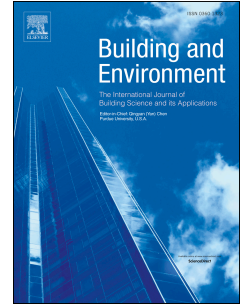# Journal Pre-proof

Sampling method for long-term monitoring of indoor environmental quality in residential buildings

Huimin Yao, Xiaojie Cheng, Shen Wei, Yuling Lv, Ang Li, Xiong Shen

Please cite this article as: Yao H, Cheng X, Wei S, Lv Y, Li A, Shen X, Sampling method for long-term monitoring of indoor environmental quality in residential buildings, *Building and Environment* (2022), doi: https://doi.org/10.1016/j.buildenv.2022.108965.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

# Sampling method for long-term monitoring of indoor environmental quality in residential buildings

Huimin Yao[1], Xiaojie Cheng[1], Shen Wei[2], Yuling Lv[1], Ang Li[1], Xiong Shen*[1]

[1] Tianjin Key Lab of Indoor Air Environmental Quality Control, School of

Environmental Science and Engineering, Tianjin University, Tianjin 300072, China

[2] The Bartlett School of Sustainable Construction, University College London (UCL),

1-19 Torrington Place, London WC1E 7HB, United Kingdom

E-mail for the corresponding author: shenxiong@tju.edu.cn

**Abstract:**

The data collected during long-term monitoring (LTM) of indoor environmental quality (IEQ) can reflect occupants' exposure to contaminants and can be used to improve thermal comfort. As there are large differences among existing guidelines for IEQ monitoring of dwellings, it is important to identify a sampling method that balances data accuracy, sample size and cost. This paper reports the major findings that developed a systematic approach to determining the sample method for IEQ monitoring. In the study, LTM was carried out in 13 naturally ventilated urban residences in Kunming, China. We proposed the continuous sampling strategy (CSS) and discrete sampling strategy (DSS). Descriptive statistics was used to evaluate the performances of both strategies, and it was found that DSS could obtain more accurate data than CSS. Next, an algorithm was developed for calculating the optimal sampling frequencies for different parameters based on the Pearson correlation coefficient. We evaluated the required number of dwellings(RND) for various parameters that satisfied the statistical

24  confidence in Kunming and other four cities of China. We found that with the increase

25  in the household number in one city, the RND will reach to a critical threshold and no

26  longer increase anymore. Using this threshold and the simple random sampling

27  principle, we also provide guidance for determining the RND for IEQ monitoring in

28  residence. It is expected that the results of this study will facilitate the selection of

29  sampling method for similar studies in the future, with reduced manpower and

30  consumption but a representative sample.

31  **Keywords**: Indoor environmental quality, Sampling strategy, Long-term monitoring,

32  Pearson correlation coefficient, Sampling frequency

33

34  **Nomenclature**

35  TVOC   Total volatile organic compound

36  $CO_2$   Carbon dioxide

37  LTM   Long-term Monitoring

38  IEQ   Indoor Environmental Quality

39  CSS   Continuous Sampling Strategy

40  DSS   Discrete Sampling Strategy

41  PCC   Pearson Correlation Coefficient

42  RND   Required number of dwellings.

43

44  **1. Introduction**

45      According to the National Human Activity Pattern Survey in the United States [1],

46    people spend about 87% of their time indoors, and this percentage is increasing. Indoor

47    air quality (IAQ) and thermal comfort are key factors in the health of building occupants

48    [2], because poor indoor environmental quality (IEQ) may lead to respiratory diseases

49    or sick building syndrome [3]. In addition, ensuring clean and comfortable indoor air is

50    an important public health goal [4]. The evaluation and analysis of IEQ commonly rely

51    on field monitoring in actual buildings. Therefore, accurate monitoring becomes

52    necessary for solving indoor pollution problems and improving people's overall well-

53    being [5].

54        The analysis of IEQ normally involves both short-term monitoring (STM) and

55    long-term monitoring (LTM). Because IEQ parameters can be affected by outdoor

56    weather (daily and seasonally), human behavior and building attributes [6-8],

57    incomplete and varying IEQ measurement results are common [9]. Because of the

58    characteristics of IEQ parameters, it is necessary to monitor the residential environment

59    for a longer time. In comparison with STM data, the data collected by LTM have several

60    advantages: 1) LTM data enable the detection of peak concentration values [10]; 2) the

61    data reflect real-time exposure to indoor pollutants [11]; and 3) the data can be used as

62    feedback signals for real-time pollutant control [7]. Therefore, LTM techniques have

63    been widely adopted to capture critical indoor environmental parameters, primarily

64    temperature, humidity, as well as concentration of carbon dioxide($CO_2$), PM2.5,

65    formaldehyde and total volatile organic compound (TVOC) [12]. In the collection of

66    data by long-term monitoring, the quality of data will determine whether the data

67    correctly reveal the basic conditions of relevant indoor parameters and represents the

68    characteristics of human exposure to pollutants [13].

69        To reduce the effort and cost entailed by IEQ-related studies, we need to consider

70    decisions about both data sampling strategies and sampling frequency [14]. Hui et al.

71    [15] used $CO_2$ as a reference to evaluate existing and proposed sampling schemes for

72    indoor pollutant concentration in terms of the necessary sampling time and sampling

73    point density and the probable errors induced at certain confidence levels of the

74    measurement. In Hong Kong, continuous sampling [16] for a measurement period of 8

75    h has been generally adopted to determine the average pollutant concentration in a

76    workplace. Since building-related contaminants normally peak in the morning in

77    workplaces, and occupant-related contaminants in the afternoon, Mui et al. [17]

78    proposed a new sampling strategy that uses the average concentration of two random

79    measurement samples. According to the results, as compared to the typical 8-h

80    measurement period of a continuous sampling method in Hong Kong, the measurement

81    time with the new method could be reduced by up to 30%. Christopher [18] analyzed a

82    rich data set and found that indoor particle events tend to be brief, intermittent, and

83    highly variable. Hence, to characterize sources of PM in indoor environments, he used

84    both continuous and time-integrated sampling instruments to simultaneously measure

85    indoor/outdoor (I/O) particle mass concentration. Based on the different characteristics

86    of indoor air parameters, previous studies have often used different monitoring and

87    sampling methods, but a unified conclusion has not been formed. Furthermore, the

88    sampling strategy and sampling frequency usually depend only on the sampling

89    precision of the instrument and the labor cost [19]. A comprehensive IEQ data sampling

90    method along with IAQ audit methodology for buildings should be established for

91    identification of indoor air problems.
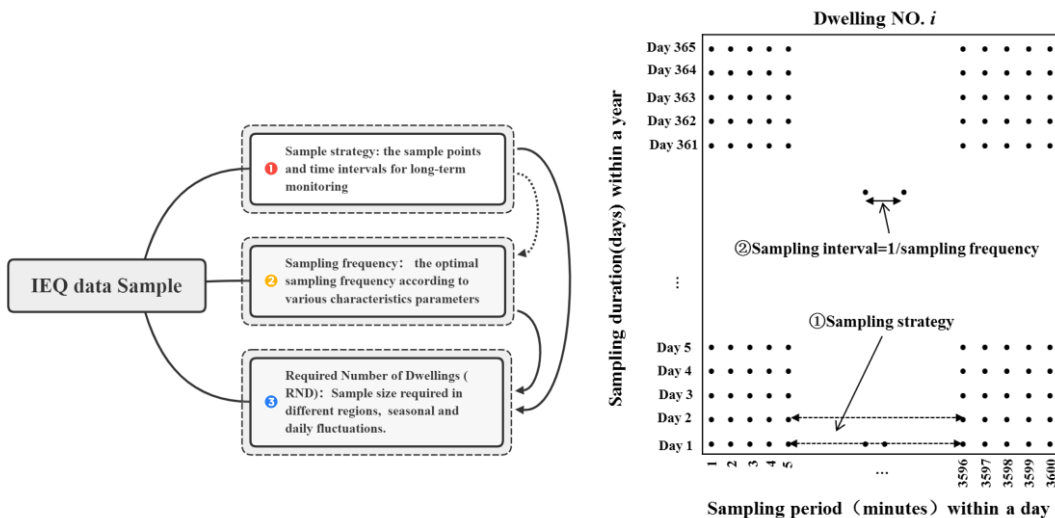
92        Both sampling frequency and sampling strategies have been addressed in existing

93    standards from different countries and regions, as summarized in Table 1. These

94    standards recommend sampling frequencies for various major IEQ parameters. For

95    example, for formaldehyde, China's GB/T18883 standard [20] recommends taking the

96    daily average. The list in Table 1 also includes different sampling periods, either real

97    time (WHO [21]) or an 8-hour average (US-EPA [22]). The sampling periods for

98    formaldehyde, PM2.5 and VOC also differ in these standards. According to the Indoor

99    Air Sampling and Evaluation Guide [22] formulated by the U.S. Environmental

100   Protection Agency, monitoring strategies for indoor air pollutants should be evaluated

101   in terms of the pollutants' exposure level, exposure time, pollutant toxicity and

102   pollutant concentration. The above standards have usually relied on previous

103   experimental results and experience in establishing the threshold value [23]. But

104   because of the different background of the experiment, different countries and regions

105   have very different requirements for the same IEQ parameter. Furthermore, according

106   to WHO, there is insufficient evidence that indoor pollutants do not cause adverse

107   effects when they fall below the thresholds in the standards [21]. Therefore, short-term

108   measurement and the threshold concentrations proposed by the standards may not be

109   enough. In order to study IEQ and ensure the health and thermal comfort of occupants,

110   we may also need to study the sampling frequency and sampling strategies of LTM.

111  Table 1 Varying regulations for threshold concentration and sampling period of IEQ parameters in
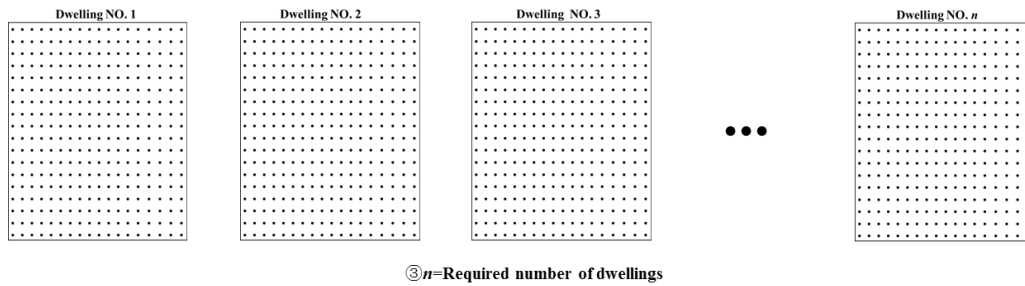
112  residential buildings

| Index | | Sampling period | Threshold concentration |
|---|---|---|---|
| **Carbon dioxide (ppm)** | | | |
| China-GB/T18883[20] | | 24-hour average | 1000 |
| China HK [24] | Good level | 8-hour average | 1000 |
| | Excellent level | 8-hour average | 800 |
| WHO-Europe[25] | | 1-hour average | 900 |
| Singapore[26] | | 8-hour average | 1000 |
| NIOSH[27] | | 8-hour average | 5000 |
| | | 15-min average | 30000 |
| Canada[28] | | real time | 3500 |
| UK[29] | | 15-min average | 15000 |
| | | 5-min average | 5000 |
| Australia[30] | | 15-min average | 30000 |
| US-EPA[22] | | real time | 800 |
| **PM2.5(µg/m³)** | | | |
| WHO[21] | | annual average | 10 |
| | | 24-hour average | 25 |
| Canada-EGR | | 1-hour average | 100 |
| | | real time | 40 |
| China- JGJ/T 309 | | 24-hour average | 75 |
| US-EPA | | 24-hour average | 65 |
| **Formaldehyde (µg/m³)** | | | |
| Canada-EGR | | real time | 120 |
| WHO | | 30-min average | 100 |
| | | real time | 200 |
| Singapore | | 8-hour average | 120 |
| NIOSH | | 15-min average | 0.1ppm |
| UK | | 15-min average | 2500 |
| Australia | | 15-min average | 2500 |
| US-EPA | | 8-hour average | 920 |
| China-GB/T18883 | | 24-hour average | 100 |
| China-GB/T50325[31] | | 24-hour average | 100 |
| China HK-IAQC | Good level | 30-min average | 100 |
| | Excellent level | 30-min average | 70 |
| **TVOC(µg/m³)** | | | |
| Singapore | | real time | 3ppm |
| China-GB/T18883 | | real time | 600 |
| China HK-IAQC | Good level | 8-hour average | 600 |
| | Excellent level | 8-hour average | 200 |
| China-GB/T50325 | | real time | 600 |
| **Temperature(℃)** | | | |
| Singapore | | | 22.5-25.5 |
| China-GB/T18883 | Summer | | 19-21 |
| | Winter | | 16-24 |
| American-ASHRAE | Summer | real time | 23-26 |
| | Winter | | 21-23 |
| Europe-AQGE | Summer | | 22-28 |
| | Winter | | 16-24 |
| China HK-IAQC | Good | 8-hour average | 25.5 |
| | Excellent | 8-hour average | 20-25.5 |
| **Relative Humidity（%）** | | | |
| Singapore | | real time | <70 |

| China-GB/T18883 | Summer | | 40-80 |
|---|---|---|---|
| | Winter | | 30-60 |
| Europe-AQGE | Summer | | 40-80 |
| | Winter | | 30-60 |
| Canada-EGR | Summer | | 30-80 |
| | Winter | | 30-55 |
| American-ASHRAE | Summer | | 50-60 |
| | Winter | | 20-30 |
| China HK-IAQC | Good | 8-hour average | <70 |

