



# Personal exposure to particulate matter in peri-urban India: predictors and association with ambient concentration at residence

Margaux Sanchez<sup>1,2,3</sup> · Carles Milà<sup>1,2,3</sup> · V. Sreekanth<sup>4</sup> · Kalpana Balakrishnan<sup>5</sup> · Sankar Sambandam<sup>5</sup> · Mark Nieuwenhuijsen<sup>1,2,3</sup> · Sanjay Kinra<sup>6</sup> · Julian D. Marshall<sup>4</sup> · Cathryn Tonne<sup>1,2,3</sup>

Received: 17 November 2018 / Revised: 11 March 2019 / Accepted: 1 May 2019  
© Springer Nature America, Inc. 2019

## Abstract

Scalable exposure assessment approaches that capture personal exposure to particles for purposes of epidemiology are currently limited, but valuable, particularly in low-/middle-income countries where sources of personal exposure are often distinct from those of ambient concentrations. We measured 2 × 24-h integrated personal exposure to PM<sub>2.5</sub> and black carbon in two seasons in 402 participants living in peri-urban South India. Means (sd) of PM<sub>2.5</sub> personal exposure were 55.1 (82.8) µg/m<sup>3</sup> for men and 58.5(58.8) µg/m<sup>3</sup> for women; corresponding figures for black carbon were 4.6(7.0) µg/m<sup>3</sup> and 6.1 (9.6) µg/m<sup>3</sup>. Most variability in personal exposure was within participant (intra-class correlation ~20%). Personal exposure measurements were not correlated ( $R_{\text{spearman}} < 0.2$ ) with annual ambient concentration at residence modeled by land-use regression; no subgroup with moderate or good agreement could be identified (weighted kappa ≤ 0.3 in all subgroups). We developed models to predict personal exposure in men and women separately, based on time-invariant characteristics collected at baseline (individual, household, and general time-activity) using forward stepwise model building with mixed models. Models for women included cooking activities and household socio-economic position, while models for men included smoking and occupation. Models performed moderately in terms of between-participant variance explained (38–53%) and correlations between predictions and measurements ( $R_{\text{spearman}}$ : 0.30–0.50). More detailed, time-varying time-activity data did not substantially improve the performance of the models. Our results demonstrate the feasibility of predicting personal exposure in support of epidemiological studies investigating long-term particulate matter exposure in settings characterized by solid fuel use and high occupational exposure to particles.

**Keywords** Black carbon · Peri-urban · Personal exposure · Exposure modeling · PM<sub>2.5</sub> · India

**Supplementary information** The online version of this article (<https://doi.org/10.1038/s41370-019-0150-5>) contains supplementary material, which is available to authorized users.

✉ Cathryn Tonne  
cathryn.tonne@isglobal.org

<sup>1</sup> Barcelona Institute for Global Health (ISGlobal), Barcelona, Spain

<sup>2</sup> Universitat Pompeu Fabra (UPF), Barcelona, Spain

<sup>3</sup> CIBER Epidemiología y Salud Pública (CIBERESP), Barcelona, Spain

<sup>4</sup> Department of Civil and Environmental Engineering, University of Washington, Seattle, WA, United States

<sup>5</sup> Department of Environmental Health Engineering, Sri Ramachandra University (SRU), Chennai, India

<sup>6</sup> Department of Non-communicable Disease Epidemiology, London School of Hygiene and Tropical Medicine, London, UK

## Introduction

The epidemiological evidence linking particulate matter with diameter <2.5 µm (PM<sub>2.5</sub>) with premature mortality and morbidity is large [1, 2]. However, most of this evidence is based on populations in high-income countries, despite indications that the majority of the attributable burden of PM<sub>2.5</sub> comes from populations from low- and middle-income countries (LMICs) [3]. The relative lack of epidemiological evidence in LMICs, especially for long-term effects of air pollution, has been highlighted previously [4].

Epidemiological studies of long-term exposure typically rely on spatial contrasts to estimate between-individual exposure. Land-use regression approaches aim to model ambient levels at residential address and have been widely used in the epidemiological literature [5, 6]. They have generally shown good performance in predicting spatial

patterns of ambient air pollution, especially within urban areas dominated by traffic sources [6].

New technology allows researchers to directly measure personal exposure and integrate it with individual characteristics, time-activity-features, and residence characteristics, moving beyond traditional estimates of ambient air pollution at residence [7–10]. Although measurements may be more accurate than models, both have limitations. Measurements are often limited to small groups and short durations. Extrapolating measurements to larger populations for epidemiological research is also challenging. One source of complexity in measuring and modeling personal exposure is the relative contribution of within-individual (temporal) and between-individual (spatial) variability in exposure, which has been stressed before [11, 12]. Understanding these sources of variability is essential to advance approaches aiming at better approximate long-term personal exposure to air pollution.

Personal exposure prediction is especially attractive for air pollution epidemiology in many LMIC settings where traffic is not necessarily the dominant source. In settings with high prevalence of cooking with solid fuels or with high occupational exposures, personal exposure is likely to differ substantially from estimates of ambient air pollution at residence. There is a need for population-level exposure estimates to enable epidemiology that reflect the complexity of LMIC settings.

We measured and analyzed 24-h integrated personal exposure to  $PM_{2.5}$  and black carbon in a relatively large sample of the general population in peri-urban South India. We specifically aimed to: (1) compare measured personal exposure with annual ambient concentration estimated at residence using land-use regression models (previously developed for the study area), and (2) develop prediction models that could be used in epidemiological analyses to predict long-term personal exposure to  $PM_{2.5}$  and black carbon.

## Methods

### Study population

We used data collected through the Cardiovascular Health Effects of Air Pollution in Telangana, India (CHAI) project nested in the Andhra Pradesh Children and Parents Study (APCAPS) cohort [13, 14]. APCAPS is a large prospective, intergenerational cohort study including ~6000 participants living near the city of Hyderabad, India. The study area consists of 28 villages each with 187–5065 households spread over 543 km<sup>2</sup> southeast of Hyderabad. Villages vary in terms of area, population size, socio-economic status, level of urbanization, and primary cooking fuel. Ethics approval was granted by the Parc de Salut Mar, Public

Health Foundation of India, National Institute of Nutrition, Sri Ramachandra University, and the European Research Council. All participants provided informed consent.

