

# Extensible database of validated biomass smoke events for health research

Ivan C. Hanigan<sup>\* 1</sup> Fay H. Johnston<sup>2</sup> Geoffrey G. Morgan<sup>3</sup> Grant J. Williamson<sup>2</sup> Farhad Salimi<sup>2</sup> Sarah B. Henderson<sup>4</sup>

<sup>1</sup>*National Centre for Epidemiology and Population Health, Australian National University*

<sup>2</sup>*Menzies School of Population Health, University of Tasmania*

<sup>3</sup>*University Centre for Rural Health, University of Sydney*

<sup>4</sup>*School of Population and Public Health, University of British Columbia*

<sup>\*</sup> Corresponding author: [ivan.hanigan@anu.edu.au](mailto:ivan.hanigan@anu.edu.au)

2016-03-01

## Abstract

**Background:** Epidemiological studies of the health effects of biomass smoke events (such as bushfires or wood-heater smoke spikes due to inversion layers) have been hindered by the lack of available datasets that explicitly list the locations and dates of pollution events from these sources. Extreme air pollution events may also be caused by dust storms, fossil fuel induced smog events or factory fires, and so validation is necessary to ensure the events are from biomass sources. This paper presents an extensible database developed by the authors to identify historical spikes in air pollution and to evaluate whether they were caused by vegetation fire smoke or by other possible sources. The ability for this database to be extended by other researchers means that new events can be added, and new information for already identified events can be described. These methods provide a systematic framework for retrospective identification of the air quality impacts of biomass smoke. In this paper, we describe the database and data acquisition methods, as well as analytical considerations when validating historical events using a range of reference types.

**Methods:** Several major urban centers and smaller regional towns in the Australian states of New South Wales, Western Australia, and Tasmania were selected as they are intermittently affected by extreme episodes of vegetation fire smoke. Air pollution data was collated and missing values were

32 imputed. Extreme values were identified and a range of sources of reference information were assessed  
33 for each date. Reference types included online newspaper archives, government and research agency  
34 records, satellite imagery and a Dust Storms database.

35 **Results:** This dataset contains validated events of extreme biomass smoke pollution across Australian  
36 cities. The authors have previously demonstrated the utility of this database in analyses of hospital  
37 admissions and mortality data for these locations to quantify the pollution-related health effects of  
38 these events.

39 **Conclusions:** The database was created using open source software and this makes the prospect for  
40 future extensions to the database possible. This is because if other scientists notice an omission or  
41 error in these data they can offer an amendment. We believe that this will improve the database and  
42 benefit the whole biomass smoke health research community.

## 43 **Epidemiological studies of outdoor air pollution**

44 For decades, researchers have studied the public health impacts of ambient outdoor air pollution,  
45 particularly from the effects of particulate and gaseous pollutants, especially associated with the  
46 combustion of coal, petroleum and biomass used for cooking (Pope & Dockery 2006). Far fewer studies  
47 have examined the effect of intermittent smoke from biomass burning, such as that which occurs in  
48 bushfires, or from woodsmoke trapped by inversion layers during winter months as wood is burned for  
49 heating (Naeher *et al.* 2007).

50 There is a gap in the epidemiological literature of health effects from ambient outdoor air pollution  
51 relating to smoke from biomass burning such as that from bushfires or woodsmoke from heating. Most  
52 literature available that focuses on biomass smoke health impacts looks at indoor pollution from cooking  
53 (Smith 1993). Particles (and perhaps noxious gases) in outdoor pollution from biomass smoke might  
54 directly influence the respiratory system through their inhalation and lodgement in the lungs. Particles  
55 may then affect the cardiovascular system after their entry into the circulatory system from the alveolae.  
56 Indirect effects on mental health and wellbeing are also plausible.

