

1 Independent validation of national satellite-based land-use regression models for nitrogen
2 dioxide using passive samplers

3

4 Luke D. Knibbs^{1*}, Craig P. Coorey², Matthew J. Bechle³, Christine T. Cowie^{4,5,6}, Mila
5 Dirgawati⁷, Jane S. Heyworth⁷, Guy B. Marks^{4,5}, Julian D. Marshall³, Lidia Morawska⁸,
6 Gavin Pereira⁹, Michael G. Hewson¹⁰

7

8 ¹ School of Public Health, The University of Queensland, Herston, QLD 4006, Australia

9 ² School of Medicine, The University of Queensland, Herston, QLD 4006, Australia

10 ³ Department of Civil and Environmental Engineering, University of Washington, Seattle,
11 WA 98195, USA

12 ⁴ South Western Sydney Clinical School, The University of New South Wales, Liverpool,
13 NSW 2170, Australia

14 ⁵ Ingham Institute of Medical Research, Liverpool, NSW 2170, Australia

15 ⁶ Woolcock Institute of Medical Research, University of Sydney, Glebe, NSW 2037, Australia

16 ⁷ School of Population Health, The University of Western Australia, Crawley, WA 6009,
17 Australia

18 ⁸ International Laboratory for Air Quality and Health, Queensland University of Technology,
19 Brisbane, QLD 4001, Australia

20 ⁹ School of Public Health, Curtin University, Perth, WA 6000, Australia

21 ¹⁰ School of Geography, Planning and Environmental Management, The University of
22 Queensland, St. Lucia, QLD 4067, Australia

23 * Corresponding author (e: l.knibbs@uq.edu.au; ph: +61 7 3365 5409; fax: +61 7 3365 5540)

24 **Abstract**

25 Including satellite observations of nitrogen dioxide (NO₂) in land-use regression (LUR)
26 models can improve their predictive ability, but requires rigorous evaluation. We used 123
27 passive NO₂ samplers sited to capture within-city and near-road variability in two Australian
28 cities (Sydney and Perth) to assess the validity of annual mean NO₂ estimates from existing
29 national satellite-based LUR models (developed with 68 regulatory monitors). The samplers
30 spanned roadside, urban near traffic (≤ 100 m to a major road), and urban background (> 100
31 m to a major road) locations. We evaluated model performance using R² (predicted NO₂
32 regressed on independent measurements of NO₂), mean-square-error R² (MSE-R²), RMSE,
33 and bias. Our models captured up to 69% of spatial variability in NO₂ at urban near-traffic
34 and urban background locations, and up to 58% of variability at all validation sites, including
35 roadside locations. The absolute agreement of measurements and predictions (measured by
36 MSE-R²) was similar to their correlation (measured by R²). Few previous studies have
37 performed independent evaluations of national satellite-based LUR models, and there is little
38 information on the performance of models developed with a small number of NO₂ monitors.
39 We have demonstrated that such models are a valid approach for estimating NO₂ exposures
40 in Australian cities.

41

42

43

44

45

46

47

48

49 **Introduction**

50 Land-use regression (LUR) is frequently used for estimating exposure to outdoor air pollution
51 in epidemiological studies. LUR models use features of the built and natural environment,
52 such as road length, impervious surfaces, and tree cover, to capture spatial variability in
53 pollutant concentrations measured at fixed locations. This allows concentrations at
54 unmeasured locations to be estimated.¹ Several recent studies have shown that the predictive
55 ability of LUR models for nitrogen dioxide (NO₂), quantified as R², increases by 2 to 15
56 percentage points when satellite-observed tropospheric NO₂ is included as a predictor
57 variable.²⁻⁷ These models aim to leverage the best attributes of satellite observations (e.g.,
58 large spatial coverage) and LUR models (e.g., local-scale predictors) to improve performance
59 and coverage compared with either technique alone.

60

61 The spatial coverage offered by satellite data makes it suitable for national or multi-national
62 applications, and satellite-based LUR models have been developed for the USA,^{2,7,8} Canada,⁶
63 Australia,⁵ Western Europe,³ and the Netherlands.⁴ A single national satellite model can offer
64 a simpler and consistent way to assign exposures to geographically dispersed study subjects
65 compared with separate non-satellite LUR models for each city, which are costly and time-
66 intensive to develop.⁹ Some national models can also offer comparable predictive ability and
67 spatial resolution to city-scale models.^{2,7}

68

69 LUR models can overfit, particularly when the number of measurement sites is small and the
70 number of potential predictors is large.¹⁰⁻¹² Validation is therefore important for assessing
71 how well they perform when applied beyond the data sets used to develop them (e.g., at the
72 home addresses of subjects in an epidemiological study).^{12,13} Numerous LUR validation
73 studies have focused on city-scale models (e.g.,^{11,14,15}). In contrast, there is little information

74 on validation of satellite-based national NO₂ models,^{2,3,7,8} especially in countries with limited
75 ground-based monitoring.⁶ Validation of these models is particularly important because they
76 are implemented at a nation-wide scale, which encompasses a wide range of land-use
77 conditions that may differ from the sites used to develop the models.

78

79 In this study, we sought to perform an independent validation of Australian national satellite-
80 based LUR models for NO₂. Through this, we wanted to determine if our models were
81 suitable for estimating residential NO₂ exposures in epidemiological studies. We also aimed
82 to add to the limited literature on satellite-based LUR evaluation by exploring the ability of
83 national models developed with a relatively small number of monitoring sites to predict NO₂
84 concentrations at sites selected to capture within-city and near-road variability.

85

86 **Experimental Materials and Methods**

87 *Models being evaluated*

88 We previously described our satellite-based LUR models for NO₂,⁵ which were developed
89 using data from 68 continuous regulatory chemiluminescence monitors throughout Australia
90 (population = 23.5 million; area = 7.7 million km²; ~0.3 NO₂ monitors/100,000 persons).
91 Two models using different satellite predictors were developed. One model included the
92 tropospheric column abundance of NO₂ molecules observed by the OMI spectrometer aboard
93 the Aura satellite as a predictor (molecules × 10¹⁵ per cm²; ‘column model’). The other model
94 included the estimated NO₂ concentration at ground-level (ppb; ‘surface model’), based on
95 applying a surface-to-column ratio from the Weather Research and Forecasting model
96 coupled with Chemistry (WRF-Chem). Using eight and nine land-use predictor variables, our
97 column and surface models respectively explained 81% (RMSE = 1.4 ppb) and 79% (RMSE

98 = 1.4 ppb) and of spatial variability in measured annual mean NO₂ in Australia during 2006-
99 11.

100

101 *Measurements used for validation*

102 In this study, we sought a data set independent of that used in our LUR models' development
103 to rigorously assess their performance. Because we had previously used most available
104 regulatory air monitoring data for development, we contacted all investigators who had
105 performed NO₂ monitoring as part of epidemiological studies between 2006 and 2014. Our
106 initial inclusion criteria were that: (a) NO₂ had been measured anywhere in Australia
107 provided that repeated, precise coordinates were collected (i.e, to 5 decimal places); (b)
108 measurements ran for at least two weeks, and; (c) a validated measurement method with
109 documented quality assurance procedures was used. We received data from five studies,
110 which, to our knowledge, represented all NO₂ monitoring that met the inclusion criteria.
111 Together, these studies included 174 measurement sites across three of Australia's six states.