In 2015, CHAI recruited a stratified (by sex and village) random sample of 402 adult participants of APCAPS. They were invited to participate in two non-consecutive 24-h monitoring sessions that included detailed measurements of self-reported time-activity patterns and particulate air pollution exposure. The first session occurred between May and July 2015 during summer season (including monsoon); the second session occurred between December 2015 and March 2016 during winter season.

Of the 402 selected participants, 81 completed one monitoring session and 278 completed two or three monitoring sessions, leading to 639 measurements of 24-h personal exposure to  $PM_{2.5}$  and black carbon. We excluded 13 measurements due to device malfunction (e.g., run time <70% of the expected 24 h or missing data) or poor compliance (the collocated accelerometer recorded no motion during monitoring) and 13 measurements due to missing covariates. We additionally excluded three measurements showing negative  $PM_{2.5}$  concentrations and 41 measurements with negative black carbon concentrations (potentially due to concentrations below the lower end of the standard curve used for correction). We therefore analyzed 610 participant-days of  $PM_{2.5}$ –569 participant-days of black carbon—24-h personal exposure, corresponding to 349 unique participants (207 men and 142 women).

### Personal exposure

Participants were asked to wear a secured backpack containing a personal exposure monitor to measure their 24-h integrated gravimetric exposure to  $PM_{2.5}$ . The inlet of the personal monitor was placed near the breathing zone on one strap of the backpack. The pump (model 224-PCMTX8, SKC Ltd, Dorset, UK) was placed inside the backpack and drew air through a sharp cut cyclone attached to a cassette containing a 37-mm filter (Emfab, 113 Pallflex<sup>®</sup>). Filters were weighed pre- and post-monitoring according to previously described protocol that follows the RTI (Research Triangle Institute) guidelines [15]. Daily  $PM_{2.5}$  concentrations were derived from filter mass after correction for mass accumulated on blank filters (session-specific correction using median blank weights based on 31 blank filters overall). Daily black carbon concentrations were derived from optical attenuation (880 nm) of the mass collected on sampled filters, using a Magee OT21 Sootscan Optical Transmissometer (Magee Scientific, Berkeley, California, USA). The factor value used for conversion was consistent with previous literature [16]. We detected a sensitivity of the OT21 output to the weight of the unexposed filter, so we corrected the attenuation factor value for filter weights.

Sessions began with a field worker setting up the monitoring equipment at the participant's house (average time: 8 a.m.) and finished the following day around the same time. Scheduling was designed to minimize disruption to participants' daily life. Participants were asked to wear the backpack during their usual activities for 24 consecutive hours. If the backpack interfered with activities (e.g., sleeping, sitting, bathing), participants were instructed to place the backpack nearby, on a stool or a chair.

### Ambient air pollution

**Background measurements:** Continuous monitoring of PM<sub>2.5</sub> was implemented from 2015 to 2016 at one site in the North of the study area [17]. Hourly PM<sub>2.5</sub> concentrations were measured using an e-BAM device (model 9800, Met One, Grants Pass, OR). Missing hourly data (22% of the monitoring sessions) were imputed using a linear regression of temperature, relative humidity, wind speed (measured at the Rajiv Gandhi International Airport, located 15.8 km from North site) and ambient PM<sub>2.5</sub> concentrations measured at the US embassy in Hyderabad, located 23.7 km from North site. The adjusted-R<sup>2</sup> of the model was 0.49 and the 10-fold cross-validation mean absolute error was 10.3 µg/m<sup>3</sup> (standard deviation of the hourly ambient time series being 21.3 µg/m<sup>3</sup>). We calculated daily average ambient PM<sub>2.5</sub>, temperature, and relative humidity to correspond with the monitoring sessions.

**Ambient concentration at residence:** We previously developed land-use regression models for the study area to estimate annual ambient concentration to PM<sub>2.5</sub> and black carbon at the residence of all participants [18]. Briefly, the PM<sub>2.5</sub> model included indicators of vegetation and urbanicity and explained 58% of the spatial variation; the black carbon model included indicators related to roads, natural spaces, and non-residential places and explained 78% of the spatial variation.

### Questionnaires

A baseline questionnaire was administered to participants by members of the field staff at most one month prior to personal exposure sampling. The questionnaire included data on general individual characteristics (e.g., age, occupation, smoking habits), usual activities (e.g., average time spent at work, cooking habits), and residence characteristics (e.g., primary stove type, kitchen type, and fuel use for cooking) of the participants. After each monitoring session, participants completed a post-monitoring questionnaire in which they were asked by member of the field staff about major sources of air pollution they had been exposed to during the session (e.g., solid fuel use for cooking, passive and active smoking, being in traffic or near open fires). The questionnaire included an hourly time-activity diary

(1-h slot with up to two main locations and activities). The questionnaire, developed at Sri Ramachandra University, has been validated in a previous study involving PM<sub>2.5</sub> measurements [19]. Questionnaires are available in the Supplementary Information.

### Statistical analysis

We performed all analyses in men and women separately as previous results showed strong difference in lifestyle and behaviors by sex in the study population [20, 21].

**Ambient-adjusted personal exposure:** For analysis, measurements of personal exposure were adjusted in order to account for the day-to-day variability driven by ambient factors (PM<sub>2.5</sub> concentration, temperature, and humidity) not related to individual characteristics. We used the background data measured by the e-BAM device located in the North of the study area. We regressed 24-h average of log-transformed ambient PM<sub>2.5</sub>, 24-h average temperature, and 24-h average relative humidity on log-transformed 24-h personal exposure (PM<sub>2.5</sub> and black carbon), with a random intercept per participant. The ambient adjustment decreased the within-participant variance component by 35 and 26%—relative to the empty model—for PM<sub>2.5</sub> and black carbon personal exposures, respectively. The relationship between the outcome and predictors was considered as linear. These predictors were uncorrelated with the variables used in the prediction models (see below) as pair-wise Pearson correlations were all <0.1. Log transformation was used to ensure normal distribution of residuals. The resulting marginal residuals were considered ambient-adjusted personal exposure and used throughout subsequent analyses.

**Variance component of personal exposure:** A linear mixed model with only a random intercept per participant (i.e., empty model) was used to partition personal exposure variability into within-participant (residual variance) and between-participant (random effect variance) variability components. We calculated intra-class correlation coefficients (ICC) i.e., the proportion of total variability attributable to between-participant variability. Subgroups analysis were performed.