113    Moreover, to obtain the average levels of residential indoor pollutants in a specific

114    region, an appropriate number of dwellings are generally required in LTM. To study

115    the risk of sick building syndrome in air-conditioned spaces, Cheung et al. [32]

116    investigated the IEQ in 8 dwellings with different building areas, orientations and

117    numbers of residents. Mentese et al. [33] carried out an LTM study in 121 dwellings

118    located in three cities/towns in Turkey, to explore the relationship between respiratory

119    diseases and indoor pollutants. Lim et al. [34] selected 25 apartments for a study of the

120    relationship between IEQ and occupant health in energy-efficient dwellings. The

121    purposes of these studies were different, but a large sample size would increase the

122    financial cost of monitoring systems [35]. Therefore, it is meaningful to discuss the

123    effect of the required number of dwellings (RND) on the results in a given region.

③*n*=Required number of dwellings

124    Figure 1 Sampling process of long-term monitoring data on residential indoor air quality

125    Hence, this study suggested that the IEQ data sample for LTM is determined by

126    the sampling strategy, sampling frequency and RND, as shown in Fig. 1. The purpose

127    of this work was to identify a suitable data sampling method for evaluation of IEQ, by

128    means of the following steps: 1) study the periodic fluctuation characteristics of IEQ

129    parameters in residential buildings, so as to propose sampling suggestions; 2) propose

130    and compare continuous and discontinuous sampling strategies for different IEQ

131    parameters; 3) design an algorithm based on the Pearson correlation coefficient to

132    optimize the sampling frequency for various parameters, combining the experimental

133    data for these parameters; and 4) calculate the RND for studying indoor air parameters

134    according to seasonal and daily fluctuations, and for five cities in different building

135    thermal zones in China. The method developed here will help researchers in the future

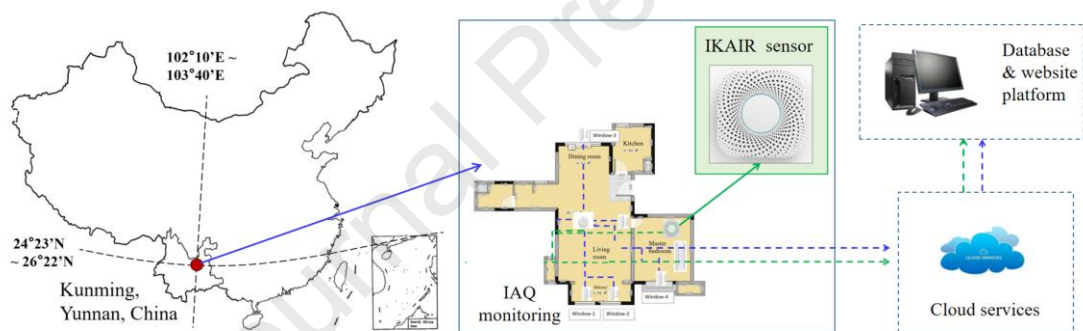136    to balance the required sample size and financial restrictions.

137    **2. Methods**

138    **2.1 Data sources and long-term monitoring (LTM) method**

139    To compare different LTM sampling strategies, this study employed IEQ-related

140    field data collected from 13 residential buildings in Kunming, China, as shown in Fig.

141    2. The study focused on dwellings in the moderate climate zone of China. The buildings

142    generally did not use heating and air conditioning systems, but only natural ventilation

143    to regulate the indoor environment[36]. In order to analyze the impact of outdoor

144    parameters on IEQ, we obtained the data of outdoor air parameters. The outdoor $CO_2$

145    concentration was measured by an IKAIR sensor that was located at the balcony of

146    dwelling. The balcony was directly exposed to the outdoor environment and was not

147    affected by indoor air conditions[8]. Outdoor PM2.5 concentrations were provided by

148    the China National Environmental Monitoring Centre, which were collected from the

149    nearest local meteorological and air quality station close to these dwellings[37]. The

150    average outdoor PM2.5 concentrations of Kunming is 39 ug/m$^3$ during the monitoring

151    period, which is one of the lowest among Chinese major cities[38]. By monitoring the

152    outdoor and indoor $CO_2$ concentration, we noticed that the concentration difference is

153    small[8]. Hence, the local outdoor pollution had less impact on IEQ. The main source

154    of indoor pollutants was from indoors [39]. The case-study building was a high-rise

155    structure, which is a common type of residential building in China [40]. The building

156    dedication year has a greater impact on the intensity and decay rate of indoor pollution

157    sources such as formaldehyde [41].

158    The long-term monitoring system for IEQ employed in this study was an integrated

159    system with various gaseous sensors that monitored $CO_2$, PM$_{2.5}$, and formaldehyde

160    concentrations, air temperature and relative humidity with a sampling period of one

161    minute. In each dwelling, we placed one IKAIR sensor in the bedroom and one in the

162    living room. The sensor is installed in the center of the monitored room, with a height

163 of 1.5 m. According to previous published research [8], we can see that there a little

164 difference in the monitoring results of various parameters between the bedrooms and

165 living rooms of residences. Thus, we did not consider the effect of sampling locations

166 in this study. The influence of sampling location on parameters is complicated, which

167 is beyond the objective of this study. Therefore, this paper will not carry out in-depth

168 discussion in this perspective. Table 2 lists the primary specifications of the sensors in

169 this system. The data collected by the sensors are dynamically sent to a data center

170 through the dwellings' Wi-Fi networks. Detailed information about the LTM system

171 can be found in Liu et al. [37].

172



173 Figure 2 Long-term monitoring in residential buildings in Kunming

174 Before and after the monitoring, sensor calibration was carried out in a 1 m³

175 integrated chamber, with a Dusttrak 8530, PB-RAE, temperature and humidity

176 measuring instrument for the detection of PM2.5, formaldehyde, temperature and

177 relative humidity, respectively. The Laskin atomization method [42] was employed to

178 generate the particles and formaldehyde using standard sources, DEHS and trimetric

179 formaldehyde. A constant temperature and humidity environmental cabin were used to

180 measure the calibrated temperature of the sensor. The temperature measurement range

181  of the environmental chamber was -70–150 ℃ with an accuracy of 0.1 ℃, and the

182  relative humidity measurement range was 20–98% with an accuracy of 0.1%. The

183  calibration process exhibited a high regression coefficient of 0.99, indicating the high

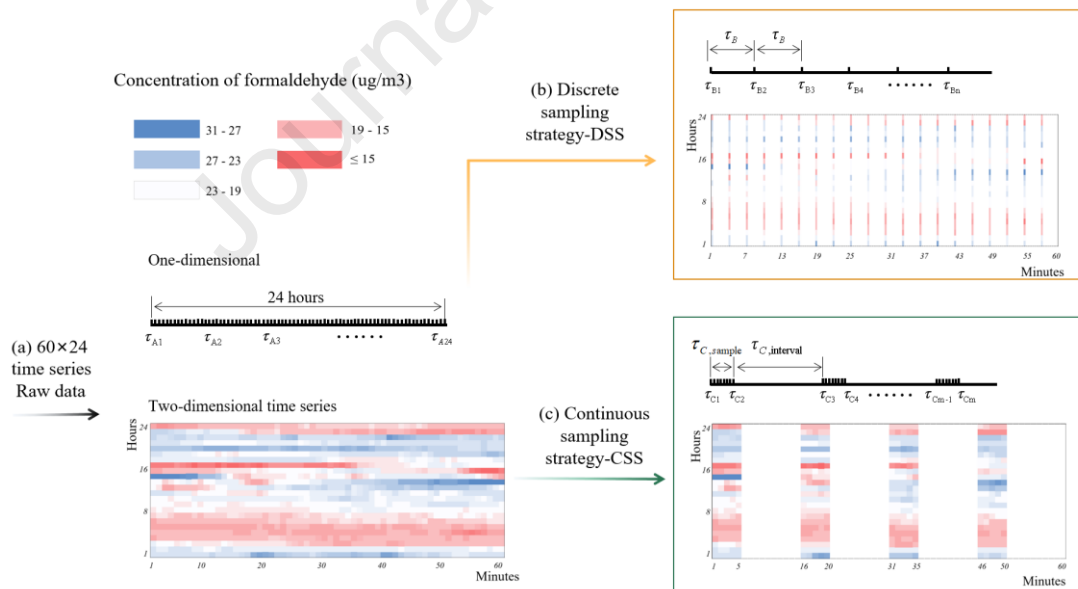184  reliability of the measurement devices used in this study.

185  Table 2 Information about sensors used for long-term monitoring of IEQ

|  | Measurement principle | Measurement range | Accuracy |
|---|---|---|---|
| PM2.5 (ug/m³) | Laser scattering | 1-1000 | ± 1 |
| Formaldehyde (ug/m³) | Electrochemistry | 0-5000 | ± 20 |
| Temperature (ºC) | Thermal resistor | -40-125 | ± 0.35 |
| Relative humidity (%) | Humidity sensitive resistor | 30-100 | ± 3 |
| Carbon dioxide (ppm) | Carbon dioxide resistor | 400-10000 | ± 30 |

186  **2.2 Sampling strategies**

187  This study proposed a continuous sampling strategy (CSS) and a discrete sampling

188  strategy (DSS). Fig. 3 explains the difference between the two strategies by using a

189  single day of formaldehyde measurements in a house as an example. The data are

190  expressed in one-dimensional time-series graph and two-dimensional stacked graph.

191  Fig. 3(a) shows the raw data obtained by measurement at a frequency of once per minute.

192  In Fig. 3, the ordinate is the 24 hours within a day and the abscissa applies the time

193  series of 60 minutes within each hour. Therefore, there are a total of $60 \times 24 = 1440$ data

194  points as presented in Fig. 3(a). Varied data points can be extracted from the raw data

195  by different sampling strategies, as shown in Fig. 3(b) and (c). For the DSS, one data

196  point of the indoor parameters is acquired at each sampling time; thus, the data points

197  collected by this method have no continuity in time. Meanwhile, as displayed in Fig.