112

113 After preliminary screening we imposed additional, more stringent, inclusion criteria for the
114 studies. Namely, we required three repeated measurements of 14 days' duration each that
115 spanned different seasons. We aimed to ensure that measurements from different studies
116 captured seasonal variation in NO₂, were of comparable duration, and able to be converted to
117 an estimated annual mean using standard methods. These criteria were informed by the well-
118 described European Study of Cohorts for Air Pollution Health Effects (ESCAPE) protocol for
119 LUR development.¹⁶ Based on this, we excluded two studies comprising 43 measurement
120 sites.

121

122 The remaining 131 sites were located in Sydney (87 sites; population = 4.9 million) and Perth
123 (44 sites; population = 2 million), the most and fourth-most populous cities in Australia,
124 respectively (Figure 1). All of the sites were located within the metropolitan area of those two
125 cities, and were selected to capture within-city and near-road variability in NO₂. All NO₂
126 measurements were performed using passive sampling techniques (Ferm-type sampler and
127 Ogawa sampler). Information on sampling dates, measurement methods, and quality
128 assurance is in Table 1.

129

130 *Conversion to annual mean NO₂*

131 Because each site was measured over two week periods in different seasons but our models'
132 predictions were for annual mean NO₂, we adjusted the measurements to an estimated annual
133 mean. We did this using the ratio of mean NO₂ measured by regulatory monitors during each
134 measurement period compared with its annual mean.^{17,18} We calculated the ratio based on
135 three separate regulatory monitors in each study area. We took that approach to improve the
136 precision of the adjusted annual mean estimate (i.e., the overall mean of adjusted
137 concentrations for each measurement period), as measured by its standard error.¹⁷ The
138 selection criteria for the regulatory monitors and the adjustment process are described in the
139 Supporting Information (pages S3-S12).

140

141 *Site classification*

142 We classified each site as either: (1) roadside (≤ 15 m to the centre of a major road), or; (2)
143 urban near traffic (not roadside, but ≤ 100 m to the centre of a major road), or; (3) urban
144 background (not roadside or urban near traffic; > 100 m to the centre of a major road). The 15
145 m distance threshold was selected to capture sites immediately influenced by vehicle
146 emissions, while the 100 m threshold was selected because it represents the approximate half-

147 life in the decay of NO₂ away from a road.¹⁹⁻²¹ Borderline sites on either side of a distance
148 threshold were manually investigated using Google Earth and Street View before assigning
149 them to a category. We assessed the sensitivity of our analyses to a halving and a doubling of
150 the distance thresholds used for classifying roadside sites (7.5 m, 30 m) and urban near traffic
151 sites (50 m, 200 m). Major roads were defined using transport hierarchy codes supplied by
152 the Public Sector Mapping Agencies.^{5,22} We also assessed the effect of changing the
153 definition of a major road on our analyses (Supporting Information, pages S22-S26).

154

155 There was only one industrial point source of NO_x within 250 m of a site, based on the
156 Australian National Pollutant Inventory.²³ The site was located 120 m from a hospital that
157 emitted a moderate amount of NO_x per year (~5000 kg), but the main source of NO₂ was
158 more likely to be traffic emissions because it was also a roadside site.

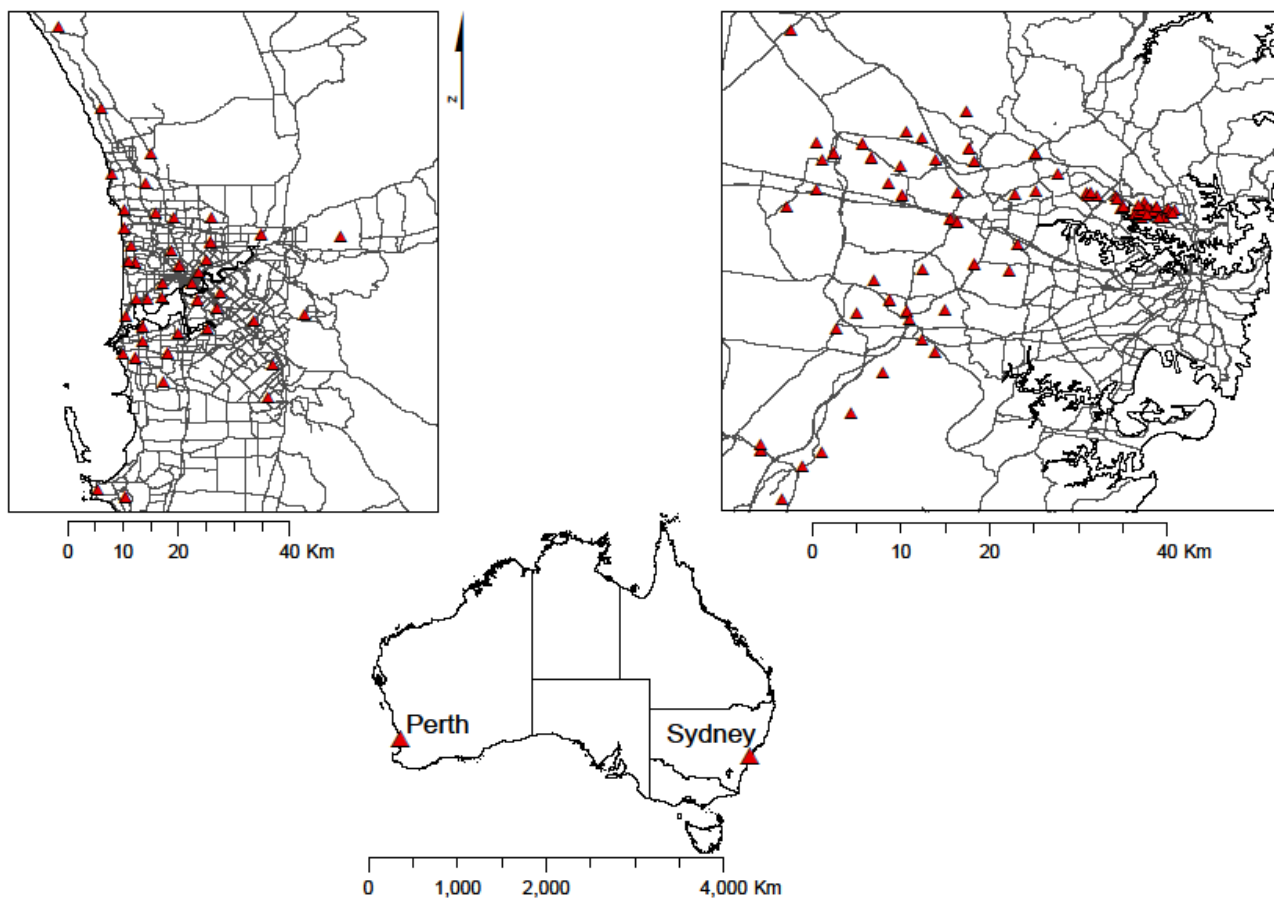
159

160

161

162

163



164

165 Figure 1. The two Australian cities (Sydney and Perth) where validation measurements were performed. The left panel shows Perth and the right
166 shows Sydney. The map shows the 123 sites used in the main analysis, denoted as red triangles. Major roads are also shown. See Figure S5 in
167 the supporting information for maps of predicted NO₂ in the study areas. The outlines were created using census data published by the
168 Australian Bureau of Statistics and roads were generated from data supplied by the Australian Public Sector Mapping Agencies.^{22,37}