**Personal exposure compared to ambient concentration at residence:** We calculated Spearman correlation coefficients between personal exposure and ambient concentration at residence modeled by land-use regression. We assessed the degree of agreement across rank quintiles of the different exposure indicators with weighted Kappa. Agreement more directly assesses whether modeled concentration at residence can be used as a proxy for measured personal exposure. Subgroups analysis were performed.

**Prediction models of personal exposure:** First, we developed a model including only time-invariant characteristics of the participant or his/her household (collected once at

baseline, complete list available in Supplementary Table S1). We identified predictors of PM<sub>2.5</sub> and black carbon following a data-driven, forward stepwise procedure. Starting from an empty mixed model with random intercept per participant, we tested each of the time-invariant candidate predictors as a fixed effect. We selected the predictor that yielded the greatest decrease in the marginal Akaike Information Criterion (AIC) and we repeated the procedure until no additional variable provided any further decrease in AIC. Within each iteration, we verified whether any of the included variables could be removed without increase in AIC and we checked multicollinearity using the Variance Inflation Factor (if >5, the predictor was excluded). After final iteration, we dropped predictors with *p*-values >0.1 according to a likelihood ratio test that compared the model with and without the predictor. We detected influential observations using adapted Cook's distance (if >4/sample size) for mixed model [22]. Predictors that were sensitive to the removal of influential observations (i.e., >20% change toward the null in the corresponding estimates) were excluded from the final model. We checked model residuals and random effects for normality and homoscedasticity. Second, we repeated the full process including time-varying characteristics as potential predictors. These characteristics related to specific activities or events that occurred during monitoring and were reported in the post-monitoring questionnaire (complete list in Supplementary Table S1). Finally, for the two models, we calculated the proportion of within- and between-participant variance of personal exposure explained by each final model as compared to the empty model (expressed in percent change) as previously used [12]. We calculated the root-mean-square error (RMSE) and Spearman correlation between averaged exposure and averaged prediction per participant, among those with ≥2 sessions.

**Evaluation of prediction models:** We evaluated the models by performing 10-fold cross-validation at the participant level. Briefly, we randomly partitioned participants in 10 similar-size subgroups—all measurements of each participant being in the same group. Each subgroup was used once as a validation dataset for the models previously developed, which were then fitted in the other nine groups. To evaluate the robustness of the selected predictors, we further performed a 10-fold cross-holdout validation [23]. Briefly, the predictor selection procedure was repeated using nine of the 10 previously partitioned subgroups and then used to predict personal exposure on the remaining subgroup, leading to 10 different model applications. We calculated the RMSE and Spearman correlation between the averaged exposure and out-of-sample predictions per participant, among participants with ≥2 sessions.

Analysis and figures were done using the statistical software R version 3.4.0 (R Foundation for Statistical Computing, Vienna, Austria) [24] using several packages [25–28].

## Results

Women were slightly older than men (mean (sd): 45 (11) vs. 41 (17), respectively) (Table 1). Most women were illiterate (80%) and engaged in manual unskilled (agricultural) work (60%). At baseline, men reported more working hours per day than women (6.7 (3.7) vs. 5.1 (3.5), respectively) and most men reported zero hours spent cooking with biomass fuel (93%). Forty percent of the participants reported biomass as the primary cooking fuel in the household. Mean number of personal exposure monitoring sessions was similar across men and women.

Personal exposure to PM<sub>2.5</sub> and black carbon was slightly higher in women than in men (Table 2). Relative to women's, men's PM<sub>2.5</sub> personal exposure showed larger variability and higher maximum values (up to 1331 µg/m<sup>3</sup>). On average, ambient concentrations, whether measured at fixed background site or modeled at residence using land-use regression, were lower than personal exposure.

**Table 1** Characteristics of the study population

		Men	Women
N unique participant		207	142
N sessions, <i>m</i> (sd)		1.8 (0.4)	1.7 (0.5)
N sessions, <i>n</i> (%)	Two	159 (77)	101 (71)
Age ≥45 years, <i>n</i> (%)		93 (45)	85 (60)
Education, <i>n</i> (%)	Illiterate	87 (42)	114 (80)
Occupation skill-level, <i>n</i> (%)	Unemployed	37 (18)	32 (22)
	Manual unskilled	72 (35)	85 (60)
	Skilled manual	84 (40)	25 (18)
Occupation type, <i>n</i> (%)	Non-manual	14 (7)	0
	Agriculture	75 (36)	75 (53)
	Industry, construction	30 (15)	4 (3)
Usual hours/day spent at work, <i>m</i> (sd)		6.7 (3.7)	5.1 (3.5)
Smoking status, <i>n</i> (%)	Non smoker	115 (55)	86 (61)
	Passive	37 (18)	56 (39)
	Active	55 (27)	0
Primary stove type, <i>n</i> (%)	Biomass	101 (49)	46 (32)
Kitchen type, <i>n</i> (%)	Separate	175 (84)	110 (77)
Usual ventilation during cooking, <i>n</i> (%)	Always	84 (41)	47 (33)
Usual hours /day spent cooking on biomass, <i>n</i> (%)	0	193 (93)	74 (52)
	1	13 (6)	51 (36)
	≥2	1 (1)	17 (12)

Unemployed category includes housewives, retired participants, and students

*m* mean, *n* number, *sd* standard deviation

**Table 2** PM<sub>2.5</sub> and black carbon exposures in the study population

		Men	Women
PM <sub>2.5</sub>	<i>N sessions</i>	367	243
	Personal exposure (24 h)	55.08 (82.78) <sup>c</sup> [6; 1331]	58.51 (58.84) <sup>d</sup> [3; 564]
	Personal exposure (48-hr average) <sup>a</sup>	53.11 (59.26) [17; 730]	60.36 (45.31) [11; 298]
	Annual ambient concentration at residence <sup>b</sup>	33.02 (2.36) [25; 37]	32.77 (2.57) [24; 37]
	Ambient fixed site (24 h)	31.90 (16.10) [13; 92]	34.16 (18.06) [13; 92]
Black carbon	<i>N sessions</i>	339	220
	Personal exposure (24 h)	4.61 (7.04) <sup>c</sup> [0; 111]	6.06 (9.63) <sup>d</sup> [0; 100]
	Personal exposure (48-hr average) <sup>a</sup>	4.72 (6.71) [0; 75]	6.18 (9.42) [0; 95]
	Annual ambient concentration at residence <sup>b</sup>	2.51 (0.19) [2; 3]	2.52 (0.24) [2; 3]

