

31 and (3) a Business Hours Design. We used Monte Carlo resampling to investigate the range of
32 possible outcomes (i.e., the resulting annual average concentration prediction) from each design
33 against the “truth”, the actual monitoring data. We found that the Balanced Design consistently
34 yielded the most accurate annual averages; Rush Hours and Business Hours Designs generally
35 resulted in comparatively more biased estimates and model predictions. Importantly, the superior
36 performance of the Balanced Design was evident when predictions were evaluated against true
37 concentrations but less detectable when predictions were evaluated against the measurements
38 from the same sampling campaign since these were themselves biased. This result is important
39 since mobile monitoring campaigns that use their own measurements to test the robustness of the
40 results may underestimate the level of bias in their results. Appropriate study design is crucial for
41 mobile monitoring campaigns aiming to assess accurate long-term exposure in epidemiologic
42 cohorts. Campaigns should aim to implement balanced designs that sample during all seasons of
43 the year, days of the week, and all or most hours of the day to produce generally unbiased, long-
44 term averages. Furthermore, differential exposure misclassification could result from unbalanced
45 designs, which may result in misleading health effect estimates in epidemiologic investigations.

46 1 Introduction

47
48 A large body of evidence links long-term exposure to air pollution to adverse health
49 effects in humans, including mortality from cardiovascular outcomes and lung cancer.¹⁻⁶ Most
50 such studies focus on criteria air pollutants such as PM_{2.5} and NO₂, in large part to use available
51 monitoring data to estimate exposures. While air pollution cohort studies may leverage data from
52 exposure assessment campaigns that supplement the regulatory network data, few focus on
53 ambient air pollution exposures that are not criteria air pollutants.⁷⁻⁹

54 Recently, mobile monitoring campaigns have been used to estimate long-term average air
55 pollutant levels.¹⁰⁻¹⁶ Mobile monitoring campaigns typically equip a vehicle with air monitors
56 and collect samples while in motion (mobile sampling) or while stopped (stationary sampling). A
57 single monitoring platform can thus be used to collect samples at many locations within a
58 relatively short period of time, making it a time and cost-efficient sampling approach. Mobile
59 campaigns are particularly well-suited for multi-pollutant monitoring of less frequently
60 monitored traffic-related air pollutants that require expensive instruments or instruments that
61 need frequent attention during the sampling period.

62 While many recent mobile monitoring campaigns have been leveraged to map air
63 pollution exposures and link them to health studies, there has not been any literature on the
64 appropriate design of a mobile monitoring campaign for application to epidemiologic cohort
65 studies. Many mobile campaigns have been short, lasting from a few weeks to months and with
66 few repeat visits to each location spanning one to three seasons.^{12,17,18} Most campaigns have
67 conducted sampling during weekday business or rush hours, ignoring the surrounding hours,
68 when air pollution concentrations can be drastically different.^{19,20} The design of these studies
69 may not be adequate for long-term exposure cohort studies.

70 The goal of this paper is to shed light on the design of a mobile monitoring campaign for
71 application to epidemiologic cohort studies. We carry out a set of simulation studies to better
72 understand the role of mobile monitoring design on the prediction of long-term average surfaces.
73 We use existing monitoring data from California to compare the primary, long-term site averages
74 when all of the data are included to subsequent analyses utilizing subsets of the data. These data
75 provide a unique opportunity to explore how temporal sampling strategies can influence the
76 resulting estimated annual-average concentration. Our analysis requires having a long-term,

77 comprehensive set of measurement data, which therefore necessitates using fixed-site
78 measurements rather than mobile measurements, to shed light on an aspect of study design for
79 mobile monitoring.

80

81 2 Methods

82

83 2.1 Data

84

85 We simulate three sampling designs (see below) using hourly observations for oxides of
86 nitrogen (NO_x) collected during 2016 from regulatory Air Quality System (AQS) sites in
87 California. NO_x was selected since it is a spatially and temporally variable traffic pollutant with
88 a strong diurnal pattern,^{11,21,22} and it is measured at many regulatory monitoring sites in
89 California, providing a large enough dataset for this analysis.²³

90 We required that NO_x observations meet various criteria to be included in this analysis.
91 First, sites needed to have readings at least 66% of the time (5,797/8,784 hourly samples; 2016
92 was a leap year). Second, sites needed to have sampling throughout the year, such that data
93 collection gaps were a maximum of 45 days long. Third, sites were required to have sampled for
94 at least 40% of the time during various two-week periods that were used in two of our
95 “common” designs (described below). This sample size ensured that we could sample during
96 these periods without replacement. Fourth, sites were required to have positive readings (> 0
97 ppb) at least 60% of the time, thus ensuring that sites had sufficient variability in their
98 concentrations and allowing us to model annual averages on the natural log scale. Finally, sites
99 in rural and industrial settings (as determined by the US EPA)²⁴ were excluded.

100

101 2.2 Sampling Designs

102

103 We conducted simulation studies to characterize the properties of three sampling designs
104 (Table 1, Supplementary Information [SI] Figure S1). Each design has a long- and a short-term
105 sampling approach. Long-term approaches use all of the data that meet each design’s definition
106 to estimate site annual averages and are analogous to traditional, fixed-site sampling approaches

107 where sampling at a given location occurs over an extended period of time. Short-term
 108 approaches only collect 28 samples per site and are analogous to mobile monitoring campaigns
 109 that collect a few repeat samples per site. (The cut-off of 28 samples reflects our preliminary
 110 analyses showing that 28 hourly NO_x samples are sufficient to estimate a site’s annual average
 111 within about 25% error or less [SI Figure S2].) Each design has multiple versions where samples
 112 are collected at slightly different times. The various design versions are intended to reflect the
 113 bias produced if only certain times are included in the measurements. We simulated each short-
 114 term sampling approach 30 times (Monte Carlo resampling), and hereafter refer to each of these
 115 simulations as a “campaign” since each represents a potential mobile monitoring study.

117 *Table 1. Simulated sampling designs used to estimate site annual averages.¹*

Design	Sampling Approach and Number of Samples²	Sampling Seasons	Sampling Days	Sampling Hours
Year-Around “Balanced” Design	Long-term (Max 8,784 (V1), 7,320 (V2), or 4,392 (V3) hourly samples per site x 1 campaign simulation)	Winter, spring, summer, fall	Mon – Sun	<u>V1</u> (All Hours) <u>V2</u> (Most Hours): 5 AM – 12 AM <u>V3</u> (Truncated Hours): 6-9 AM, 1-5 PM, 8-10 PM
	Short-term (28 hourly samples per site x 30 campaign simulations)	7 samples per season	5/7 weekday samples; 2/7 weekend samples	Random hours according to V1, V2, or V3
Two-Season Weekday “Rush Hours” Design⁴	Long-term (Max 160 hourly samples per site x 1 campaign simulation)	<u>V4-5</u> : winter & summer (2-wk period per season)	Mon – Fri	7-10 AM, 3-6 PM

	simulation)	<u>V6-7</u> : spring & fall (2-wk period per season)		
	Short-term (28 hourly samples per site x 30 campaign simulations)	14 samples per season		Random Rush Hours according to V4-5 or V6-7
Two-Season Weekday “Business Hours” Design⁴	Long-term	<u>V4-5</u> : winter & summer (2-wk period per season)	Mon – Fri	9 AM – 5 PM
	(Max 180 hourly samples per site x 1 campaign simulation)	<u>V6-7</u> : spring & fall (2-wk period per season)		
	Short-term	14 samples per season		Random Business Hours according to V4-5 or V6-7
	(28 hourly samples per site x 30 campaign simulations)			

118 ¹ There are three archetypal sampling designs, each with long- and short-term sampling
 119 approaches and multiple versions. Long-term approaches are analogous to traditional, fixed-site
 120 sampling, while short-term approaches are analogous to mobile monitoring campaigns. Short
 121 names for the sampling design appear in quotes.