198 3(b), formaldehyde concentration data points were sampled every three minutes. As

199 shown in Table 1, the DSS exhibits significant differences across existing IEQ

200 standards, and currently there is no conclusive reference for the appropriate DSS. For

201 the CSS, data were sampled continuously. This continuity refers to the continuity of

202 sampling of the raw data, so data are sampled at a certain sampling frequency to obtain

203 multiple data points, as shown in Fig. 3(c). Therefore, the data represent IEQ conditions

204 within a short period. As shown in Table 1, the CSS has commonly been employed in

205 existing IEQ standards to determine the mean concentrations of indoor air pollutants,

206 with duration of 8 hours, 24 hours or 1 year. Within these sampling durations, however,

207 there are different frequencies for the data sampling. Nevertheless, the evaluation of

208 IEQ was based on the average of the sampled data during the monitoring period.

209

Figure 3 Principles of the DSS and CSS for the sampling of raw data

211 **2.3 Data preprocessing method**

212 Different indoor air parameters have different dimensions and orders of magnitude.

213  To evaluate and compare the sampling strategies for different IEQ parameters in

214  parallel and to improve the credibility of the research results, it is necessary to

215  preprocess the raw data obtained from LTM [43]. Commonly used data preprocessing

216  procedures include missing value filling and data standardization processing [44]. Data

217  standardization refers to scaling of the data to a small specific interval. After

218  standardization, the data is transformed into a dimensionless pure value, so that

219  comprehensive evaluation and analysis can be carried out. The z-score standardization

220  method can be used. The processed data conform to the standard normal distribution;

221  that is, the mean value is 0 and the standard deviation is 1, with the standardized form

222  of data calculated by Eq. 1,

$$A_i(i=1,2,\ldots\ldots n) = \frac{A_i - \mu_A}{\sigma_A} \tag{1}$$

$$\mu_A = \frac{1}{n}\sum\nolimits_{i=1}^{n} A_i$$

$$\sigma_A = \sqrt{\sum\nolimits_{i=1}^{n}(A_i - \mu_A)^2}$$

223  where matrix $A_i(i=1,2,\cdots\cdots n)$ represents the raw data time series for LTM of IEQ

224  parameters; $i=1,2,\cdots\cdots n$ represent the time scale dynamically measured IEQ

225  parameters in the time series; $\mu_A$ and $\delta_A$ are the mean value and the standard deviation,

226  respectively, of the raw data.

227  **2.4 Evaluation method for various sampling strategies**

228    The Pearson correlation coefficient (PCC) is one of the most widely used

229  relationship measures [45] and is a statistical metric for the strength and direction of a

230    linear relationship between two random variables[46]. Based on this attribute of PCC,

231    we can use it to measure the correlation between the samples and the raw data to

232    evaluate the performance of the two sampling strategies. In this study, there are two

233    random variable matrices: raw data $A_i(i=1,2, \cdots n)$ and samples $B_i (i=1,2, \cdots n)$. The

234    definition matrix $B_i(i=1,2, \cdots n)$ represents the sampled data obtained by sampling the

235    raw data using the DSS or CSS. The PCC of matrices $A_i$ and $B_i$ is formally defined as

236    the product of the covariance of two random variables divided by their standard

237    deviations (which acts as a normalization factor). If each variable has n scalar

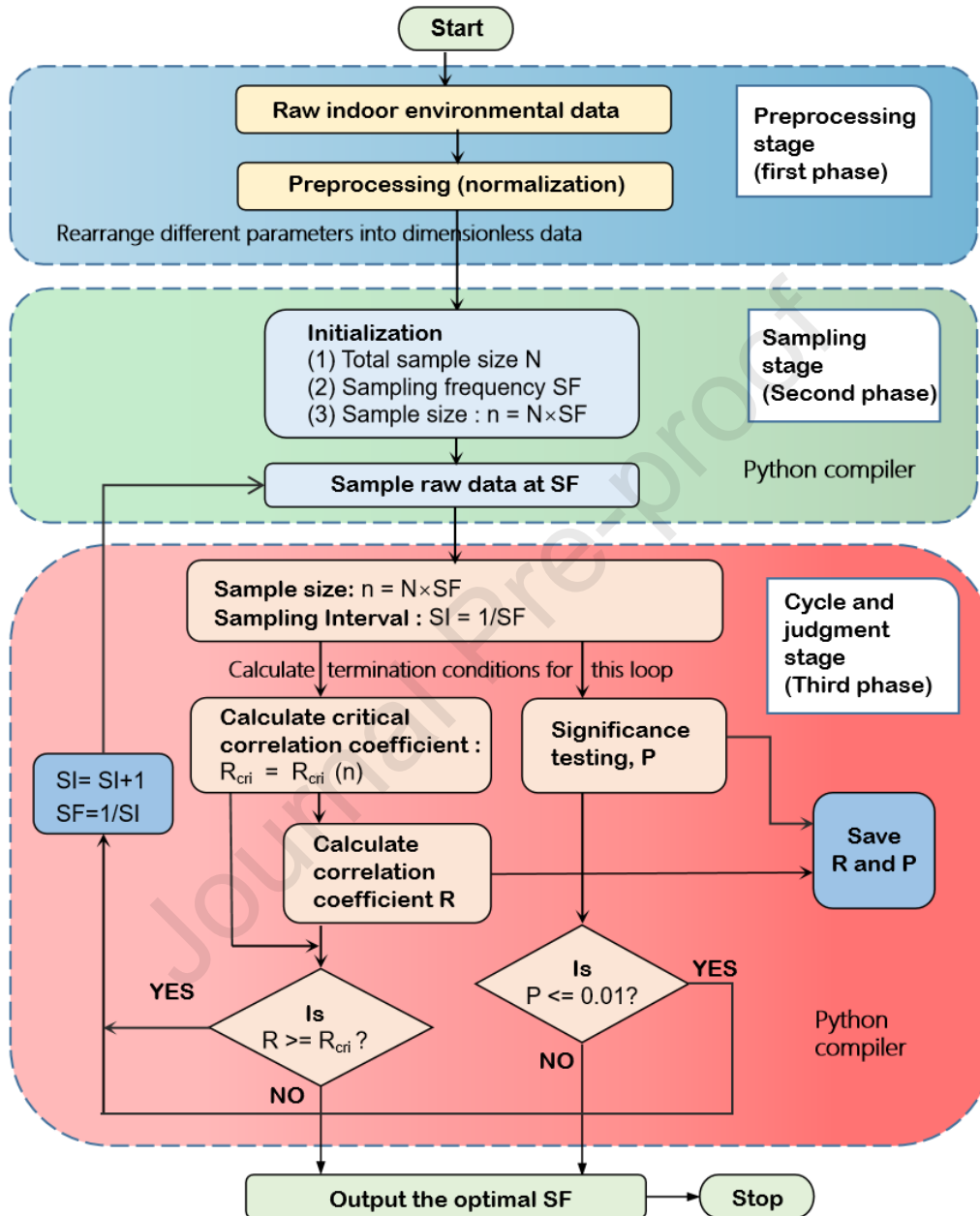238    observations, then the PCC of a certain IEQ parameter can be expressed by Eq. 2 [46],

$$R(A,B) = \frac{\text{cov}(A,B)}{\delta_A \delta_B} = \frac{1}{n-1} \sum_{i=1}^{n} \left( \frac{A - \mu_A}{\delta_A} \right) \left( \frac{B - \mu_B}{\delta_B} \right) \tag{2}$$

239    Here, $\mu_A$ and $\delta_A$ are the mean and standard deviation of $A_i(i=1,2, \ldots n)$,

240    respectively, and $\mu_B$ and $\delta_B$ are the mean and the standard deviation of $B_i (i=1,2, \ldots n)$,

241    respectively. The correlation coefficient $R(A, B)$ ranges from -1 to 1. The closer the

242    absolute value of $R(A, B)$ to 1, the stronger the correlation between the two random

243    variables.

244    **2.5 Algorithm to determine sampling frequency**

245    The quality of various sampling strategies can be assessed on the basis of the source

246    data by means of correlation analysis. Since correlation analysis determines the optimal

247    sampling frequency from the marginal value between correlation and frequency

248    increase [47], it has been widely used to compare different data sampling strategies [45].

249    In this study, the correlation analysis was implemented in Python with the algorithm

250    shown in Fig. 4.



251
252              Figure 4 Algorithm for obtaining the optimal sampling frequency by correlation analysis

253    The algorithm reaches the optimal sampling frequency *SF* by continuously

254    increasing the sampling interval *SI*. Each cycle uses the new *SI* as an intermediate

255    variable to iterate the value of SF. The iteration uses PCC parameter *R* between the

256 sampled data and the source data as the indicator variable. It updates the SI values until

257 iteration stops when R is less than the critical Pearson correlation coefficient (CPCC)

258 or passes the significance test P = 0.01. We looked up the CPCC value when $P = 0.01$,

259 and the polynomial regression curve of the CPCC data was fitted in this study as shown

260 in Fig. 4. The fitted regression curve accurately represents the condition of the loop

261 termination in the algorithm.

$$R_{cri}=7.98\times10^{-19}\times n^6 - 2.8\times10^{-15}\times n^5 + 4.1\times10^{-12}\times n^4 \\ -3.06\times10^{-9}\times n^3 + 1.26\times10^{-5}\times n^2 - 0.003n + 0.45 \tag{3}$$

262 In Eq. 3, $R_{cri}$ and $n$ represent the CPCC and the final sampling size, respectively.

263 After testing, the polynomial regression curve represented the critical correlation

264 coefficient with a degree of fit equal to 0.936. To avoid accidental error, at least 10

265 iterations were performed unless the $R$ values converged. The effect of data

266 dimensionality reduction was also analyzed after the optimum had been reached.

267 **2.6 Method for determining RND**

268 The sampled RND is an important factor in the time, cost and effort required for

269 data collection. In the determination of the appropriate sampled RND $n$ for a specific

270 region, the choice of $n$ from the total RND $n = 13$ is crucial in longitudinal studies to

271 ensure that the results are representative for that region. Based on simple random

272 sampling, $n$ depends on the degree of overall difference, the allowable error, the

273 confidence level and the adopted sampling method [48]. After these four items have

274 been confirmed, $n$ can be obtained by referring to the value that produces the minimum

275 sampling error. In the simple random sampling method, $n$ can be determined from Eq.

276     4,

$$n = \frac{z^2 \sigma^2 N}{\Delta_{\bar{x}}^2 N + z^2 \sigma^2} \tag{4}$$

277     where $N$ is the total RND in the region; $\bar{x}$ is the overall average of the results for $n$

278     dwellings; $\Delta_{\bar{x}}$ is the limited error of the $\bar{x}$; $\sigma$ is the overall standard deviation; and $z$ is

279     the degree of probability, which is directly connected to the degree of confidence with

280     a probability function of $F(z)$. In simple random sampling, $F(z)$ refers to the normal

281     distribution function. When $N$ is sufficiently large ($n/N \leq 5\%$), the RND = $n$ can be
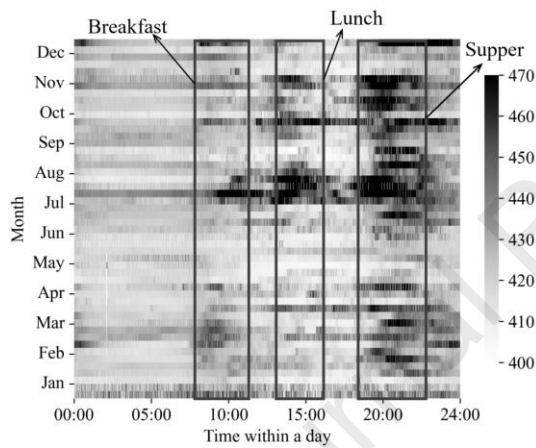
282     determined by simplifying Eq. 4 into Eq. 5,

$$n = \frac{z^2 s^2}{\Delta_{\bar{x}}^2}, \text{ if } n/N \leq 5\% \tag{5}$$

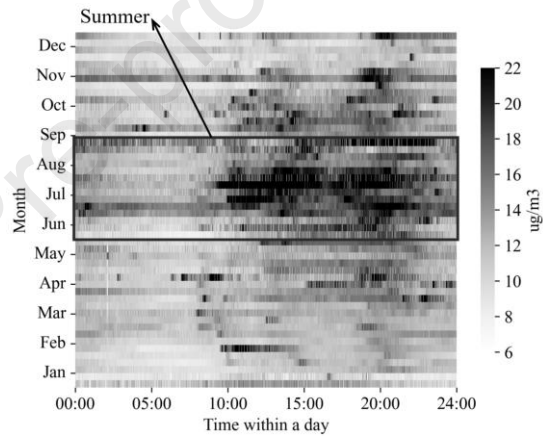283     **3. Results and discussion**

284     **3.1 Heat map analysis of IEQ characteristics**

285     Fig. 5 shows the indoor environment in the living room of one dwelling in the case-

286     study building, with a monitoring period from 1st January to 12th December. In general,

