

Article

Not peer-reviewed version

A Data Fusion Approach to Assessing the Contribution of Wildland Fire Smoke to Fine Particulate Matter in California

[Hongjian Yang](#)^{*}, Sofia Ruiz-Suarez, Brian Reich, Yawen Guan, [Ana Rappold](#)

Posted Date: 10 July 2023

doi: 10.20944/preprints202307.0596.v1

Keywords: Bayesian analysis; calibration; citizen science; spatiotemporal methods; spectral analysis



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Article

A Data Fusion Approach to Assessing the Contribution of Wildland Fire Smoke to Fine Particulate Matter in California

Hongjian Yang ^{1,*}, Sofia Ruiz-Suarez ², Brian J Reich ³, Yawen Guan ⁴ and Ana G Rappold ⁵

¹ North Carolina State University; hyang23@ncsu.edu

² University of Toronto, National University of Comahue; sofia.ruizsuarez@utoronto.ca

³ North Carolina State University; bjreich@ncsu.edu

⁴ University of Nebraska - Lincoln; yawen.guan@unl.edu

⁵ US Environmental Protection Agency; Rappold.Ana@epa.gov

* Correspondence: hyang23@ncsu.edu

Abstract: The escalating frequency and severity of global wildfires necessitate an in-depth understanding and monitoring of wildfire smoke impacts, specifically its contribution to fine particulate matter (PM_{2.5}). We propose a data-fusion method to study wildfire contribution to PM_{2.5} using satellite-derived smoke plume indicators and PM_{2.5} monitoring data. Our study incorporates two types of monitoring data, the high-quality but sparse Air Quality System (AQS) stations and the abundant but less accurate PurpleAir (PA) sensors that are gaining popularity among citizen scientists. We propose a multi-resolution spatiotemporal model specified in the spectral domain to calibrate the PA sensors against accurate AQS measurements, and leverage the two networks to estimate wildfire contribution to PM_{2.5} in California in 2020 and 2021. A Bayesian approach is taken to incorporate all uncertainties and our prior intuition that the dependence between networks, as well as the accuracy of PA network, vary by frequency. We find that 1% to 3% increase in PM_{2.5} concentration due to wildfire smoke, and that leveraging PA sensors improves accuracy.

Keywords: Bayesian analysis; calibration; citizen science; spatiotemporal methods; spectral analysis

1. Introduction

Airborne particles are a serious environmental health risk globally, contributing in excess of 7 million premature deaths each year [1]. Fine particulate matter (PM_{2.5}, particles with a diameter of less than 2.5 micrometers) has been causally linked to cardiovascular morbidity and mortality [2] and are therefore regulated under the provisions of the Clean Air Act [3] to protect human health and wellbeing. As a result, the emissions of PM_{2.5} from many anthropogenic sources, such as transpiration and industry, have been on a steady decline [4] and wildfires have become the single largest source [5], potentially offsetting reduction in emissions from other sources.

High concentrations of fine particles and gasses found in smoke have also produced alarming impacts on health [6,7]. During peak wildfire seasons, smoke exposure can exacerbate health problems, causing a spike in emergency department visits [8]. In an epidemiological study of health impacts, [9] estimated that 2.2 % of annual respiratory health burden, or 92 ED visits per 100,000 people, is attributed to ambient PM and that wildfire days account for over 15% of that burden. However, providing a definite answer as to how much of particle pollution can be attributed to wildfires remains a challenging problem because instruments measure a total ambient concentration which is composed of natural, anthropogenic, and wildfire sources.

Previous research has studied the contribution of wildfires on PM_{2.5} concentrations by integrating remote sensing data on the location and extent of smoke plumes and PM_{2.5} readings from Air Quality System (AQS) monitors deployed by the Environmental Protection Agency (EPA). These studies revealed that wildfires contribute to 40% of unhealthy days and substantially increase PM_{2.5} concentrations [10,11]. Wildfire smoke impacts are dynamic and often affect areas without a monitoring

station, as AQS monitors have limited spatial coverage due to the high cost and difficulty in installation. It is important to make air quality information available to the public quickly during wildfires, therefore AQS alone provides insufficient data source for monitoring wildfire emissions.

The increased incidence of days with poor air quality due to wildfires has created a demand and public interest for monitoring particulate pollution. Perhaps the most prevalent sensors are PurpleAir (PA), which are installed by members of the public, providing a near real-time monitoring of $PM_{2.5}$ with extensive spatial coverage [12]. However, it is known that PA sensors are less reliable compared to the AQS, and thus correction to the sensor readings is needed [13,14]. [15] developed a correction equation using meteorological conditions including relative humidity and temperature, as both measurements affect the accuracy of the instrument; however, this calibration is developed for a US-wide correction and without smoke impacts. Another simple linear correction model under smoke impacted conditions was proposed in [16]. As the sensor performance can be affected by geographic and environmental conditions, it is more reasonable to relax the assumption of a constant spatially varying bias, but rather capture the spatiotemporally varying bias.

Smoke impacts are dynamic and often affect areas without a monitoring station, as AQS monitors have limited spatial coverage. Previous studies have either separated anthropogenic PM from smoke emissions using chemical transport models or by subtracting out historically observed averages [17]. However, neither approach provides a definite answer as to how much of particle pollution can be attributed to wildfires. Data fusion is a widely used method that integrates information from different types of sensors to provide a robust and complete description of a process of interest [18,19]. It has been used extensively to estimate spatially and temporally resolved air quality surfaces. For example, [20–23] use data fusion method to study the complex relationship between monitoring data and outputs from Community Multi-Scale Air Quality (CMAQ), a deterministic chemical transport model. [24] combines observations from two noisy datasets to predict the true aerosol process. More recently, several researchers have exploited the usefulness of low-cost sensors such as Purple Air to map air quality and quantify the uncertainty of estimation [25–27]. Most similar to our approach is Stein et al (2005), who also use a spectral transformation in time and spatial processes to capture dependence between stations for a single fixed monitoring network [28]. We extend this approach to handle multiple data networks.

This study aims to provide an estimate of wildfire contribution on air quality by supplementing the remotely sensed smoke plume indicators with PurpleAir data. We propose a multi-resolution Bayesian approach fusing information from both AQS monitors and PA sensors to estimate the contribution to $PM_{2.5}$ caused by wildfires. We apply a Discrete Fourier Transform (DFT) to account for temporal correlation, transforming the data from the time domain to the frequency domain, and model the spatial correlation in the frequency domain. To quantify the relative increase in $PM_{2.5}$ concentrations due to wildfires, we propose regression and matching estimators, as discussed in Section 2.3. Our findings will not only enhance understanding of the relationship between wildfires and air pollution but also inform policy and decision-making related to wildfire management, public health, and climate change impacts.

The remainder of this paper is structured as follows: In Section 2, we present an overview of the Hazard Mapping System (HMS) and its smoke plume data. We also delve into the specifics of the AQS and PurpleAir sensors, contrasting their features and functionality. The latter part of Section 2 elucidates our use of the Discrete Fourier Transform (DFT) and the data fusion model, with detailed insights into our Markov Chain Monte Carlo (MCMC) approach. Section 3 presents our findings, including the estimated smoke plume parameters and the calculated contribution of wildfires to $PM_{2.5}$ concentrations, as determined by both regression and matching estimators. We also offer a comparison of various models, illustrating the effectiveness of the PurpleAir sensors and the benefits of our proposed data fusion method. Finally, in Section 4, we discuss the implications of our findings, the limitations of our study, and potential directions for future research.

2. Materials and Methods

2.1. Data sources and exploratory analysis

Our analysis incorporates data from three distinct sources: satellite-derived smoke plume indicators obtained through the National Oceanic and Atmospheric Administration's Hazard Mapping System (HMS), $PM_{2.5}$ measurements from Air Quality System (AQS) monitoring stations, and $PM_{2.5}$ readings from PurpleAir (PA) monitoring stations. We collect hourly data and average them to daily level from each source for 2020 and 2021 fire seasons, spanning July 1 to October 31. The original $PM_{2.5}$ readings from both AQS and PA stations are right-skewed so we apply log-transformation to all $PM_{2.5}$ readings.

