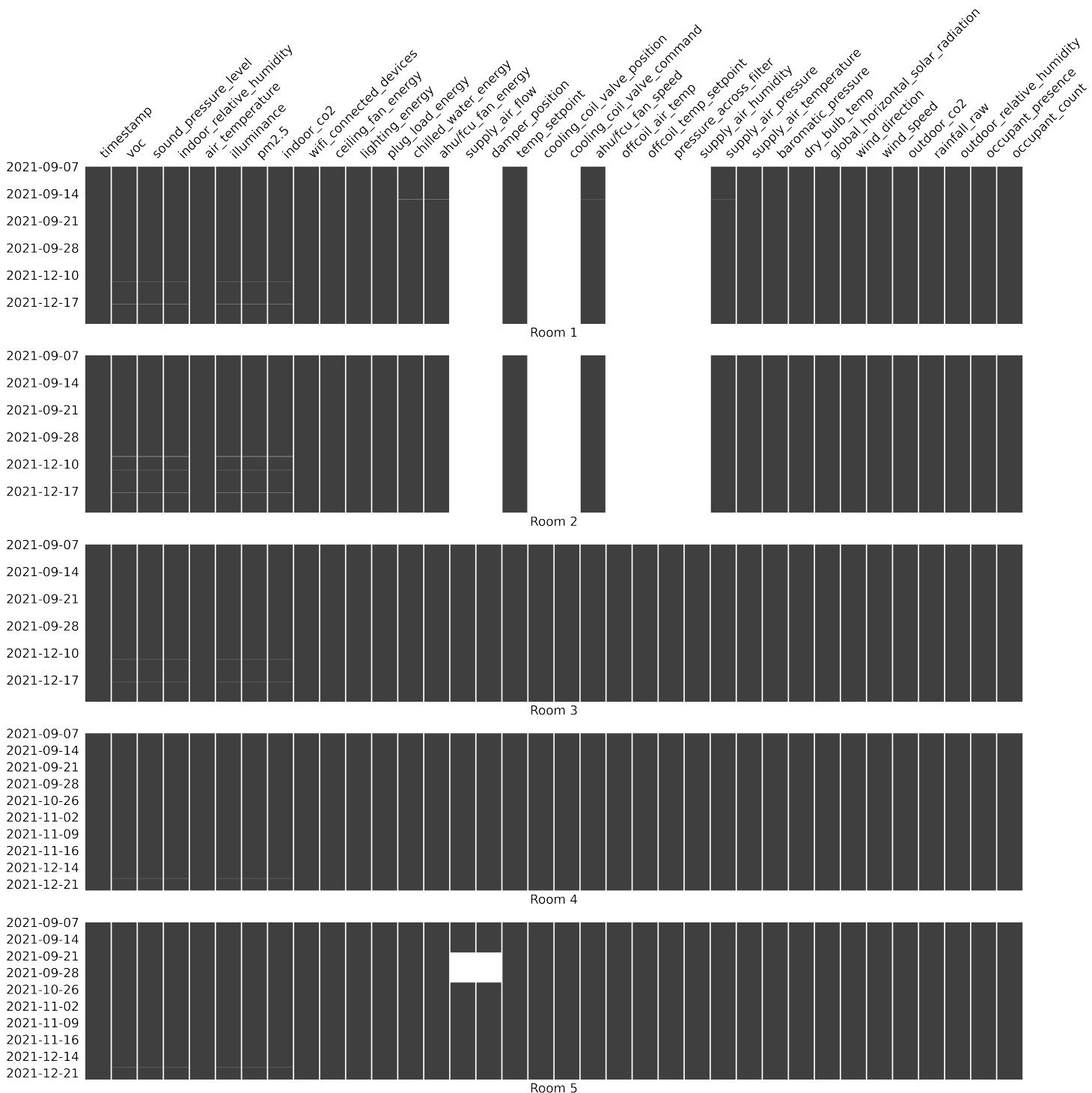


Figure 2. The amount of missing data in each column and their temporal relationship for each room.