122 ² Long-term approaches have 1 campaign simulation (each includes all of the available data that
 123 meet that design’s criteria), while short-term approaches have 30 campaign simulations (each
 124 with 28 samples). Maximum hourly samples per site varied because some sites had missing
 125 readings. Year 2016 was a leap year.

126 ⁴ See SI Table S1 for each version’s exact sampling periods

127

128 The Year-Around “Balanced” Design represents an “ideal” sampling scheme: sampling is
 129 conducted during all seasons, days of the week, and all or most hours of the day. Version 1
 130 collects samples during all hours of the day. Versions 2-3 reduce the sampling hours to reflect
 131 the logistical constraints of executing an extensive campaign: samples occur during most hours
 132 of the day (5 AM – 12 AM only; “Version 2”) or during 6-9 AM, 1-5 PM and 8-10 PM

133 (“Version 3”). Estimates from the long-term Balanced Design Version 1 are analogous to what
134 might be collected from a traditional, year-around, fixed-site sampling scheme. For simplicity,
135 we interchangeably refer to these as the “true” estimates or the “gold standard” hereafter, though
136 we acknowledge that some error exists (e.g., due to missing hours or instrument accuracy).

137 The Two-Season Weekday “Rush Hours” and “Business Hours” Designs reflect common
138 designs in the literature. Samples are collected either during summer and winter (Versions 4-5)
139 or spring and fall (Versions 6-7). Sampling for each version occurs on weekdays during a two-
140 week period each relevant season (See SI S1 for each version’s exact sampling periods).
141 Sampling is restricted to the hours of 7-10 AM and 3-6 PM (Rush Hours Design) or 9 AM – 5
142 PM (Business Hours Design). The short-term approach collects 14 random samples during each
143 season.

144

145 2.3 Prediction Models

146

147 We estimated site annual averages from the data collected during each campaign. We
148 log-transformed these before using them as the outcome variable in partial least squares (PLS)
149 regression models, which summarized hundreds of geographic covariate predictors (e.g., land
150 use, road proximity, and population density; see SI Table S2 for the covariates considered) into
151 two PLS components (using the `pls` function in the `pls` package in R). We evaluated the
152 performance of each campaign using ten-fold cross-validated (CV) predictions on the native
153 scale, incorporating re-estimation of the PLS components in each fold. The cross-validation
154 groups were randomly selected and, importantly, fixed across all campaigns to allow for
155 consistent model performance comparisons across design versions.

156 To best understand the role of design, we present results for annual average estimates,
157 predictions, and model performance statistics. In descriptive analyses, we compare design-
158 specific annual average estimates and predictions to the gold standard (long-term Balanced
159 Design Version 1). We compare predicted site concentrations against predictions from the gold
160 standard since epidemiologic air pollution studies often rely on predicted exposure, and the gold
161 standard prediction represents the best possible prediction of annual-average concentrations that
162 a study could hope to achieve. We complement this approach with model assessment evaluations
163 of design-specific site predictions against two different references: an assessment against the true

164 averages, and a traditional model assessment evaluation against the respective design-specific
165 annual average estimates. The traditional assessment compares the predicted exposures to the
166 observed site measurements from which they were derived. This allows us to document the
167 quantities that would normally be available from modeling the data measured from any specific
168 campaign. We summarize the model performance in terms of cross-validated mean squared error
169 (MSE)-based R^2 (R^2_{MSE}), regression-based R^2 (R^2_{reg}), and root mean squared error (RMSE).
170 R^2_{MSE} reflects bias as well as variation around the one-to-one line. R^2_{reg} is based on the best fit
171 line between the measurements and predictions, which adjusts for bias and slopes different than
172 one, and is defined as the squared correlation between a measurement and a prediction. See SI
173 Equations 1-3 for definitions.

174 We repeated these analyses for nitrogen dioxide (NO₂) and nitrogen monoxide (NO),
175 adding a two ppb constant to all of the hourly NO readings before log-transforming to eliminate
176 negative and zero concentration readings.

177 All analyses were conducted in R (v 3.6.2, using RStudio v 1.2.5033).²⁵

178

179 3 Results

180

181 We included the 69 of 105 California AQS sites that met our data criteria (Figure 1, SI
182 Figure S3). Sites were located in both urban and suburban settings, in residential and commercial
183 areas.

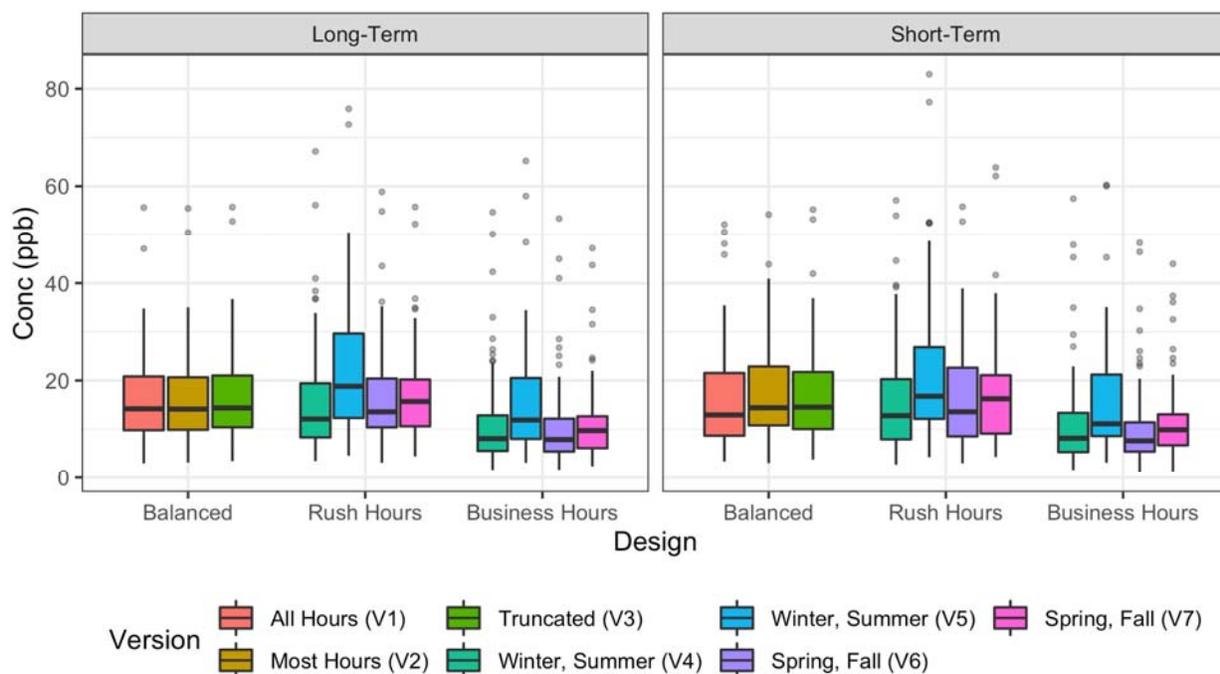
195 3.2 Annual Average Estimates

196

197 Across the 69 monitor locations, measured annual average concentrations (long-term
198 Balanced Design Version 1), had a median (IQR) of 14 (10 - 21) ppb and ranged from 3-56 ppb.
199 The short-term and long-term sampling approaches resulted in similar distributions of annual
200 averages for different design versions. Figure 2 shows the long-term and a single short-term
201 approach for each design. Overall, the long-term and short-term approach for each design
202 version had very similar distributions. All of the Balanced Design versions resulted in only slight
203 differences in their medians and IQRs. The Rush Hours Design versions generally resulted in
204 slightly higher annual averages than the true averages, with some versions being more variable
205 and having somewhat different distributions. The Business Hours Design versions resulted in
206 annual averages that were generally lower than the true averages and less variable than the Rush
207 Hours Design versions. See SI Table S5 for summary statistics. SI Figure S7 shows annual
208 average estimates for all campaigns and pollutants.