169

170 Table 1. Details of each of the three sampling campaigns used for validation.

	Perth 1	Sydney 1	Sydney 2
Year	2012	2006-2008	2013-2014
<i>n</i> sites	44	40	47
Site selection	Following ESCAPE protocol ¹⁶	Selected to represent the expected variability of NO ₂ in the study area	Following ESCAPE protocol ¹⁶
Sample height	1.5-2 m above ground	2.2 m above ground	2.3-2.4 m above ground
Duration per sample	14 days per sample	14 days per sample	14 days per sample
Timing of samples	1 sample in each of summer, autumn and winter	1 sample in each of summer, winter and spring in each year during 2006-8 (subset of 11-13 sites also sampled in autumn)	1 sample in each of summer (2013) autumn (2014), and winter (2014)
Measurement method	Ogawa sampler ²⁴	Ferm-type sampler ²⁵	Ogawa sampler ²⁴

Analysis	Spectrophotometry based on Saltzman method	Spectrophotometry based on Saltzman method	Spectrophotometry based on Saltzman method
Quality assurance	co-located with chemiluminescence monitors, field blanks, duplicates for each sample	co-located with chemiluminescence monitors, field blanks, duplicates for one in five samples	co-located with chemiluminescence monitors, field blanks, duplicates for one in five samples
Limit of detection	2.0 ppb	0.5 ppb	2.0 ppb
Reference	Dirgawati et al. ²⁶	Rose et al. ²⁷	Not yet published

171

172

173

174 *Model predictions*

175 We used our satellite-based LUR models to predict annual mean NO₂ concentrations at each
176 site. Surface and column model predictions were determined for the year in which the
177 validation measurements were performed. Where measurements were done across more than
178 one year, we averaged the predicted NO₂ concentrations to match the measurement period.
179 Measurements from two campaigns (2012 and 2013-14) were performed outside the 2006-11
180 period used to develop our models. We obtained updated satellite column and surface
181 estimates of NO₂ for those years using our previous methods,⁵ and applied them using our
182 existing models. We used all other LUR predictor variables unmodified, based on the
183 assumption that they were unlikely to change substantially over 1-3 years.

184

185 We excluded validation sites that had values of one or more LUR predictor variables that
186 were outside the range observed at the 68 regulatory monitoring sites used for model
187 development. We did this to prevent unrealistic predictions, based on the approach of Wang
188 et al.^{9,12} Eight sites were excluded, leaving a total of 123 available for validation. We
189 assessed the effect of excluding those sites on our results by comparing them to results with
190 the sites included.

191

192 *Validation*

193 We used standard methods to validate our LUR models,^{12,28} and summarized their
194 performance using an independent validation R² (predicted NO₂ regressed on independent
195 measurements of NO₂), the regression slope and 95% confidence intervals, RMSE (absolute
196 and percentage scale), and bias (absolute and fractional). The R² we calculated is analogous
197 to a hold-out validation R² (HV-R²),¹¹⁻¹³ except our validation data were a set of unrelated,
198 independent measurements, rather than a subset of model development sites held out for

199 validation. As such, we refer to our validation metric as R^2 rather than HV- R^2 . We performed
200 standard diagnostics on the normality of residuals and their variance. We assessed the spatial
201 correlation of residuals using Moran's I .

202

203 Because R^2 is based on the correlation between validation measurements and model
204 predictions, it does not reflect their absolute agreement. Therefore, we also calculated a
205 mean-square-error R^2 (MSE- R^2) that took absolute values into account.^{10,12,28} MSE-
206 R^2 indicates how well the relationship between measurements and predictions follows a 1:1
207 line; its derivation is described extensively elsewhere.^{10,12,28,29} Using both R^2 and MSE- R^2
208 can identify LUR model predictions that are well-correlated with measurements but have
209 poor absolute agreement.¹⁰ Unlike R^2 , MSE- R^2 can have negative values if the average of
210 measurements leads to a lower MSE than the predictions.^{10,12,28,29}

211

212 We evaluated LUR model predictions for the entire validation set, by site classification, and
213 by each of the three validation measurement campaigns. We used R version 3.2.2 for all
214 analyses (R Project for Statistical Computing, Vienna, Austria).

215

216 **Results**

217 *NO₂ concentrations*

218 There were 8,177 days of NO₂ measurements performed in total across the 123 validation
219 sites during 2006-2014. Measured NO₂ concentrations adjusted to annual means are
220 summarized in Table 2. Higher concentrations were observed at roadside sites, followed by
221 urban near traffic sites, then urban background sites, and concentrations were higher at sites
222 in Sydney than those in Perth (Table 2). The concentrations we used for validation were
223 slightly higher than those used to develop the LUR models (Table S4). The effects of

224 changing the definitions used to classify sites on concentration percentiles were minor (Table
225 S5).

226

227 *Site classification*

228 There were 25 roadside sites, 18 urban near traffic sites, and 80 urban background sites using
229 the standard classification criteria. There was a greater proportion of roadside sites and a
230 smaller proportion of urban background sites used for validation compared with LUR model
231 development, particularly in Perth (Tables S6, S7). However, the percentiles of LUR
232 predictors at validation sites were comparable to the model development sites overall (Table
233 S8). Changing the definitions used to classify sites led to moderate changes in the number of
234 sites in each category (Table S9).

235

236 *Model validation*

237 Table 3 presents key validation statistics. The surface and column models captured 58%
238 ($\text{MSE-R}^2 = 51\%$) and 55% ($\text{MSE-R}^2 = 52\%$), respectively, of spatial variability in annual
239 mean NO_2 at the 123 validation sites overall (Figures 2a and 2b). The figures show some
240 evidence of increasing variance of errors with increasing NO_2 concentrations, but plots of
241 predicted NO_2 against residuals did not indicate overt violation of homoscedasticity (Figures
242 S1-S2).

243

244 Table 2. Percentiles of annual NO₂ concentrations (ppb) measured at validation sites. * Any negative concentrations following subtraction of
 245 field blank values were randomly assigned a value between zero and the limit of detection (2.0 ppb) in the Perth study (see Dirgawati et al. ²⁶).

Location	Min.	5th	25th	50th	75th	95th	Max.	Mean	S.D.
overall (<i>n</i> = 123)	0.4	2.9	5.9	8.5	11.2	14.6	19.3	8.6	3.7
roadside (<i>n</i> = 25)	5.1	5.6	8.0	11.0	13.1	17.6	19.3	11.0	3.9
urban near traffic (<i>n</i> = 18)	4.8	4.9	5.8	9.5	11.5	14.8	16.5	9.3	3.6
urban background (<i>n</i> = 80)	0.4	2.8	5.2	8.2	10.0	12.4	15.3	7.8	3.3
Sydney (<i>n</i> = 80)	3.9	5.8	8.2	9.9	11.9	16.5	19.3	10.2	3.1
Perth (<i>n</i> = 43)*	0.4	1.3	4.3	5.1	7.1	11.0	11.5	5.7	2.8

246

247

248

249

250

251 Table 3. Validation statistics for the surface and column models. RMSE = root-mean-square error; FB = fractional bias. Other abbreviations are
 252 defined in the main text.