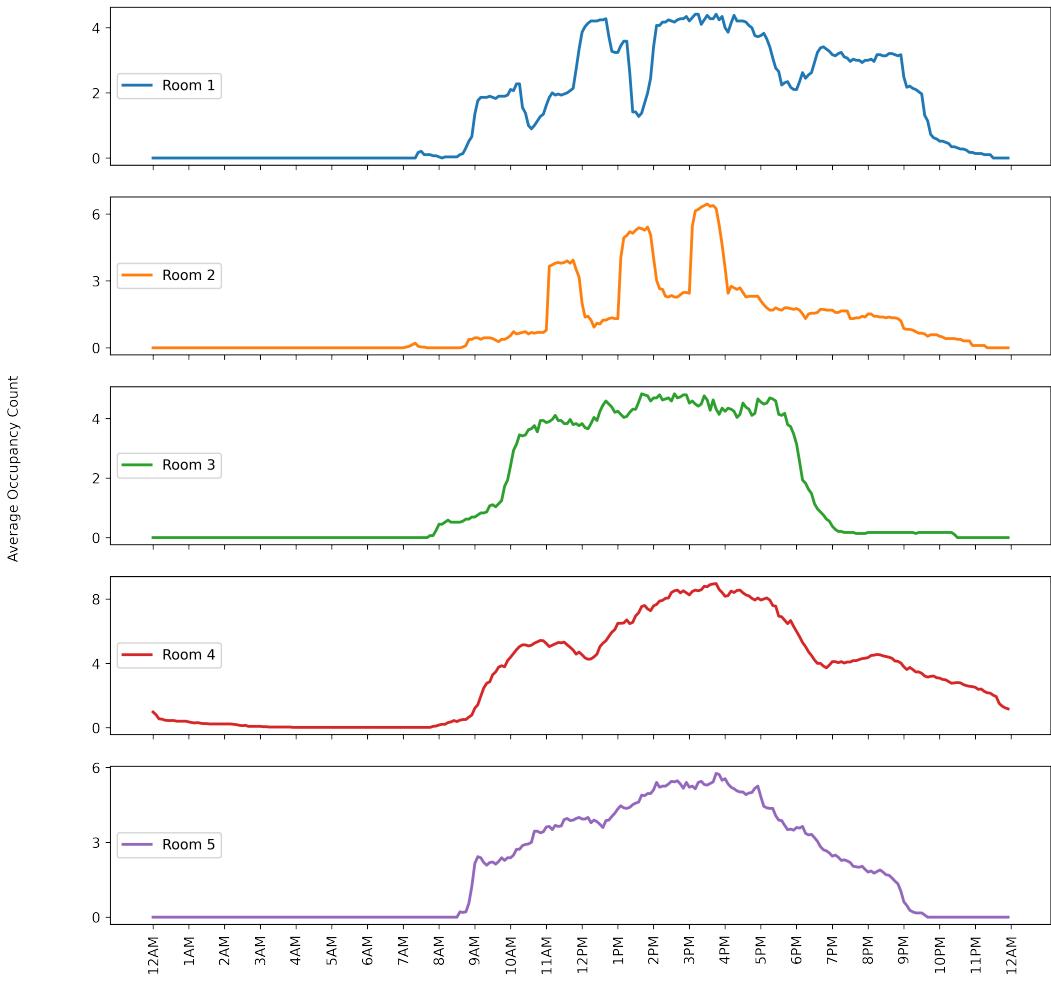


Figure 3. Average occupant count for each room on an average day.