2.1.1. Satellite-derived smoke plume indicators

Exposure to wildland fire smoke is assessed using smoke plume indicators supplied by the Hazard Mapping System (HMS) [29]. This automated data product integrates observations from multiple polar and geostationary satellites to generate daily polygons representing smoke plume extents in near-real time. Distinct polygons are provided for low-, medium-, and high-density plumes. Figure 1 shows the HMS smoke plumes for September 20, 2021. These smoke plume indicators tend to underestimate the actual intensity of smoke, as they primarily rely on satellite imagery with an approximate spatial resolution spanning several miles [30]. Additionally, smoke visibility is limited to daytime hours, resulting in a significant underestimation of smoke levels during the night.

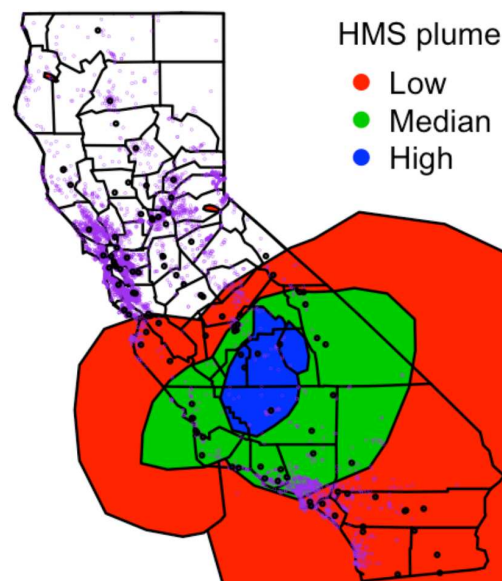


Figure 1. HMS smoke plume density on September 20, 2021 (shaded regions) the locations of PurpleAir (purple dots) and Air Quality System (black dots) monitoring stations.

2.1.2. AQS monitoring stations

The Air Quality System (AQS) monitoring stations, deployed by the US Environmental Protection Agency (EPA) and state, local, and tribal air pollution control agencies, provide precise $PM_{2.5}$ measurements. However, their distribution is spatially sparse due to the high cost and complexity associated with their installation and maintenance. The black dots in Figure 1 depict the distribution of AQS monitoring sites across California in 2021. Figure 2 illustrates the distribution of $PM_{2.5}$ levels, aggregated across stations, by smoke plume intensity.

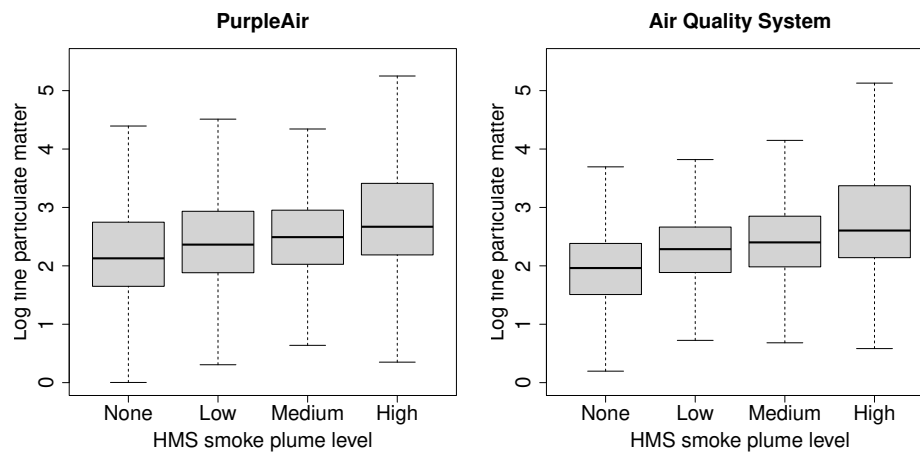


Figure 2. Distribution of log $PM_{2.5}$ ($\mu g/m^3$) by smoke plume level for PurpleAir (PA) and Air Quality System (AQS) stations. Four smoke plume levels from left to right are: no smoke, low, medium, and high plume density.

2.1.3. PurpleAir sensors

PurpleAir (PA) sensors are low-cost monitoring devices deployed by individuals and organizations for continuous ambient air pollutant tracking. Despite their affordability and ease of installation, PA sensors offer less accurate $PM_{2.5}$ readings and are significantly influenced by environmental factors, such as temperature and humidity [15]. We use bias corrected data for all analyses. However, this initial bias correction based on [15] may be insufficient because it only depends on a linear trend in temperature and humidity and is constant across space and time. Therefore, our Bayesian data fusion model adds a more flexible spatiotemporal bias correction term.

In 2021, more than 7,800 outdoor PA sensors were operational in California. The purple dots in Figure 1 represent the distribution of PA monitoring sites throughout the state. We have included only those PA stations that reported fewer than 18 missing days during the fire seasons, resulting in a total of 1,080 for 2020 and 712 PA stations for 2021. Figure 2 displays the distribution of $PM_{2.5}$, aggregated across stations for 2021, by smoke plume intensity. A similar pattern is observed in both PA stations and AQS stations where $PM_{2.5}$ measurements escalate in the presence of a smoke plume.

Figure 3 investigates the relationship between AQS stations and their corresponding nearby PA sites. Each point in Figure 3 symbolizes one AQS station. For every AQS station, we pinpoint the closest PA station and calculate the sample correlation between the daily $PM_{2.5}$ measurement series for these two locations. We can see that the correlation is generally high for nearby stations and decreases with distance, suggesting that PA data will be a useful supplement to the spatial model.

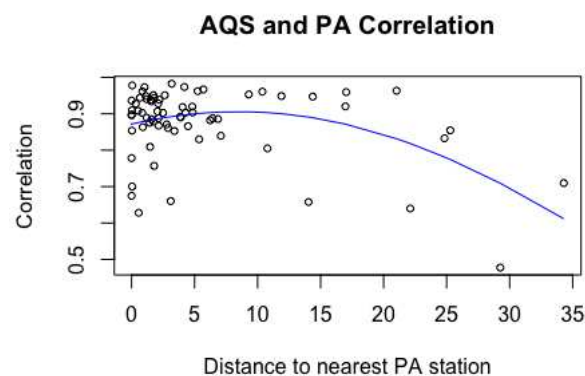


Figure 3. Sample correlation between each AQS stations the nearest PA stations versus the distance (km) between the two stations.

2.2. Statistical model

We propose a multi-resolution Bayesian model for modeling AQS and PA measurements jointly in the spectral domain. Let $Y_{1t}(\mathbf{s})$ and $Y_{2t}(\mathbf{s})$ be AQS and PA measurements, respectively, for spatial location \mathbf{s} at time (day) $t \in \{1, \dots, n_t\}$, and $\mathbf{X}_t(\mathbf{s}) = \{X_{0t}(\mathbf{s}), \dots, X_{pt}(\mathbf{s})\}$ be a corresponding vector of covariates with $X_{0t}(\mathbf{s}) = 1$ for the intercept. The $p = 5$ covariates are temperature, relative humidity and indicators of low, medium and high density smoke plumes at site \mathbf{s} and day t . We note that temperature and relative humidity are standardized to have mean zero and variance one and that the AQS and PA measurements are not taken at the same spatial locations.

The observations are decomposed as $Y_{jt}(\mathbf{s}) = Z_{jt}(\mathbf{s}) + \varepsilon_{jt}(\mathbf{s})$ for $j \in \{1, 2\}$, where $Z_{1t}(\mathbf{s})$ and $Z_{2t}(\mathbf{s})$ are spatiotemporal processes, and $\varepsilon_{jt}(\mathbf{s}) \stackrel{\text{indep}}{\sim} \text{Normal}(0, \tau_j^2)$ is error. The time span of our data is relatively short, therefore, it is reasonable to assume the spatiotemporal processes are stationary within the modeling period. We will apply Fourier transformation to the spatiotemporal processes $Z_{jt}(\mathbf{s})$ with respect to time to remove the temporal dependence. The resulting spectral processes $Z_{jl}^*(\mathbf{s})$ capture periodicity, are independent over frequency $\{\omega_l, l = 1, \dots, n_t\}$ and spatially correlated. For time series observed at equal time intervals, we can apply the discrete Fourier transformation (DFT). The spectral processes at frequency ω_l is

$$Z_{jl}^*(\mathbf{s}) = \sum_{t=1}^{n_t} \exp(-it\omega_l) Z_{jt}(\mathbf{s}) \quad (1)$$

and measures the variation in $Z_{jt}(\mathbf{s})$ at frequency ω_l . Terms with small ω_l (low frequency) represent long-term trends such as month-to-month averages and terms with large ω_l (high frequency) represent short-term trends such as day-to-day variation. Let $\{Z_{j1}^*(\mathbf{s}), \dots, Z_{jn_t}^*(\mathbf{s})\}$ be the unique real components of the DFT of $\{Z_{j1}(\mathbf{s}), \dots, Z_{jn_t}(\mathbf{s})\}$ at frequencies $\{\omega_1, \dots, \omega_{n_t}\}$ with $\omega_1 \leq \dots \leq \omega_{n_t}$.