209

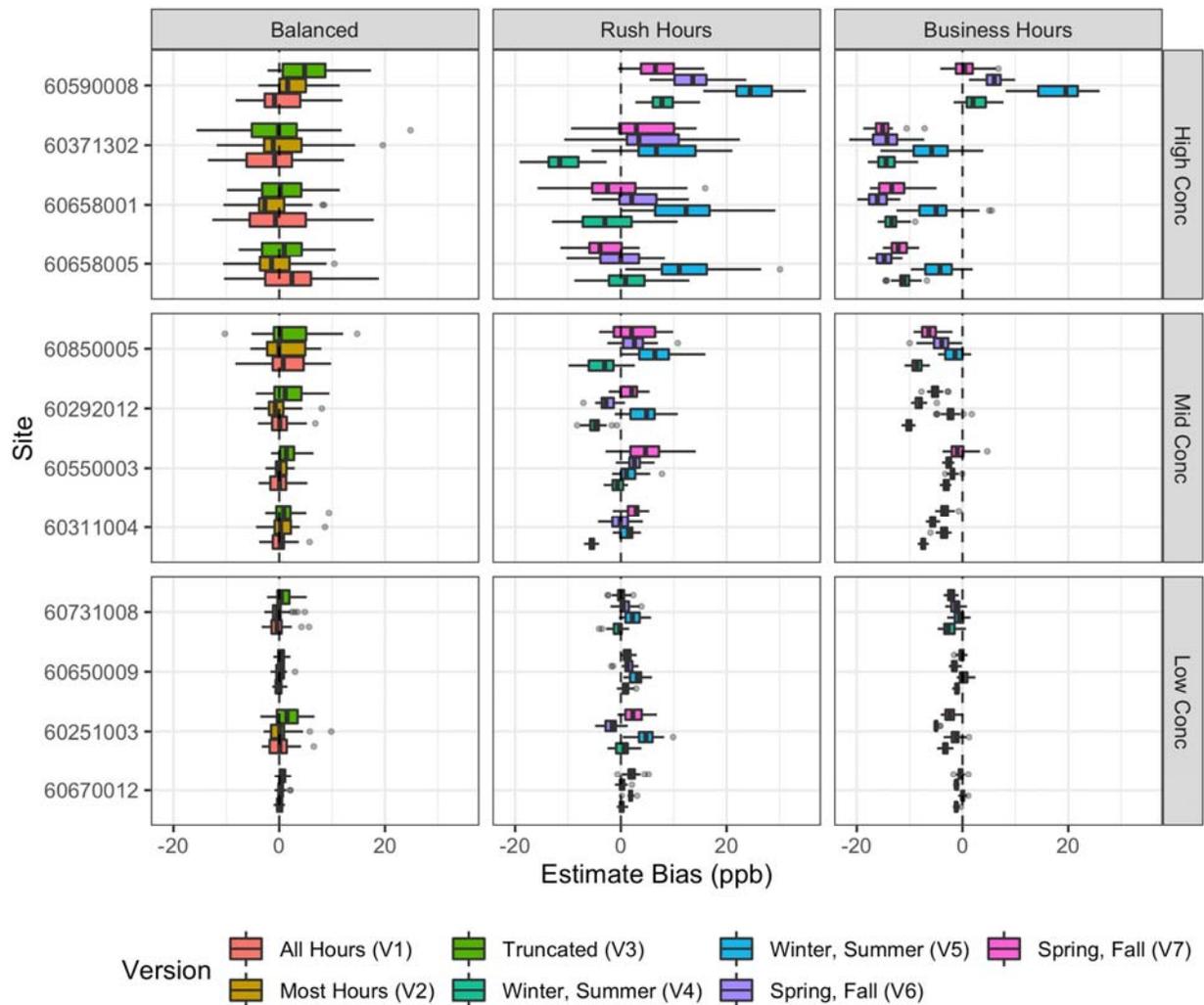
210



211
212 *Figure 2. Distribution of NO_x annual averages (N=69 sites) from different design versions. Showing the*
213 *one campaign for each long-term approach and one example campaign for each short-term approach.*

214
215 Figure 3 shows the site-specific distributions of annual averages across designs for short-
216 term approaches relative to the true averages for a stratified random sample of 12 sites. Sites are
217 stratified by whether their true mean concentration was in the low (<25th percentile), middle
218 (25th-75th percentile) or high (>75th percentile) concentration category. The variation of averages
219 across campaigns increases with concentration in all designs. Site-specific averages are similar to
220 the true averages for all Balanced Design versions while there were multiple sites from the
221 Business Hours Design versions with averages systematically lower. The Rush Hours Design
222 versions also had many biased averages, although the direction of the bias varied by site and
223 design version. SI Figure S8 shows these biases for all sites.