57 Epidemiological studies that investigate the relationship between health and air pollution exposures  
58 have primarily used time-series methods that study variations of some health outcomes such as deaths  
59 or hospitalisations from specific disease groups (Peng & Dominici 2008). These outcomes are usually  
60 monitored by day across whole cities, and relationships with atmospheric variables estimated in  
61 regression models. These typically focus on daily levels of ambient air pollution measured by a network  
62 of monitoring sites scattered across a city, time matched to the health outcomes on the same day or a  
63 few days after. In general biomass smoke forms only a small part of the mixture of pollutants in the  
64 air, however when a bushfire or inversion layer event occurs there is often a concomitant spike in the  
65 pollution levels primarily composed of biomass smoke. There is then the ability to study statistical  
66 associations between these pollution spikes and the health outcomes around those days. Anomalous  
67 levels of pollution can be arbitrarily defined using a threshold such as the 95th percentile and these  
68 might be assumed to be biomass smoke days, however there are other events that might cause such as  
69 spike such as dust storms, factory fires or even sea salt being driven by certain wind events. There is a  
70 need then to validate the dates on which events are ascribed in any correlational study of pollution  
71 spikes and health that claims the high levels are due to biomass smoke.

## 72 **The development of this biomass smoke events database**

73 This open and extensible database was developed by the authors to identify historical spikes in particulate  
74 matter concentrations and to evaluate whether they were caused by vegetation fire smoke or by other  
75 means. A summary of the protocol for developing this database and a summary of the data we collated  
76 is published already as a descriptive paper (Johnston *et al.* 2011). This paper describes how the  
77 database has been extended to be able to be distributed in an open, extensible format that allows the  
78 research community to add to the history of these events.

## 79 **System design**

80 The system is described in Figure 1. The procedure starts with the online database and web interface  
81 that is maintained by the Data Manager (DM) in our group. The DM extracts a snapshot of the database