Surface model	R²	β (95% CI)	MSE-R²	RMSE (ppb)	RMSE (%)	Bias (ppb)	FB (-)
overall (<i>n</i> = 123)	0.58	0.69 (0.61, 0.78)	0.51	2.6	29.6	-0.8	-0.10
roadside (<i>n</i> = 25)	0.36	0.55 (0.29, 0.81)	-0.18	4.1	37.5	-2.5	-0.26
urban near traffic (<i>n</i> = 18)	0.71	0.97 (0.70, 1.24)	0.60	2.2	23.9	-0.2	-0.03
urban background (<i>n</i> = 80)	0.68	0.74 (0.65, 0.84)	0.66	1.9	24.6	-0.5	-0.06
urban near traffic + urban background (<i>n</i> = 98)	0.69	0.80 (0.71, 0.89)	0.66	2.0	24.5	-0.4	-0.06
Column model							
overall (<i>n</i> = 123)	0.55	0.64 (0.55, 0.72)	0.52	2.5	29.5	-0.6	-0.07
roadside (<i>n</i> = 25)	0.29	0.47 (0.21, 0.74)	-0.13	4.0	36.7	-2.1	-0.21
urban near traffic (<i>n</i> = 18)	0.70	0.91 (0.65, 1.17)	0.64	2.1	22.8	0.1	0.01
urban background (<i>n</i> = 80)	0.64	0.67 (0.57, 0.76)	0.64	2.0	25.3	-0.2	-0.03
urban near traffic + urban background (<i>n</i> = 98)	0.66	0.73 (0.65, 0.82)	0.65	2.0	24.8	-0.2	-0.02

253 The surface model captured 71% ($\text{MSE-R}^2 = 60\%$) and 68% ($\text{MSE-R}^2 = 66\%$) of spatial
254 variability at urban near traffic and urban background sites, respectively. The column model
255 captured 70% ($\text{MSE-R}^2 = 64\%$) and 64% ($\text{MSE-R}^2 = 64\%$), respectively. When we combined
256 urban near traffic and urban background sites but excluded the 25 roadside sites, the surface
257 and column models captured 69% ($\text{MSE-R}^2 = 66\%$) and 66% ($\text{MSE-R}^2 = 65\%$) of spatial
258 variability at the remaining 98 sites, respectively (Figures 3a and 3b). The RMSE and bias of
259 both models was reduced compared with the analysis that included roadside sites. The surface
260 and column models captured 36% ($\text{MSE-R}^2 = -18\%$) and 29% ($\text{MSE-R}^2 = -13\%$),
261 respectively, of spatial variability at roadside sites.

262

263 *Prediction bias and RMSE*

264 Both models modestly but consistently under-predicted annual mean NO_2 , and the column
265 model predicted NO_2 with slightly less bias than the surface model (Table 3). The absolute
266 bias of both models was less than -0.5 ppb for most analyses. Fractional bias was mostly less
267 than -0.10. The absolute RMSE was very similar across both models; approximately 2 ppb
268 (~25% in relative terms). Residuals had an approximately normal distribution and constant
269 variance across all analyses (Figures S1-S4). There was no evidence of spatial correlation
270 among residuals (Table S10).

271

272 *Sensitivity of results*

273 Moving the distance thresholds used to classify roadside and urban near traffic sites led to
274 similar results to the main analysis (Table S9). Likewise, changing the classification of major
275 roads did not substantially alter the results (Table S9). The results of validation stratified by
276 each of the three measurement campaigns are presented in the Supporting Information (Table
277 S11). The predictive ability of both models was lower than that observed when the data were

278 pooled across all sampling campaigns. Including the eight sites that had predictors outside the
279 range used to develop the models resulted in comparable R^2 values, but lower MSE- R^2 values
280 (Table S12). That finding supported the decision to exclude the sites.

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

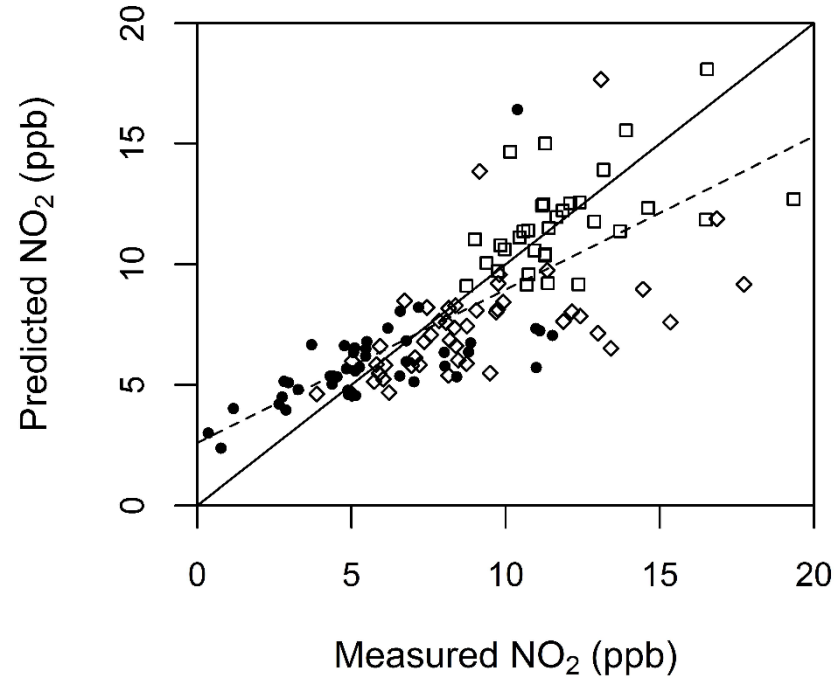
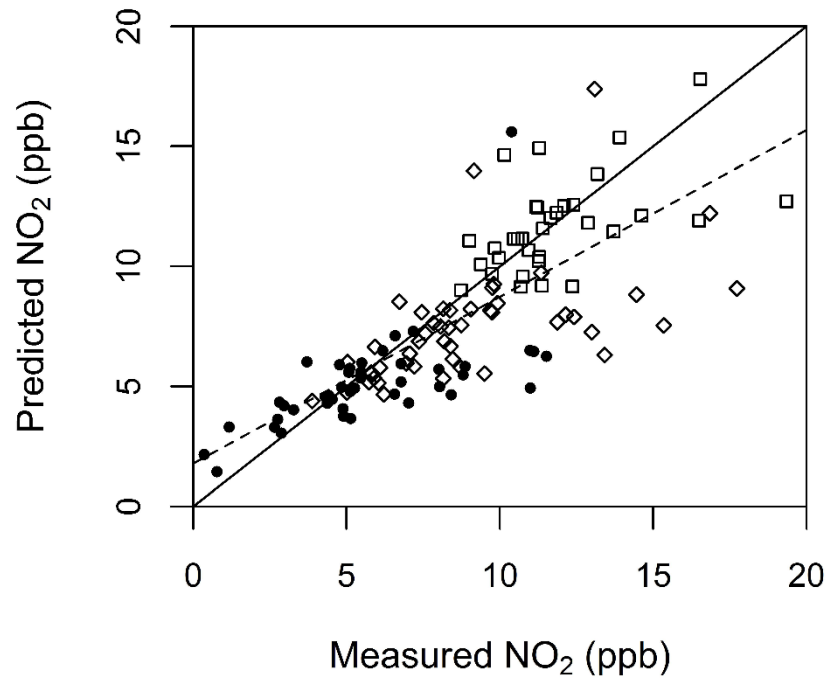
297