Figures are mean (sd) [min;max]. Concentrations are expressed in µg/m<sup>3</sup>

<sup>a</sup>Only participants with two sessions

<sup>b</sup>Estimated by land-use regression

<sup>c</sup>Corresponding geometric means (geometric standard deviations) were 41.36 (1.95) for PM<sub>2.5</sub> and 3.21 (2.40) for black carbon

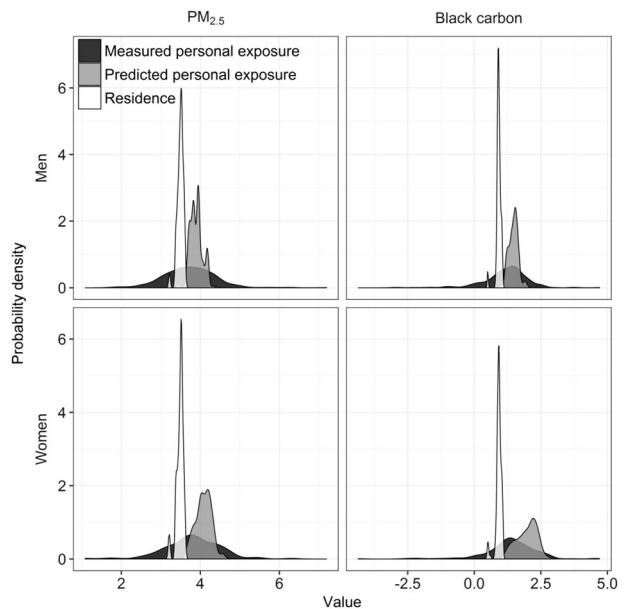
<sup>d</sup>Corresponding geometric means (geometric standard deviations) were 44.78 (2.07) for PM<sub>2.5</sub> and 3.78 (2.82) for black carbon

### Variance components of personal exposure

Log-transformed measured personal exposure to PM<sub>2.5</sub> and black carbon showed much higher within-participant than between-participant variance, resulting in ICC of 0% in men and of 12–18% in women. Adjusting for daily ambient factors decreased the within-participant variance, resulting in higher ICC for men (18–20%) and slightly higher ICC for women (21–22%, Supplementary Table S2). The greatest between-participant variability was observed among participants without separate kitchen (36–39% for PM<sub>2.5</sub> and 56–60% for black carbon) and actively smoking men (32% for PM<sub>2.5</sub> and 46% for black carbon). Patterns of ICC according to subgroups were not always the same across PM<sub>2.5</sub> and black carbon or across men and women.

### Measured personal compared to ambient concentration at residence

Figure 1 compares the probability distribution of measured personal exposure and ambient concentration at residence of PM<sub>2.5</sub> and black carbon in the study population. Distributions were overlapping but measured personal exposure showed a much wider distribution. Very weak correlations were



**Fig. 1** Probability density of measured and predicted personal exposure and annual ambient concentration at residence in men and women. Values are natural log-transformed for clarity. Personal exposure predicted with model including time-invariant predictors. Annual ambient concentrations at residence estimated by previously developed land-use regression model

observed between measured personal exposure and ambient concentration at residence of the same pollutant (Spearman correlation coefficients between  $-0.18$  and  $0.06$ ). Figures were similar when considering averaged personal exposure among participants with  $\geq 2$  sessions (between  $-0.16$  and  $0.09$ ). In men, averaged personal exposure to black carbon was more correlated with residential PM<sub>2.5</sub> than with residential black carbon ( $0.25$  and  $0.07$ , respectively).

Agreement between rank quintiles of averaged personal exposure and ambient concentration at residence of the same pollutant was poor (weighted Kappas  $< 0.09$ ). Poor agreement between rank quintiles was consistent across population subgroups ( $\leq 0.33$ ). The highest level of agreement, though still poor, was found for black carbon in the subgroup of women living close to Hyderabad. The direction of the discrepancy (personal exposure being higher or lower than ambient concentration at residence) varied with subgroups (Supplementary Fig. S1-A, Supplementary Fig. S1-B). For example, for most participants living close to Hyderabad, PM<sub>2.5</sub> concentrations were ranked lower for personal exposure than for ambient concentration at residence, while the reverse was observed for most participants with biomass as primary stove type.

### Predictive models of personal exposure

Figure 1 compares the probability distribution of measured and predicted personal exposures. Distributions were

**Table 3** Prediction models of personal exposure to PM<sub>2.5</sub> and black carbon using time-invariant predictors

	Regression equation	R <sup>2</sup> <sub>within</sub>	R <sup>2</sup> <sub>between</sub>	RMSE	R <sub>spearman</sub>	
Men (n = 367)	Empty model	(ref)	(ref)	–	–	
	+Time-invariant predictors	$\alpha + \alpha_{\text{participant}} + \epsilon_{\text{session}}$ $\alpha + \alpha_{\text{participant}} - 21\%$ if non-smoking <sup>a</sup> – 24% if passive smoking <sup>a</sup> + 30% if construction or industry job – 12% per 1 h spent cycling + 16% if unemployed + 4% if skilled manual occupation + 30% if unskilled manual occupation + $\epsilon_{\text{session}}$	0	53%	0.43	0.31
(n = 339)	+Time-invariant predictors	$\alpha + \alpha_{\text{participant}} + 38\%$ if office or shop job + 27% if biomass primary stove + 4% per 1 $\mu\text{g}/\text{m}^3$ of ambient PM <sub>2.5</sub> <sup>b</sup> at residence + $\epsilon_{\text{session}}$	1%	20%	0.68	0.30
Women (n = 243)	Empty model	(ref)	(ref)	–	–	
	+Time-invariant predictors	$\alpha + \alpha_{\text{participant}} + 13\%$ per 1 h spent cooking with biomass + 22% if biomass primary stove + 2% if unemployed household head + 23% if unskilled manual household head + $\epsilon_{\text{session}}$	0	38%	0.53	0.42
(n = 220)	+Time-invariant predictors	$\alpha + \alpha_{\text{participant}} + 62\%$ if biomass primary stove + 20% per 1 h spent cooking with biomass – 26% if motorcycle household ownership + $\epsilon_{\text{session}}$	1%	57%	0.70	0.50

Models based on personal exposure measurements previously adjusted for ambient factors (see methods). Time-invariant predictors collected at baseline questionnaire.