224



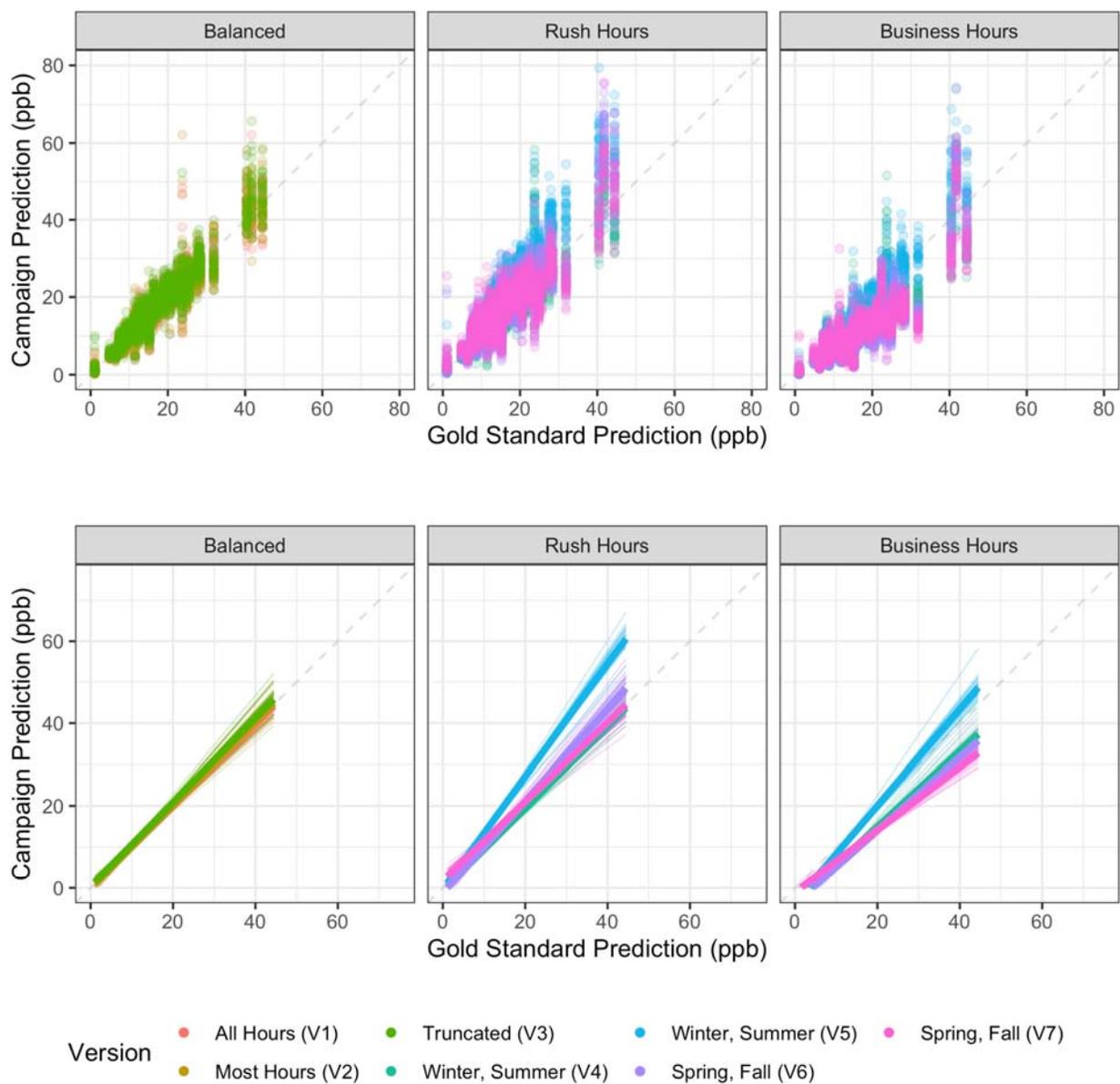
225
 226 *Figure 3. Site-specific NO_x measurement biases for short-term designs (N = 30 campaigns) as compared*
 227 *to the true annual average at that site (long-term Balanced Design Version 1). Showing a stratified*
 228 *random sample of 12 sites, stratified by whether their true concentration was in the low (<25th*
 229 *percentile), middle (25th-75th percentile) or high (>75th percentile) concentration category and arranged*
 230 *within each stratum with lower concentration sites being closer to the bottom.*

231
 232 **3.3 Model Predictions**

233
 234 The PLS model of the true annual average had a root mean square error (RMSE) of 7.2
 235 ppb and a mean square error-based coefficient of determination (R^2_{MSE}) of 0.46.

236 We compared PLS model predictions from each short-term design to the gold standard
 237 model predictions. SI Figure S9 shows the relative standard deviations of predictions by design
 238 version, with 1 indicating that design predictions have the same standard deviation as the gold

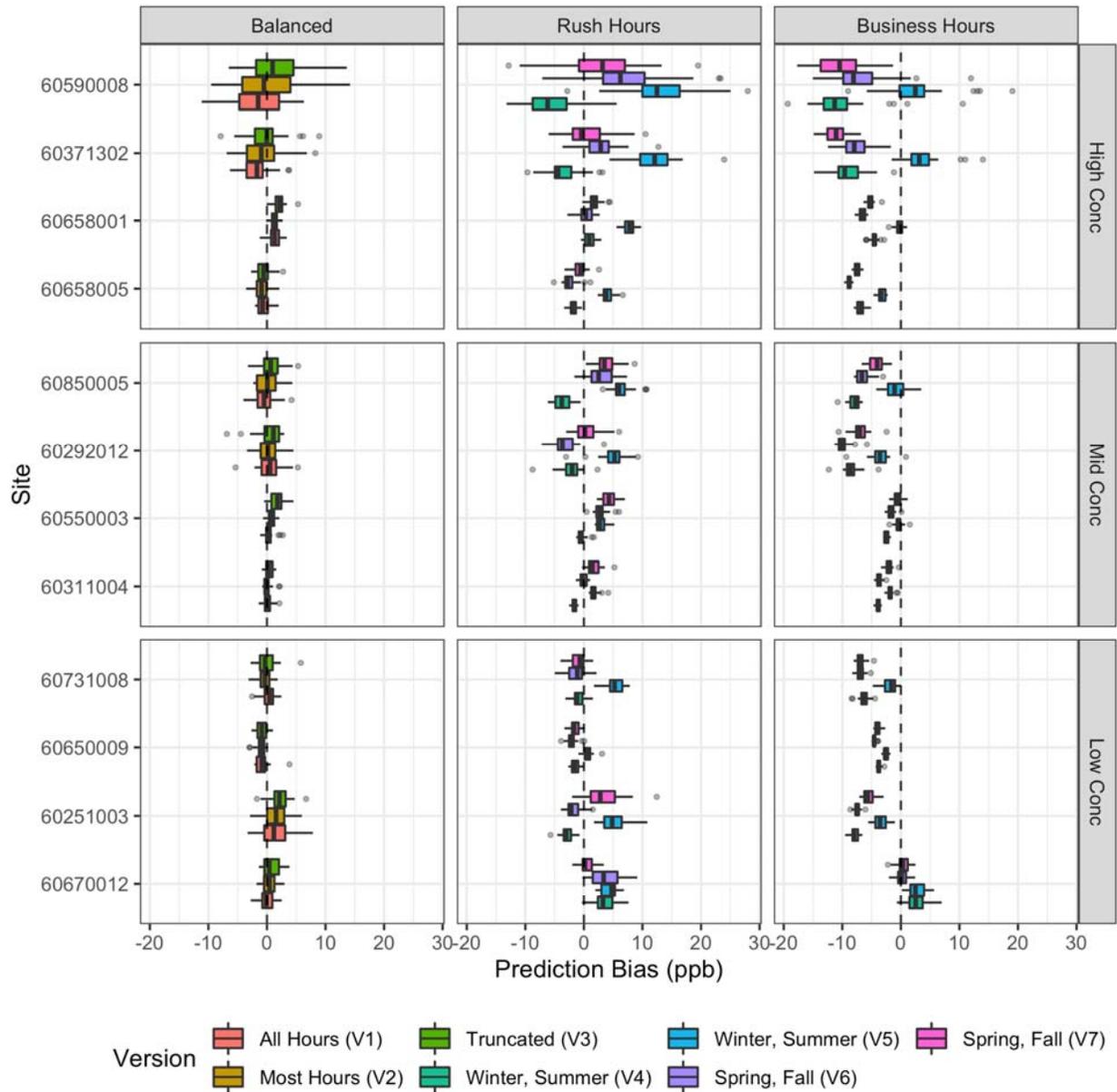
239 standard model predictions. Overall, the Balanced Design predictions have similar variability to
240 those of the gold standard (range: 0.87-1.28), the Rush Hours Design predictions are more
241 variable (range: 0.90-1.74), and the Business Hours Design predictions are mixed: some less and
242 some more variable (range: 0.73-1.54). Figure 4 displays these comparisons as scatterplots and
243 best fit lines. The scatterplots show that there are a few sites, some of which have high leverage,
244 that have variable predictions in all designs. From the best fit lines, we observe that all short-
245 term Balanced Design versions resulted in the most accurate predictions on average, as indicated
246 by their overlapping general trends along the one-to-one line. The Rush Hours Design versions
247 were more likely to have a positive general trend, while the Business Hours Design versions
248 were more likely to have a negative general trend, indicating, for example, that higher
249 concentrations were more likely to be over- or under-estimated, respectively. However, there
250 was heterogeneity in this overall pattern across the Rush and Business Hours Design versions.
251 Furthermore, there was additional heterogeneity across individual campaigns. The SI contains
252 comparable figures comparing design predictions to the gold standard and additional figures for
253 NO and NO₂ (SI Figures S10-S13).
254