287     IEQ parameters would have certain changes in characteristics with time [49]. These

288     characteristics are related to residents' behavior, environmental conditions and outdoor

289     meteorological parameters. A heat map was produced for this dwelling, including

290     annual indoor $CO_2$, formaldehyde, $PM_{2.5}$, TVOC, relative humidity, temperature. This

291     map displays the concentration fluctuation characteristics of the measured data for each

292     parameter during the overall monitoring period. The vertical axis represents the first

293     week of the study. To display the characteristics of the parameters in different quarters

294     more clearly, months were used as the vertical axis labels. In addition, to show the data

295     characteristics of each parameter at different times during one day, the data for a 24-
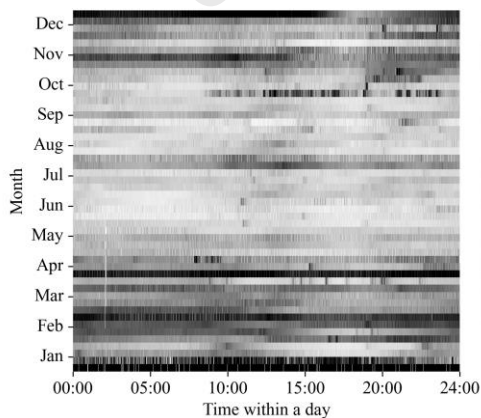
296    hour period was used to represent the indoor environmental quality for this week. The

297    24-hour data were the average of the continuous measurement data for the seven days

298    of the week; therefore, they indicate the parameter characteristics for the week. The

299    horizontal coordinate is in hours, from 0 to 24 hours. Analysis of the heat map allows

300    the effects of residents' behavior, environmental conditions and outdoor parameters on

301    the long-term measurement data to be identified. Next, the sample duration can be

302    calculated, the sampling frequency that best represents the indoor air quality can be

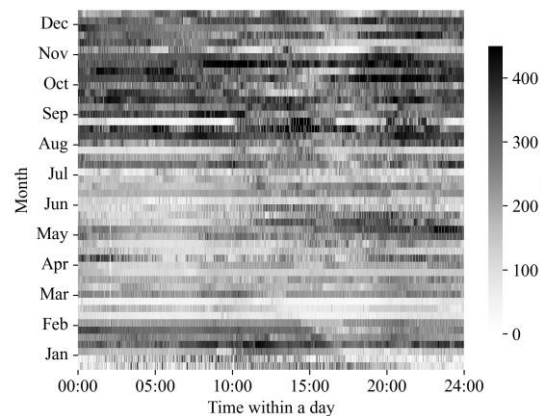303    determined, and the appropriate sampling strategy can be proposed.
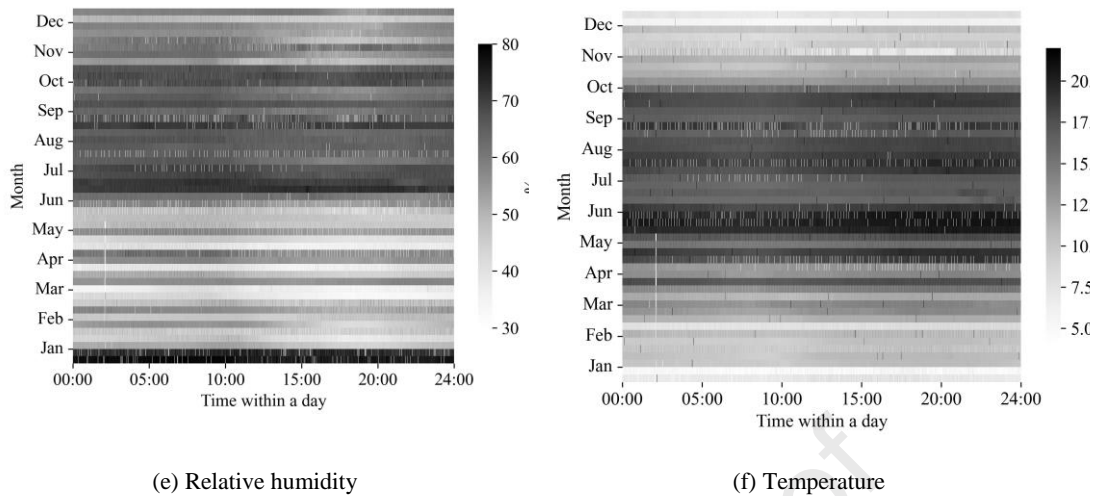


(a) Carbon dioxide                      (b) Formaldehyde

© PM2.5                                (d) TVOC

(e) Relative humidity                    (f) Temperature

Figure 5 Heat map of measured carbon dioxide, formaldehyde, PM2.5, TVOC, relative humidity and air temperature for the 2017–2018 period in a dwelling in Kunming, China

According to previous research[8, 50], indoor $CO_2$ emission sources mainly involve human breathing and fuel combustion in kitchen. Hence, it is generally believed that these two sources of $CO_2$ are present only when the residents are at home [51]. $CO_2$ is an indicator of the general level of air pollution related to the presence of humans indoors and can therefore reflect the level of human exposure [52]. It can be seen in Fig. 5(a) that the concentration of $CO_2$ changed with time. In the course of a day, both the value and the gradient of the concentration varied greatly. As shown in Fig. 5(a), the $CO_2$ concentration value was much higher during the periods from 11:00 to 15:00 and from 19:00 to 21:00 than during the rest of the day. The value between 7:00 and 10:00 was also higher than the daily average. The value was significantly lower during the period from 0:00 to 7:00. From this measurement, it can be concluded that the outdoor $CO_2$ concentration was quite stable at about 400 ppm, so the change in $CO_2$ concentration indoors was driven mainly by indoor sources. The peak $CO_2$

319  concentration was maintained for the most part in the 460–480 ppm range, which meets

320  the requirements of the Chinese IEQ standard (GB 18883). According to the peak

321  concentration of $CO_2$ on the heat map, cooking behavior increased with the indoor $CO_2$

322  concentration during meal times, especially at lunch and dinner. Therefore, attention

323  should be paid to sampling periods with high $CO_2$ concentrations, especially during

324  cooking times. Large variations in the average value and the gradient of the $CO_2$

325  concentration can reflect the occupants' behavior, e.g., the fact that the residents are at

326  home. Hence, obtaining more detailed information requires high-frequency sampling.

327  Meanwhile, changes in the indoor formaldehyde concentration throughout the year

328  are shown in Fig. 5(b). The peak concentration of formaldehyde within a day mainly

329  appeared at daytime from 10:00 to 21:00, and the concentration of formaldehyde is

330  relatively lower at night. Moreover, the gradient of formaldehyde in different seasons

331  within a year is also large. Among them, the peak value of formaldehyde mainly

332  appeared in summer (June, July, and August). As shown in Fig. 5(b) and (f), it can be

333  deduced that the change of temperature will cause the change of formaldehyde

334  concentration. The results is similar to another research, which reported the

335  formaldehyde emission from decoration and furniture materials surfaces is closely

336  correlated with air temperature [53]. In some cases, high formaldehyde concentrations

337  are difficult to predict, so the best observation times for formaldehyde concentration

338  can be determined by monitoring the periods with high air temperature.

339  As shown in Fig. 5(c) and (d), PM2.5 and TVOC concentrations did not fluctuate

340  significantly in the course of a day. However, the variations from one week to another

341     were large. In Fig. 5(c), it can be seen that the PM2.5 concentration had relatively high

342     values from November to April, when the outdoor PM2.5 concentration was high. Fig.

343     5(d) indicates that the concentration of indoor TVOCs was relatively high in January,

344     February, and from August to December, but low in other months. It is generally

345     believed that indoor TVOCs are emitted mainly by building and decoration materials

346     [54] and are independent of air temperature and humidity; therefore, the main factor in

347     the concentration is the air change rate of the dwelling. In autumn and winter, occupants

348     tend to close windows and doors so that the air exchange rate is lower than in other

349     seasons, which may have caused the increase in TVOC concentration. The indoor

350     PM2.5 concentration exhibited a similar trend to that of the TVOCs, as can be seen in

351     Fig. 5(c). The PM2.5 and TVOC concentrations did not vary greatly in the course of a

352     day, possibly because the occupants were nonsmokers or did not smoke inside the

353     dwellings [55]. Above all, when detecting PM2.5 and TVOC concentrations, it is

354     recommended that the effects of the season and the fluctuations within months and

355     seasons be taken into account.

356         It should be noted that temperature and humidity vary greatly from season to season.

357     As shown in Fig. 5(e), the relative humidity in July, August, September and October is

358     high (65–75%), and the relative humidity values from February to June and in

359     November and December are low. It can be seen in Fig. 5(f) that the distribution of the

360     temperature heat map is also related only to the month. May through October are the

361     peak months, with a temperature range of 16—22 °C, and the temperature in the other

362     months is lower, ranging from 5°C to 13 °C. The dwelling in this study was located in

the moderate climate zone of China, with small annual temperature differences and discrete dry and wet seasons. Summer and autumn in this region are rainy, and the outdoor humidity is relatively high. Therefore, on the heat map, the measured relative humidity inside residential buildings in summer and autumn is relatively high, and the temperature in summer is also relatively high. However, the temperature and humidity did not fluctuate significantly in the course of a day. Therefore, it is recommended that the seasonal characteristics be taken into account when sampling. Since a numerical gradient was observed only from week to week, the sampling of temperature and relative humidity can be performed at a weekly frequency.

Fig. 5 presents all the data collected for the dwelling throughout the year. The figure also shows detailed characteristics as follows: 1) according to periodic variation characteristics of $CO_2$, formaldehyde, PM2.5, TVOC, relative humidity and air temperature, the cyclical behavior of the occupants in residential dwellings can be preliminarily observed; and 2) there are large differences between IEQ parameters, so the sampling strategy should be formulated according to the time fluctuation characteristics of the parameters.