$\alpha$  and  $\alpha_{\text{participant}}$  are overall and participant-specific random intercept, respectively. Estimates ( $\beta$ ) are transformed to express percent change in personal exposure using  $100 \times [\exp(\beta) - 1]$ . R<sub>spearman</sub> represents Spearman correlation coefficient between averaged measured and averaged predicted values per participant, among those with two measurements. RMSE represents root-mean-squared error between averaged measured and averaged predicted values per participant, among those with two measurements. R<sup>2</sup><sub>within</sub> and R<sup>2</sup><sub>between</sub> represent the proportion of variance explained relative to the empty model

<sup>a</sup>Active smoking was used as reference category

<sup>b</sup>Ambient concentrations at residence were modeled by land-use regression (see methods)

overlapping with similar means but the distributions for predicted exposures were much narrower.

In women, predictors of PM<sub>2.5</sub> and black carbon exposure mostly related to cooking activities, whether time-invariant (Table 3) or time-varying (Supplementary Table S3). For example, 24-h average PM<sub>2.5</sub> personal exposure increased by 13% for each hour spent cooking with biomass, as reported at baseline. Other predictors likely reflected the socio-economic status of the household (occupation of the household head, vehicle ownership, and time spent in vehicle). Predictors selected during validation process were highly consistent (Table 4). Models with time-invariant predictors explained 38% (PM<sub>2.5</sub>) and 57% (black carbon) of the between-woman variability in personal exposure but explained no within-participant variability (Table 3). Correlations between measurements and predictions were moderate (0.42–0.50) and decreased during validation process, particularly for PM<sub>2.5</sub> (from 0.42 to 0.12, Table 4). The inclusion of time-varying variables increased the explained between-participant variance by 26% for PM<sub>2.5</sub> and 10% for black carbon (Supplementary Table S4) but it improved none of the other metrics considered (within-participant variability, RMSE, and correlations between predictions and measurements).

Predictors of PM<sub>2.5</sub> and black carbon exposure were more diverse for men compared to women. Predictors associated with increased personal exposure primarily related to occupation (time-invariant or time-varying) and smoking (Table 3, Supplementary Table S3). Non-smoking was associated with a 21% decrease in PM<sub>2.5</sub> personal exposure as compared to active smoking, but was not a predictor for black carbon. Annual ambient PM<sub>2.5</sub> at residence was associated with an increase in personal exposure for black carbon. Black carbon personal exposure increased by 4% for each hour spent working during the monitoring session (Supplementary Table S3). Time-invariant predictors explained 53% (PM<sub>2.5</sub>) and 20% (black carbon) of the between-man variability in personal exposure; correlation coefficients between predicted and measured values were low (~0.30, Table 3). Similar predictors were selected during validation process but correlation coefficients halved (Table 4). Inclusion of time-varying predictors did not improve the model performance metrics (between- or within-participant variance, RMSE, and correlation coefficients) for either PM<sub>2.5</sub> or black carbon in men (Supplementary Table S3, Supplementary Table S4).

## Discussion

We analyzed personal exposure to PM<sub>2.5</sub> and black carbon in a relatively large sample of the general population of peri-urban South India. Personal exposures to PM<sub>2.5</sub> and

**Table 4** Cross-validation of the prediction models of personal exposure to PM<sub>2.5</sub> and black carbon using time-invariant predictors

	10-fold cross-validation		10-fold cross-holdout validation		Most-selected predictors (% of inclusion)
	RMSE	R <sub>spearman</sub>	RMSE	R <sub>spearman</sub>	
Men	0.45	0.24	0.47	0.14	Usual time cycling (80%), construction or industry job (70%), primary occupation skill-level (60%), smoking status (40%)
Black carbon	0.69	0.25	0.72	0.12	Ambient PM <sub>2.5</sub> <sup>a</sup> at residence (50%), primary stove type (50%), office or shop job (40%)
Women	0.55	0.33	0.59	0.12	Primary stove type (60%), ambient black carbon <sup>a</sup> at residence (40%), occupation of household head (40%), manual unskilled job (30%)
Black carbon	0.73	0.43	0.75	0.39	Usual time cooking with biomass (60%), separate kitchen (50%)

Time-invariant predictors obtained from baseline questionnaire. R<sub>spearman</sub> represents Spearman correlation coefficient between averaged measured and averaged predicted values per participant. RMSE represents root-mean-squared error between averaged measured and averaged predicted values per participant. Only predictors included in ≥30% of the models are shown <sup>a</sup>Ambient concentrations at residence were modeled by land-use regression (see methods)

black carbon were, on average, higher than and relatively poorly correlated with annual ambient concentrations at residence. Personal exposure variability was substantially larger within participant than between participant. Predictors of personal exposure to PM<sub>2.5</sub> and black carbon included cooking activities (women), occupation (men), and smoking (men). Prediction models explained a moderate amount of between-person variability in measured personal exposure, except for black carbon in men where model performance was poor.

We observed larger within-participant variance compared to between-participant variance in measured personal exposure for both PM<sub>2.5</sub> and black carbon, even after temporal adjustment for daily ambient factors. Compared to published results for other populations in LMICs, our results showed lower between-participant variation (Supplementary Table S5). McCracken et al. reported an ICC of 33% in children and 29% in adult women for personal exposure to CO in Guatemala [12, 29]. Dionisio et al. reported an ICC of 39% for personal CO for children in The Gambia [30]. Several studies from high-income countries have reported higher between-participant variation in PM<sub>2.5</sub> personal exposure [31–34], but not all [35]. The relatively low ICC observed in our study may be due to the limited number of measurements (2 days in two different seasons). However, a nested panel study within this population with up to 6 measurement-days per person throughout the year observed similar between-participant variability in PM<sub>2.5</sub> personal exposure [36]. The large within-participant variability of personal exposure we observed for women may reflect the variability generally observed in cooking stove emissions [37], a major source of particulate matter for women. This high temporal variability in personal exposure likely reduces the precision of these measurements as surrogates for long-term exposure in the study population.