255
 256 *Figure 4. Scatterplots and best fit lines of cross-validated short-term predictions for 30 campaigns vs the*
 257 *gold standard predictions for NO_x. Thin transparent lines are individual campaigns, colored by design*
 258 *version; thicker lines are the overall version trend. (One prediction is excluded for clarity from the Rush*
 259 *Hours Version 4 scatterplot at x=24 ppb, y=109 ppb [site 60731016] but included in the line plots.)*

260
 261 Figure 5 shows site-specific comparisons of predictions across 30 short-term campaigns
 262 relative to the gold standard predictions for a stratified random sample of 12 sites in order to
 263 characterize relative bias (see SI Figure S14 for all sites). Overall, the short-term Balanced
 264 Design predictions had a median (IQR) bias of 0.2 (-1 – 1.4) ppb relative to the gold standard

265 predictions (see SI Table S7 for details). All Balanced Design predictions were very similar to
266 the gold standard predictions, though some sites frequently had larger biases. The Rush Hours
267 and Business Hours Design versions were more likely to consistently produce biased site
268 predictions, with a median (IQR) bias of 1.2 (-1.2 – 4) ppb and -3.8 (-6.6 – -1.4) ppb,
269 respectively. While the Rush Hours Design versions generally resulted in higher predictions
270 across sites (with some inconsistency across versions for a few sites), the Business Hours Design
271 versions resulted in predictions that were both lower and higher than the gold standard
272 predictions across sites. There were also a few sites that tended to have more biased and/or more
273 variable predictions relative to the gold standard across all designs. We observed similar patterns
274 when looking at estimate (rather than prediction) biases (See Figure 3, SI Figure S8).
275



276

277 *Figure 5. Site-specific NO_x prediction biases for short-term designs (N = 30 campaigns) as compared to*
 278 *the gold standard predictions (long-term Balanced Design Version 1). Showing a stratified random*
 279 *sample of 12 sites, stratified by whether true concentrations were in the low (Conc < 0.25), middle (0.25*
 280 *≤ Conc ≤ 0.75) or high (Conc > 0.75) concentration quantile and arranged within each stratum with lower*
 281 *concentration sites closer to the bottom.*

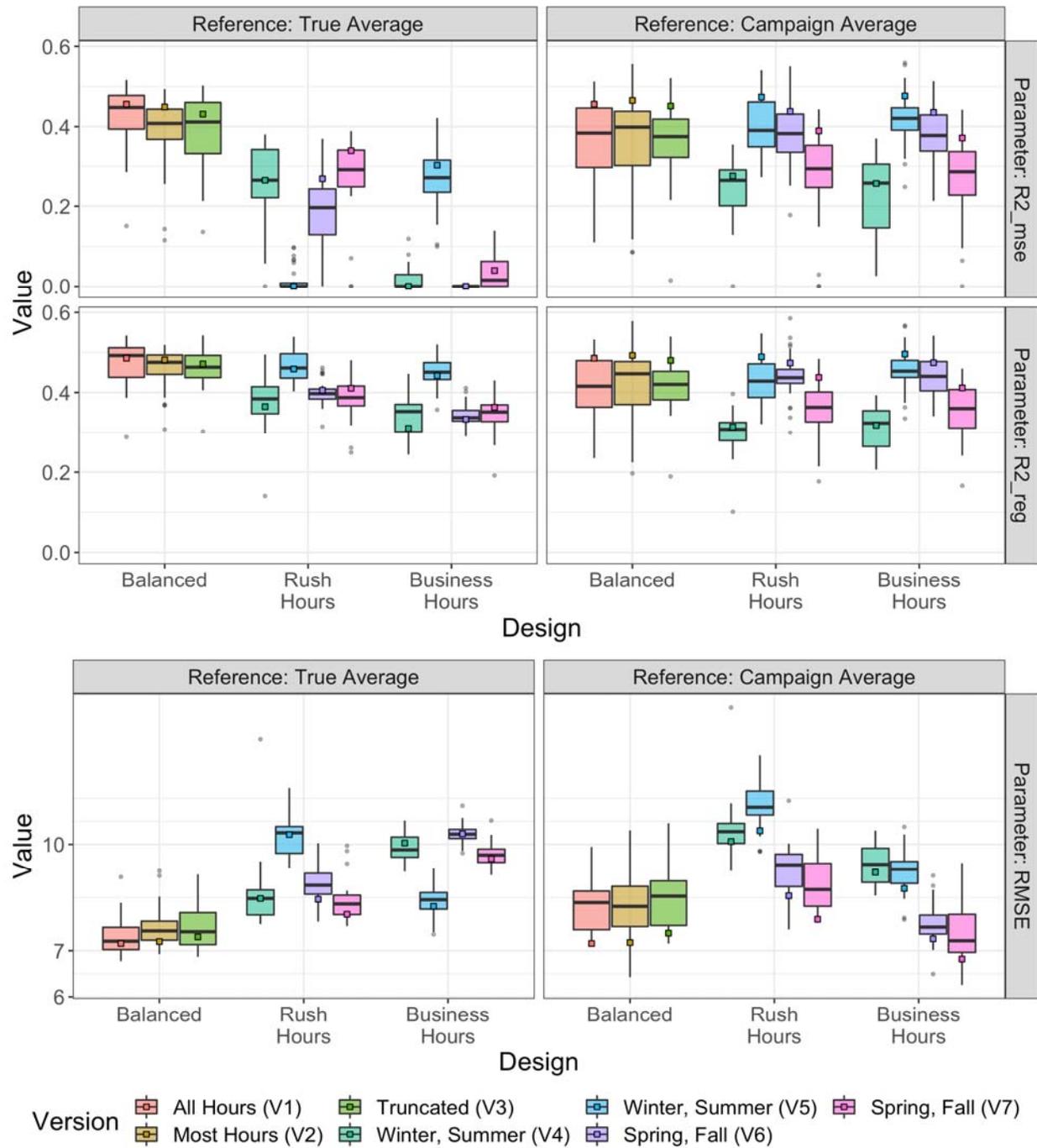
282

283 3.4 Model Assessment

284

285 Figure 6 shows the out-of-sample prediction performances relative to the observations
286 from the true averages (left column) and the specific design (right column), for both the long-
287 term and short-term approaches. The boxplots quantify the distribution of performance statistics
288 across all 30 short-term campaigns while the squares show the performance for the long-term
289 approach of the same design version. When assessed against the true averages, all the Balanced
290 Design versions generally perform better than either the Rush Hours or Business Hours Design
291 versions with higher $CV R^2_{MSE}$ and $CV R^2_{reg}$, and lower $CV RMSE$ estimates. This is particularly
292 apparent for the long-term approach. Furthermore, within design the performance for the long-
293 term approach is better than the majority of the short-term campaigns. There is considerable
294 heterogeneity in performance across the Rush Hours and Business Hours Design versions. In
295 contrast, when assessed against observations from the same design, as would typically be done in
296 practice, the role of sampling design on prediction performance is not as evident. The superior
297 performance of the Balanced Design is not as apparent, and some of the Rush Hours and
298 Business Hours Design versions appear to perform better. There are also a few campaigns that
299 show poor performance, even under the Balanced Design. SI Figure S15-S16 show similar
300 results for NO_2 and NO , with NO showing more variability and some lower performing statistics.