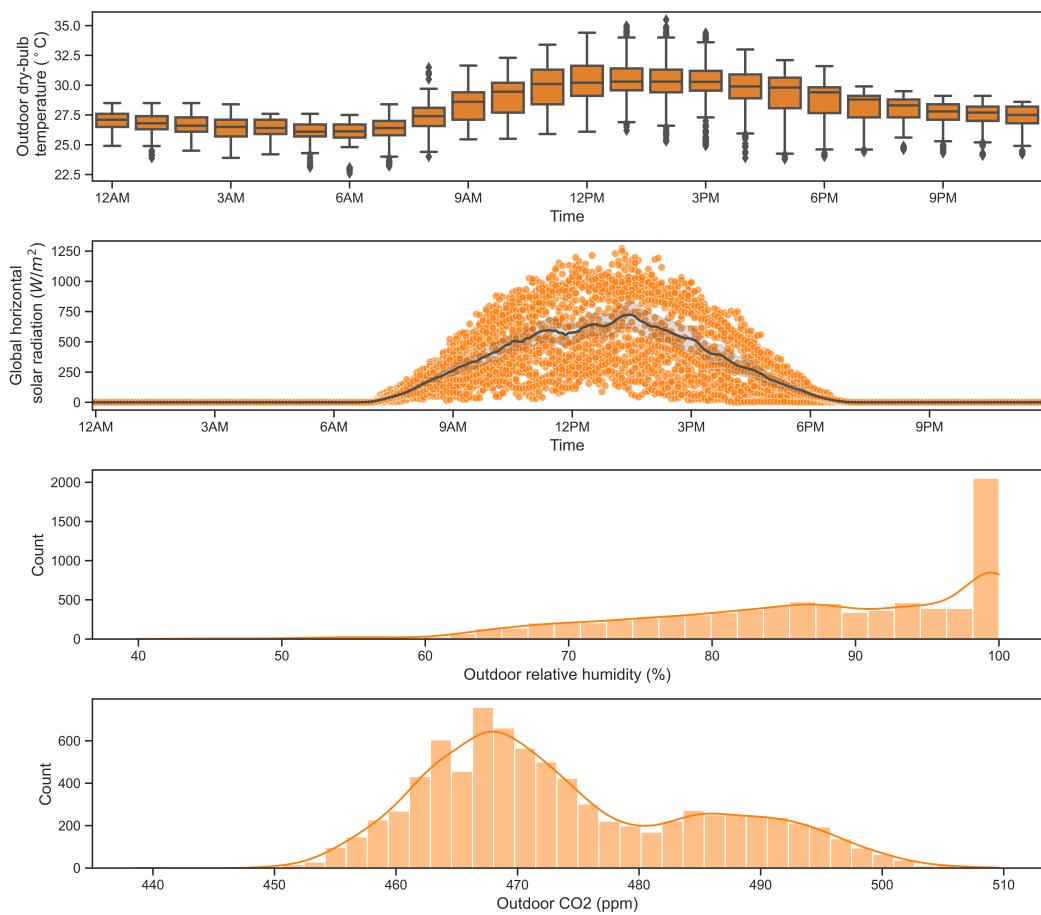


Figure 4. Data visualisations for outdoor dry-bulb temperature, relative humidity and CO₂ levels.

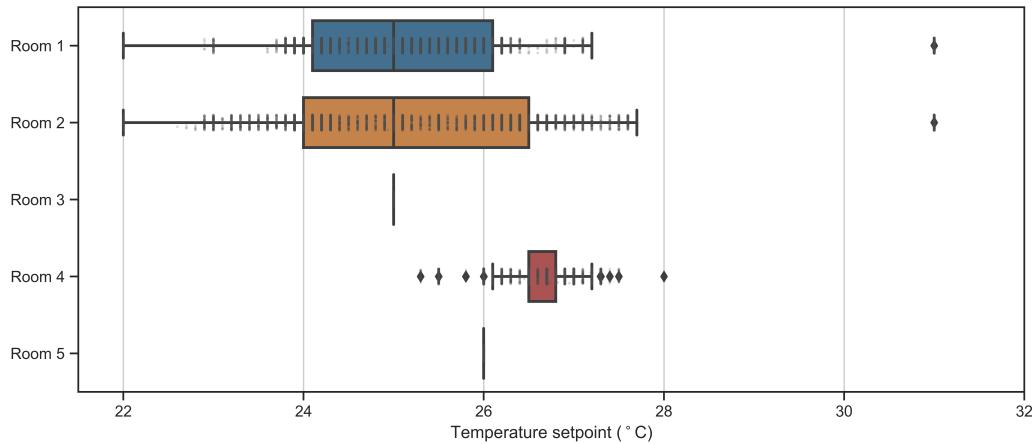


Figure 5. Distributions of room temperature setpoints for each room.