**3.2 Evaluations of different sampling strategies**

In the study, we evaluated two sampling strategies: the discrete sampling strategy (DSS) and the continuous sampling strategy (CSS). With the use of continuously measured data for six parameters, PM2.5, formaldehyde, TVOC, $CO_2$, temperature and relative humidity, in a residence in Kunming throughout the year, sampling with different measurement intervals was performed by means of the DSS and CSS.
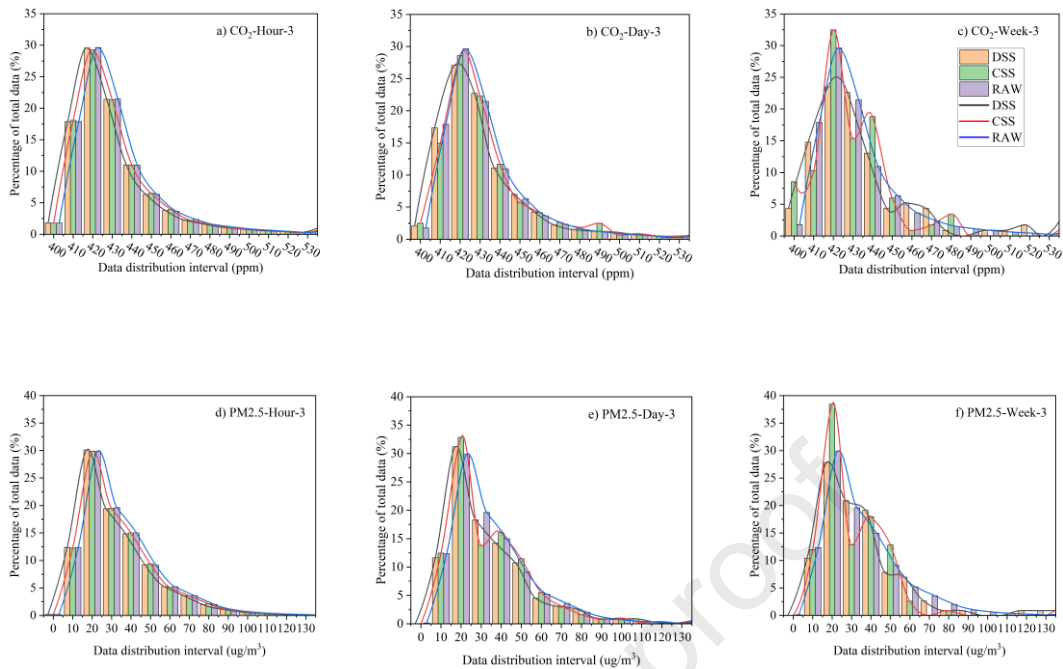
385      The sampling period was set to hours, days and weeks. And in each sampling period,

386      *i* samples were chosen from the raw data. The sampling method is abbreviated as

387      Sampling Strategies (SS) - Sampling period-*i*. Different sampling strategies will extract

388      the same amount of data during the same sampling period. Therefore, the final sample

389      sizes obtained by the two sampling strategies were the same.

390      Statistical analysis could easily and intuitively process an array of time series data.

391      it plays an important role in comparing the representativeness of data obtained under

392      different sampling strategies. In this study, frequency distribution histogram analysis

393      and descriptive statistics analysis were performed on different parameters. Considering

394      that Pearson Correlation Coefficient (PCC) can be used to calculate the correlation

395      between two arrays, by comparing the PCC calculated by raw data and the sampled

396      data, we can obtain the optimal sampling strategy.

397    **3.2.1 Effects of different sampling strategies on data distribution**

398      Using the CSS and DSS, we sampled the annual data for PM2.5 and $CO_2$ in the

399      case-study dwelling. In order to compare the distribution of the sampled data with the

400      raw data, a histogram of the frequency distribution under different sampling frequencies

401      was plotted, as shown in Fig. 6. The IEQ data were sampled at frequencies of hours,

402      days and weeks.

403    Figure 6 Density of distribution of PM2.5 and $CO_2$ concentrations at various frequencies under the

404    DSS and CSS in a single dwelling in Kunming, China; histograms (a), (b) and (c) refer to $CO_2$, and (d),

405    (e) and (f) refer to PM2.5

406    The density distribution exhibited a positive trend with the sampling frequency, as

407    shown in Fig. 6. According to Fig. 6(a) and (d), when $CO_2$ and PM2.5 were sampled at

408    a frequency of once every 3 hours, the density distribution under both the CSS and DSS

409    was close to that of the data source. When the sampling frequency was once every 3

410    hours, the density distributions under the CSS and DSS have little difference. With a

411    sampling frequency of 3 days as shown in Fig. 6(b) an©(e), the density distribution

412    under the CSS differed greatly from that under the DSS. When the sampling frequency

413    was reduced to once every three weeks as shown in Fig. 6(c) and (f), the distributions

414    of the DSS and CSS differed more strongly. Moreover, the DSS generally yielded a

415    better comparison with the raw data, especially for the peak values. The data sampled

416    under the CSS fluctuated less, with a more scattered distribution. The density

417    distribution in general prohibits a lognormal distribution. From the above, it can be

418    concluded that the sampling frequency should be on the order of hours or days, but no

419    less frequently than week for $CO_2$ and PM2.5.

420    **3.2.2 Impact of different sampling strategies on descriptive statistics**

421    The study investigated not only the average value but also the standard deviation,

422    coefficient of variation, partial peak, and kurtosis value. The average value indicates

423    the overall level of the sample, and the standard deviation and coefficient of variation

424    reflect the degree of dispersion of the data. The skewness and kurtosis value measure

425    the degree of skewness and flatness of the distribution, respectively. The maximum and

426    minimum values indicate the data range. It is recommended that a sampling strategy be

427    selected that has closer agreement with the raw data as well as fewer data points.

428    Table 3 Comparison of DSS and CSS by descriptive statistics for PM2.5 and $CO_2$

| Sampling strategy | Mean | Standard deviation | Correlation of variance | Partial peak | Kurtosis | Minimum | Maximum | Data points |
|---|---|---|---|---|---|---|---|---|
| PM2.5 | | | | | | | | |
| Raw data | 30.4 | 24.33 | 0.8 | 3.4 | 28.6 | 0 | 852 | 384694 |
| DSS-Hour-3 | 30.43 | 24.55 | 0.81 | 3.41 | 25.5 | 0 | 489 | 19235 |
| CSS-Hour-3 | 29.8 | 23.56 | 0.79 | 2.73 | 13.61 | 0 | 243 | 19236 |
| DSS-Day-3 | 30.86 | 24.39 | 0.79 | 2.73 | 13.43 | 1 | 229 | 802 |
| CSS-Day-3 | 30.29 | 23.74 | 0.78 | 2.74 | 13.58 | 1 | 213 | 802 |
| DSS-Week-3 | 32.21 | 25.58 | 0.79 | 2.92 | 13.4 | 2 | 189 | 115 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| CSS-Week-3 | 29.57 | 31.85 | 1.08 | 4.39 | 22.49 | 1 | 213 | 117 |
| Carbon dioxide | | | | | | | | |
| Raw data | 427.81 | 27.59 | 0.06 | 5.59 | 133.26 | 400 | 2087 | 384694 |
| DSS-Hour-3 | 427.9 | 29.42 | 0.07 | 11.9 | 526.91 | 400 | 2070 | 19236 |
| CSS-Hour-3 | 427.91 | 27.55 | 0.06 | 4.79 | 63.57 | 400 | 1121 | 19235 |
| DSS-Day-3 | 428.29 | 27.23 | 0.06 | 4.53 | 45.51 | 400 | 800 | 802 |
| CSS-Day-3 | 428.79 | 24.71 | 0.06 | 1.91 | 4.26 | 400 | 548 | 804 |
| DSS-Week-3 | 434.13 | 44.62 | 0.1 | 5.35 | 39.8 | 400 | 800 | 115 |
| CSS-Week-3 | 426.57 | 24.35 | 0.06 | 2.48 | 8.65 | 400 | 538 | 117 |

429      In previous studies, the maximum, minimum and average values have been used

430      as evaluation criteria. The mean values differed by less than 5.9% between the DSS and

431      CSS. The standard deviation and coefficient of variation for PM2.5 varied greatly with

432      sampling frequency under the DSS. The accuracy of the monitoring instrument of $CO_2$

433      is ±30 ppm, yet the standard deviation of the data is lower than the accuracy. Therefore,

434      we do not study the average of $CO_2$. According to a comprehensive analysis of peak

435      and kurtosis values for $CO_2$ and PM2.5, the DSS exhibited closer agreement with the

436      raw data. The maximum value also indicates a closer agreement under the DSS than the

437      CSS. The number of data points is another important factor in the choice of sampling

438      frequency.

439      We noticed that the peak concentration of $CO_2$ is an important IEQ parameter.

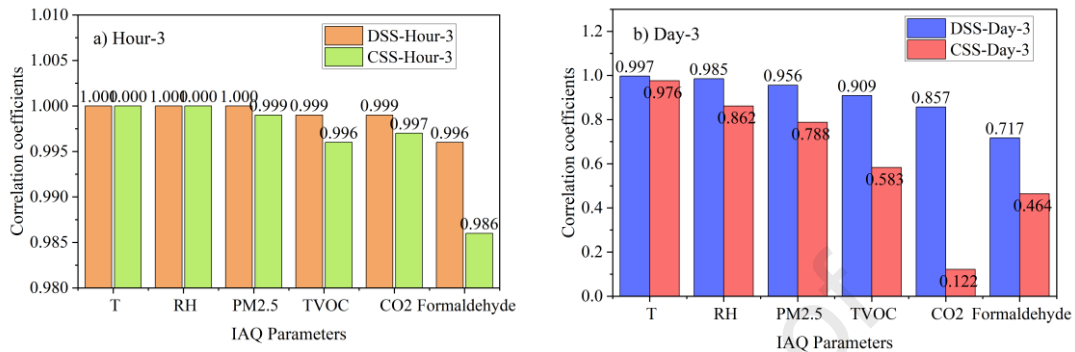440      Since if the maximum of $CO_2$ is high, more fresh air needed to be supplied to the house

441 to improve the IAQ. According to Table 1, the threshold peak concentration of CO2 in

442 China is 1000 ppm, while in WHO is 900 ppm. Table 3 shows the maximum obtained

443 by DSS is much higher than CSS and is closer to the raw data. In other words, the DSS

444 data are more likely to indicate that the house needs fresh air. Thereby, it is

445 recommended to use DSS-Hour-3 for the sampling of $CO_2$.

446 In short, although the statistical measures of mean, standard deviation, and

447 coefficient of variation are not regular, analysis of the data for other IEQ parameters is

448 still needed. However, the descriptive statistical characteristics of the DSS samples are

449 considered to be closer than CSS in terms of partial peaks, kurtosis values and sample

450 ranges.

451 **3.2.3 Correlation analysis of different sampling strategies**

452 The DSS and CSS were used to sample PM2.5, temperature(T), relative

453 humidity(RH), formaldehyde, $CO_2$ and TVOCs in a residence in Kunming throughout

454 the year, and the Pearson correlation coefficient(PCC) between the sample and the raw

455 data was compared for the two strategies. The PCC reflects the correlation between

456 samples. The calculation of the correlation coefficient requires that the variables have

457 the same sample size. Because the data distributions of DSS-Week-$i$ and CSS-Week-$i$

458 differ greatly from the raw data according to the conclusion of section 3.2.1, this

459 comparison only involves the sampling strategies of DSS(CSS)-Day-$i$ and CSS(CSS)-

460 Hour-$i$. The limits set by China-GB/T18883, China-JGJ/T309, and China-GB/T50325

461 for residential PM2.5, formaldehyde and $CO_2$ are all daily averages. Therefore, we

462 calculated the daily average for the data obtained with different sampling frequencies

463    and different sampling strategies to unify the sample size. We were then able to analyze

464    the correlation coefficient between the daily average value and the raw data.



465    Figure 7 Comparison of DSS and CSS in terms of correlation coefficients for six IEQ

466    parameters: (a) DSS-Hour-3, CSS-Hour-3; (b) DSS-Day-3, CSS-Day-3

467    Fig. 7 displays the results for the DSS and CSS with sampling frequency of once

468    every three hours or once every three days, referred to as Hour-3 and Day-3,

469    respectively. The Pearson correlation coefficient of the DSS, namely, DSS-Hour-3 and

470    DSS-Day-3, for the six IEQ parameters indicates a closer relationship with the raw data

471    than does the PCC of the CSS. The linear relationship between the sampled data under

472    the DSS and the raw data is obvious. Therefore, the samples obtained under the DSS

473    more accurately reflect the characteristics of IEQ.