301



302

303 *Figure 6. Model performances (MSE-based R^2 , Regression-based R^2 , and RMSE), as determined by each*
 304 *campaign's cross-validated predictions relative to: a) the true averages (long-term Balanced Version 1),*
 305 *and b) its respective campaign averages. Boxplots are for short-term approaches (30 campaigns), while*
 306 *squares are for long-term approaches (1 campaign).*

307

308 3.5 NO and NO₂

309

310 We found similar results for NO and NO₂ (See the SI).

311

312 4 Discussion

313

314 In this paper we have used existing regulatory monitoring data to deepen our
315 understanding of the importance of mobile monitoring study design for application to
316 epidemiologic cohort studies. Others have shown that short-term data can be used to estimate
317 long-term averages.^{11,12} What has been missing from the literature until now, however, is the
318 impact of mobile monitoring study design on the accuracy and precision of long-term exposure
319 estimates and model predictions, particularly when the goal is to produce predictions for an
320 epidemiologic study. Our results indicate that for designs with a sufficient number of short-term
321 samples (about 28 or more), the design rather than the sampling approach (short- vs long-term)
322 has the largest impact on the estimated long-term averages. We focus the rest of this discussion
323 on the short-term approaches for each design, which resemble mobile monitoring, though the
324 long-term approaches produced similar results.

325 In terms of specific design, we found that all of the Balanced Design versions resulted in
326 similar annual averages as the true averages (long-term Balanced Version 1), while the Rush
327 Hours and Business Hours Design versions were more likely to result in more biased and more
328 or less variable annual average estimates. Specifically, the Rush Hours Design was more likely
329 to overestimate, while the Business Hours Design was more likely to underestimate site
330 averages. This result was likely because the Balanced Design captured much of NO_x's temporal
331 variability by allowing for samples to be collected during each season, day of the week, and all
332 or most times of the day, all periods during which meteorology and traffic activity patterns
333 impact air pollution concentrations (SI Figure S4-S6). The Rush Hours Design, on the other
334 hand, was restricted to two sampling seasons and was more likely to sample during high
335 concentration times of day and days of the week. The Business Hours Design had similar
336 limitations though it was more likely to sample during low concentration times.

337 We found a similar pattern with the predictions: similar predictions across all Balanced
338 Design versions, while most of versions in the Rush Hours tended to overpredict and those in the
339 Business Hours tended to underpredict. However, this varied by design version, suggesting that
340 the particular four weeks of sampling are an important source of heterogeneity in the results. The
341 predictions were more variable for all Rush Hours Design versions and one Business Hours
342 Design version (SI Figure S9). One Business Hours Design version was less variable, while two
343 versions were about the same relative to the gold standard predictions.

344 The similarity in annual averages and predictions across all of the Balanced Design
345 versions suggests that campaigns with slightly reduced sampling hours (for example, due to
346 logistical constraints) should to a large degree still produce unbiased annual averages at most
347 sites. On the other hand, campaigns that follow more temporally restricted sampling designs such
348 as the Rush Hours and Business Hours Designs may produce systematically biased results, with
349 the degree and direction of error being heavily impacted by the sampling window that happens to
350 be selected.

351 At the site level, we saw that while any individual study campaign had the potential to
352 produce biased estimates and predictions, the Rush Hours and Business Hours Designs were
353 more likely to do so than the Balanced Design. The direction and magnitude of the bias varied by
354 site and depended upon the sampling design and the typical seasonal, day of week, and time of
355 day patterns of pollution at that site. This suggests that a simple correction factor (e.g., the ratio
356 of the true annual average concentration relative to the resulting concentration from a given
357 design) is unlikely to appropriately adjust for bias at the site level. Given that higher
358 concentration sites were more likely to have greater degrees of bias and variation (Figure 4 –
359 Figure 5), non-balanced designs may misrepresent some sites more than others and lead to
360 differential exposure misclassification in epidemiologic studies. Thus, while non-balanced
361 design may be appropriate for non-epidemiologic purposes including characterizing the spatial
362 impact of traffic related air pollutants during peak hours for urban planning and policy purposes,
363 these could be misleading in epidemiologic applications.

364 In this study we were able to evaluate prediction model performance against the true
365 annual average exposure as well as against the observations typically available for model
366 performance assessment. Performance assessment against the true averages indicates that the
367 Balanced Design is clearly the best, and that there is little degradation in performance across

368 design versions. This means that it is possible to design high quality mobile monitoring studies
369 that accommodate some measure of logistical feasibility, for example, by not requiring sampling
370 in the middle of the night. In contrast, the performance of the Rush Hours and Business Hours
371 Designs is comparatively worse, indicating that the logistically appealing approach that samples
372 only four weeks during two seasons, during daytime hours, and only during weekdays is
373 inadequate for providing high quality estimates of annual averages. Further, the performance of
374 these designs varies considerably and unpredictably depending upon the specific pair of two-
375 week periods that were selected for sampling. Additionally, comparison of the two R^2 estimates
376 (R^2_{MSE} and R^2_{reg}) indicates that not all of their poor performance is due to the inability to predict
377 the same value as the truth (R^2_{MSE}), but due to systematic bias in the design.

378 Further, it is notable that the standard approach to model assessment, comparing model
379 predictions to observations collected during the sampling campaign, doesn't clearly reveal the
380 superior performance of the Balanced Design or the inherent flaws of the Rush Hours and
381 Business Hours Designs. In fact, some of the Rush Hours and Business Hours Design versions
382 perform better than the Balanced Design when evaluated against the campaign's observations.
383 This is because the evaluation doesn't take into account that the observations are biased because
384 of the sampling design.

385 It is notable that the performance of our short-term campaigns was fairly consistent with,
386 though generally slightly worse than, the performance observed in the long-term campaign for
387 each design version (Figure 6). However, occasionally there was an "unlucky" short-term
388 campaign with meaningfully poorer performance than the other campaigns of the same design.
389 This is true even for the Balanced Design versions where 1-2 of the 30 campaigns (~3-6%) had
390 notably worse performance as quantified by R^2 . It may be possible that this result is driven by a
391 few high-leverage outlier sites that impact the prediction model performance. In practice, mobile
392 monitoring study investigators are likely to investigate high-leverage sites and address their
393 influence in their prediction modeling.

394 Our study focused on short-term campaigns with 28 repeat samples per site. We did not
395 consider campaigns with fewer or more visits. As evident in SI Figure S2, the percent error in
396 estimating the annual average from fewer than 25 visits skyrockets, suggesting that site estimates
397 will be considerably noisier in mobile campaigns with few repeat visits, regardless of the study
398 design. Prediction model performance is thus likely to decrease as the number of visits per site

399 decrease. Logistically, it is also difficult to achieve balance in sampling over time across season,
400 day of week, and time of day with fewer than 28 samples per site. Furthermore, we note that this
401 study focused on a few generalizable, common designs in the literature, though many other
402 approaches have been taken. We expect that the variety of mobile campaign designs that have
403 been implemented will all produce slightly different results.