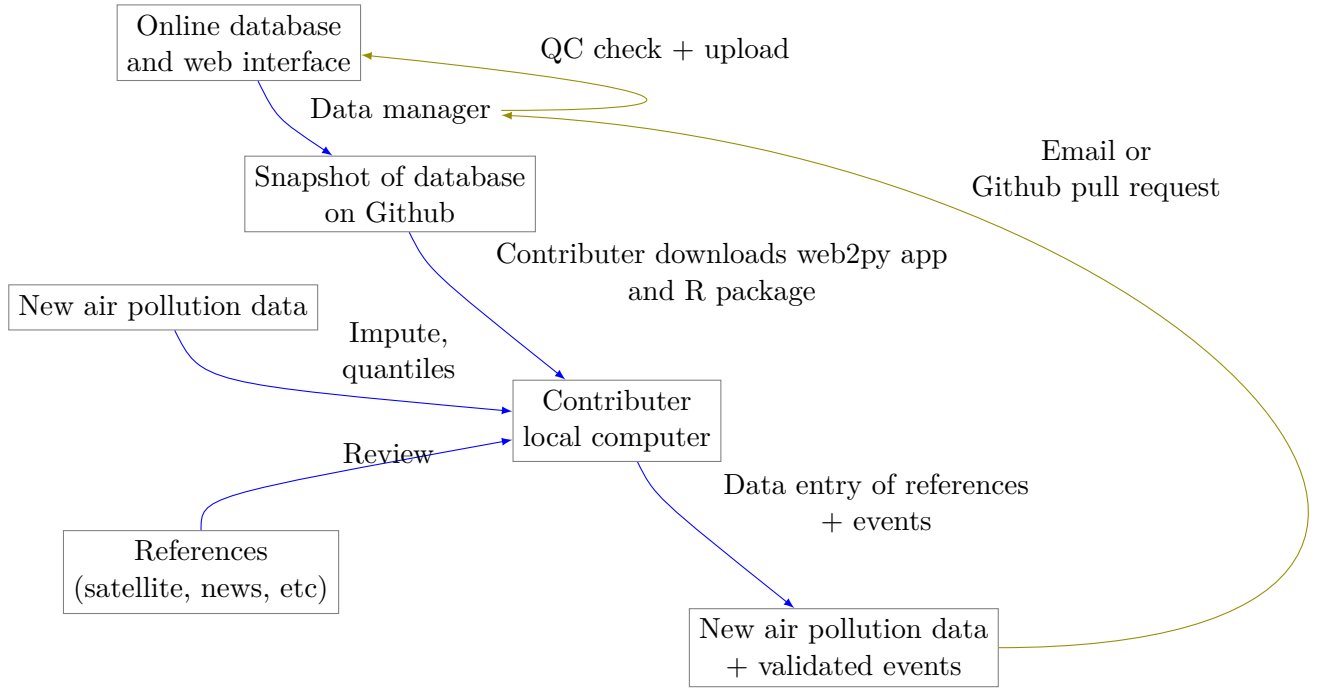


Figure 1: Schematic diagram of the online database and offline processes for extending the database

82 (with a specific version identifier from the Git version control system) and makes a ‘standalone’ version  
 83 available on Github. This standalone version uses web2py so that it is capable of being downloaded and  
 84 run on any operating system used by other computers. Contributors may download that version and  
 85 use it as a local database. The R package is also available on Github, and contains functions that may  
 86 be used to impute any missing data gaps using the APHEA procedure (Katsouyanni *et al.* 1996) as per  
 87 the study protocol. The contributor needs to have new air pollution data available, and access to the  
 88 required reference materials for validation. The R package is used to compute the quantiles of the new  
 89 extended time-series of imputed pollution data, to identify events above the 95th percentile threshold  
 90 that has been set to define ‘extreme events’. The contributor uses the web2py data entry forms to add  
 91 the information that is used to meet the validation criteria. Once they complete their review of all  
 92 events they notify the DM either with email or by using the Github ‘pull request’ feature. The DM  
 93 performs Quality Control (QC) checks and then uploads the new data to the online database. The  
 94 procedure then starts again and a new version is loaded into the Github repository with descriptions of  
 95 the additional changes that have been incorporated.

## 96 General overview of protocols

97 For each location in the original study there were up to 13 years (between 1994 and 2007) of daily air  
98 quality data measured as Particulate Matter (PM) less than 10 $\mu$ m ( $PM_{10}$ ) or less than 2.5  $\mu$ m ( $PM_{2.5}$ )  
99 in aerodynamic diameter were examined. Air pollution data were provided by government agencies in  
100 the states of Western Australia, New South Wales, and Tasmania. Daily averages for each site were  
101 calculated excluding days with less than 75% of hourly measurements. In Sydney and Perth, where  
102 data were collected from several monitoring stations, the missing daily site-specific PM concentrations  
103 were imputed using available data from other proximate monitoring sites in the network. The daily  
104 city-wide PM concentrations were then estimated following the protocol of the Air Pollution and Health:  
105 a European Approach studies (Atkinson *et al.* 2001).

106 First a ‘filling-in’ procedure was used to improve data completeness. It entailed the substitution of the  
107 missing daily values with a weighted average, using the weights of the missing sites 3-month average  
108 proportional to the network average. The weights are calculated against the values from the rest of the  
109 monitoring stations. The pollutant measures from all stations providing data were then averaged to  
110 provide single, city-wide estimates of the daily levels of the pollutants

111 For each city, all days in which PM<sub>10</sub> or PM<sub>2.5</sub> exceeded the 95th percentile were identified over the  
112 entire time series. These extreme values were termed ‘events’. A range of sources was ex- amined  
113 to identify the cause of particulate air pollution events, including electronic news archives, Internet  
114 searches for other reports, government and research agencies, satellite imagery and a Dust Storms  
115 database. Also examined were remotely sensed aerosol optical thickness (AOT) data to provide further  
116 information about days for which the other methods did not.