Measured personal exposure to PM<sub>2.5</sub> and black carbon showed neither correlation nor agreement (between quintiles) with annual ambient concentration at residence modeled by land-use regression. No subgroup with moderate or good level of agreement could be identified. However, some subgroups were identified as having notably low agreement, for example for PM<sub>2.5</sub>: male smokers, males with non-manual occupation, and women with biomass primary stove. Previous studies from high-income countries have found modest correlation coefficient between measured personal exposure and long-term ambient (measured or modeled) concentrations of PM<sub>2.5</sub> [38–40] or black carbon [41, 42], but generally higher than what we observed. Although methods may differ across studies (e.g., population characteristics, modeling method, monitoring time and period), the literature generally supports the use of modeled long-term ambient levels as a surrogate of personal exposure in areas where ambient sources (e.g., traffic) are major

contributors of personal exposure. This is not the case in our study area, in which biomass cooking fuel, smoking, and occupational exposures are dominant sources of personal exposure. The high variability observed in personal exposure due to these diverse sources likely contributes to the poor agreement with ambient concentrations at residence. Land-use regression models aim to estimate spatial variability in ambient concentrations, while personal exposure is a mixture of temporal and spatial variations in ambient and non-ambient levels. They likely capture different components of the true exposure and reflect the contribution from different sources [43]; they are therefore both potentially relevant in epidemiological studies.

To predict personal exposure, we developed an empirical, data-driven model based on a stepwise predictor selection using individual data, residence characteristics, and general time-activity data. Consistently with literature in LMICs, women's personal exposure to  $PM_{2.5}$  and black carbon were driven by cooking activities or habits (use of biomass, time spent cooking, and ventilation) and socioeconomic position of the household (use of biomass, separate kitchen, and motorcycle ownership) [44–47]. Occupational exposure was an important contributor of  $PM_{2.5}$  personal exposure in men but we were unable to identify specific occupational tasks relevant for personal exposure, possible due to the wide variability in occupation types (e.g. industry or agriculture) and related tasks throughout the year. In the study area, some of the selected predictors of personal exposure (use of biomass, motorcycle ownership, occupation type, and smoking status) correlated strongly with the urbanization level. Urbanization level could thus contribute to total personal exposure through these interrelated aspects. Overall, our models performed moderately well (except the black carbon model for men) in terms of between-participant variance explained (38–53%) and correlations between predictions and measurements (Spearman coefficients from 0.3–0.5). These results are comparable with several previous publications. A model for personal CO exposure in Guatemalan children used *a priori* selected predictors and explained 47% of between-participant variability [12]. Using a backward stepwise procedure, a  $PM_{2.5}$  model explained 43% of between-participant variability in pregnant Canadian women [32]. Personal  $PM_{2.5}$  prediction model explained 74% of variability in summer but only 5% in winter in another study based on a Canadian population [34]. The  $PM_{10}$  personal exposure predicted in an Irish population correlated more strongly (Pearson coefficients between 0.55 and 0.84) with measurements than did ours [48]. Specific activities or locations can have a large impact on personal exposure through peaks of exposure e.g., commuting, smoking, and using a kerosene lamp [46, 49, 50]. More time-resolved data could thus help to explain variability in personal exposure.

Yet, in the present study, the addition of more detailed, time-varying variables about activities performed during the monitoring did not improve the overall performance or predictive ability of our models—except for the addition of cooking activities performed during monitoring in women's model for  $PM_{2.5}$  exposure. This overall lack of improvement might be a consequence of the coarse time resolution of the time-activity questionnaire and its inability to capture short-term activities or tasks that may be related to peak exposures [20, 21]. Previous analyses in the study population showed better performance of the diary for women than for men; a result possibly related to the rather homogeneous activities (related to cooking) performed by women. This could explain why women's prediction models for  $PM_{2.5}$  explained more variability after adding the diary data; however, the predictive ability remained similar.

Our results support the feasibility of using predicted long-term personal exposure for epidemiological studies in LMIC contexts. Our aim was to capture total personal exposure and not just the component of personal exposure due to ambient concentrations, which could be estimated using other approaches. Ambient concentrations appear to contribute little to total exposures for our study population in peri-urban India as compared to local sources (biomass burning, occupational exposures, and smoking). The prediction models showed better out-of-sample predictive ability for personal exposure (in terms of correlation with measurements) than annual concentrations at residence modeled by land-use regression. The variance and interquartile range of the predicted values were also larger, potentially increasing the statistical power if used as exposure in an epidemiological analysis. A main advantage of the continuous predicted values is that they facilitate estimating an exposure-response function, a clear advantage over categorical indicators based on cooking fuel, which have been previously used for epidemiological studies in LMICs [51]. The literature in settings with prevalent biomass cooking fuel use is largely limited to women and children, shedding little light on exposure levels and health effects in men.

The limitations of the personal exposure predictions warrant consideration. First, the majority of the total variability in measured exposure was within participant, which the models did not explain beyond adjustment for ambient factors ( $PM_{2.5}$  concentrations and meteorology). More time-resolved or detailed activity data could have improved the performance of our models (regarding within-participant variability) as such data could capture peaks in daily exposure. However, in the context of epidemiological studies of long-term exposure, between-participant variability is more relevant, for which most of our models had moderate performance. For other research objectives focusing on more time-resolved exposures (e.g., hourly concentrations), improved time-activity data are likely to be important as we previously demonstrated in a



nested panel study using highly time-resolved activity information derived from wearable cameras [21]. Second, we had available only 2 days of personal exposure measurements, which may not be sufficiently representative of long-term exposure. Additional repeated measurements throughout a year would likely have provided a better reflection of long-term exposure. Nonetheless, in a nested panel study with up to six repeated measurements, the ICC remained low [36], suggesting a very large number of repeated measurements would be required for a notable improvement. Third, we did not have data to validate the prediction models in an independent external dataset. The models appeared robust during the extensive evaluation process but the evidence of applicability to other populations cannot be ensured. However, as the study population represents a stratified random sample, the prediction models could be applied to the general population of the study area. Finally, the uncertainties introduced when deriving black carbon measurements might be an explanation for the poorer predictive ability of the prediction models for black carbon as compared to  $PM_{2.5}$ .

Our results provide valuable insights into the limited agreement between measured personal exposure and estimates of annual ambient concentration at residence in a LMIC setting, where exposure is not dominated by sources correlated with land use. This has important implications for epidemiology in this and similar settings, as estimates of exposure to ambient pollution are likely to capture only a small fraction of true personal exposure. Our results demonstrate the feasibility of combining personal exposure measurements with questionnaire data on usual activities to generate estimates of particulate matter exposure for a relatively large population. These estimates appear to capture aspects of exposure independent of what is captured by land-use regression. Epidemiological studies of long-term exposure to particulate air pollution in LMIC settings will require exposure assessment approaches that consider both sources of ambient concentrations and of personal exposure.