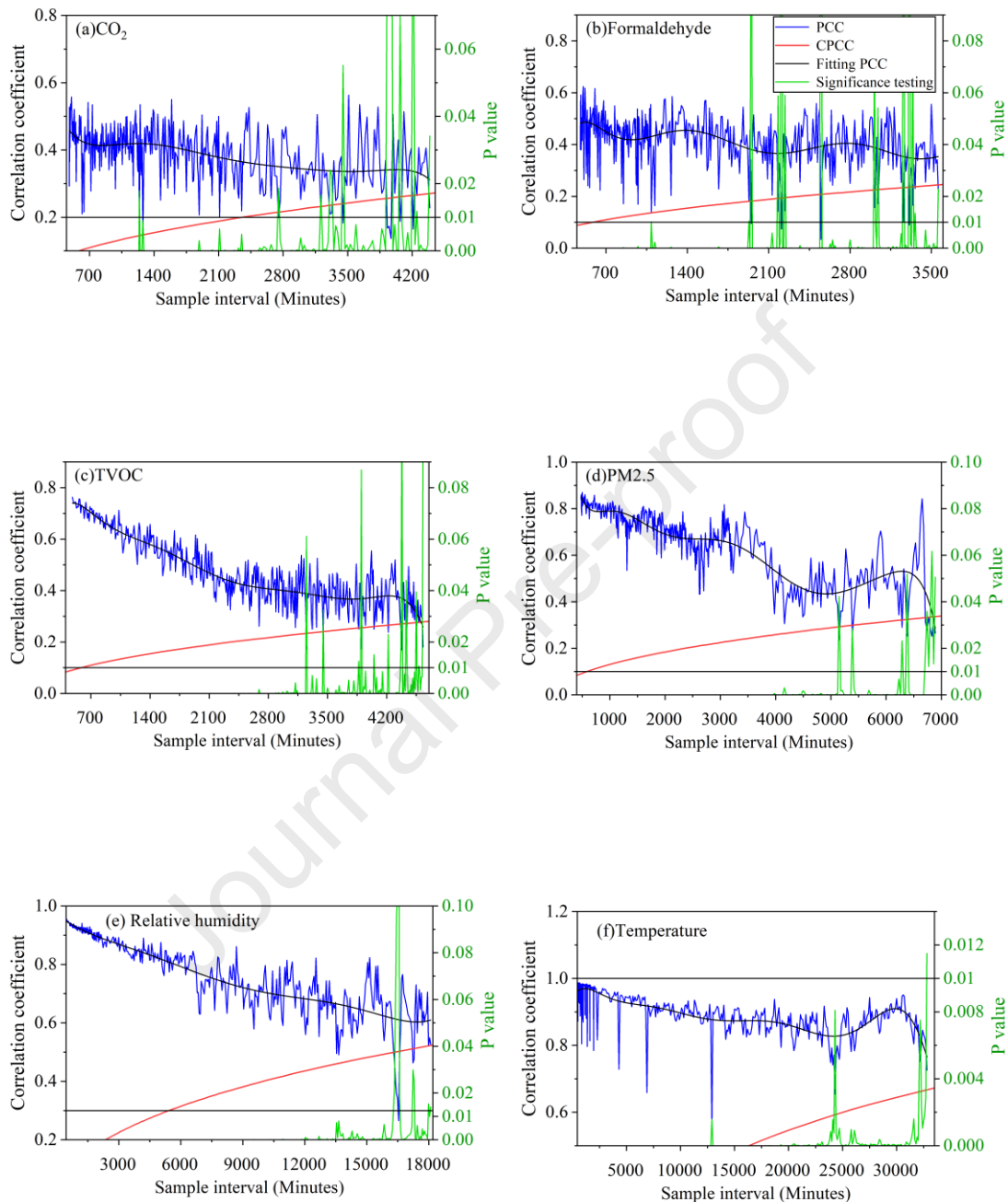474    **3.3 Optimal sampling frequency under DSS**

475    According to the descriptive statistics and correlation coefficients for the two

476    strategies, the sampled data obtained under the DSS are superior to the CSS data in

477    reflecting the characteristics of the raw data. As discussed in Section 3.1, the values of

478    the IEQ parameters fluctuated periodically with the days and weeks. In order to

479    understand the differences among various parameters, it is necessary to further study

480    the sampling frequency and determine the characteristics of the data.

481

482



483    Figure 8 PCC and significance test for sampling frequencies under DSS for dix indoor air

484    parameters:(a)Carbon dioxide, (b)Formaldehyde, (c)TVOC, (d)PM2.5, (e)Relative

485    humidity,(f)Temperature. Legends are indicated in (b):CPCC means critical Pearson correlation

486    coefficient; PCC, CPCC  refers to the coordinates on the left; Significance test results refer to the

487    coordinates on the right

488    The calculated results for PCC under the DSS are shown in Fig. 8, where the

489    sample interval was equal to the inverse of the sampling frequency. As the sampling

490    interval increased, the PCC of the sample gradually decreased, but the CPCC gradually

491    increased. The fitted curves converged at the point where the sampling interval reached

492    its optimum. At the same time, we took into account the significance of the results.

493    When the $P$ value was less than 0.01, we considered the test to be unqualified and the

494    sample data to no longer be correlated with the raw data. When the significance test

495    failed, the two curves might not intersect, and the iteration stopped.

496    Fig. 8(a) shows the PCC and CPCC data for $CO_2$ and the significance test results

497    with the change in sampling interval. The optimal sampling interval obtained by the

498    algorithm was 4390 minutes; that is, the sampling period was 3.04 days, which accounts

499    for only a single data point every three days. Regardless of accidental errors, the

500    correlation coefficients at certain frequencies also reached a critical value (intersection

501    point). After the determination of correlation coefficients, the significance test was

502    performed. According to the test result, the null hypothesis that the sampled data is not

503    correlated with the raw data should be rejected. Therefore, the optimal sampling

504    frequency of $CO_2$ is at least once every 3 days.

505    The calculated sampling frequency for formaldehyde in Fig. 8(b) is once every 2.47

506    days. Thus, the formaldehyde in this dwelling should be sampled at least once every 2

507    days. According to the calculation result for TVOCs in Fig. 8(c), the sampling interval

508    is 4630 minutes and the sampling frequency is once every 3.2 days, which means that

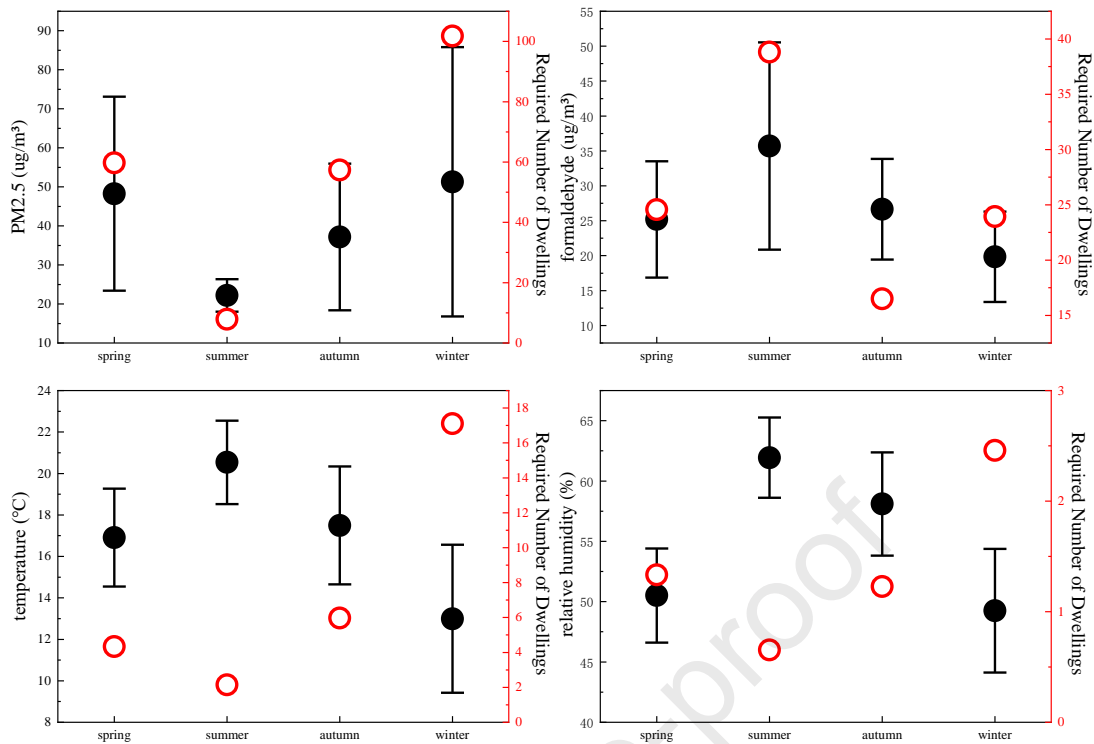509    the optimal sampling frequency for TVOCs is at least once every 3 days. Meanwhile,

510    the sampling interval for PM2.5 in Fig. 8(c) is 6890 min with an optimal sampling

511    frequency of at least once every 4 days. The sampling intervals for relative humidity

512    and temperature in Fig. 8(e) and (f) are 18090 and 32800 minutes, respectively, with

513    optimal sampling frequencies of at least once every 12 and 22 days.

514        In many studies, the IEQ parameters have all been monitored with the same

515    frequency by integrated sensors that are similar to the ones in this study. In cases in

516    which six IEQ parameters are measured simultaneously with the same sampling

517    frequency, the frequency that is selected should be the one that is the lowest among all

518    the parameters. In the dwelling used in the present study, the lowest frequency was that

519    for formaldehyde. The PCC and significance test are therefore highly recommended for

520    the determination of optimal sampling frequency, not only for individual sensors but

521    also when integrated sensors are used.

522    **3.4 RND for seasonal and daily average data**

523        For studies of different durations, the number of samples to be monitored is often

524    different. In this section, residential sample sizes required for the study of indoor air

525    parameters are proposed in accordance with seasonal and daily fluctuations of these

526    parameters in different dwellings [13, 37, 56]. In terms of seasonal averages, different

527    numbers of dwellings are required for each season. As shown in Fig 9, the seasonal

528    RND ranges from 10 to 100 for PM2.5, from 15 to 40 for formaldehyde, from 2 to 17

529    for temperature, and from 1 to 3 for relative humidity. Therefore, PM2.5 requires the

530    largest RND for long-term monitoring.

Figure 9 RND for determining seasonal averages for indoor PM2.5, formaldehyde, air temperature and

relative humidity in Kunming, where black dots and red circles represent the change of averages and

RNDs, respectively

The RND for PM2.5, temperature, and relative humidity exhibited similar trends

with seasonal changes. The seasonal RND is the largest in winter and the smallest in

summer. Meanwhile, it is equivalent in spring and autumn. The seasonal RND is

affected neither by the sample mean and the individual variance of the sample, nor by

the ratio of the sample variance to the sample mean.