## 117 Detailed data preparation and validation methods

### 118 Step 1: Source air pollution data

119 Step 1.0 Source air pollution data. Both time series observations and spatial data regarding site  
120 locations.

121 Step 1.1. NSW data downloaded from an online data server. Site locations (Lat and Long) obtained  
122 from website.

123 Step 1.2. WA data sent on CD from contacts at the WA Government Department, these were hourly  
124 data as provided. Cleaned so as only days with >75% of hours are used. Licence puts restrictions on our  
125 right to provide to a third party. Therefore those observed and imputed data are not included, only the  
126 events.

127 Step 1.3. Tasmanian data sent via email from contact at the Department, these were daily data.

128 Step 1.4. All data combined and Quality Control checked in the PostGIS database.

### 129 Step 2. Define spatial extent for cities

130 The cities and towns were selected based on the aims of the health study to investigate Cardio-respiratory  
131 disease and air pollution from biomass smoke events. These were Albany, Albury, Armidale, Bathurst,  
132 Bunbury, Busselton, Geraldton, Gosford-Wyong, Hobart, Illawarra, Launceston, Newcastle, Perth,  
133 Sydney, Tamworth and Wagga Wagga.

134 The spatial extent of each city and town was devised by intersecting Australian Bureau of Statistics  
135 Statistical Local Areas (SLAs) from the various Census editions. These boundaries were set so give the  
136 best possible representation of hospital admissions from the population.

137 Air pollution monitoring sites were then selected on the basis of their proximity to these populations.

### 138 **Step 3. Imputation to fill in gaps in the time-series and calculate a network average**

139 In cities where data were collected from several monitoring stations, the missing daily site-specific PM  
140 concentrations were imputed using available data from other proximate monitoring sites in the network.  
141 The daily city-wide PM concentrations were then estimated following the protocol of the Air Pollution  
142 and Health: a European Approach studies (Katsouyanni *et al.* 1996).

143 Step 3.1. Prepare Data. First it was necessary to find the minimum date that the series of continuous  
144 observations can be considered to start. In the Australian datasets the initial observations could not be  
145 used because they were sometimes only one day per week, only during a particular season or of poor  
146 quality due to teething problems with equipment and procedures. Then it was necessary to identify  
147 missing dates. Get a list of the sites to include – that is with more than 70% observed over the time  
148 period (as defined after assessing min and max dates of period).

149 Step 3.2. Loop over each station individually and calculate a daily network average of all the other  
150 non-missing sites (ie an average of all stations except the focal station of that iteration in the loop).

151 Step 3.3. Calculate three monthly seasonal mean of these non-missing stations. Calculate a three-month  
152 seasonal mean for MISSING site. Estimate missing days at missing sites. The missing value was  
153 replaced by the mean level of the remaining stations, multiplied by a factor equal to the ratio of the  
154 seasonal (centred three month) mean for the missing station, over the corresponding mean from the  
155 stations available on that particular day.

156 Step 3.4. Join all sites for city wide averages and fill any missing days at the site-level with average of  
157 the days immediately before and after the missing days (only when this is below a threshold).

158 Step 3.5 Take the average of all sites per day for city wide averages.

159 Step 3.6. Fill any missing days at the city-wide level with the average of before and after (if this is less  
160 than 5% of days).

## 161 **Step 4. Validate events and identify the causes**

162 Select any events with PM10 or PM2.5 greater than 95 percentile. Manually validate events using online  
163 newspaper archives, government and research agency records, satellite imagery and other sources (such  
164 as a Dust Storm database). Enter the information for each event into the custom built data entry forms.  
165 For any events with references for multiple types of source, assess the likelihood of any single source  
166 being the dominant source. Double check any remaining 99th percentile dates with no references.

## 167 **Step 5. Insert contributed pollution and validated events, and downstream dissem-** 168 **ination**

- 169 • To close the loop the data are then inserted back into the DB.