404 In putting these results in context, it is important to recognize that in this simulation study
405 we are using existing regulatory monitoring data that has been through extensive quality
406 assurance and quality control processes to approximate the data that might be collected by
407 mobile monitoring campaigns. For instance, we used hourly averages to approximate much
408 shorter-term sampling durations (e.g., a few minutes or less) that would be collected during a
409 mobile campaign. Shorter duration sampling will affect the noise in the data, to an amount that
410 depends on the environment (e.g., temporal patterns in the concentrations of the pollutant being
411 measured) and the instruments. (For comparison, however, our additional evaluations of minute-
412 level data suggest that the decrease in percent error in going from two-minute to hour-long
413 samples is at most a few percent.) Further, our study took place throughout California, a large,
414 geographically diverse area with varying climate profiles.²⁶ This may explain the moderate
415 model performance of the gold standard campaign. While such a large sampling domain would
416 be challenging for a real-world mobile monitoring campaign, the overall conclusions of this
417 study are likely generalizable given that traffic-related air pollution concentrations in generally
418 exhibit temporal patterns that vary by location. While we observed moderate model
419 performances, campaigns with smaller, less spatially heterogeneous study areas may see higher
420 performances.

421 There are several other differences between our study and a typical mobile monitoring
422 study. By definition, mobile monitoring campaigns collect non-stationary (mobile)
423 measurements, which are subject to jostling inherent with on-road sampling, even if some of the
424 sampling occurs while the vehicle is stationary. Further, mobile samples may be less precise than
425 what we observed from fixed, regulatory monitoring sites due to differences in instrumentation,
426 instrument quality, and maintenance. Mobile sampling platforms are more likely to be
427 immediately near another vehicle (e.g., in a traffic queue while stopped at a traffic signal) than a
428 fixed-site monitor. Another distinction is that while we sampled measurements within sites at
429 random, mobile campaigns typically sample from sites along a fixed route or in a designated

430 area. This induces some spatial correlation in the mobile monitoring results that is not part of our
431 simulations. Furthermore, we did not consider the importance of the distribution of sampling
432 locations in this study, which is particularly relevant when the exposure assessment goal is an
433 epidemiologic application. Selecting sites that are representative of the target cohort's residence
434 locations will ensure the spatial compatibility assumption is met, which is an important way to
435 reduce the role of exposure measurement error in epidemiologic inference.²⁷ This consideration
436 is especially relevant for mobile monitoring near major sources (e.g., airports, marine activity,
437 and industry),^{11,12,28-34} which may or may not represent a study cohort of interest.

438 Our evaluation focused on NO_x, NO, and NO₂, which are quickly and moderately
439 decaying air pollutants (concentrations reach background levels approximately 400-600 m from
440 roadway sources).²¹ Campaigns that measure these pollutants may be more susceptible to
441 sampling design than campaigns that measure less spatially- and/or temporally-variable
442 pollutants such as PM_{2.5}.²¹ We selected NO_x, NO, and NO₂ because these are often measured in
443 short-term campaigns, and data for these pollutants are more widely available. Non-criteria
444 pollutants, for example ultrafine particulates (UFP), however, have also received increasing
445 attention in recent years given their emerging link to adverse health effects.^{7,35-37} Still, high-
446 quality information about their spatial distribution is essentially absent, and most studies have
447 implemented short-term mobile sampling approaches³⁴ that may not be temporally balanced and
448 potentially be misleading.

449 An important next step in this work is to understand whether the differences in exposure
450 estimates that we observed across study designs have a meaningful impact on epidemiologic
451 inferences. This is of particular interest considering that year-around, balanced designs are
452 resource-intensive and rare, while shorter, more convenient campaigns are more common in the
453 literature. More research is needed to better understand how and whether unbalanced mobile
454 monitoring campaigns may contribute high quality exposure assessments for epidemiology.
455 Regardless of design, we expect that the predictions from all of the campaigns will result in both
456 classical-like and Berkson-like error.^{27,38-40} Specifically, the predictions capture only part of the
457 true long-term exposure (Berkson-like error), while the parameters in the prediction model are
458 inherently noisy (classical-like error). However, these measurement error methods have not to
459 date considered exposure assessment study design, beyond considering the importance of spatial
460 compatibility, i.e., that distribution of monitoring locations is the same as the distribution of

461 participant locations. Our work suggests that deeper understanding of the role of exposure
462 assessment design on epidemiologic inference is an important area of research.

463

464 4.1 Conclusions and Recommendations for Mobile Monitoring Campaigns

465

466 Mobile monitoring study design should be an important consideration for campaigns
467 aiming to assess long-term exposure in an epidemiologic cohort. Given the temporal trends in air
468 pollution, campaigns should implement balanced designs that sample during all seasons of the
469 year, days of the week, and hours of the day in order to produce unbiased long-term averages.
470 Nonetheless, restricting the sampling hours in balanced designs, for example due to logistical
471 considerations, will still generally produce unbiased estimates at most sites. On the other hand,
472 unbalanced sampling designs like those often seen in the literature are more likely to produce
473 biased long-term estimates, with some sites being more biased than others. And while
474 predictions from these restricted designs may at times perform similarly to balanced designs (or,
475 more problematically, may erroneously *appear* to perform similarly when evaluated against
476 measurements which are themselves biased samples), this performance may strongly depend on
477 the exact sampling period chosen and may thus be difficult or impossible to anticipate prior to
478 conducting a new sampling campaign. Furthermore, the differential exposure misclassification
479 that may result from these designs may be problematic in epidemiologic investigations. Finally,
480 studies that implement unbalanced sampling designs may have hidden exposure misclassification
481 given that both the observations and model predictions may be systematically incorrect. By
482 implementing a balanced sampling design, campaigns can thus increase their likelihood of
483 capturing accurate long-term exposure averages.

484

485 5 Funding

486

487 This work was funded by the Adult Changes in Thought – Air Pollution (ACT-AP) Study
488 (National Institute of Environmental Health Sciences [NIEHS], National Institute on Aging
489 [NIA], R01ES026187), and BEBTEH: Biostatistics, Epidemiologic & Bioinformatic Training in
490 Environmental Health (NIEHS, T32ES015459).

491 Research described in this article was conducted under contract to the Health Effects
492 Institute (HEI), an organization jointly funded by the United States Environmental Protection
493 Agency (EPA) (Assistance Award No. CR-83998101) and certain motor vehicle and engine
494 manufacturers. The contents of this article do not necessarily reflect the views of HEI, or its
495 sponsors, nor do they necessarily reflect the views and policies of the EPA or motor vehicle and
496 engine manufacturers.

497 6 References

498

499 1. Hoek, G. *et al.* Long-term air pollution exposure and cardio- respiratory mortality: a review.

500 *Environ. Health* **12**, 43 (2013).

501 2. Kampa, M. & Castanas, E. Human health effects of air pollution. *Environ. Pollut.* **151**, 362–

502 367 (2007).

503 3. R uckerl, R., Schneider, A., Breitner, S., Cyrys, J. & Peters, A. Health effects of particulate

504 air pollution: a review of epidemiological evidence. *Inhal. Toxicol.* **23**, 555–592 (2011).

505 4. Schwartz, J. Air pollution and daily mortality: a review and meta analysis. *Environ. Res.* **64**,

506 36–52 (1994).

507 5. Chen, H., Goldberg, M. & Villeneuve, P. A systematic review of the relation between long-

508 term exposure to ambient air pollution and chronic diseases. *Rev. Environ. Health* **23**, 243–

509 298 (2008).

510 6. Pope, C. A., Dockery, D. W. & Schwartz, J. Review of epidemiological evidence of health

511 effects of particulate air pollution. *Inhal. Toxicol.* **7**, 1–18 (1995).