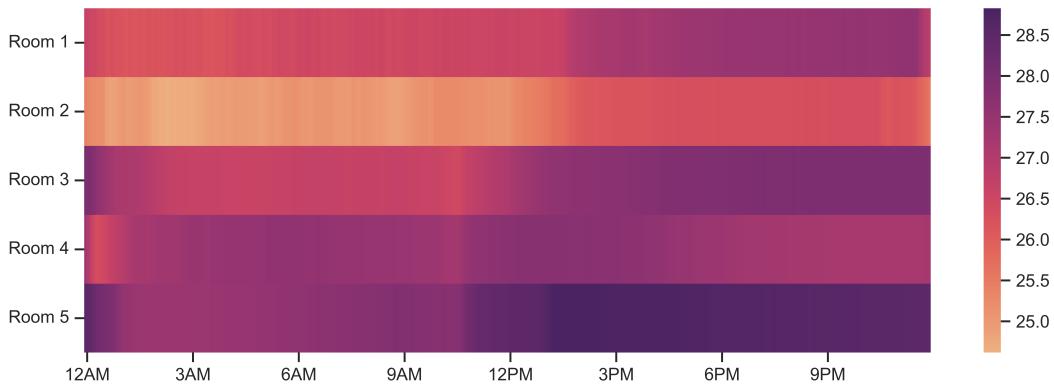


Figure 6. Average room air temperature ($^{\circ}\text{C}$) for each room.

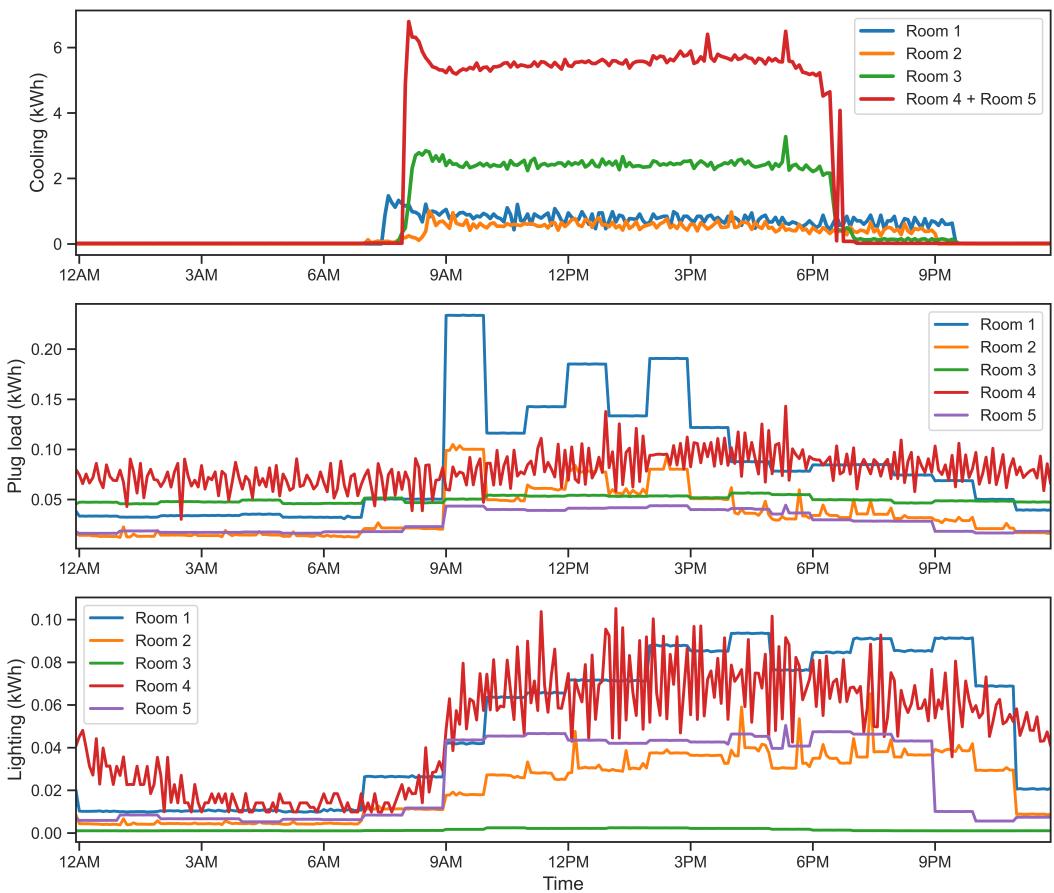


Figure 7. Average daily energy consumption of space cooling, plug load, and lighting for each room. For cooling, it should be noted that Room 3, Room 4 and Room 5 are conditioned by AHUs, which are connected to multiple rooms. In this case, Room 4 and Room 5 share the same AHU and therefore have the same cooling consumption as reflected in the figure.