**Acknowledgements** The research leading to these results received funding from the European Research Council under ERC Grant Agreement number 336167 for the CHAI Project. The third wave of data collection and village socio-demographic surveys for the APCAPS study were funded by the Wellcome Trust (Grant 084674/Z). C.T. was funded through a Ramón y Cajal fellowship (RYC-2015-17402) awarded by the Spanish Ministry of Economy and Competitiveness. We thank all participants and study teams of the APCAPS and CHAI studies.

## Compliance with ethical standards

**Conflict of interest** The authors declare that they have no conflict of interest.

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

## References

1. World Health Organization. Review of evidence on health aspects of air pollution - REVIHAAP Project. 2013. <http://www.euro.who.int/en/health-topics/environment-and-health/air-quality/publications/2013/review-of-evidence-on-health-aspects-of-air-pollution-revihaap-project-final-technical-report>.
2. U.S. EPA. Integrated Science Assessment for Particulate Matter. 2009 EPA/600/R-08/139F.
3. Cohen AJ, Brauer M, Burnett R, Anderson HR, Frostad J, Estep K, et al. Estimates and 25-year trends of the global burden of disease attributable to ambient air pollution: an analysis of data from the Global Burden of Diseases Study 2015. *Lancet*. 2017;389:1907–18.
4. Tonne C. A call for epidemiology where the air pollution is. *Lancet Planet Heal*. 2017;1:e355–e356.
5. Özkaynak H, Baxter LK, Dionisio KL, Burke J. Air pollution exposure prediction approaches used in air pollution epidemiology studies. *J Expo Sci Environ Epidemiol*. 2013;23:566–72.
6. Marshall JD, Nethery E, Brauer M. Within-urban variability in ambient air pollution: comparison of estimation methods. *Atmos Environ*. 2008;42:1359–69.
7. Lane KJ, Levy JI, Scammell MK, Patton AP, Durant JL, Mwamburi M, et al. Effect of time-activity adjustment on exposure assessment for traffic-related ultrafine particles. *J Expo Sci Environ Epidemiol*. 2015;25:506–16.
8. Deffner V, Küchenhoff H, Maier V, Pitz M, Cyrys J, Breitner S, et al. Personal exposure to ultrafine particles: two-level statistical modeling of background exposure and time-activity patterns during three seasons. *J Expo Sci Environ Epidemiol*. 2016; 26:17–25.
9. Glasgow ML, Rudra CB, Yoo E-H, Demirbas M, Merriman J, Nayak P, et al. Using smartphones to collect time-activity data for long-term personal-level air pollution exposure assessment. *J Expo Sci Environ Epidemiol*. 2016;26:356–64.
10. Steinle S, Reis S, Sabel CE. Quantifying human exposure to air pollution-Moving from static monitoring to spatio-temporally resolved personal exposure assessment. *Sci Total Environ*. 2013;443:184–93.
11. Isaacs K, McCurdy T, Glen G, Nysewander M, Erickson A, Forbes S, et al. Statistical properties of longitudinal time-activity data for use in human exposure modeling. *J Expo Sci Environ Epidemiol*. 2013;23:328–36.
12. McCracken JP, Schwartz J, Bruce N, Mittleman M, Ryan LM, Smith KR. Combining individual- and group-level exposure information. *Epidemiology*. 2009;20:127–36.
13. Tonne C, Salmon M, Sanchez M, Sreekanth V, Bhogadi S, Sambandam S, et al. Integrated assessment of exposure to  $PM_{2.5}$  in South India and its relation with cardiovascular risk: design of the CHAI observational cohort study. *Int J Hyg Environ Health*. 2017;220:1081–8.
14. Kinra S, Radha Krishna K, Kuper H, Rameshwar Sarma K, Prabhakaran P, Gupta V, et al. Cohort Profile: Andhra Pradesh Children and Parents Study (APCAPS). *Int J Epidemiol*. 2014;43:1417–24.
15. Balakrishnan K, Sambandam S, Ramaswamy P, Ghosh S, Venkatesan V, Thangavel G, et al. Establishing integrated rural-urban cohorts to assess air pollution-related health effects in pregnant women, children and adults in Southern India: an overview of objectives, design and methods in the Tamil Nadu Air Pollution and Health Effects (TAPHE) s. *BMJ Open*. 2015;5:e008090–e008090.
16. Zanatta M, Gysel M, Bukowiecki N, Müller T, Weingartner E, Areskoug H, et al. A European aerosol phenomenology-5: Climatology of black carbon optical properties at 9 regional background sites across Europe. *Atmos Environ*. 2016;145:346–64.