Formaldehyde is distinct from PM2.5, temperature, and relative humidity. The

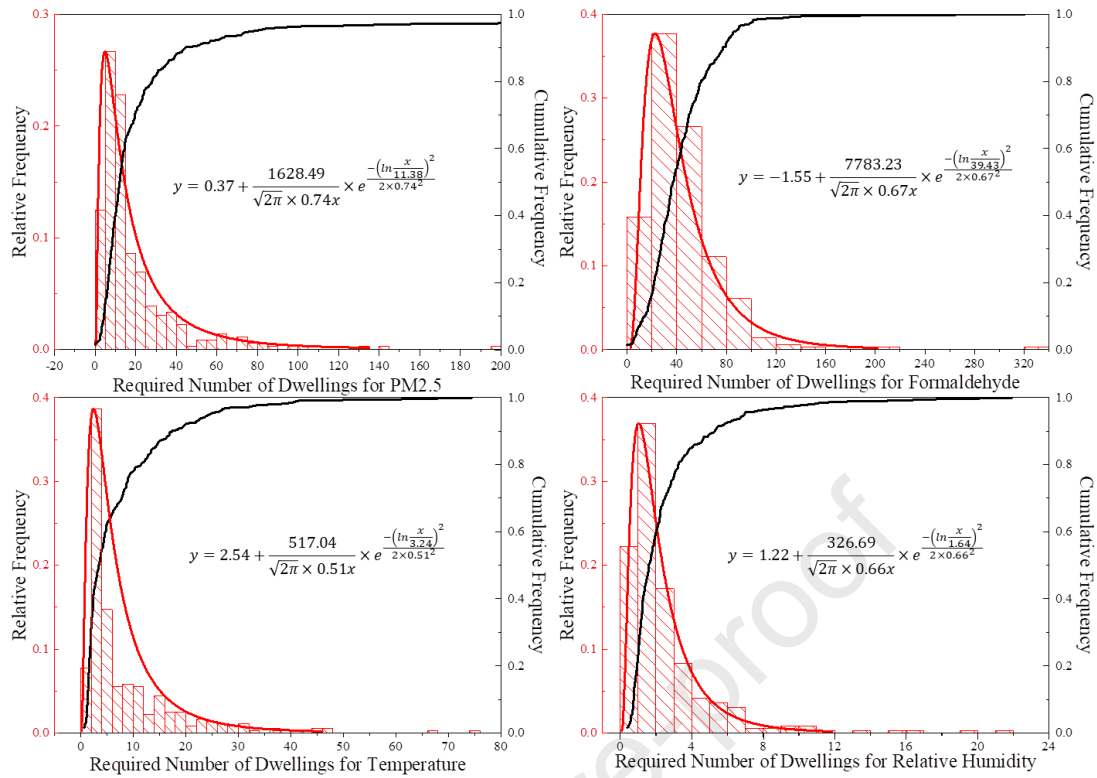change in seasonal RND with the season for PM2.5, temperature, and relative humidity

in the 13 dwellings was close to each other. Formaldehyde has the largest RND in

summer because of a large difference among dwellings in the strength of the

formaldehyde emission source. The variations in temperature and humidity in summer

545    have a significant effect on the release of formaldehyde from decorative products. This

546    finding agreed well with those of related studies [57-59], which reported formaldehyde

547    emissions at high levels when temperature and humidity were high.

548        Fig. 10 displays a distribution histogram of the daily RND for PM2.5,

549    formaldehyde, temperature and relative humidity [13, 60]. The RND satisfies the

550    lognormal distribution throughout the year. The mean daily RND for PM2.5,

551    formaldehyde, temperature and relative humidity are 11, 39, 3 and 2, respectively, and

552    the variance of the RND is 0.74, 0.67, 0.51 and 0.66. The annual, seasonal, monthly or

553    daily RND for temperature and humidity is less than the RND for PM2.5 and

554    formaldehyde. Indoor temperature and relative humidity varied only slightly among

555    dwellings because they were affected mainly by the outdoor climate conditions. In

556    contrast, indoor PM2.5 and formaldehyde varied greatly because the strength of indoor

557    pollutant sources differed from one dwelling to another.

558

559    Figure 10 RND for daily average of indoor PM2.5, formaldehyde, air temperature and relative

560        humidity. Fitting formulas are logarithmic normal distributions of relative frequency

561        Based on the cumulative frequency of the RND (see Table 4), when the confidence

562    level is 90%, the RND for PM2.5, formaldehyde, temperature and relative humidity are

563    45, 77, 20, and 5, respectively. When the confidence is 95%, the RND for PM2.5,

564    formaldehyde, temperature and relative humidity are 78, 87, 28, and 7, respectively.

565    Therefore, the greater the confidence level, the larger the RND. Hence, when

566    calculating RND, it is necessary to take into consideration either the confidence level

567    or the number of data points.

568    Table 4 Relationship between RND and degree of confidence

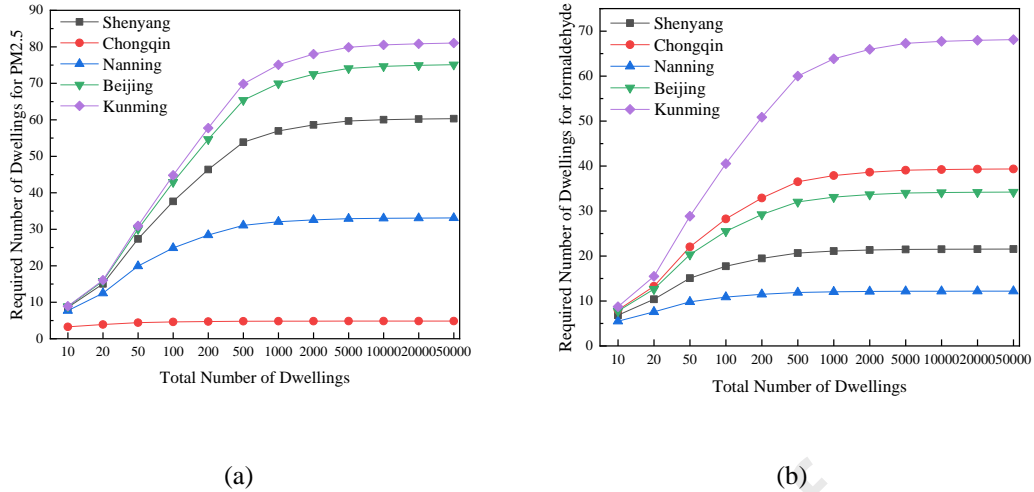| Confidence level | PM2.5 (ug/m³) | Formaldehyde （ug/m³） | Temperature (°C) | Relative humidity (%) |
|---|---|---|---|---|

| | | | | |
|---|---|---|---|---|
| 90% | 45 | 77 | 20 | 5 |
| 95% | 78 | 87 | 28 | 7 |

569

**3.5 RND in various thermal zones in China**

It can be seen in Fig. 11 that different thermal zones have different RND for studying the levels of PM2.5 and formaldehyde in a region. The RND is affected mainly by the coefficient of variation of the study parameters. The larger the coefficient of variation, the greater the difference in the concentration of pollutants between dwellings, and the larger the RND. Based on the current monitoring data, the difference in formaldehyde between Chongqing and Nanning is greater than that for PM2.5, while the difference in formaldehyde among Kunming, Beijing, and Shenyang is less than that for PM2.5. At the same time, a larger coefficient of variation leads to a larger RND. In Kunming, Beijing and Shenyang, the PM2.5 varies greatly among dwellings, while in Chongqing and Nanning, the difference in PM2.5 is small. For PM2.5, the RND is 60 in Kunming, while the RND in Beijing and Shenyang is less than 35, and in Chongqing and Nanning it is less than 5. For formaldehyde, the RND in each city is less than 20.

(a)                                                        (b)

Figure 11 Relationship between RND and total RND *N* for indoor PM2.5 (a) and formaldehyde (b) in

Shenyang, Chongqing, Nanning, Beijing and Kunming

Fig. 11 shows the relationship between the total RND *N* and RND for PM2.5 and

formaldehyde in five large cities in China. It can be seen that as the total RND increases,

the RND increases as well, but when the total RND exceeds a certain value, the RND

no longer increases. In this context, the threshold is called the critical total RND. If the

number of dwellings in the region is estimated to be more than the critical total RND,

*N* is considered to be equal to the critical total RND. In the five cities of Kunming,

Beijing, Nanjing, Chongqing and Shenyang, the critical total RND *N* was

approximately 5000 for the calculation of RND with $\Delta_{\bar{x}}=0.1\bar{x}$ and *F(t)* = 87%.

**4. Conclusion**

This study has proposed a series of systematic methods for the sampling of indoor

environmental parameters in urban residential dwellings. This paper starts with the

sampling strategies and sampling frequency for various indoor environmental

parameters in a single residence, and finally proposes the total required number of

599 dwellings for cities in different thermal zones of China. Compared with sampling

600 methods that correspond to the thresholds of various standards, the approach in this

601 study is more systematic. In order to focus on the periodic characteristics of various

602 parameter fluctuations, this study used statistical methods, thereby enhancing the

603 representativeness and credibility of the samples. The relevant findings can be

604 summarized as follows.

605 1) A heat map revealed obvious daily and weekly fluctuations in $CO_2$, TVOCs,

606 formaldehyde, PM2.5, air temperature and relative humidity. The $CO_2$ concentration

607 reaches a peak during meal times. High concentrations of formaldehyde occur in

608 periods with high temperature. PM2.5 exhibits a high concentration in autumn and

609 winter. The distributions of other parameters display various seasonal fluctuations.

610 2) Because of the variations in environmental parameters inside residential

611 dwellings, two sampling strategies, the continuous sampling strategy (CSS) and the

612 discrete sampling strategy (DSS), were compared by means of descriptive statistical

613 analysis. In general, the DSS performs better than the CSS in terms of accuracy.

614 3) Using the DSS as the sampling strategy, this study proposed an algorithm to

615 calculate the optimal sampling frequencies for different indoor environmental

616 parameters. We found that the optimal sampling frequencies for the concentration of

617 TVOCs, PM2.5, $CO_2$, formaldehyde, relative humidity and temperature are 3 days, 4

618 days, 3 days, 2 days, 12 days and 22 days, respectively. The algorithm can effectively

619 extract the periodic fluctuation characteristics of different indoor environmental

620 parameters, thus providing more representative indoor environmental quality data and

621     effectively reducing sampling costs.

622     4) A method was proposed for determining the required number of dwellings (RND)

623     in different thermal zones. Based on simple random sampling, the RND for studying

624     different indoor environmental parameters was calculated. The results of the study can

625     be confirmed with the first and third conclusions; that is, the degree of a parameter's

626     fluctuation determines the sample size that is needed to accurately reflect the

627     characteristics of the data. The required number of dwellings depends on the coefficient

628     of variation of the sampled data. PM2.5 and formaldehyde had greater RND values than

629     did temperature and humidity. The RND satisfies the lognormal distribution throughout

630     the year.

631     However, it must be admitted that this study has certain limitations. First, the

632     methodology of this study is mainly based on statistical principles and cannot be used

633     trace the sources of pollutants. It can only make certain assumptions based on the

634     fluctuations of parameters. Second, the raw data is limited to that from residential

635     dwellings in Kunming, China. Nevertheless, the series of sampling strategies and

636     algorithms proposed in this paper are applicable to other types of sampling processes

637     and to residential environment research in other regions. Since the proposed algorithm

638     can effectively extract the periodic fluctuation characteristics of different indoor

639     environmental parameters and reduce sampling costs, we expect that the method will

640     provide valuable assistance to future IEQ long-term monitoring studies.

641

## 5. Acknowledgements

## References

[1] N.E. Klepeis, W.C. Nelson, W.R. Ott, J.P. Robinson, W.H.J.J.E.A.E.E. Engelmann, The National Human Activity Pattern Survey (NHAPS): a resource for assessing exposure to environmental pollutants, 11(3) (2000) 231-252.

[2] H.-J. Oh, N.-N. Jeong, J.-R. Sohn, J. Kim, Personal exposure to indoor aerosols as actual concern: Perceived indoor and outdoor air quality, and health performances, Building and Environment 165 (2019).

[3] S. Mentese, N.A. Mirici, M.T. Otkun, C. Bakar, E. Palaz, D. Tasdibi, S. Cevizci, O. Cotuker, Association between respiratory health and indoor air pollution exposure in Canakkale, Turkey, Building and Environment 93 (2015) 72-83.

[4] F. Wu, D. Jacobs, C. Mitchell, D. Miller, M.H.J.E.H.P. Karol, Improving Indoor Environmental Quality for Public Health: Impediments and Policy Recommendations, 115(6) (2007) 953-957.

[5] G. de Gennaro, P.R. Dambruoso, A. Di Gilio, V. Di Palma, A. Marzocca, M. Tutino, Discontinuous and Continuous Indoor Air Quality Monitoring in Homes with Fireplaces or Wood Stoves as Heating System, Int J Environ Res Public Health 13(1) (2015) 78.

[6] S.C. Doll, E.L. Davison, B.R. Painting, Weatherization impacts and baseline indoor environmental quality in low income single-family homes, Building and Environment 107 (2016) 181-190.

[7] J. Langevin, P.L. Gurian, J. Wen, Tracking the human-building interaction: A longitudinal field study of occupant behavior in air-conditioned offices, Journal of Environmental Psychology 42 (2015) 94-115.