298

299

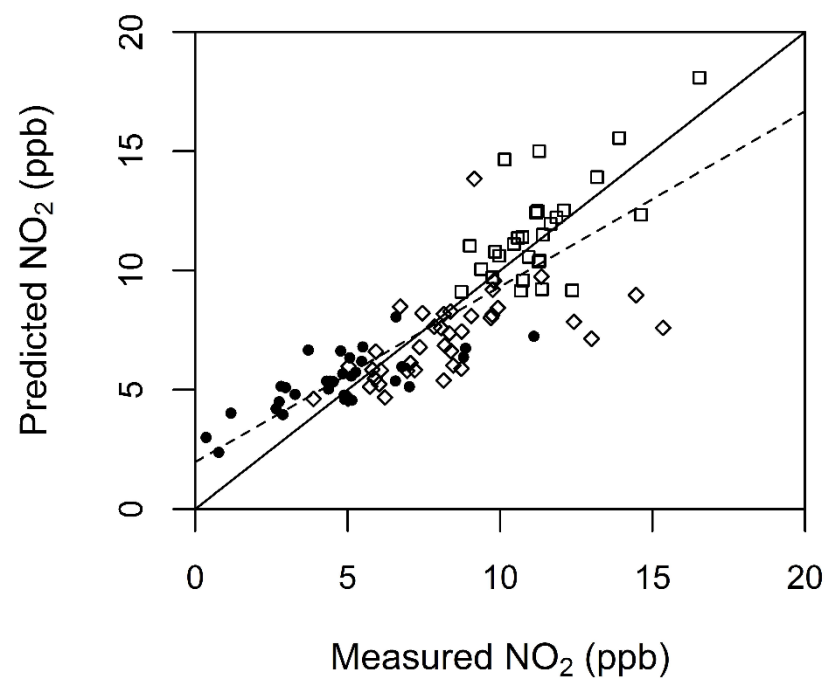
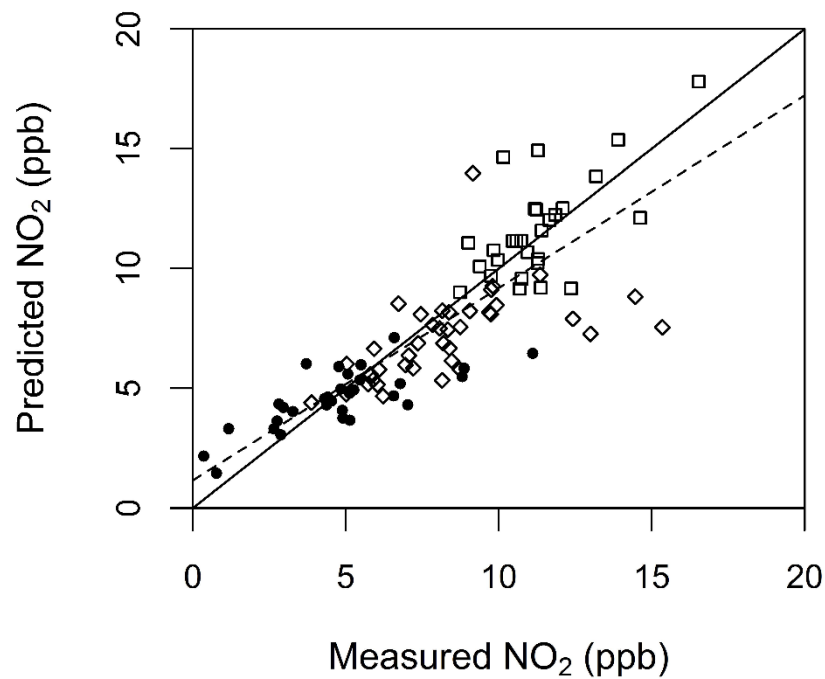
300

301



302

303 Figure 2a (left) and 2b (right). Measured *vs.* predicted annual mean NO₂ at 123 validation sites (roadside, urban near traffic, and urban
 304 background combined) for the surface (2a) and column (2b) models. The dashed line is the line of best fit (see Table 3 for fit statistics). The solid
 305 line is the line of agreement. Symbols denote different measurement campaigns: solid circles = Perth 1; hollow squares = Sydney 1; hollow
 306 diamonds = Sydney 2.



307

308 Figure 3a (left) and 3b (right). Measured vs. predicted annual mean NO₂ at 98 urban near traffic and urban background validation sites combined
 309 for the surface (3a) and column (3b) models. The dashed line is the line of best fit (see Table 3 for fit statistics). The solid line is the line of
 310 agreement. Symbols denote different measurement campaigns: solid circles = Perth 1; hollow squares = Sydney 1; hollow diamonds = Sydney 2.

311 **Discussion**

312 *Key results and comparison to other studies*

313 Validation of LUR models with data not used in their development is the optimum method
314 for quantifying how well they perform.¹² In this study, we used a large independent set of
315 NO₂ measurements in two Australian cities ($n = 123$ sites) that was not available at the time
316 of model development to assess the ability of our national satellite-based LUR models ($n =$
317 68 sites) to capture within and near-road variability. We previously used five-fold cross-
318 validation with five replications to validate our models due to the scarcity of long-term
319 regulatory NO₂ data in Australia.⁵ The model R² was 79% (RMSE = 19%) and 81% (RMSE
320 = 19%), respectively, for the surface and column models. Here, we found that our surface and
321 column models explained 69% (RMSE = 25%) and 66% (RMSE = 25%), respectively, of
322 spatial variation in measured annual mean NO₂ at urban near traffic and urban background
323 validation sites combined ($n = 98$).

324

325 Excluding roadside sites, which are discussed in a separate section below, we observed a
326 decrease in R² from model development to independent validation of between 10 and 15
327 percentage points. Bechle et al.⁷ assessed their satellite-based LUR for NO₂ in the USA by
328 varying the proportion of sites held out from 10 to 95%. With approximately 70 sites for
329 development and 300 sites for validation, both the model build R² (median ~80%) and
330 decrease in R² when validated (approximately 10 percentage points) were consistent with
331 what we observed here and in our previous study.⁵ Our results also agree with those reported
332 by Wang et al.¹² for a Dutch national, non-satellite LUR for NO₂ developed with 70 sites.

333

334 The R² decrease we found was less than that described by Hystad et al.⁶ for their Canadian
335 national satellite-based LUR for NO₂. They found an average decrease from model

336 development to independent validation at 618 sites of 34 percentage points (73% vs. 39%).
337 Because of the diverse siting of validation sites in their study, its results are more comparable
338 with our overall validation results at 123 sites (i.e., including roadside sites). In that analysis,
339 we observed a decrease in R^2 of 21 and 26 percentage points for the surface and column
340 models, respectively. The smaller reduction in R^2 in this study might reflect the reduced
341 number of sites we used for validation, or the standard criteria we used for repeat
342 measurements and annual adjustment at validation sites to capture seasonal variation in NO_2 ,
343 which Hystad et al.⁶ did for some, but not all, of their sites. It might also reflect that their
344 model had fewer variables (4 predictors vs. 8 and 9 predictors in our models) and was not
345 geared towards detecting emissions attributable to heavy industry and biomass combustion,
346 which the authors noted may have affected their results.