17. Kumar MK, Sreekanth V, Salmon M, Tonne C, Marshall JD. Use of spatiotemporal characteristics of ambient PM 2.5 in rural South India to infer local versus regional contributions. *Environ Pollut*. 2018;239:803–11.
18. Sanchez M, Ambros A, Milà C, Salmon M, Balakrishnan K, Sambandam S, et al. Development of land-use regression models for fine particles and black carbon in peri-urban South India. *Sci Total Environ*. 2018;634:77–86.
19. Balakrishnan K, Sambandam S, Ramaswamy P, Mehta S, Smith KR. Exposure assessment for respirable particulates associated with household fuel use in rural districts of Andhra Pradesh, India. *J Expo Sci Environ Epidemiol*. 2004;14:S14–S25.
20. Sanchez M, Ambros A, Salmon M, Bhogadi S, Wilson R, Kinra S, et al. Predictors of daily mobility of adults in Peri-Urban South India. *Int J Environ Res Public Health*. 2017;14:783.
21. Salmon M, Milà C, Bhogadi S, Addanki S, Madhira P, Muddepaka N, et al. Wearable camera-derived microenvironments in relation to personal exposure to PM 2.5. *Environ Int*. 2018;117:300–7.
22. Snijders TAB, Berkhof J Diagnostic Checks for Multilevel Models. In: *Handbook of Multilevel Analysis*. New York: Springer New York, pp 141–75.
23. Wang M, Brunekreef B, Gehring U, Szpiro A, Hoek G, Beelen R. A New Technique for Evaluating Land-use Regression Models and Their Impact on Health Effect Estimates. *Epidemiology*. 2016;27:51–6.
24. R Core Team. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. 2015. <https://www.r-project.org/>.
25. Nieuwenhuis R, Te Grotenhuis M, Pelzer B. Influence.ME: tools for detecting influential data in mixed effects models. *R J*. 2012;4:38–47.
26. Wickham H. tidyverse: Easily Install and Load ‘Tidyverse’ Packages. R package version 1.0.0. 2016. <https://cran.r-project.org/package=tidyverse>.
27. Wickham H. ggplot2: Elegant Graphics for Data Analysis. R package version 2.2.0. 2009. <http://ggplot2.org>.
28. Bates D, Maechler M, Bolker B, Walker S. Fitting linear mixed-effects models using lme4. *J Stat Softw*. 2015;67:1–48.
29. McCracken JP, Schwartz J, Diaz A, Bruce N, Smith KR. Longitudinal relationship between personal CO and personal PM2.5 among women cooking with woodfired cookstoves in Guatemala. *PLoS One*. 2013;8:e55670.
30. Dionisio KL, Howie SRC, Dominici F, Fornace KM, Spengler JD, Donkor S, et al. The exposure of infants and children to carbon monoxide from biomass fuels in The Gambia: a measurement and modeling study. *J Expo Sci Environ Epidemiol*. 2012;22:173–81.
31. Lee K, Bartell SM, Paek D. Interpersonal and daily variability of personal exposures to nitrogen dioxide and sulfur dioxide. *J Expo Sci Environ Epidemiol*. 2004;14:137–43.
32. Nethery E, Teschke K, Brauer M. Predicting personal exposure of pregnant women to traffic-related air pollutants. *Sci Total Environ*. 2008;395:11–22.
33. Nethery E, Leckie SE, Teschke K, Brauer M. From measures to models: an evaluation of air pollution exposure assessment for epidemiological studies of pregnant women. *Occup Environ Med*. 2008;65:579–86.
34. MacNeill M, Wallace L, Kearney J, Allen RW, Van Ryswyk K, Judek S, et al. Factors influencing variability in the infiltration of PM2.5 mass and its components. *Atmos Environ*. 2012;61:518–32.
35. Johannesson S, Gustafson P, Molnár P, Barregard L, Sällsten G. Exposure to fine particles (PM2.5 and PM1) and black smoke in the general population: personal, indoor and outdoor levels. *J Expo Sci Environ Epidemiol*. 2007;17:613–24.
36. Milà C, Salmon M, Sanchez M, Ambrós A, Bhogadi S, Sreekanth V, et al. When, where, and what? characterizing personal PM2.5 exposure in Periurban India by integrating GPS, wearable camera, and ambient and personal monitoring data. *Environ Sci Technol*. 2018;52:13481–90.
37. Carter E, Archer-Nicholls S, Ni K, Lai AM, Niu H, Secrest MH, et al. Seasonal and diurnal air pollution from residential cooking and space heating in the Eastern Tibetan Plateau. *Environ Sci Technol*. 2016;50:8353–61.
38. Sarnat JA, Brown KW, Schwartz J, Coull BA, Koutrakis P. Ambient gas concentrations and personal particulate matter exposures. *Epidemiology*. 2005;16:385–95.
39. Michikawa T, Nakai S, Nitta H, Tamura K. Validity of using annual mean particulate matter concentrations as measured at fixed site in assessing personal exposure: an exposure assessment study in Japan. *Sci Total Environ*. 2014;466–7:673–80.
40. Miller KA, Spalt EW, Gasset AJ, Curl CL, Larson TV, Avol E, et al. Estimating ambient-origin PM2.5 exposure for epidemiology: observations, prediction, and validation using personal sampling in the Multi-Ethnic Study of Atherosclerosis. *J Expo Sci Environ Epidemiol*. 2019;29:227–37.
41. Nieuwenhuijsen MJ, Donaire-Gonzalez D, Rivas I, de Castro M, Cirach M, Hoek G, et al. Variability in and agreement between modeled and personal continuously measured black carbon levels using novel smartphone and sensor technologies. *Environ Sci Technol*. 2015;49:2977–82.
42. Montagne D, Hoek G, Nieuwenhuijsen M, Lanki T, Pennanen A, Portella M, et al. Agreement of land use regression models with personal exposure measurements of particulate matter and nitrogen oxides air pollution. *Environ Sci Technol*. 2013;47:8523–31.
43. Wilson WE, Brauer M. Estimation of ambient and non-ambient components of particulate matter exposure from a personal monitoring panel study. *J Expo Sci Environ Epidemiol*. 2006;16:264–74.
44. Balakrishnan K, Parikh J, Sankar S, Padmavathi R, Srividya K, Venugopal V, et al. Daily average exposures to respirable particulate matter from combustion of biomass fuels in rural households of Southern India. *Environ Health Perspect*. 2002;110:1069–75.
45. Pant P, Habib G, Marshall JD, Peltier RE. PM 2.5 exposure in highly polluted cities: a case study from New Delhi, India. *Environ Res*. 2017;156:167–74.
46. Van Vliet EDS, Asante K, Jack DW, Kinney PL, Whyatt RM, Chillrud SN, et al. Personal exposures to fine particulate matter and black carbon in households cooking with biomass fuels in rural Ghana. *Environ Res*. 2013;127:40–8.
47. Nayek S, Padhy PK. Daily personal exposure of women cooks to respirable particulate matters during cooking with solid bio-fuels in a rural community of West Bengal, India. *Aerosol Air Qual Res*. 2017;17:245–52.
48. McCreddin A, Alam MS, McNabola A. Modelling personal exposure to particulate air pollution: an assessment of time-integrated activity modelling, Monte Carlo simulation and artificial neural network approaches. *Int J Hyg Environ Health*. 2015;218:107–16.
49. Meng QY, Spector D, Colome S, Turpin B. Determinants of indoor and personal exposure to PM2.5 of indoor and outdoor origin during the RIOPA study. *Atmos Environ*. 2009;43:5750–8.
50. Rivas I, Donaire-Gonzalez D, Bousso L, Esnaola M, Pandolfi M, de Castro M, et al. Spatiotemporally resolved black carbon concentration, schoolchildren’s exposure and dose in Barcelona. *Indoor Air*. 2016;26:391–402.
51. McCracken JP, Wellenius GA, Bloomfield GS, Brook RD, Tolunay HE, Dockery DW, et al. Household air pollution from solid fuel use. *Glob. Heart*. 2012;7:223–34.