512 7. Weichenthal, S. *et al.* Within-city Spatial Variations in Ambient Ultrafine Particle

513 Concentrations and Incident Brain Tumors in Adults. *Epidemiology* **31**, (2020).

514 8. Weichenthal, S. *et al.* Long-term exposure to ambient ultrafine particles and respiratory

515 disease incidence in in Toronto, Canada: a cohort study. *Environ. Health* **16**, 64 (2017).

516 9. Downward, G. S. *et al.* Long-term exposure to ultrafine particles and incidence of

517 cardiovascular and cerebrovascular disease in a prospective study of a Dutch cohort.

518 *Environ. Health Perspect.* **126**, 127007 (2018).

- 519 10. Hankey, S. & Marshall, J. D. Land Use Regression Models of On-Road Particulate Air
520 Pollution (Particle Number, Black Carbon, PM_{2.5}, Particle Size) Using Mobile Monitoring.
521 *Environ. Sci. Technol.* **49**, 9194–9202 (2015).
- 522 11. Apte, J. S. *et al.* High-Resolution Air Pollution Mapping with Google Street View Cars:
523 Exploiting Big Data. *Environ. Sci. Technol.* **51**, 6999–7008 (2017).
- 524 12. Hatzopoulou, M. *et al.* Robustness of Land-Use Regression Models Developed from Mobile
525 Air Pollutant Measurements. *Environ. Sci. Technol.* **51**, 3938–3947 (2017).
- 526 13. Patton, A. P. *et al.* Spatial and temporal differences in traffic-related air pollution in three
527 urban neighborhoods near an interstate highway. *Atmos. Environ.* (2014)
528 doi:10.1016/j.atmosenv.2014.09.072.
- 529 14. Van den Bossche, J. *et al.* Mobile monitoring for mapping spatial variation in urban air
530 quality: Development and validation of a methodology based on an extensive dataset. *Atmos.*
531 *Environ.* (2015) doi:10.1016/j.atmosenv.2015.01.017.
- 532 15. Kerckhoffs, J. *et al.* Comparison of ultrafine particle and black carbon concentration
533 predictions from a mobile and short-term stationary land-use regression model. *Environ. Sci.*
534 *Technol.* **50**, 12894–12902 (2016).
- 535 16. Xie, X. *et al.* A Review of Urban Air Pollution Monitoring and Exposure Assessment
536 Methods. *ISPRS International Journal of Geo-Information* vol. 6 (2017).
- 537 17. Weichenthal, S. *et al.* A land use regression model for ambient ultrafine particles in
538 Montreal, Canada: A comparison of linear regression and a machine learning approach.
539 *Environ. Res.* **146**, 65–72 (2016).

- 540 18. Minet, L., Gehr, R. & Hatzopoulou, M. Capturing the sensitivity of land-use regression
541 models to short-term mobile monitoring campaigns using air pollution micro-sensors.
542 *Environ. Pollut.* **230**, 280–290 (2017).
- 543 19. Batterman, S., Cook, R. & Justin, T. Temporal variation of traffic on highways and the
544 development of accurate temporal allocation factors for air pollution analyses. *Atmos.*
545 *Environ.* **107**, 351–363 (2015).
- 546 20. Saha, P. K. *et al.* Quantifying high-resolution spatial variations and local source impacts of
547 urban ultrafine particle concentrations. *Sci. Total Environ.* **655**, 473–481 (2019).
- 548 21. Karner, A. A., Eisinger, D. S. & Niemeier, D. A. Near-roadway air quality: Synthesizing the
549 findings from real-world data. *Environ. Sci. Technol.* **44**, 5334–5344 (2010).
- 550 22. Riley, E. A. *et al.* Multi-pollutant mobile platform measurements of air pollutants adjacent to
551 a major roadway. *Atmos. Environ.* **98**, 492–499 (2014).
- 552 23. US EPA. Air Quality System (AQS). *US Environmental Protection Agency*
553 <https://www.epa.gov/aqs> (2019).
- 554 24. US EPA. AirData Pre-Generated Data Files. *US Environmental Protection Agency*
555 https://aqs.epa.gov/aqsweb/airdata/download_files.html (2019).
- 556 25. R Core Team. R: A Language and Environment for Statistical Computing. *R Foundation for*
557 *Statistical Computing* <https://www.r-project.org> (2019).
- 558 26. Li, L. *et al.* Ensemble-based deep learning for estimating PM_{2.5} over California with
559 multisource big data including wildfire smoke. *Environ. Int.* **145**, 106143 (2020).
- 560 27. Szpiro, A. A. & Paciorek, C. J. Measurement error in two-stage analyses, with application to
561 air pollution epidemiology. *Environmetrics* (2013) doi:10.1002/env.2233.

- 562 28. Dodson, R. E., Houseman, E. A., Morin, B. & Levy, J. I. An analysis of continuous black
563 carbon concentrations in proximity to an airport and major roadways. *Atmos. Environ.* **43**,
564 3764–3773 (2009).
- 565 29. Riley, E. A. *et al.* Correlations between short-term mobile monitoring and long-term passive
566 sampler measurements of traffic-related air pollution. *Atmos. Environ.* **132**, (2016).
- 567 30. Austin, E. *et al.* *Mobile Observations of Ultrafine Particles: The MOV-UP study report.*
568 (2019).
- 569 31. Hudda, N., Gould, T., Hartin, K., Larson, T. V. & Fruin, S. A. Emissions from an
570 international airport increase particle number concentrations 4-fold at 10 km downwind.
571 *Environ. Sci. Technol.* **48**, 6628–6635 (2014).
- 572 32. Lack, D. A. & Corbett, J. J. Black carbon from ships: a review of the effects of ship speed,
573 fuel quality and exhaust gas scrubbing. *Atmospheric Chem. Phys.* **12**, (2012).
- 574 33. Kozawa, K. H., Fruin, S. A. & Winer, A. M. Near-road air pollution impacts of goods
575 movement in communities adjacent to the Ports of Los Angeles and Long Beach. *Atmos.*
576 *Environ.* **43**, 2960–2970 (2009).
- 577 34. Riffault, V. *et al.* Fine and Ultrafine Particles in the Vicinity of Industrial Activities: A
578 Review. *Crit. Rev. Environ. Sci. Technol.* **45**, 2305–2356 (2015).
- 579 35. Kilian, J. & Kitazawa, M. The emerging risk of exposure to air pollution on cognitive decline
580 and Alzheimer ' s disease e Evidence from epidemiological and animal studies. *Biomed. J.*
581 **41**, 141–162 (2018).
- 582 36. Lane, K. J. *et al.* Association of modeled long-term personal exposure to ultrafine particles
583 with inflammatory and coagulation biomarkers. *Environ. Int.* **92–93**, 173–182 (2016).

- 584 37. US EPA. Integrated science assessment (ISA) for particulate matter (final report, Dec 2019).
585 *US Environ. Prot. Agency* (2019).
- 586 38. Gryparis, A., Paciorek, C. J., Zeka, A., Schwartz, J. & Coull, B. A. Measurement error
587 caused by spatial misalignment in environmental epidemiology. *Biostatistics* **10**, 258–274
588 (2009).
- 589 39. Szpiro, A. A., Sheppard, L. & Lumley, T. Efficient measurement error correction with
590 spatially misaligned data. *Biostatistics* **12**, 610–623 (2011).
- 591 40. Sheppard, L. *et al.* Confounding and exposure measurement error in air pollution
592 epidemiology. *Air Qual. Atmosphere Health* **5**, 203–216 (2012).
593