347

348 *Relevance of LUR validation to epidemiological studies*

349 LUR models that have higher out-of-sample R^2 (i.e., between 3 and 16 percentage points
350 lower than model R^2) introduce substantially less attenuation in health effect estimates (from
351 1% to 14%).³⁰ The attenuation due to models with lower out-of-sample R^2 (i.e., between 16
352 and 74 percentage points lower than model R^2) ranges from 9 to 57%, depending on the
353 number of predictors and sites used to develop the model.³⁰ In the present study, we observed
354 a relatively modest decrease in R^2 from model build through to validation at urban near
355 traffic and urban background sites (10 to 15 percentage points), which was consistent
356 with that in other comparable studies, as outlined above.

357

358 Recent work has shown that LUR models with higher independent validation R^2 values
359 produce larger effect estimates than those with lower R^2 values when applied to the
360 association between NO_2 and forced viral capacity (FVC) in children.¹³ Model performance

361 evaluated using leave-one-out-cross-validation (LOOCV) had a much weaker correlation
362 with effect estimates, which underscores the importance of independent validation to
363 determine the utility of LUR models in health studies.¹³ Our results demonstrate that the
364 national satellite-based LUR models can be used to estimate with reasonable accuracy the
365 annual mean NO₂ exposures of people living in the metropolitan parts of Australia.

366

367 The absolute agreement between pollutant measurements and LUR model predictions is
368 important when models are used to assign exposures in epidemiological studies.¹⁰ Because
369 we aimed to determine if our models were fit for this purpose, we assessed absolute
370 agreement using MSE-R². We observed between one and three percentage points difference
371 in R² and MSE-R² values for urban near traffic and urban background sites combined, and
372 between three and seven percentage points for all sites combined. The differences we found
373 was mostly comparable to those reported by Wang et al.¹² and Basagãna et al.¹⁰ in their
374 European studies. The consistency we observed between R² and MSE-R² demonstrates that in
375 addition to being correlated, predicted and measured NO₂ also showed similar absolute
376 agreement.

377

378 Improving the accuracy of LUR model predictions does not always improve health effect
379 estimates.^{29,31} This has been demonstrated when the variability in an LUR predictor is smaller
380 at the measurement sites used to develop the model than the locations to which it will be
381 applied. In turn, this leads to an increase in classical-like measurement error associated with
382 estimating the predictor, which increases bias in the effect estimate compared with a model
383 that has a lower R² but less classical error.²⁹ Such findings illustrate that careful attention
384 needs to be paid to the characteristics of the sites used to develop LUR models versus those
385 they are applied to. In this study, we demonstrated that the percentiles of predictors at

386 validation sites were well-matched to the model development sites (Table S8), and both sets
387 of sites were generally consistent with the ~350,000 census block centroids across Australia
388 (Table S8, ⁵). This suggests that our models can be applied to a range of geographic settings
389 within Australia.

390

391 *Surface vs. column model performance*

392 Our surface and column models had similar R^2 , MSE- R^2 and RMSE values (Table 3), which
393 agrees with our original model development results.⁵ The column model had slightly lower
394 absolute and fractional bias compared with the surface model. We previously reported that
395 column models are a more straightforward and less time-consuming approach, which do not
396 require the simulation of surface-to-column ratios that the surface model does.⁵ Since then,
397 Bechle et al.⁷ also found that models using tropospheric NO₂ columns performed slightly
398 better than those using surface estimates in a national LUR for the USA. The validation we
399 have described here confirms that column-based NO₂ LUR models for Australia offer a
400 simpler alternative to surface-based models.

401

402 *Performance at roadside sites*

403 The predictive ability of our models at roadside sites ($n = 25$) was markedly reduced and
404 prediction error increased compared with urban near traffic and urban background sites. The
405 R^2 at roadside sites was 36% (RMSE = 4.1 ppb [38%]) and 29% (RMSE = 4.0 ppb [37%])
406 for the surface and column models, respectively, indicating some correlation between
407 roadside measurements and predictions. The MSE- R^2 values were negative in both cases,
408 indicating poor absolute agreement and that the mean of measurements performed better than
409 model predictions in terms of MSE. Both models under-predicted at roadside locations, with
410 bias of -2.5 ppb and -2.1 ppb for the surface and column models, respectively.

411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435

Our satellite-based LUR models were developed using ambient regulatory monitors, which are deliberately sited away from hotspots like roads. Although the roadside sites used for validation had predictors within the range observed at ambient sites, there was a higher proportion of roadside sites in the validation compared with development data; 20% versus 3%, respectively (Table S6). This is a likely explanation for the lower predictive performance at roadside sites. Also, our models were developed for all of Australia and did not include traffic density data because they are not available nationally. We instead used road length data, and the lower predictive ability at roadside sites is probably partially due to the difficulty associated with capturing the variability in NO₂ associated with complex, highly trafficked locations.³²

We previously geocoded the residential addresses of 15,000 Australian women randomly selected from Australia's universal healthcare database. We found that the median distance to a major road was 296 m in that cohort, where 84% of women lived in the major cities and inner regional areas of Australia.³³ Moreover, 5.7% of women lived within our definition of a roadside location (≤ 15 m from a major road), while 8.5% of women lived ≤ 30 m from a major road. Here, we were mainly interested in the ability of our models to predict at a typical residential address in an epidemiological study, most of which are unlikely to be located immediately proximate to a major road. Our models' performance at roadside locations is therefore less influential on decisions about implementing them in health studies.^{12,32}

436 *Limitations*

437 Our study has some important limitations. The validation data we used came from two
438 Australian cities, Sydney and Perth, while the models we sought to validate had national
439 coverage. Sydney and Perth combined (6.9 million people) account for 29% of the Australian
440 population, but it is possible that our validation sites may be less representative of other
441 areas. However, the values of LUR model predictors at our validation sites were largely
442 consistent with those at ~350,000 Australian census block centroids across the country (Table
443 S8), suggesting that the sites are appropriate for validating a national model. Our sites were
444 all located in the metropolitan part of the two cities, which means that validation was not
445 possible in rural and remote parts of Australia. Over 70% of Australians live in major cities,
446 and more than 85% of the population live in urban areas, making Australia one of the world's
447 most urbanized countries.³⁴ We therefore focused our models' validation on the locations
448 where they will be applied most frequently.

449

450 Although our LUR models were developed using continuous regulatory chemiluminescence
451 monitors we validated them using data from Ferm-type and Ogawa passive samplers.
452 However, these methods have consistently been shown to correlate and agree well for the two
453 week measurement periods we used.^{25,35,36}

454

455 Our main analysis only included validation sites that had predictors within the range used to
456 develop our satellite-based LUR models. We did this to prevent unreasonably high or low
457 predictions.^{9,12} This means that the predictive performance we observed holds for situations
458 where the predictors are within the models' development range.^{10,12} Options for assigning
459 exposures to out-of-range sites in epidemiological studies have been discussed by Wang et
460 al.¹²

461 In summary, we capitalized on the availability of a large number of NO₂ measurements
462 performed in Australia using standard passive sampling methods, which were not available at
463 the time we built our LUR models. We used almost double the number of sites to validate our
464 models ($n = 123$) as we used to develop them ($n = 68$). Our results add to the scant literature
465 on independent validation of national satellite-based LUR models for NO₂, particularly those
466 developed using a relatively small ground-based monitoring network. Our models captured
467 up to 69% of spatial variability in annual mean NO₂ at independent urban near traffic and
468 urban background validation sites, and up to 58% at all validation sites (including roadside
469 sites). Our findings indicate that satellite-based LUR models provide a valid, consistent, and
470 cost-effective method for assigning NO₂ exposures, even when the number of sites available
471 to develop them is limited. Based on the results, we will use the models to estimate
472 residential NO₂ concentrations in a national study of children's respiratory health.

473

474 **Supporting Information**

475 Tables S1-S12: adjustment to annual mean NO₂; NO₂ concentration percentiles; LUR model
476 development and validation site information; percentiles of predictors at development and
477 validation sites and census block centroids; site classification effects; spatial correlation
478 results; validation results by sampling campaign; effects of excluding sites.

479

480 Figures S1-S5: predicted NO₂ vs. residuals; Q-Q plots of residuals; predicted NO₂ in 2008
481 for Sydney and Perth.

482

483 **Acknowledgements**

484 L.D.K. acknowledges an NHMRC Early Career (Australian Public Health) Fellowship
485 (APP1036620). G.P. acknowledges a Sidney Sax Fellowship (APP1052236) and project

486 grants (APP1099655 and APP1047263) from the NHMRC. J.H. acknowledges an NHMRC
487 project grant (APP1003589). C.T.C. acknowledges funding from the Clean Air Research
488 Programme through the Commonwealth Department of Environment, Water, Heritage and
489 the Arts for NO₂ sampling during 2006-2008. Please contact the corresponding author to
490 obtain LUR model predictions for research purposes.

491

492 **References**

- 493 1. Hoek, G.; Beelen, R.; de Hoogh, K.; Vienneau, D.; Gulliver, J.; Fischer, P.; Briggs, D.,
494 A review of land-use regression models to assess spatial variation of outdoor air
495 pollution. *Atmos. Environ.* **2008**, *42*, (33), 7561-7578.
- 496 2. Novotny, E. V.; Bechle, M. J.; Millet, D. B.; Marshall, J. D., National Satellite-Based
497 Land-Use Regression: NO₂ in the United States. *Environ. Sci. Technol.* **2011**, *45*, (10),
498 4407-4414.
- 499 3. Vienneau, D.; de Hoogh, K.; Bechle, M. J.; Beelen, R.; van Donkelaar, A.; Martin, R.
500 V.; Millet, D. B.; Hoek, G.; Marshall, J. D., Western European Land Use Regression
501 Incorporating Satellite- and Ground-Based Measurements of NO₂ and PM₁₀. *Environ.*
502 *Sci. Technol.* **2013**, *47*, (23), 13555-13564.
- 503 4. Hoek, G.; Eeftens, M.; Beelen, R.; Fischer, P.; Brunekreef, B.; Boersma, K. F.;
504 Veefkind, P., Satellite NO₂ data improve national land use regression models for
505 ambient NO₂ in a small densely populated country. *Atmos. Environ.* **2015**, *105*, 173-
506 180.
- 507 5. Knibbs, L. D.; Hewson, M. G.; Bechle, M. J.; Marshall, J. D.; Barnett, A. G., A
508 national satellite-based land-use regression model for air pollution exposure assessment
509 in Australia. *Environ. Res.* **2014**, *135*, 204-211.

- 510 6. Hystad, P.; Setton, E.; Cervantes, A.; Poplawski, K.; Deschenes, S.; Brauer, M.; van
511 Donkelaar, A.; Lamsal, L.; Martin, R.; Jerrett, M.; Demers, P., Creating National Air
512 Pollution Models for Population Exposure Assessment in Canada. *Environ. Health*
513 *Perspect.* **2011**, *119*, (8), 1123-1129.
- 514 7. Bechle, M. J.; Millet, D. B.; Marshall, J. D., National Spatiotemporal Exposure Surface
515 for NO₂: Monthly Scaling of a Satellite-Derived Land-Use Regression, 2000–2010.
516 *Environ. Sci. Technol.* **2015**, *49*, (20), 12297-12305.
- 517 8. Young, M. T.; Bechle, M. J.; Sampson, P. D.; Szpiro, A. A.; Marshall, J. D.; Sheppard,
518 L.; Kaufman, J. D., Satellite-Based NO₂ and Model Validation in a National Prediction
519 Model Based on Universal Kriging and Land-Use Regression. *Environ. Sci. Technol.*
520 **2016**, *50*, (7), 3686-3694.
- 521 9. Wang, M.; Beelen, R.; Bellander, T.; Birk, M.; Cesaroni, G.; Cirach, M.; Cyrus, J.; de
522 Hoogh, K.; Declercq, C.; Dimakopoulou, K.; Eeftens, M.; Eriksen, K. T.; Forastiere, F.;
523 Galassi, C.; Grivas, G.; Heinrich, J.; Hoffmann, B.; Ineichen, A.; Korek, M.; Lanki, T.;
524 Lindley, S.; Modig, L.; Molter, A.; Nafstad, P.; Nieuwenhuijsen, M. J.; Nystad, W.;
525 Olsson, D.; Raaschou-Nielsen, O.; Ragettli, M.; Ranzi, A.; Stempfelet, M.; Sugiri, D.;
526 Tsai, M. Y.; Udvardy, O.; Varro, M. J.; Vienneau, D.; Weinmayr, G.; Wolf, K.; Yli-
527 Tuomi, T.; Hoek, G.; Brunekreef, B., Performance of multi-city land use regression
528 models for nitrogen dioxide and fine particles. *Environ. Health Perspect.* **2014**, *122*,
529 (8), 843-9.
- 530 10. Basagaña, X.; Rivera, M.; Aguilera, I.; Agis, D.; Bouso, L.; Elosua, R.; Foraster, M.; de
531 Nazelle, A.; Nieuwenhuijsen, M.; Vila, J.; Künzli, N., Effect of the number of
532 measurement sites on land use regression models in estimating local air pollution.
533 *Atmos. Environ.* **2012**, *54*, 634-642.

- 534 11. Johnson, M.; Isakov, V.; Touma, J. S.; Mukerjee, S.; Özkaynak, H., Evaluation of land-
535 use regression models used to predict air quality concentrations in an urban area.
536 *Atmos. Environ.* **2010**, *44*, (30), 3660-3668.
- 537 12. Wang, M.; Beelen, R.; Eeftens, M.; Meliefste, K.; Hoek, G.; Brunekreef, B., Systematic
538 Evaluation of Land Use Regression Models for NO₂. *Environ. Sci. Technol.* **2012**, *46*,
539 (8), 4481-4489.
- 540 13. Wang, M.; Brunekreef, B.; Gehring, U.; Szpiro, A.; Hoek, G.; Beelen, R., A New
541 Technique for Evaluating Land-use Regression Models and Their Impact on Health
542 Effect Estimates. *Epidemiology* **2016**, *27*, (1), 51-6.
- 543 14. Beelen, R.; Voogt, M.; Duyzer, J.; Zandveld, P.; Hoek, G., Comparison of the
544 performances of land use regression modelling and dispersion modelling in estimating
545 small-scale variations in long-term air pollution concentrations in a Dutch urban area.
546 *Atmos. Environ.* **2010**, *44*, (36), 4614-4621.
- 547 15. de Nazelle, A.; Aguilera, I.; Nieuwenhuijsen, M.; Beelen, R.; Cirach, M.; Hoek, G.; de
548 Hoogh, K.; Sunyer, J.; Targa, J.; Brunekreef, B.; Künzli, N.; Basagaña, X., Comparison
549 of performance of land use regression models derived for Catalunya, Spain. *Atmos.*
550 *Environ.* **2013**, *77*, 598-606.
- 551 16. ESCAPE Study Manual, ENV.2007.1.2.2.2. European cohort on air pollution,
552 2008; http://www.escapeproject.eu/manuals/ESCAPE-Study-manual_x007E_final.pdf
- 553 17. Hoek, G.; Meliefste, K.; Cyrus, J.; Lewné, M.; Bellander, T.; Brauer, M.; Fischer, P.;
554 Gehring, U.; Heinrich, J.; van Vliet, P.; Brunekreef, B., Spatial variability of fine
555 particle concentrations in three European areas. *Atmos. Environ.* **2002**, *36*, (25), 4077-
556 4088.

- 557 18. Lewné, M.; Cyrus, J.; Meliefste, K.; Hoek, G.; Brauer, M.; Fischer, P.; Gehring, U.;
558 Heinrich, J.; Brunekreef, B.; Bellander, T., Spatial variation in nitrogen dioxide in three
559 European areas. *Sci. Total Environ.* **2004**, *332*, (1–3), 217-230.
- 560 19. Roorda-Knape, M. C.; Janssen, N. A. H.; de Hartog, J.; Van Vliet, P. H. N.; Harssema,
561 H.; Brunekreef, B., Traffic related air pollution in city districts near motorways. *Sci.*
562 *Total Environ.* **1999**, *235*, (1–3), 339-341.
- 563 20. Gilbert, N. L.; Woodhouse, S.; Stieb, D. M.; Brook, J. R., Ambient nitrogen dioxide
564 and distance from a major highway. *Sci. Total Environ.* **2003**, *312*, (1–3), 43-46.
- 565 21. Pleijel, H.; Pihl Karlsson, G.; Binsell Gerdin, E., On the logarithmic relationship
566 between NO₂ concentration and the distance from a highroad. *Sci. Total Environ.* **2004**,
567 *332*, (1–3), 261-264.
- 568 22. Public Sector Mapping Agencies, Transport and topography product description, v3.6
569 2013;[https://www.pasma.com.au/sites/default/files/transport_and_topography_product_d](https://www.pasma.com.au/sites/default/files/transport_and_topography_product_description.pdf)
570 [escription.pdf](https://www.pasma.com.au/sites/default/files/transport_and_topography_product_description.pdf)
- 571 23. National Pollutant Inventory (Australia), 2016; <http://www.npi.gov.au/>
- 572 24. Ogawa and Co.; NO, NO₂, NO_x and SO₂ Sampling Protocol Using The Ogawa
573 Sampler, 2006; [http://ogawausa.com/wp-content/uploads/2014/04/prono-](http://ogawausa.com/wp-content/uploads/2014/04/prono-noxno2so206.pdf)
574 [noxno2so206.pdf](http://ogawausa.com/wp-content/uploads/2014/04/prono-noxno2so206.pdf)
- 575 25. Ayers, G. P.; Keywood, M. D.; Gillett, R.; Manins, P. C.; Malfroy, H.; Bardsley, T.,
576 Validation of passive diffusion samplers for SO₂ and NO₂. *Atmos. Environ.* **1998**, *32*,
577 (20), 3587-3592.
- 578 26. Dirgawati, M.; Barnes, R.; Wheeler, A. J.; Arnold, A.-L.; McCaul, K. A.; Stuart, A. L.;
579 Blake, D.; Hinwood, A.; Yeap, B. B.; Heyworth, J. S., Development of Land Use

- 580 Regression models for predicting exposure to NO₂ and NO_x in Metropolitan Perth,
581 Western Australia. *Environ. Modell. Softw.* **2015**, *74*, 258-267.
- 582 27. Rose, N.; Cowie, C.; Gillett, R.; Marks, G. B., Validation of a Spatiotemporal Land Use
583 Regression Model Incorporating Fixed Site Monitors. *Environ. Sci. Technol.* **2011**, *45*,
584 (1), 294-299.
- 585 28. Gulliver, J.; de Hoogh, K.; Hansell, A.; Vienneau, D., Development and Back-
586 Extrapolation of NO₂ Land Use Regression Models for Historic Exposure Assessment
587 in Great Britain. *Environ. Sci. Technol.* **2013**, *47*, (14), 7804-7811.
- 588 29. Szpiro, A. A.; Paciorek, C. J.; Sheppard, L., Does more accurate exposure prediction
589 necessarily improve health effect estimates? *Epidemiology* **2011**, *22*, (5), 680-5.
- 590 30. Basagaña, X.; Aguilera, I.; Rivera, M.; Agis, D.; Foraster, M.; Marrugat, J.; Elosua, R.;
591 Kunzli, N., Measurement error in epidemiologic studies of air pollution based on land-
592 use regression models. *Am. J. Epidemiol.* **2013**, *178*, (8), 1342-6.
- 593 31. Beckerman, B. S.; Jerrett, M.; Serre, M.; Martin, R. V.; Lee, S.-J.; van Donkelaar, A.;
594 Ross, Z.; Su, J.; Burnett, R. T., A Hybrid Approach to Estimating National Scale
595 Spatiotemporal Variability of PM_{2.5} in the Contiguous United States. *Environ. Sci.*
596 *Technol.* **2013**, *47*, (13), 7233-7241.
- 597 32. Dijkema, M. B.; Gehring, U.; van Strien, R. T.; van der Zee, S. C.; Fischer, P.; Hoek,
598 G.; Brunekreef, B., A Comparison of Different Approaches to Estimate Small-Scale
599 Spatial Variation in Outdoor NO₂ Concentrations. *Environ. Health Perspect.* **2011**,
600 *119*, (5), 670-675.
- 601 33. Lazarevic, N.; Dobson, A. J.; Barnett, A. G.; Knibbs, L. D., Long-term ambient air
602 pollution exposure and self-reported morbidity in the Australian Longitudinal Study on
603 Women's Health: a cross-sectional study. *BMJ Open* **2015**, *5*, e008714.

- 604 34. Australian Bureau of Statistics, Australian Historical Population Statistics,
605 2014; <http://www.abs.gov.au/ausstats/abs@.nsf/mf/3105.0.65.001>
- 606 35. Henderson, S. B.; Beckerman, B.; Jerrett, M.; Brauer, M., Application of Land Use
607 Regression to Estimate Long-Term Concentrations of Traffic-Related Nitrogen Oxides
608 and Fine Particulate Matter. *Environ. Sci. Technol.* **2007**, *41*, (7), 2422-2428.
- 609 36. Eeftens, M.; Beelen, R.; Fischer, P.; Brunekreef, B.; Meliefste, K.; Hoek, G., Stability
610 of measured and modelled spatial contrasts in NO(2) over time. *Occup. Environ. Med.*
611 **2011**, *68*, (10), 765-70.
- 612 37. Australian Bureau of Statistics, 2011, Digital Mesh Block Boundaries (Australia),
613 viewed 05 May2016, <[http://www.abs.gov.au/AUSSTATS/abs@.nsf/DetailsPage](http://www.abs.gov.au/AUSSTATS/abs@.nsf/DetailsPage/1270.0.55.001July%202011)
614 [/1270.0.55.001July%202011](http://www.abs.gov.au/AUSSTATS/abs@.nsf/DetailsPage/1270.0.55.001July%202011)>