## 170 **Availability and requirements**

171 Lists the following:

- 172 • Project name: BiosmokeValidatedEvents
- 173 • Project home page: <http://swish-climate-impact-assessment.github.io/BiosmokeValidatedEvents/>
- 174 • Operating system(s): R package is platform independent. Data Entry forms are Microsoft  
175 Windows.
- 176 • Programming language: R and SQL
- 177 • Other requirements: PostgreSQL (PostGIS is desirable)
- 178 • License: CC BY 4.0
- 179 • Any restrictions to use: amendments of errors of omission or commission are invited but will be  
180 vetted before insertion into the master database.



## 181 Availability of supporting data

## 182 Air pollution data provided

183 The NSW Air pollution data are available to download from <http://www.environment.nsw.gov.au/AQMS/search.htm>

## 185 Data derived

186 The data set supporting the results of this article are available in the repository from the website  
187 [http://swish-climate-impact-assessment.github.io/biomass\\_smoke\\_events\\_db](http://swish-climate-impact-assessment.github.io/biomass_smoke_events_db)

188 We have applied the license under Creative Commons - Attribution 4.0. This allows others to copy,  
189 distribute and create derivative works provided that they credit the original source.

190 Users should cite the Johnston 2011 Journal of the Air & Waste Management Association as the validation  
191 protocol and the Database itself as: Hanigan, IC., Johnston, FH., Morgan, GG., and contributors[\*].  
192 (2015). The Validated Bushfire Smoke Events Database. [https://gislibrary-extreme-weather.anu.edu.](https://gislibrary-extreme-weather.anu.edu.au/biomass_smoke_events)  
193 [au/biomass\\_smoke\\_events](https://gislibrary-extreme-weather.anu.edu.au/biomass_smoke_events)

## 194 References

195 Atkinson, R.W., Anderson, R.H., Sunyer, J., Ayres, J., Baccini, M., Vonk, J.M., Boumghar, A.,  
196 Forastiere, F., Forsberg, B., Touloumi, G., Schwartz, J. & Katsouyanni, K. (2001). Acute Effects of  
197 Particulate Air Pollution on Respiratory Admissions. *American Journal of Respiratory and Critical*  
198 *Care Medicine*, 164(10),: 1860–1866.

199 Johnston, F., Hanigan, I., Henderson, S., Morgan, G., Portner, T., Williamson, G. & Bowman, D.  
200 (2011). Creating an integrated historical record of extreme particulate air pollution events in Australian  
201 cities from 1994 to 2007. *Journal of the Air Waste Management Association*, 61(4),: 390–398.

202 Katsouyanni, K., Schwartz, J., Spix, C., Touloumi, G., Zmirou, D., Zanobetti, A., Wojtyniak, B., Vonk,  
203 J.M., Tobias, A., Ponka, A., Medina, S., Bacharova, L. & Anderson, H.R. (1996). Short term effects

204 of air pollution on health: a European approach using epidemiologic time series data: the APHEA  
 205 protocol. *Journal of Epidemiology & Community Health*, 50(Suppl 1),: S12–S18.

206 Naeher, L.P., Brauer, M., Lipsett, M., Zelikoff, J.T., Simpson, C.D., Koenig, J.Q. & Smith, K.R. (2007).  
 207 Woodsmoke health effects: A review. *Inhalation Toxicology*, 19(1),: 67–106.

208 Peng, R.D. & Dominici, F. (2008). *Statistical methods for environmental epidemiology with R. A case*  
 209 *study in air pollution and health*. Springer Science & Business Media, New York, USA.

210 Pope, C.A. & Dockery, D.W. (2006). Health Effects of Fine Particulate Air Pollution: Lines that  
 211 Connect. *Journal of the Air & Waste Management Association*, 56(6),: 709–742.

212 Smith, K. (1993). Fuel combustion, air pollution exposure, and health: The situation in developing  
 213 countries. *Annual Review of Energy and the Environment*, 18: 529–66.