The spectral processes $Z_{jl}^*(\mathbf{s})$ are dependent across $j = 1, 2$, as they represent the two networks measuring the same underlying PM_{2.5} process. They are also spatially dependent processes as locations nearby may exhibit similar periodicity. We model the cross network dependence and spatial dependence for each ω_l as

$$Z_{1l}^*(\mathbf{s}) = U_l(\mathbf{s}) \quad \text{and} \quad Z_{2l}^*(\mathbf{s}) = A_l U_l(\mathbf{s}) + V_l(\mathbf{s}), \quad (2)$$

where spatial process $U_l(\mathbf{s})$ is the true PM_{2.5} concentration. The PA stations are assumed to be measuring a biased and noisy version of the true PM_{2.5} with discrepancy $V_l(\mathbf{s})$. The coefficient A_l controls the dependence across networks. Both the bias $V_l(\mathbf{s})$ and cross-dependence A_l vary by ω_l to allow for a multi-resolution calibration of the two networks. We model A_l linearly as $A_l = \beta_{A0} + \beta_{A1} \cdot \omega_l$, where β_{A0} and β_{A1} are unknown coefficients. This allows the correlation between the processes to vary stochastically with frequency. For example, if PA is more reliable for long-term trends than day-to-day variation, then we expect larger (smaller) correlation between networks for small (large) ω_l .

The true process $U_l(\mathbf{s})$ and discrepancy term $V_l(\mathbf{s})$ are both regressed onto the covariates. Since we are developing a model in the spectral domain, we will also apply DFT to each covariate in $\mathbf{X}_t(\mathbf{s})$ with respect to time and denote this as $\mathbf{X}_j^*(\mathbf{s}) = \{X_{0l}^*(\mathbf{s}), \dots, X_{pl}^*(\mathbf{s})\}$. Define the covariates for the true process U_l as $\mathbf{X}_{ul}^*(\mathbf{s}) = \mathbf{X}_j^*(\mathbf{s})$, containing all five covariates, and define $\mathbf{X}_{vl}(\mathbf{s}) = \{X_{0l}^*(\mathbf{s}), X_{1l}^*(\mathbf{s}), X_{2l}^*(\mathbf{s})\}$ to include only temperature and relative humidity for bias correction [15]. We model $U_l(\mathbf{s})$ and $V_l(\mathbf{s})$ as independent (with each other and over l) Gaussian processes with means $E\{U_l(\mathbf{s})\} = \mathbf{X}_{ul}^*(\mathbf{s})\beta_u$ and $E\{V_l(\mathbf{s})\} = \mathbf{X}_{vl}(\mathbf{s})\beta_v$, variances $\text{Var}\{U_l(\mathbf{s})\} = \sigma_{ul}^2$ and $\text{Var}\{V_l(\mathbf{s})\} = \sigma_{vl}^2$, and spatial correlations $\text{Cor}\{U_l(\mathbf{s}), U_l(\mathbf{s}')\} = \exp(-\|\mathbf{s} - \mathbf{s}'\|/\rho_u)$ and $\text{Cor}\{V_l(\mathbf{s}), V_l(\mathbf{s}')\} = \exp(-\|\mathbf{s} - \mathbf{s}'\|/\rho_v)$.

The regression coefficients $\beta_u = (\beta_{u0}, \dots, \beta_{up})^T$ control the effects of the covariates on the true PM_{2.5} process U . Although we specify the model in the spectral domain, the DFT is a linear operator

and thus the covariates can be interpreted as usual in the spatial domain since the mean AQS response is

$$E\{Y_{1t}(\mathbf{s})\} = \mathbf{X}_t(\mathbf{s})\boldsymbol{\beta}_u.$$

Therefore, $\boldsymbol{\beta}_u$ is of primary interest. In particular, the components of $\boldsymbol{\beta}_u$ that correspond to the smoke plume indicators are used to summarize the wildland fire contribution to $\text{PM}_{2.5}$.

The regression coefficients $\boldsymbol{\beta}_v = (\beta_{v0}, \beta_{v1}, \beta_{v2})^T$ control the effect of the covariates on the discrepancy term V , and thus the contribution of the covariates to the PA bias. By allowing the covariance parameters σ_{ul}^2 and σ_{vl}^2 to vary by frequency (l), we allow for a different degree of dependence between the networks at different temporal scales, with

$$\text{Cor}\{Z_{1l}^*(\mathbf{s}), Z_{2l}^*(\mathbf{s})\} = \frac{A_l}{\sqrt{A_l^2 + \sigma_{vl}^2 / \sigma_{ul}^2}}. \quad (3)$$

The prior for the variance components is

$$\sigma_{ul}^2 \sim \text{InvGamma}(a_{ul}, b_{ul}) \quad \text{and} \quad \sigma_{vl}^2 \sim \text{InvGamma}(a_{vl}, b_{vl}),$$

where the hyperparameters are modelled as log-linear in frequency, e.g., $\log(a_{ul}) = \gamma_{au1} + \gamma_{au2} \cdot \omega_l$ the prior captures the intuition that the variance is higher in month-to-month variation than day-to-day variation, and the correlation between two sources vary over frequencies.

2.3. Quantifying the wildland fire contribution

To estimate the $\text{PM}_{2.5}$ contribution from wildfire, given the estimated parameters above, we consider two metrics based on either regression or matching. For the regression metric, let $\mathbf{X}_t^0(\mathbf{s})$ be the covariate vector with three plume indicators fixed at zero. For the matching estimator, define $\mathcal{P}(\mathbf{s})$ as the set of days for which site \mathbf{s} is in a smoke plume (any density) and $\bar{\mathcal{P}}(\mathbf{s})$ as the set of non-plume days. We match each plume day with a non-plume day with similar meteorology and time period. Let $\mathcal{A}_t(\mathbf{s}) = \bar{\mathcal{P}}(\mathbf{s}) \cap \{t - 30, \dots, t + 30\}$ be the set of non-plume days within 30 days of plume day t . For each plume day, we selected the matching day $m_t(\mathbf{s})$ as

$$m_t(\mathbf{s}) = \arg \min_{d \in \mathcal{A}_t(\mathbf{s})} |\text{temp}_t(\mathbf{s}) - \text{temp}_d(\mathbf{s})| + w|\text{humidity}_t(\mathbf{s}) - \text{humidity}_d(\mathbf{s})|,$$

where w above is a scaling factor adjusting the magnitude of humidity and temperature, we set $w = 1$ so that the best matching station has equal weights on temperature and humidity. Then at site \mathbf{s} the estimated contribution from wildland fires per day are

1. Regression estimator: $\delta_1(\mathbf{s}) = \frac{1}{n_t} \sum_{t=1}^{n_t} \{\mathbf{X}_t(\mathbf{s}) - \mathbf{X}_t^0(\mathbf{s})\} \boldsymbol{\beta}_u$
2. Matching estimator: $\delta_2(\mathbf{s}) = \frac{1}{n_t} \sum_{t \in \mathcal{P}(\mathbf{s})} \{Z_{1t}(\mathbf{s}) - Z_{1t'}(\mathbf{s})\}$ for $t' = m_t(\mathbf{s})$.

In the matching estimator, $Z_{1t}(\mathbf{s})$ is the true $\text{PM}_{2.5}$, the transformed pairs of $Z_{1l}^*(\mathbf{s})$ in (2) obtained by inverse DFT, and thus this estimator accounts for spatiotemporal bias and correlation. Since the analysis is on the log-scale, we plot $\exp\{\delta_1(\mathbf{s})\}$ and $\exp\{\delta_2(\mathbf{s})\}$ which estimate the multiplicative effect, i.e., $\exp\{\delta_1(\mathbf{s})\} = 1.05$ corresponds to a 5% increase in $\text{PM}_{2.5}$ in the presence of a smoke plume.

2.4. Computational Algorithm

To complete the Bayesian model, we specify uninformative prior distributions for the model parameters. The regression coefficients have Gaussian priors $\boldsymbol{\beta}_u, \boldsymbol{\beta}_v \sim \text{Normal}(\mathbf{0}, c^2 \mathbf{I}_{p+1})$. The variance parameters have conjugate priors $\tau_j^2 \sim \text{InvGamma}(a, b)$. The hyperparameters have Gaussian priors $\gamma_{au1}, \gamma_{au2}, \gamma_{av1}, \gamma_{av2} \sim \text{Normal}(0, c^2)$. To give uninformative priors we set $a = b = 0.01$ and $c = 10$. Due to poor convergence, the dependence parameters β_{A0} and β_{A1} were fixed based on cross-validation to minimize mean squared prediction error for AQS stations.

The main computational bottleneck of spatial modeling is manipulating spatial covariance matrices to estimate the range parameters ρ_u and ρ_v . Given the large size of the air pollution dataset, a reasonable simplification is to estimate the range parameters using variogram and then assume they are fixed for the purpose of fitting the final model. The estimated spatial range from variograms are $\rho_u = 177$ and $\rho_v = 111$ kilometers.

Given the range parameters are fixed, the remaining parameters are estimated using Markov Chain Monte Carlo (MCMC) methods. In particular, we perform Gibbs sampling steps for most parameters and Metropolis sampling for some hyperparameters. Missing values are addressed using Bayesian multiple imputation (refer to Appendix A for more information). We generate 8,000 posterior samples and discard the first 5,000 as burn-in. Convergence is monitored using trace plots. A simulation study is included in the Appendix to verify the algorithm produces reliable parameter estimates. Convergence plots for several representative parameters are shown in Figure A1 in the Appendix for the real data analysis.

3. Results

3.1. Summary of the fitted model

Table 1 gives the estimates of the regression coefficients for both the true process β_u and bias correction term β_v . All three smoke plume levels positively affect PM_{2.5} concentrations, with high smoke plumes having the greatest impact, followed by medium and low smoke plumes. These results are consistent between 2020 and 2021. The bias correction terms, however, are not significant. Given that PurpleAir readings have already been corrected as per [15] using temperature and relative humidity, it is reasonable that these variable do not explain trends in bias. We note that our model does include more general spatiotemporal bias correction in $V_l(s)$ and including this bias term leads to improved results, as discussed below.

Table 1. Posterior mean (95% interval) for the model parameters. The regression coefficients are given separately for the true PM_{2.5} process (β_u) and bias correction (β_v). A “****” indicates that the 95% interval excludes zero.

2020 fire season		
Parameter	True PM _{2.5}	Bias correction
Temperature	0.115 (0.106,0.125)***	-0.002 (-0.009,0.005)
Humidity	0.064 (0.048,0.080)***	0.012 (-0.002,0.035)
Plume – Low	0.007 (0.003,0.011)***	/
Plume – Medium	0.022 (0.012,0.032)***	/
Plume – High	0.049 (0.033,0.065)***	/
2021 fire season		
Parameter	True PM _{2.5}	Bias correction
Temperature	0.006 (0.004,0.008)***	0.006 (-0.003,0.015)
Humidity	0.000 (-0.001,0.001)	-0.011 (-0.026,0.003)
Plume – Low	0.011 (0.001,0.021)***	/
Plume – Medium	0.018 (0.007,0.029)***	/
Plume – High	0.041 (0.031,0.051)***	/

Figure 4 plots the estimated wildland fire contribution both years and both metrics. The estimated wildland fire contribution ranges from a 1-3% increase in PM_{2.5}, depending on the location. Both metrics yield similar estimates of contribution and spatial patterns. The impact of wildfires varies across the state and years. In 2020, both Northern and Central California experienced significant wildfire impacts, while only Northern California faced major effects in 2021. This is in line with the fact that 2020 had the highest frequency of wildfires across all states, whereas 2021 witnessed a single, massive wildfire in Northern California [31].

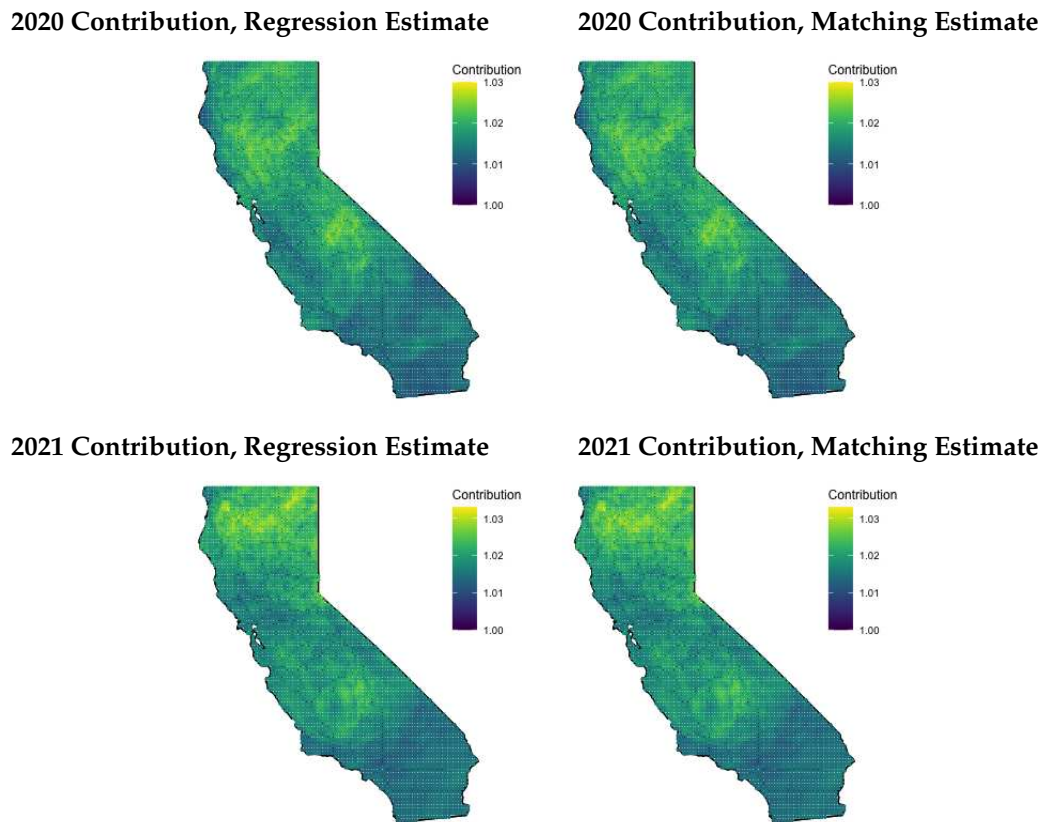


Figure 4. Smoke contribution to PM_{2.5}. Contributions are exponentiated to reflect actual percentage contribution. For example, a contribution of 1.02 means wildfire roughly contributes to 2% increase in PM_{2.5}.

In addition to covariate effects, the data-fusion model provides an evaluation of the concordance between AQS and PA stations. Equation 3 defines the correlation between the two networks as a function of the spectral frequency, ω_l . Figure 5 plots the correlation between AQS and PA by period, i.e., $1/\omega_l$. For example, period 7 (30) corresponds to variation that occurs on a weekly (monthly) scale. Figure 5 shows that the correlation between AQS and PA stations increases from short-term, such as day-to-day variation, to long-term, such as month-to-month variation. In the short-term, the correlation is lower since the readings are taken at different spatial locations and are subject to small scale variability. Over the long run, the correlation is higher as both sources estimate ambient unbiased PM_{2.5} readings.

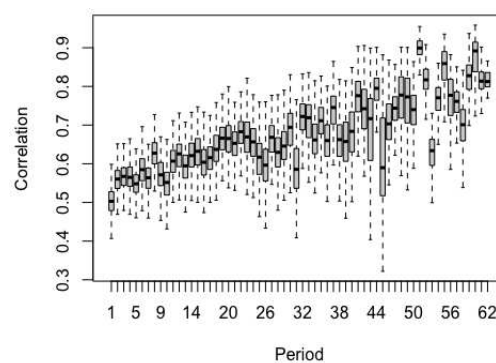


Figure 5. Posterior distribution of the correlation between AQS and PA by period. Small periods capture short-term variation, such as day-to-day variation, while large periods capture long-term variation, such as monthly trends.

3.2. Model Comparisons

To assess the effectiveness of integrating additional PurpleAir readings, we compared the proposed data-fusion model (“Data fusion”) with two simpler alternatives. The first uses only AQS data (“AQS only”) and discards the PA data (i.e., sets $A_l = 0$ for all l). The second naively (“Naive”) combines AQS and PA data and treats them as a single source without spatiotemporal bias adjustment (i.e., sets $A_l = 1$ and $V_l(\mathbf{s}) = 0$ for all \mathbf{s} , and includes an indicator variable in the regression term, β_u , to distinguish two types of data).

The estimated parameters for each model, along with the corresponding posterior standard deviations, are presented in Table 2. Clearly, incorporating PurpleAir monitors significantly reduces the posterior standard deviation. For many of the parameters the reduction in uncertainty is striking, with the standard deviation being 2-4 times smaller for the data-fusion model. Also, with the AQS-only model, only high smoke plumes exhibit a significant contribution due to a higher standard deviation. In contrast, when merging AQS and PurpleAir data, both medium and high smoke plume levels show significant contributions.

Table 2. Posterior mean (standard deviation) for the model parameters β_u for the CA data using the proposed data fusion model, the model that uses only AQS data, and the naive data fusion model that ignores bias in the PA data. A “***” indicates that the 95% interval excludes zero.

2020 fire season			
Parameter	Data fusion	AQS Only	Naive
Temperature	0.115 (0.005)***	0.105 (0.024)***	-0.418 (0.066)***
Humidity	0.064 (0.008)***	0.086 (0.022)***	-1.125 (0.052)***
Plume - Low	0.007 (0.002)***	0.005 (0.012)	0.107 (0.078)
Plume - Medium	0.022 (0.005)***	0.020 (0.014)	0.271 (0.052)***
Plume - High	0.049 (0.008)***	0.042 (0.016)***	0.637 (0.079)***
2021 fire season			
Parameter	Data fusion	AQS Only	Naive
Temperature	0.006 (0.001)***	0.015 (0.003)***	-0.014 (0.006)***
Humidity	0.000 (0.000)	0.008 (0.002)***	-0.039 (0.003)***
Plume - Low	0.011 (0.004)***	-0.001 (0.014)	-0.330 (0.032)***
Plume - Medium	0.018 (0.004)***	0.023 (0.016)	0.230 (0.074)***
Plume - High	0.041 (0.005)***	0.054 (0.017)***	0.980 (0.071)***

Furthermore, to verify that our proposed methodologies not only improve parameter estimation but also lead to accurate PM_{2.5} predictions, we performed a 5-fold cross-validation for the three models. Their performance was compared based on three key metrics: Root mean squared error, 95% prediction coverage, and prediction variance. Detailed descriptions of the models and their results can be found in the Appendix. We find that the AQS-only and data fusion model produce fairly similar out-of-sample prediction accuracy, therefore the main benefit of including the PurpleAir data is reducing uncertainty in parameter estimates. Also, the Naive model gives a 50% larger RMSE and low coverage, emphasizing the need for a careful data fusion approach.

4. Discussion

In this study, we examine the impact of wildland fires on PM_{2.5} concentrations in California during the fire seasons of 2020 and 2021. To do this efficiently, we combine remotely-sensed smoke-plume indicators with AQS and PA measurement networks. To model the spatiotemporal correlation of PM_{2.5} concentration and relationship between AQS and PA monitors, we first transform the data from spatial domain to frequency domain, and then use a data-fusion approach to model spatial correlations while accounting for biases in the PA data. Furthermore, we use a Bayesian approach to compute posterior distributions of the quantities of interest to fully characterize uncertainty.

We find that including PurpleAir monitors significantly increases the precision of the estimated contribution of wildland fire smoke to total $PM_{2.5}$. Using only AQS data we find that medium and high smoke plume levels significantly contribute to $PM_{2.5}$ concentration with standard deviations as large as 0.017, and the data fusion approach that supplements AQS with PA data gives similar parameter estimation, with standard deviation as small as 0.004. Moreover, the data fusion model also estimates a significant low smoke plume level contribution. However, since $PM_{2.5}$ concentration is relatively smooth across space and AQS stations are evenly distributed across the state, incorporating PurpleAir readings does not improve prediction performance even for the data-fusion approach. Comparing prediction performance does reveal that simple data fusion model such as the model that ignores bias in the PA data gives inferior prediction results. With our model, all three smoke plume levels demonstrate a significant contribution to $PM_{2.5}$ concentration, and the impact varies across different regions depending on the year. This study highlights the value of utilizing both AQS and PurpleAir data in understanding the impact of wildfires on air quality and informs future monitoring and management efforts.

There are some limitations of our current work. First, as mentioned above, the satellite-derived smoke plume levels might underestimate the actual smoke level, which may lead to underestimation of wildfires' contribution to $PM_{2.5}$ [30]. Second, due to computational limitations and poor MCMC convergence, we fixed the spatial correlation range parameters for both AQS and PA monitors and parameters that control the relationships between AQS and PA data. The analysis would more fully quantify uncertainty if we are able to implement a fully Bayesian analysis.

We have taken a purely statistical approach to estimating the contribution of wildland fires on ambient air pollution. An area of future work is to incorporate numerical models to simulate the process. Dispersion models, e.g., HySPLIT [32], combine the location and size of fires and meteorological conditions in a mathematical model to track particulate matter emanating from a fire. Of course, numerical models also have bias and other limitations [33], but combining their output within our statistical framework would likely further refine our estimates. Further, instead of using one range parameter for all frequencies, it is possible to get variogram estimates of ranges over frequencies. Also, to extend the current work, we can estimate the contribution over the entire U.S. continent, and investigate how different areas suffer from wildfires, although more computationally efficient models are required to analyze all US data.

Author Contributions: Conceptualization, AR and BR; methodology, HY, YG, BR; validation, HY, SRS; formal analysis, HY; data curation, HY, SRS.; writing-original draft preparation, HY; writing—review and editing, SRS, BR, YG, AR. All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded by the National Institutes of Health (R01ES031651-01) and the National Science Foundation (DMS2152887).

Data Availability Statement: AQS data is a publicly available dataset, which is part of this study. This data can be found on [EPA website](#). Purple Air data is a 3rd party data and restrictions apply to the availability of these data. Data was obtained from Purple Air and are available from [PurpleAir API](#) with the permission of Purple Air. HMS smoke plume data is publicly available and can be downloaded at Office of Satellite and Product Operations [website](#). The codes to download and analyze data in this paper is available at <https://github.com/hyang199723/PAFusion>.

Conflicts of Interest: The authors declare no conflict of interest.

Appendix A. Data Cleaning

Before the data was subjected to a statistical model, we implemented several pre-processing steps on the PurpleAir data and standardized certain covariates to meet the assumption of normality. PurpleAir stations feature two independent channels, Channel A and Channel B, both of which measure ambient $PM_{2.5}$ independently. To achieve a more accurate estimation of the actual ambient $PM_{2.5}$ concentration, we discarded readings where the measurement difference exceeded $200 \mu g/m^3$ and constant high $PM_{2.5}$ readings over $2000 \mu g/m^3$. Subsequently, the mean reading from Channel A and Channel B was considered as the PurpleAir measurement.

Most PurpleAir stations measure temperature (in Fahrenheit) and humidity (as a percentage). Given the spatially smooth changes in temperature and humidity, we employed a 10-nearest-neighbor approach to impute stations with missing temperature and humidity values. Both temperature and humidity were standardized before fitting the model as specified in the analysis.

Appendix B. MCMC algorithm

Assume the n_1 AQS monitors are at spatial locations $\mathbf{s}_1, \dots, \mathbf{s}_{n_1}$ and the n_2 PA monitors are located at $\mathbf{s}_{n_1+1}, \dots, \mathbf{s}_{n_s}$ for $n_s = n_1 + n_2$. The observations can be written as the vectors $\mathbf{Y}_{1t} = [Y_{1t}(\mathbf{s}_1), \dots, Y_{1t}(\mathbf{s}_{n_1})]^T$, $\mathbf{Y}_{2t} = [Y_{2t}(\mathbf{s}_{n_1+1}), \dots, Y_{2t}(\mathbf{s}_{n_s})]^T$ and $\mathbf{Y}_t = (\mathbf{Y}_{1t}^T, \mathbf{Y}_{2t}^T)^T$. Similarly, for frequency l let \mathbf{Y}_{jl}^* , \mathbf{U}_{jl} and \mathbf{V}_{jl} be vectors of length n_j and \mathbf{Y}_l^* , \mathbf{U}_l and \mathbf{V}_l be vectors of length n_s , analogous to \mathbf{Y}_t . The covariate matrices of size $n_j \times p$ are denoted \mathbf{X}_{jl}^* and \mathbf{X}_l^* is the $n_s \times p$ matrix that stacks \mathbf{X}_{1l}^* and \mathbf{X}_{2l}^* . Then the model in the spectral domain is

$$\mathbf{Y}_{1l}^* = \mathbf{U}_l + \mathbf{E}_{1l} \quad \text{and} \quad \mathbf{Y}_{2l}^* = A_l \mathbf{U}_l + \mathbf{V}_{2l} + \mathbf{E}_{2l}$$

where $\mathbf{E}_{jl} \stackrel{\text{indep}}{\sim} \text{Normal}(\mathbf{0}, \tau_j^2 \mathbf{I}_{n_j})$. Using this notation, the spatial models are defined by $E(\mathbf{U}_{jl}) = \mathbf{X}_{jl}^* \boldsymbol{\beta}_u$, $E(\mathbf{V}_{jl}) = \mathbf{X}_{jl}^* \boldsymbol{\beta}_v$, $\text{Cov}(\mathbf{U}_{jl}, \mathbf{U}_{kl}) = \sigma_{ul}^2 \boldsymbol{\Sigma}_{ujk}$ and $\text{Cov}(\mathbf{V}_{jl}, \mathbf{V}_{kl}) = \sigma_{vl}^2 \boldsymbol{\Sigma}_{vjk}$. The full $n_s \times n_s$ spatial correlation matrices are denoted $\boldsymbol{\Sigma}_u$ and $\boldsymbol{\Sigma}_v$.

Each MCMC iteration we impute missing data and update the error variance parameters in the spatial domain, and then update all remaining parameters in the spectral domain. The missing values are simply drawn from the univariate normal distribution

$$Y_{jt} | \text{rest} \sim \text{Normal}(Z_{jt}(\mathbf{s}), \tau_j^2),$$

independently over j and t . The error variances are drawn from full conditional distribution $\tau_1^2 | \text{rest} \sim \text{InvGamma}[n_1 n_t / 2 + a, \sum_{i=1}^{n_1} \sum_{t=1}^{n_t} \{(Y_{1t}(\mathbf{s}_i) - Z_{1t}(\mathbf{s}_i))^2 / 2 + b\}]$ and $\tau_2^2 | \text{rest} \sim \text{InvGamma}[n_2 n_t / 2 + a, \sum_{i=n_1+1}^{n_s} \sum_{t=1}^{n_t} \{(Y_{2t}(\mathbf{s}_i) - Z_{2t}(\mathbf{s}_i))^2 / 2 + b\}]$.

After imputation in the spatial domain, the data are complete and can be projected into the spectral domain where they are independent over time. The spatial processes are updated as

$$\begin{aligned} \mathbf{U}_l | \text{rest} &\sim \text{Normal} \left\{ \boldsymbol{\Omega}_{ul} \left(\mathbf{T} \mathbf{A}_l^1 (\mathbf{Y}_l^* - \mathbf{V}_l) + \frac{1}{\sigma_{ul}^2} \boldsymbol{\Sigma}_u^{-1} \mathbf{X}_l^* \boldsymbol{\beta}_u \right), \boldsymbol{\Omega}_{ul} \right\} \\ \mathbf{V}_{2l} | \text{rest} &\sim \text{Normal} \left\{ \boldsymbol{\Omega}_{vl} \left(\frac{1}{\tau_2^2} (\mathbf{Y}_{2l}^* - A_l \mathbf{U}_{2l}) + \frac{1}{\sigma_{vl}^2} \boldsymbol{\Sigma}_{v22}^{-1} \mathbf{X}_{2l}^* \boldsymbol{\beta}_v \right), \boldsymbol{\Omega}_{vl} \right\} \end{aligned} \quad (\text{A1})$$

where \mathbf{A}_l^k is diagonal with first n_1 elements equal one and the remaining n_2 elements equal A_l^k , \mathbf{T} is diagonal with first n_1 elements equal τ_1^{-2} and the remaining n_2 elements equal τ_2^{-2} , \mathbf{V}_l is the vector with n_1 zeros followed by \mathbf{V}_{2l} , $\boldsymbol{\Omega}_{ul}^{-1} = \mathbf{T} \mathbf{A}_l^2 + \frac{1}{\sigma_{ul}^2} \boldsymbol{\Sigma}_u^{-1}$ and $\boldsymbol{\Omega}_{vl}^{-1} = \frac{1}{\tau_2^2} \mathbf{I}_{n_2} + \frac{1}{\sigma_{vl}^2} \boldsymbol{\Sigma}_{v22}^{-1}$.

The regression coefficients and bias parameters are updated as

$$\begin{aligned} \boldsymbol{\beta}_u | \text{rest} &\sim \text{Normal} \left\{ \mathbf{P}_u \left(\sum_{l=1}^{n_t} \frac{1}{\sigma_{ul}^2} \mathbf{X}_l^{*T} \boldsymbol{\Sigma}_u^{-1} \mathbf{U}_l \right), \mathbf{P}_u \right\} \\ \boldsymbol{\beta}_v | \text{rest} &\sim \text{Normal} \left\{ \mathbf{P}_v \left(\sum_{l=1}^{n_t} \frac{1}{\sigma_{vl}^2} \mathbf{X}_{2l}^{*T} \boldsymbol{\Sigma}_{v22}^{-1} \mathbf{V}_{2l} \right), \mathbf{P}_v \right\} \end{aligned} \quad (\text{A2})$$

where $\mathbf{P}_u^{-1} = \sum_{l=1}^{n_t} \frac{1}{\sigma_{ul}^2} \mathbf{X}_l^{*T} \boldsymbol{\Sigma}_u^{-1} \mathbf{X}_l^* + \frac{1}{c^2} \mathbf{I}_p$ and $\mathbf{P}_v^{-1} = \sum_{l=1}^{n_t} \frac{1}{\sigma_{vl}^2} \mathbf{X}_{2l}^{*T} \boldsymbol{\Sigma}_{v22}^{-1} \mathbf{X}_{2l}^* + \frac{1}{c^2} \mathbf{I}_p$. The remaining hyperparameters are updated as

$$\begin{aligned} \sigma_{ul}^2|_{\text{rest}} &\sim \text{InvGamma}\left(\frac{n_s}{2} + a_{ul}, \frac{(\mathbf{U}_l - \mathbf{X}_l^* \boldsymbol{\beta}_u)^T \boldsymbol{\Sigma}_u^{-1} (\mathbf{U}_l - \mathbf{X}_l^* \boldsymbol{\beta}_u)}{2} + b_{ul}\right) \\ \sigma_{vl}^2|_{\text{rest}} &\sim \text{InvGamma}\left(\frac{n_2}{2} + a_{vl}, \frac{(\mathbf{V}_{2l} - \mathbf{X}_{2l}^* \boldsymbol{\beta}_v)^T \boldsymbol{\Sigma}_{v22}^{-1} (\mathbf{V}_{2l} - \mathbf{X}_{2l}^* \boldsymbol{\beta}_v)}{2} + b_{vl}\right). \end{aligned} \tag{A3}$$

Finally, γ_{au1} , γ_{au2} , γ_{av1} and γ_{av2} are updated using a Metropolis step with Gaussian candidate distribution tuned to give acceptance rate around 0.4.

Appendix C. Simulation Results

We conduct a simulation study to demonstrate the reliability of the MCMC algorithm. The regression parameters, $\boldsymbol{\beta}_u$ and $\boldsymbol{\beta}_v$, are fixed at the mean of the 2021 model output in Table 1. We generate a total number of 80 AQS stations and 500 Purple Air stations with 60 time steps. The spatial locations are randomly sampled from the region $(0, 15)^2$. The data was generated in the frequency domain using the following equations:

$$Y_{1l}(\mathbf{s}) = U_l(\mathbf{s}) + \epsilon_1(\mathbf{s}) \quad \text{and} \quad Y_{2l}(\mathbf{s}) = A_l U_l(\mathbf{s}) + V_{2l}(\mathbf{s}) + \epsilon_2(\mathbf{s}). \tag{A4}$$

The variables U_l and V_l are drawn from Gaussian processes as described in (2). The range parameters are set to $\rho_u = 2$ and $\rho_v = 4$. The error variances of $\epsilon_1(\mathbf{s})$ and $\epsilon_2(\mathbf{s})$ are set to 1.6 and 3.6, respectively. The values of A_l are fixed at the best A_l selected from the real data which is $A_l = 0.2$.

To simulate realistic smoke plume frequencies, we assigned percentages to represent the occurrence of low, medium, and high smoke plume levels. Specifically, 20% of the days corresponded to low smoke plume levels, 15% to medium levels, and 10% to high levels. Temperature and humidity values were randomly generated from standard normal distributions.

The covariates were initially generated in the time domain and then transformed to the frequency domain. The values of σ_{ul} form a decreasing sequence ranging from 50 to 10, with larger values assigned to lower frequencies. Similarly, σ_{vl} follows a decreasing sequence from 40 to 10. Finally, the values of $\boldsymbol{\beta}_u$ and $\boldsymbol{\beta}_v$ are the mean values from Table 1.

We generate 50 datasets from this model. For each simulated dataset, we fit the model with ρ_u , ρ_v and A_l fixed at the true values and generate 8000 MCMC iterations and discard the first 5000 as burn-in. Since our main interest is in the covariate effects, for each dataset we record the effective sample size of the MCMC algorithm [34] and the posterior mean estimator and 95% posterior interval.

For each dataset and each parameter, we compute the posterior mean, standard deviation and 95% interval and measure MCMC convergence using the effective sample size. The average of the posterior means, standard deviations and effective samples sizes, and the empirical coverage of 95% intervals are shown in Table A1. The posterior means show small bias, the coverage is near the nominal level and the effective sample size coefficients indicate reasonable convergence.

Table A1. True value used for the fixed effects for the true PM_{2.5} ($\boldsymbol{\beta}_u$) and bias ($\boldsymbol{\beta}_v$) to simulate data and the average (SD) over the 50 datasets of the posterior mean estimators ("Ave post mean"), coverage of 95% posterior intervals and average (SD) effective sample size based on 3000 MCMC iterations.

Type	Covariate	True value	Average post mean	Coverage	ESS
PM _{2.5}	Temperature	0.118	0.117 (0.013)	100%	420.23 (0.14)
	Humidity	0.064	0.069 (0.022)	96%	307.27 (0.10)
	Plume - Low	0.007	0.006 (0.132)	100%	875.99 (0.29)
	Plume - Medium	0.022	0.020 (0.037)	98%	376.91 (0.13)
	Plume - High	0.049	0.050 (0.176)	100%	480.22 (0.16)
Bias	Temperature	-0.002	0.003 (0.019)	92%	168.75 (0.06)
	Humidity	0.012	0.009 (0.041)	96%	176.97 (0.06)

Appendix D. Cross-Validation Results

We compare methods using a 5-fold cross-validation using data from 2021. We randomly split the AQS stations into five folds. For each fold, we build predictive models based on the other AQS stations and all PA stations and make predictions at the test sites. We compare the proposed data-fusion model (“Data fusion”) with the AQS only and Naive models described in the main text. For all models, we fix the spatial range parameters (ρ_u and ρ_v) based on the variogram analysis of the full dataset. The cross-dependence parameter A_l is fixed at 0.2.

The results in Table A2 show that the performance of the AQS-only analysis is fairly similar to the proposed data-fusion approach, with slightly smaller prediction mean squared error and larger average prediction variance. Therefore, carefully including the additional PA data mainly reduces the prediction variance. However, naively including the PA data gives much higher prediction errors and low coverage.

Table A2. Root mean squared error (“RMSE”), coverage of 95% prediction intervals (“Coverage”) and average prediction variance (“Ave Var”) for the cross-validation study comparing the proposed data fusion model to models that ignore PA data (“AQS only”) and includes PA data without bias correction (“Naive”).

Model	RMSE	Coverage	Ave Var
Data Fusion	0.42	0.89	0.13
AQS only	0.40	0.91	0.16
Naive	0.66	0.73	0.18

Appendix E. Real Data Convergence

We display several representative trace plots of the data fusion model to verify the convergence of our MCMC algorithm for the 2021 CA analysis. After burn-in, the MCMC chains appear to have converged.

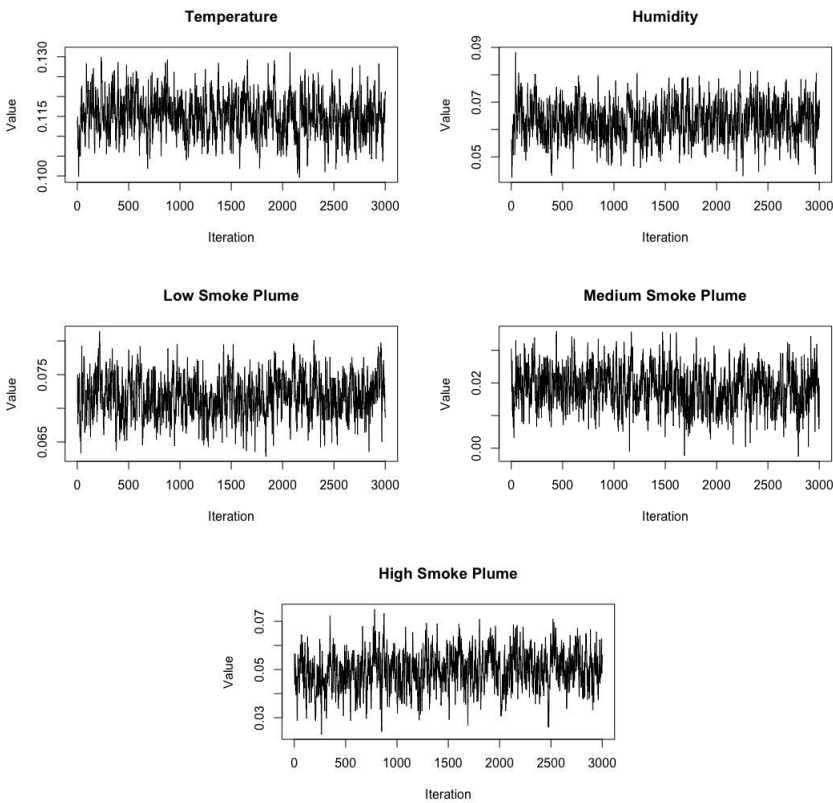


Figure A1. Trace plots of parameters of interest (β_u) for the 2021 California data analysis.

References

1. Dennekamp, M.; Abramson, M.J. The effects of bushfire smoke on respiratory health. *Respirology* **2011**, *16*, 198–209.
2. Dennekamp, M.; Straney, L.D.; Erbas, B.; Abramson, M.J.; Keywood, M.; Smith, K.; Sim, M.R.; Glass, D.C.; Del Monaco, A.; Haikerwal, A.; others. Forest fire smoke exposures and out-of-hospital cardiac arrests in Melbourne, Australia: a case-crossover study. *Environmental Health Perspectives* **2015**, *123*, 959–964.
3. Melnick, R.S. *Regulation and the courts: The case of the Clean Air Act*; Brookings Institution Press, 2010.
4. Sager, L.; Singer, G. Clean identification? The effects of the Clean Air Act on air pollution, exposure disparities and house prices **2022**.
5. McClure, C.D.; Jaffe, D.A. US particulate matter air quality improves except in wildfire-prone areas. *Proceedings of the National Academy of Sciences* **2018**, *115*, 7901–7906.
6. Johnston, F.H.; Henderson, S.B.; Chen, Y.; Randerson, J.T.; Marlier, M.; DeFries, R.S.; Kinney, P.; Bowman, D.M.; Brauer, M. Estimated global mortality attributable to smoke from landscape fires. *Environmental Health Perspectives* **2012**, *120*, 695–701.
7. Rappold, A.G.; Stone, S.L.; Cascio, W.E.; Neas, L.M.; Kilaru, V.J.; Carraway, M.S.; Szykman, J.J.; Ising, A.; Cleve, W.E.; Meredith, J.T.; others. Peat bog wildfire smoke exposure in rural North Carolina is associated with cardiopulmonary emergency department visits assessed through syndromic surveillance. *Environmental Health Perspectives* **2011**, *119*, 1415–1420.
8. Haikerwal, A.; Akram, M.; Sim, M.R.; Meyer, M.; Abramson, M.J.; Dennekamp, M. Fine particulate matter (PM 2.5) exposure during a prolonged wildfire period and emergency department visits for asthma. *Respirology* **2016**, *21*, 88–94.
9. Thilakarathne, R.; Hoshiko, S.; Rosenberg, A.; Hayashi, T.; Buckman, J.R.; Rappold, A.G. Wildfires and the changing landscape of air pollution-related health burden in California. *American Journal of Respiratory and Critical Care Medicine* **2023**, *207*, 887–898.
10. Larsen, A.E.; Reich, B.J.; Ruminski, M.; Rappold, A.G. Impacts of fire smoke plumes on regional air quality, 2006–2013. *Journal of Exposure Science & Environmental Epidemiology* **2018**, *28*, 319–327.
11. Matz, C.J.; Egyed, M.; Xi, G.; Racine, J.; Pavlovic, R.; Rittmaster, R.; Henderson, S.B.; Stieb, D.M. Health impact analysis of PM_{2.5} from wildfire smoke in Canada (2013–2015, 2017–2018). *Science of The Total Environment* **2020**, *725*, 138506.
12. Barkjohn, K.K.; Gantt, B.; Clements, A.L. Development and application of a United States-wide correction for PM 2.5 data collected with the PurpleAir sensor. *Atmospheric Measurement Techniques* **2021**, *14*, 4617–4637.
13. Tryner, J.; L'Orange, C.; Mehaffy, J.; Miller-Lionberg, D.; Hofstetter, J.C.; Wilson, A.; Volckens, J. Laboratory evaluation of low-cost PurpleAir PM monitors and in-field correction using co-located portable filter samplers. *Atmospheric Environment* **2020**, *220*, 117067.
14. Wallace, L.; Bi, J.; Ott, W.R.; Sarnat, J.; Liu, Y. Calibration of low-cost PurpleAir outdoor monitors using an improved method of calculating PM_{2.5}. *Atmospheric Environment* **2021**, *256*, 118432.
15. Barkjohn, K.; Gantt, B.; Clements, A. Development and Application of a United States wide correction for PM_{2.5} data collected with the PurpleAir sensor. *Atmos. Meas. Tech. Discuss.* <https://doi.org/10.5194/amt-2020-413> **2020**.
16. Holder, A.L.; Mebust, A.K.; Maghran, L.A.; McGown, M.R.; Stewart, K.E.; Vallano, D.M.; Elleman, R.A.; Baker, K.R. Field evaluation of low-cost particulate matter sensors for measuring wildfire smoke. *Sensors* **2020**, *20*. doi:10.3390/s20174796.
17. Kosmopoulos, G.; Salamalikis, V.; Pandis, S.; Yannopoulos, P.; Bloutsos, A.; Kazantzidis, A. Low-cost sensors for measuring airborne particulate matter: Field evaluation and calibration at a South-Eastern European site. *Science of The Total Environment* **2020**, *748*, 141396.
18. Durrant-Whyte, H.; Henderson, T.C. Multisensor data fusion. *Springer Handbook of Robotics* **2016**, pp. 867–896.
19. Luo, R.C.; Kay, M.G. A tutorial on multisensor integration and fusion. IECON'90: 16th Annual Conference of IEEE Industrial Electronics Society. IEEE, 1990, pp. 707–722.
20. Reich, B.J.; Chang, H.H.; Foley, K.M. A spectral method for spatial downscaling. *Biometrics* **2014**, *70*, 932–942.

21. Warren, J.L.; Miranda, M.L.; Tootoo, J.L.; Osgood, C.E.; Bell, M.L. Spatial distributed lag data fusion for estimating ambient air pollution. *The Annals of Applied Statistics* **2021**, *15*, 323.
22. Friberg, M.D.; Zhai, X.; Holmes, H.A.; Chang, H.H.; Strickland, M.J.; Sarnat, S.E.; Tolbert, P.E.; Russell, A.G.; Mulholland, J.A. Method for fusing observational data and chemical transport model simulations to estimate spatiotemporally resolved ambient air pollution. *Environmental Science & Technology* **2016**, *50*, 3695–3705.
23. Friberg, M.D.; Kahn, R.A.; Holmes, H.A.; Chang, H.H.; Sarnat, S.E.; Tolbert, P.E.; Russell, A.G.; Mulholland, J.A. Daily ambient air pollution metrics for five cities: Evaluation of data-fusion-based estimates and uncertainties. *Atmospheric Environment* **2017**, *158*, 36–50.
24. Nguyen, H.; Cressie, N.; Braverman, A. Spatial statistical data fusion for remote sensing applications. *Journal of the American Statistical Association* **2012**, *107*, 1004–1018.
25. Gressent, A.; Malherbe, L.; Colette, A.; Rollin, H.; Scimia, R. Data fusion for air quality mapping using low-cost sensor observations: Feasibility and added-value. *Environment International* **2020**, *143*, 105965.
26. Datta, A.; Saha, A.; Zamora, M.L.; Buehler, C.; Hao, L.; Xiong, F.; Gentner, D.R.; Koehler, K. Statistical field calibration of a low-cost PM_{2.5} monitoring network in Baltimore. *Atmospheric Environment* **2020**, *242*, 117761.
27. Lin, Y.C.; Chi, W.J.; Lin, Y.Q. The improvement of spatial-temporal resolution of PM_{2.5} estimation based on micro-air quality sensors by using data fusion technique. *Environment International* **2020**, *134*, 105305.
28. Stein, M.L. Statistical methods for regular monitoring data. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **2005**, *67*, 667–687.
29. National Oceanic and Atmospheric Administration. Hazard Mapping System Fire and Smoke Product. Available at <https://www.ospo.noaa.gov/Products/land/hms.html> (10/15/2022).
30. O'Dell, K.; Ford, B.; Fischer, E.V.; Pierce, J.R. Contribution of wildland-fire smoke to US PM_{2.5} and its influence on recent trends. *Environmental Science & Technology* **2019**, *53*, 1797–1804.
31. California Department of Forestry and Fire Protection. Top 20 Largest California Wildfires. Available at <https://www.fire.ca.gov/our-impact/statistics>.
32. Draxler, R.; Rolph, G. HYSPLIT (HYbrid Single-Particle Lagrangian Integrated Trajectory) model access via NOAA ARL READY website (<http://ready.arl.noaa.gov/HYSPLIT.php>), NOAA Air Resources Laboratory. *Silver Spring, MD* **2010**, 25.
33. Su, L.; Yuan, Z.; Fung, J.C.; Lau, A.K. A comparison of HYSPLIT backward trajectories generated from two GDAS datasets. *Science of the Total Environment* **2015**, *506*, 527–537.
34. Geyer, C.J. Introduction to Markov Chain Monte Carlo. *Handbook of Markov Chain Monte Carlo* **2011**, 20116022, 45.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.