[8] T. Deng, X. Shen, X. Cheng, J. Liu, Investigation of window-opening behaviour and indoor air quality in dwellings situated in the temperate zone in China, Indoor and Built Environment (2020).

[9] Y. Geng, B. Lin, J. Yu, H. Zhou, W. Ji, H. Chen, Z. Zhang, Y. Zhu, Indoor environmental quality of green office buildings in China: Large-scale and long-term measurement, Building and Environment 150 (2019) 266-280.

[10] Y. Men, J. Li, X. Liu, Y. Li, K. Jiang, Z. Luo, R. Xiong, H. Cheng, S. Tao, G. Shen, Contributions of internal emissions to peaks and incremental indoor PM2.5 in rural coal use households, Environ Pollut 288 (2021) 117753.

[11] J. Palmisani, A. Di Gilio, M. Viana, G. de Gennaro, A. Ferro, Indoor air quality evaluation in oncology units at two European hospitals: Low-cost sensors for TVOCs, PM2.5 and CO2 real-time monitoring, Building and Environment 205 (2021).

[12] X. Dai, J. Liu, X. Li, L. Zhao, Long-term monitoring of indoor CO2 and PM2.5 in Chinese homes: Concentrations and their relationships with outdoor environments, Building and Environment 144 (2018) 238-247.

680    [13] K. Huang, J. Song, G. Feng, Q. Chang, B. Jiang, J. Wang, W. Sun, H. Li, J. Wang, X. Fang, Indoor
681    air quality analysis of residential buildings in northeast China based on field measurements and
682    longtime monitoring, Building and Environment 144 (2018) 171-183.
683    [14] A. Tunyagi, T. Dicu, K. Szacsvai, B. Papp, G. Dobrei, C. Sainz, A. Cucoş, Automatic system for
684    continuous monitoring of indoor air quality and remote data transmission under SMART_RAD_EN
685    Project, Studia Universitatis Babeş-Bolyai Ambientum 62(2) (2017) 71-80.
686    [15] P.S. Hui, L.T. Wong, K.W. Mui, Sampling strategies of indoor air quality assessment for offices,
687    Facilities 25(5/6) (2007) 179-184.
688    [16] A.K.J.A.T. Persily, Evaluating building IAQ and ventilation with indoor carbon dioxide, 103(Pt
689    2) (1996).
690    [17] K.W. Mui, L.T. Wong, P.S. Hui, A New Sampling Approach for Assessing Indoor Air Quality,
691    Indoor and Built Environment 15(2) (2016) 165-172.
692    [18] C.M. Long, H.H. Suh, P. Koutrakis, Characterization of indoor particle sources using continuous
693    mass and size monitors, J Air Waste Manag Assoc 50(7) (2000) 1236-50.
694    [19] E. Asadi, M.C. da Silva, J.J. Costa, A systematic indoor air quality audit approach for public
695    buildings, Environ Monit Assess 185(1) (2013) 865-75.
696    [20] China's Ministry of Health, GB/T18883. Indoor Air Quality Standard, 2002.
697    [21] World Health Organization(WHO), WHO Air quality guidelines for particulate matter, ozone,
698    nitrogen dioxide and sulfur dioxide: global update 2005, 2020.
699    [22] United States Environmental Protection Agency(EPA), Typical Indoor Air Pollutants. , 2009.
700    [23] M. Mannan, S.G. Al-Ghamdi, Indoor Air Quality in Buildings: A Comprehensive Review on the
701    Factors Influencing Air Pollution in Residential and Commercial Structure, Int J Environ Res Public
702    Health 18(6) (2021).
703    [24] Indoor Air Quality Management Group, A Guide on Indoor Air Quality Certification Scheme
704    for Offices and Public Places., 2019.
705    [25] World Health Organization(WHO), Air Quality Guidelines for Europe, 2000.
706    [26] Institute of Environmental Epidemiology; Ministry of the Environment: Singapore, Guidelines
707    for Good Indoor Air Quality in Office Premises, 1996, pp. 1-47.
708    [27] National Institute for Occupational Safety and Health(NIOSH), NIOSH Pocket Guide to
709    Chemical Hazards (NPG). 2004.
710    [28] Health Canada, Residential Indoor Air Quality Guideline, 2007.
711    [29] Health and Safety Executive, EH40/2005 Workplace Exposure Limits, 2011.
712    [30] The National Health and Medical Research Council, Goals for Maximum Permissible Levels of
713    Pollutants in Indoor Air. In Interim National Indoor Air Quality Goals; The National Health and
714    Medical Research Council: Melbourne, Australia, 1996.
715    [31] China's Ministry of Construction, GB-50325-2020. Standard for indoor environmental
716    pollution control of civil building engineering, 2020.
717    [32] P.K. Cheung, C.Y. Jim, Impacts of air conditioning on air quality in tiny homes in Hong Kong,
718    Sci Total Environ 684 (2019) 434-444.
719    [33] S. Mentese, N.A. Mirici, T. Elbir, E. Palaz, D.T. Mumcuoğlu, O. Cotuker, C. Bakar, S. Oymak, M.T.
720    Otkun, A long-term multi-parametric monitoring study: Indoor air quality (IAQ) and the sources
721    of the pollutants, prevalence of sick building syndrome (SBS) symptoms, and respiratory health
722    indicators, Atmospheric Pollution Research 11(12) (2020) 2270-2281.
723    [34] A.Y. Lim, M. Yoon, E.H. Kim, H.A. Kim, M.J. Lee, H.K. Cheong, Effects of mechanical ventilation

724    on indoor air quality and occupant health status in energy-efficient homes: A longitudinal field
725    study, Sci Total Environ 785 (2021) 147324.

726    [35] S. Abraham, X. Li, A Cost-effective Wireless Sensor Network System for Indoor Air Quality
727    Monitoring Applications, Procedia Computer Science 34 (2014) 165-171.

728    [36] China's Ministry of Housing and Urban-Rural Development, thermal Design Code for Civil
729    Buildings.(GB 50176-93), China building industry press, 2016.

730    [37] J. Liu, X. Dai, X. Li, S. Jia, J. Pei, Y. Sun, D. Lai, X. Shen, H. Sun, H. Yin, K. Huang, H. Tan, Y. Gao,
731    Y. Jian, Indoor air quality and occupants' ventilation habits in China: Seasonal measurement and
732    long-term monitoring, Building and Environment 142 (2018) 119-129.

733    [38] Y. Li, Y.-h. Chiu, L.C. Lu, Energy and AQI performance of 31 cities in China, Energy Policy 122
734    (2018) 194-202.

735    [39] Y. Zhou, H. Xiao, H. Guan, N. Zheng, Z. Zhang, J. Tian, L. Qu, J. Zhao, H. Xiao, Chemical
736    composition and seasonal variations of PM2.5 in an urban environment in Kunming, SW China:
737    Importance of prevailing westerlies in cold season, Atmospheric Environment 237 (2020).

738    [40] Y. Wang, D. Mauree, Q. Sun, H. Lin, J.L. Scartezzini, R. Wennersten, A review of approaches to
739    low-carbon transition of high-rise residential buildings in China, Renewable and Sustainable
740    Energy Reviews 131 (2020).

741    [41] J. Pei, Y. Yin, J. Liu, Long-term indoor gas pollutant monitor of new dormitories with natural
742    ventilation, Energy and Buildings 129 (2016) 514-523.

743    [42] C.J. Kahler, B. Sammler, J. Kompenhans, Generation and control of tracer particles for optical
744    flow investigations in air, EXPERIMENTS IN FLUIDS 33(6) (2002) 736-742.

745    [43] L. Lei, W. Chen, Y. Xue, W. Liu, A comprehensive evaluation method for indoor air quality of
746    buildings based on rough sets and a wavelet neural network, Building and Environment 162 (2019).

747    [44] S. García, J. Luengo, F. Herrera, Tutorial on practical tips of the most influential data
748    preprocessing algorithms in data mining, Knowledge-Based Systems 98 (2016) 1-29.

749    [45] H. Zhou, Z. Deng, Y. Xia, M. Fu, A new sampling method in particle filter based on Pearson
750    correlation coefficient, Neurocomputing 216 (2016) 208-215.

751    [46] W.A.N. Joseph Lee Rodgers, Thirteen Ways to Look at the Correlation Coefficient., The
752    American Statistician 42 (1988) 59-66.

753    [47] A. Houghton, C. Castillo-Salgado, Analysis of correlations between neighborhood-level
754    vulnerability to climate change and protective green building design strategies: A spatial and
755    ecological analysis, Building and Environment 168 (2020).

756    [48] J.P. Gao, L.P. Ma, Statistics, Capital Economy and Trade University Publishing Company, Beijing
757    (2004) 160-190.

758    [49] S.C. Sekhar, S.E. Goh, Thermal comfort and IAQ characteristics of naturally/mechanically
759    ventilated and air-conditioned bedrooms in a hot and humid climate, Building and Environment
760    46(10) (2011) 1905-1916.

761    [50] M.M.M. Abdel-Salam, Outdoor and indoor factors influencing particulate matter and carbon
762    dioxide levels in naturally ventilated urban homes, Journal of the Air & Waste Management
763    Association 71(1) (2021) 60-69.

764    [51] N. Jacek, G. Taseusz, K.Z. Bozena, L. Jerzy, Indoor Air Quality (IAQ), Pollutants, Their Sources
765    and Concentration Levels, Building and Environment Vol. 27, No. 3 (1992) No. 3, pp. 339-356.

766    [52] Y. Zhao, A. Li, R. Gao, P. Tao, J. Shen, Measurement of temperature, relative humidity and
767    concentrations of CO, CO2 and TVOC during cooking typical Chinese dishes, Energy and Buildings

768   69 (2014) 544-561.

769   [53] C. Huang, W. Liu, J. Cai, X. Wang, Z. Zou, C. Sun, Household formaldehyde exposure and its

770   associations with dwelling characteristics, lifestyle behaviours, and childhood health outcomes in

771   Shanghai, China, Building and Environment 125 (2017) 143-152.

772   [54] A. Stamatelopoulou, D.N. Asimakopoulos, T. Maggos, Effects of PM, TVOCs and comfort

773   parameters on indoor air quality of residences with young children, Building and Environment 150

774   (2019) 233-244.

775   [55] Z. Li, Q. Wen, R. Zhang, Sources, health effects and control strategies of indoor fine particulate

776   matter (PM2.5): A review, Sci Total Environ 586 (2017) 610-622.

777   [56] Y. Zhao, H. Sun, D. Tu, Effect of mechanical ventilation and natural ventilation on indoor

778   climates in Urumqi residential buildings, Building and Environment 144 (2018) 108-118.

779   [57] M. Guo, X. Pei, F. Mo, J. Liu, X. Shen, Formaldehyde concentration and its influencing factors

780   in residential homes after decoration at Hangzhou, China, Journal of Environmental Sciences 25(5)

781   (2013) 908-915.

782   [58] T. Godish, J. Rouch, Mitigation of residential formaldehyde contamination by indoor climate

783   control, American Industrial Hygiene Association journal 47 (1986) 792-797.

784   [59] C.R. Frihart, J.M. Wescott, T.L. Chaffee, K.M. Gonner, Formaldehyde Emissions from Urea-

785   Formaldehyde- and No-Added-Formaldehyde-Bonded Particleboard as Influenced by

786   Temperature and Relative Humidity, Forest Products Journal 62 (2012) 551-558.

787   [60] L. Zhao, J. Liu, J. Ren, Impact of various ventilation modes on IAQ and energy consumption in

788   Chinese dwellings: First long-term monitoring study in Tianjin, China, Building and Environment

789   143 (2018) 99-106.

790

**Highlights**

- A systematic approach was presented for IEQ data sampling.

- Sampling strategy, frequency and required number of dwellings were studied.

- Discrete sampling strategy achieves better performance.

- An algorithm was proposed to calculate sampling frequency of IEQ parameters.

- Required number of dwellings depends on the variable coefficient of parameters.

**Declaration of interests**

☒ The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

☐ The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: