



A tuneable framework for outlier detection in PM_{2.5} air sensor networks during wildland fire smoke events

Stuart J. Illson¹, Karoline K. Barkjohn²

¹School of Environmental and Forestry Sciences, University of Washington, Seattle, 98195, United States of America

5 ²Office of State Air Partnerships, U.S. Environmental Protection Agency, Durham, NC 27711, United States of America

Correspondence to: Stuart Illson (illson@uw.edu)

Abstract. In recent years the use of air sensors has rapidly expanded across North America to measure fine particulate matter (PM_{2.5}), particularly in response to increasing air quality impacts from wildland fire. With the benefit of enhanced spatial and temporal coverage, the scientific community and the public have come to rely on sensor networks as valuable sources of air quality information. With an increasing variety of sensor devices being deployed, there is a need to validate and harmonize PM_{2.5} data between different device types. While significant attention has been given to calibration and correction equations to improve the accuracy of a given sensor's measurement, there is a need to develop tractable and generalizable methods of identifying malfunctioning or unreliable sensors, given the maintenance, siting, and operation of many of these devices is unknown. In this paper, we propose a method of identifying outlier PM_{2.5} sensors, defined as those whose measurements deviate strongly from other local measurements due to hardware faults or to hyper-local environmental conditions that are not representative of typical ambient air quality conditions. While detecting outliers during typical conditions is a fairly straightforward task, detecting outliers during smoke events is challenging due to real, erratic shifts in PM_{2.5} concentrations. Here, we present a novel method of detecting outliers within sensor networks by combining measures from information theory and machine learning. We first define a tuneable, rule-based detection function that balances the Shannon entropy of a local network against the information content of an individual sensor's measurement. We then use this function, together with additional information-theoretic and short-term temporal features, to train a gradient-boosted decision tree for automated outlier detection. Hourly PM_{2.5} measurements from various device types were collected for 11 unique smoke events across North America in 2024 and 2025, and a stratified sample of sensor data were randomly perturbed to simulate 5 commonly seen faults. In each of these cases, we assessed each method's ability to detect the simulated faults. We demonstrate that either of these methods, while trained on a semi-synthetic dataset, can act as a useful data validation procedure when applied to both real-time air quality reporting and retrospective analysis.

1 Introduction

Smoke from wildland fire affects millions in North America each year, contributing to thousands of deaths and a multitude of adverse health outcomes (Reid et al., 2016; Gould et al., 2024). The fine particulate matter (PM_{2.5}, particles with aerodynamic



30 diameters smaller than 2.5 μm) released from wildland fire smoke is of particular concern, as exposure is strongly associated
with a number of negative health effects (Schwartz et al., 1996; Pope et al., 2002; Brook et al., 2010). Populations in North
America are being exposed to substantially higher levels of smoke-related $\text{PM}_{2.5}$ now than they were two decades ago, with
wildfire smoke contributing up to one-quarter of annual-average $\text{PM}_{2.5}$ exposure across the United States and erasing about
25% of the air-quality gains made since 2000 (Burke et al., 2021; Burke et al., 2023). Recent projections indicate that increases
35 in wildfire smoke could result in tens of thousands of additional deaths annually in the United States by mid-century, with
associated economic damages on the order of \$600 billion USD annually (Qiu et al., 2025). Adequately evaluating these health
impacts requires estimation of population-level exposure to $\text{PM}_{2.5}$ at a fine spatiotemporal resolution using ground-based air
quality observations.

While many parts of North America, particularly large population centres, are covered by regulatory $\text{PM}_{2.5}$ monitoring
40 networks, these stations are sparse in smoke-prone regions. This limits the ability to accurately assess exposure levels in the
communities most frequently affected (Larkin, 2019). Additionally, these observational gaps make it difficult to resolve the
full extent and progression of smoke events in space and time, as smoke events are episodic and highly variable in duration,
magnitude, and spatial extent. With increasing frequency of smoke events, a large number of publicly owned, outdoor $\text{PM}_{2.5}$
sensors have been deployed across North America in recent years, with over 18,000 on the AirNow Fire and Smoke Map in
45 2025 (<https://fire.airnow.gov>, last accessed 7 April, 2025). This represents roughly a fourfold increase since 2019, with ~4,000
sensors accounted for in the continental United States in 2019 (Jaffe et al., 2020). These sensors are commonly referred to as
“low-cost”, as they typically cost a fraction of their regulatory-grade counterparts, contributing to their widespread deployment.
Sensor networks have shown promise in increasing the resolution and accuracy of population-level exposure models, furthering
understanding of the distribution of air quality impacts (Clements et al., 2017; Morawska et al., 2018; Bi et al., 2020; deSouza
50 et al., 2020; Keyes et al., 2023; Raysoni et al., 2023).

While expanded monitoring across the United States has furthered our ability to assess wildfire smoke exposure,
lower-cost sensors come with their own set of challenges. Sensors are being developed by a variety of manufacturers, each
with unique hardware designs, firmware configurations, and proprietary calibration strategies, making data harmonization
across networks a growing challenge. Off-the-shelf sensor accuracy has been shown to vary depending on particle size,
55 chemical composition, hygroscopicity and local meteorological conditions such as temperature and relative humidity
(Manikonda et al., 2016; Jayaratne et al., 2018; Zamora et al., 2019; Mehadi et al., 2020). While substantial attention has been
given to increasing the accuracy of sensor measurements through calibration and correction equations to account for source
pollutant characteristics and environmental factors (Zimmerman et al., 2018; Malings et al., 2020; Barkjohn et al., 2021; Raheja
et al., 2023), significant challenges remain. Sensors vary widely in build quality, stability, and data reliability, and may be
60 installed or deployed to locations that may bias their measurements or reduce their performance. Devices may fail unexpectedly
due to reduced hardware quality or lack of protection from environmental elements, and the operation and maintenance
practices applied to a given sensor can vary widely across deployments. Common faults included power loss, communication
failure, data corruption, hardware malfunction, and firmware error (Barkjohn et al., 2025a). When aggregating air quality data,



65 faulty sensors often manifest as statistical outliers within a network, either as persistently low or high extreme values, flatlined readings, abrupt baseline shifts, or unusual divergence from collocated measurements. Even if the device is performing as intended, publicly operated sensors may produce aberrant measurements if they are sited or used with the intent to capture hyper-local conditions of interest and therefore may not reflect ambient air quality conditions at large. Localized emissions from sources such as woodstoves, individuals moving their outdoor sensors indoors during smoke events, or other personal-use scenarios may impact public health messaging and confound network-scale analyses. As sensor networks continue to expand, enabling broader spatiotemporal coverage, it will be increasingly important for users of sensor data to implement automated fault detection and normalization strategies that can identify these outliers without sacrificing data integrity.

In this work, we explore a generalizable framework for outlier detection, primarily by inferring sensor validity using a local network of nearby devices (which may include both regulatory-grade PM_{2.5} monitors and sensors) that we may consider representative of similar localized conditions. Our objective is to develop a method that can be applied during wildfire smoke events, and over the wide range of landscapes that may be impacted by smoke. The degree of difficulty in this task is compounded by the unequal spatial distribution of sensors (deSouza and Kinney, 2021) and the highly variable scale and magnitude of smoke events, which can produce gradients of PM_{2.5} concentrations that may be sharp, patchy, and transient over small distances (Jaffe et al., 2020). These sharp gradients may be further exacerbated by terrain influences, where topographic features can channel, pool, and disperse smoke in ways that decouple nearby sensors. To overcome these challenges, a useful definition of representativeness must be agnostic of sensor manufacturer, scale with varying network densities, and account for the timing and variability of PM_{2.5} concentration shifts across networks. Additionally, methods must account for the fact that many sensors are intentionally deployed in pollutant hotspots (Madhwal et al., 2024; Sablan et al., 2024) where elevated concentrations or abrupt spatial variability may be real and meaningful rather than indicative of faulty measurements (Barkjohn et al., 2025a). Thus, a robust outlier detection framework will need to maintain tolerance for disagreement within the network during smoke intrusion into a region, while maintaining sensitivity to correctly identify true faults.

Recent work has explored network approaches for validating sensor measurements, including graph-based methods (Ferrer-Cid et al., 2022). Efforts such as this underscore a growing interest in using the collective behaviour of a sensor network to identify anomalous measurements. However, such approaches depend on either gradual measurement variation or a historical relationship between network members, both of which may break down during wildfire smoke events. In this work, we build on this general idea of network-informed detection but adapt it to the realities of wildfire smoke by designing an approach that tolerates local disagreement, incorporates temporal variability, and remains robust to the shifting structure of inter-sensor relationships during smoke events.

By using a dataset of PM_{2.5} observations during 11 smoke events across North America in 2024 and 2025, we simulate common fault types in existing sensor networks and optimize an algorithm to identify them. Using measures of entropy and information from the field of information theory, we define and optimize a rule-based equation that is tuned for identifying outliers at an hourly resolution. We then expand upon this concept and train a gradient-boosted decision tree model using our



outlier detection function in the feature set. The results from these models, and their practical implications when applied to continental-scale sensor networks, are discussed.

2 Study design and data collection

100 2.1 Study Design Overview

This study proceeds in two phases. The first focuses on method development, model training, and benchmarking using a semi-synthetic dataset with known sensor faults (Sections 3-4). The second applies a subset of these methods to 9 months of continental-scale operational data (Section 5).

For phase one, we collected ground-based $PM_{2.5}$ observations during 11 wildfire smoke events that occurred in North
105 America in 2024 and 2025. Data were collected from several types of $PM_{2.5}$ measurement devices including regulatory permanent monitors, temporarily deployed reference-grade monitors, state-developed sensors, and publicly deployed sensors that were active during the smoke events. Air quality devices were grouped into networks at several spatial distance thresholds, creating lookup tables between any given device and its neighbouring devices. For each smoke event, a random, stratified sample of devices was selected to have its data perturbed to simulate a fault. This included both monitors and sensors, and for
110 the purpose of our analysis, these devices were treated the same. Five types of faults were introduced to each device, perturbing the data of only one unit at a time during a given event (see Section 2.5). In total, 176 devices were used to simulate potential outliers, resulting in 880 different scenarios to evaluate the ability of our methods to detect perturbations.

In phase two, we expanded beyond the scope of the 11 short-duration smoke events used in parameter tuning and model training to understand the broader utility of our methods. We applied eight of the outlier detection methods generated
115 on the semi-synthetic dataset to 9 months of air quality data (January 1, 2025 through September 30, 2025) from 9 sites in the United States and one site in Canada. We manually evaluated the accuracy of the different detection methods to determine whether the outliers flagged were faulty devices, hyper-local events, or misclassified (i.e., well-performing sensors).

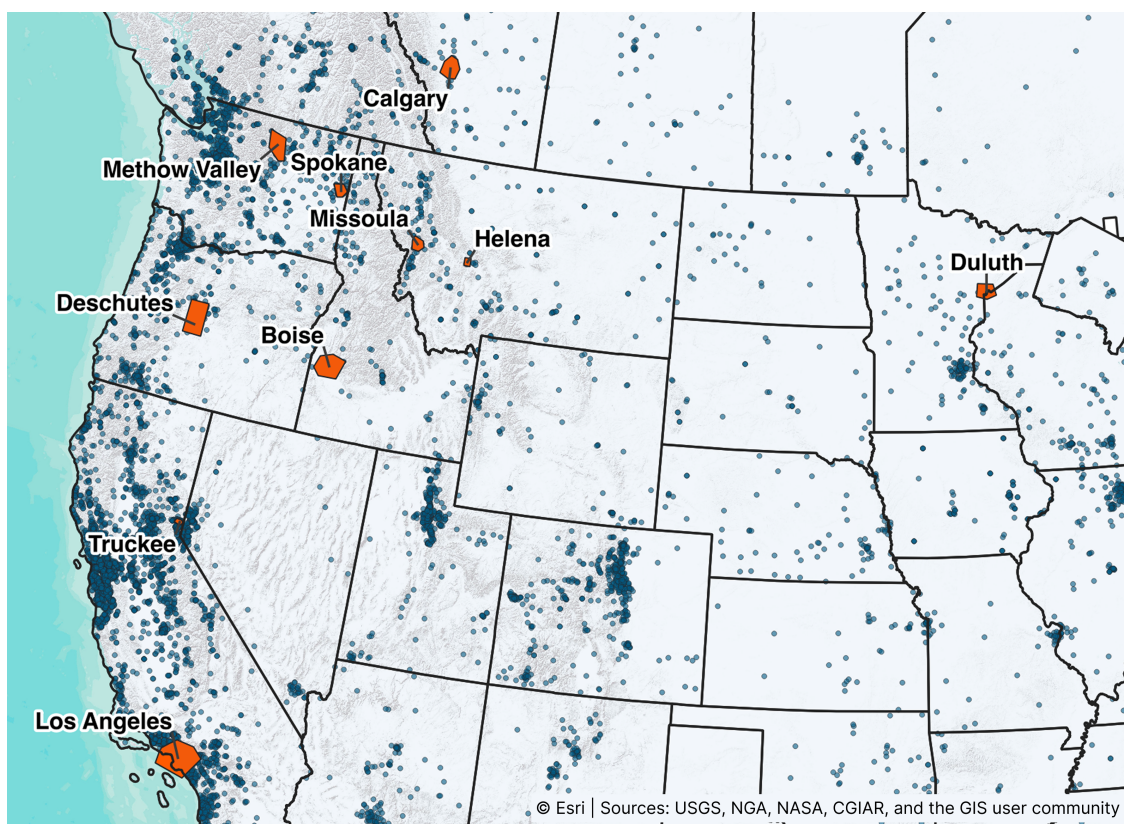
2.2 Smoke Event Selection

We define a smoke event as a geographic region impacted by elevated $PM_{2.5}$ concentrations as a direct result from wildland
120 fire. These events have a clear start and end time, delineated by lower, ambient $PM_{2.5}$ concentrations. The sites used in this study were selected from known events that occurred during 2024 – 2025 that were measured by sufficient $PM_{2.5}$ monitors and sensors to capture the spatial and temporal extent of the smoke. There were 10 geographic sites used, and 11 total events, with one geographic area used for both an event caused by regional wildfire (WF) smoke, and another event due to a local prescribed fire (Rx).

125 The 10 geographic sites used are: Boise, Idaho (ID); Calgary, Alberta (AB), Canada; Deschutes County, Oregon (OR); Duluth, Minnesota (MN); Helena, Montana (MT); Los Angeles, California (CA); Methow Valley, Washington (WA);



Missoula, MT; Spokane, WA; and Truckee, CA, and can be seen in **Fig 1**. A summarized list of the smoke events including geographic area, date range in which the event occurred, and the number of PM_{2.5} devices used is provided in **Table 1**.



130

Figure 1: Map showing the location and geographic extent of the 10 geographic sites used, and PM_{2.5} monitors and sensors with recorded observations during the study period. (Basemap sources: Esri, USGS, NGA, NASA, CGIAR, NOAA | Powered by Esri)



135 **Table 1.** Site names, smoke source (wildfire, WF or prescribed, Rx), spatial extent, and the date ranges used in capturing the smoke event. Includes the number of monitors (Permanent or Temporary) and the number of sensors (PurpleAir, SensOR, or SensWA), and the number of devices (includes both monitors and sensors) perturbed to simulate outliers.

Site Name	Area (km ²)	Smoke Event Dates	Smoke Source	Number of Monitors	Number of Sensors	Number of Simulated Outliers
Boise, ID	4,053	08/28/2024 – 09/12/2024	WF	4	46	9
Calgary, AB (Canada)	2,330	07/20/2024 – 07/31/2024	WF	4	24	10
Deschutes County, OR	4,627	05/14/2024 – 05/17/2024	Rx	8	152	16
Deschutes County, OR	4,627	09/03/2024 – 09/15/2024	WF	4	204	22
Duluth, MN	1,865	08/11/2024 – 08/16/2024	WF	3	13	6
Helena, MT	379	07/20/2024 – 07/31/2024	WF	1	15	8
Los Angeles, CA	7,957	01/06/2025 – 01/12/2025	WF	7	553	58
Methow Valley, WA	2,790	08/10/2024 – 08/21/2024	WF	1	55	10
Missoula, MT	1,006	08/31/2024 – 09/12/2024	WF	2	36	7
Spokane, WA	974	11/05/2024 – 11/12/2024	Rx	7	55	13
Truckee, CA	261	11/07/2024 – 11/10/2024	Rx	1	84	17

140 Consideration was given to selecting a diverse range of sites and smoke events. Rather than focusing solely on regions with the most frequent smoke events, sites were chosen to represent a range of conditions, including urban and rural settings, impacts from both WF and Rx smoke, localized versus regional smoke events, topographically complex versus flat regions, dense versus sparse networks of sensors, and differing environmental regions. A detailed site and event characterization that covers terrain, land cover, sensor network composition, and smoke source attribution for each event is provided as a characterization document in the companion repository (Illson, 2026b).

145 2.3 Air Quality Data

2.3.1 Data Sources

150 Air quality data were collected for four different categories of devices: Permanent monitors, temporary monitors, state-developed sensors, and publicly owned sensors. Each device type is briefly explained in the following sub-sections. Data were obtained from three sources. First, data from permanent monitors, some temporary monitors, and state-developed sensors were obtained from the U.S. Environmental Protection Agency’s (EPA) AirNow (<https://airnow.gov>) database, which provides air quality data for over 2,500 monitors and sensors across the United States, Canada, and parts of Mexico. Second, data were accessed using the AirNow API (<https://docs.airnowapi.org/>) and processed by the open-source R packages AirMonitor and



AirMonitorIngest, which were developed by Mazama Science (<https://github.com/MazamaScience/AirMonitor>,
<https://github.com/pnwairfire/AirMonitorIngest>).

155 Additionally, temporary monitor data were obtained from the Western Regional Climate Center's (WRCC,
<https://wrcc.dri.edu/>) archive of temporary monitors. These monitors are available on loan to state, local, and Tribal air quality
or public health organizations and to the Interagency Wildland Fire Air Quality Response Program (IWFAQRP) for use by
Air Resource Advisors (ARAs). Data from the WRCC archive were also downloaded and processed by the AirMonitor and
AirMonitorIngest R package. PurpleAir sensor data were obtained via the PurpleAir API (<https://api.purpleair.com/>) and were
160 accessed using a custom Python library to download and perform quality assurance/quality control (QA/QC). Since we limited
the number of devices considered in our analysis and the time ranges of our events were sufficiently short, this dataset was
manually validated to remove sensors that were known to be faulty or otherwise confound our results.

2.3.2 Permanent Monitors

The permanent monitors used in this study are currently operated and maintained by State, local, or Tribal air monitoring
165 agencies in the United States. Most permanent monitors were federal equivalent methods (FEM) including Met One Beta
Attenuation Monitor (BAM) 1020 (Grants Pass, OR, USA), Met One BAM 1022 (Grants Pass, OR, USA), and Teledyne T640
(San Diego, CA, USA). Some non-FEM methods including radiance research M903 nephelometers (Seattle, WA, USA) with
heated inlets and other automated surrogate PM_{2.5} measurements were used as well. A few sites did not report the monitor type
to AirNow.

170 2.3.3 Temporary Monitors

The temporary monitors are typically Met One Environmental Beta Attenuation monitors (EBAM; Grants Pass, OR, USA) or
Met One Environmental Smoke Attenuation Mass monitors (ESAM; Grants Pass, OR, USA), which have shown to have
reasonable agreement with the Federal Reference Method (FRM) network (Schweizer et al., 2016). These devices are deployed
to areas specifically to monitor smoke events and are operated, maintained, and quality controlled by air quality experts.

175 2.3.4 State-developed Sensors

SensWA sensors were developed and are operated by the Washington Department of Ecology (Washington State Department
of Ecology, 2024b). These sensors are used to supplement existing permanent monitors and as temporary monitors during
wildfires or other emergent air quality events, along with other uses. Each device contains two Sensirion SPS30 sensors
(Sensirion, Stäfa, Switzerland). The Sensirion PM₁ (particulate matter with diameters of 1 μm or less) output is used to estimate
180 PM_{2.5} concentrations as it is better correlated with reference PM_{2.5} measurements across Washington (Washington State
Department of Ecology, 2024a) and past literature has shown that Sensirion SPS30 based sensors do not accurately size
particles (Molina Rueda et al., 2023). Washington Department of Ecology corrects sensors using regional and seasonal
corrections. There are 10 regions across Washington state that are expected to have similar particle properties due to similar



185 meteorology and sources and BAMs are used to generate corrections for each region during both the summer, when wildfire smoke impacts the state, and non-summer periods (typically October – July 5th). A three-part equation is applied to the PM₁ data to account for the nonlinear relationship between the SensWA and BAMs. SensWA are corrected using a median regression from 20-100 $\mu\text{g}/\text{m}^3$, median regression above 100 $\mu\text{g}/\text{m}^3$ (if enough data is available), and a quadratic equation below 20 $\mu\text{g}/\text{m}^3$ (Washington State Department of Ecology, 2024a).

190 SensOR sensors were developed and are operated by the Oregon Department of Environmental Quality (State of Oregon Department of Environmental Quality, 2023). These sensors are used to communicate EPA’s Air Quality Index (AQI) health information to the public and have a variety of partner hosts (e.g., schools) (State of Oregon Department of Environmental Quality, n.d.). SensOR sensors contain two Plantower PMS5003ST sensors (Plantower, Beijing, People's Republic of China). A fan draws air through the enclosure and through a heated aluminium inlet. Oregon State’s data quality objectives for these sensors are $\geq 75\%$ completeness and concentrations should be within $\pm 20\%$ of the FRM or FEM data. Site specific linear corrections against nephelometer backscattering are used. At high backscattering levels, a wildfire smoke correction is used annually from June 1 through Oct 31. The wildfire smoke correction is different than the site-specific correction that accounts for winter woodsmoke. The sensors are monitored by Oregon for issues in the comparison between duplicate sensors and other data issues monthly with quarterly site QC checks (Johnson, 2024).

2.3.5 PurpleAir Sensors

200 The PurpleAir sensor network is a crowdsourced network of air sensors deployed by a variety of users around the world. Each sensor contains two Plantower sensors (Plantower, Beijing, People’s Republic of China) measuring PM_{2.5}. PurpleAir cf_atm measurements are retrieved from the PurpleAir API as 10-minute averages from both the A and B measurement channels. The sub-hourly measurements are rolled up into a single, channel-averaged, hourly value when data are present for at least 4 of the 6 10-minute intervals within the hour. The hourly values are corrected with a 5-piece nonlinear equation for cf_atm that includes training data from 0-1500 $\mu\text{g}/\text{m}^3$ (Barkjohn et al., 2025b). Data are excluded if the two PM_{2.5} measurement channels show substantial disagreement, defined as both an absolute difference greater than 5 $\mu\text{g}/\text{m}^3$ and a percent difference greater than 70% (Barkjohn et al., 2022). Only records meeting at least one of these agreement criteria were retained. Manual validation removed some sensors from the dataset if a review of the time series suggested the sensor was broken (e.g., intermittently reporting, persistently reporting high or low values).

210 2.4 Network Generation

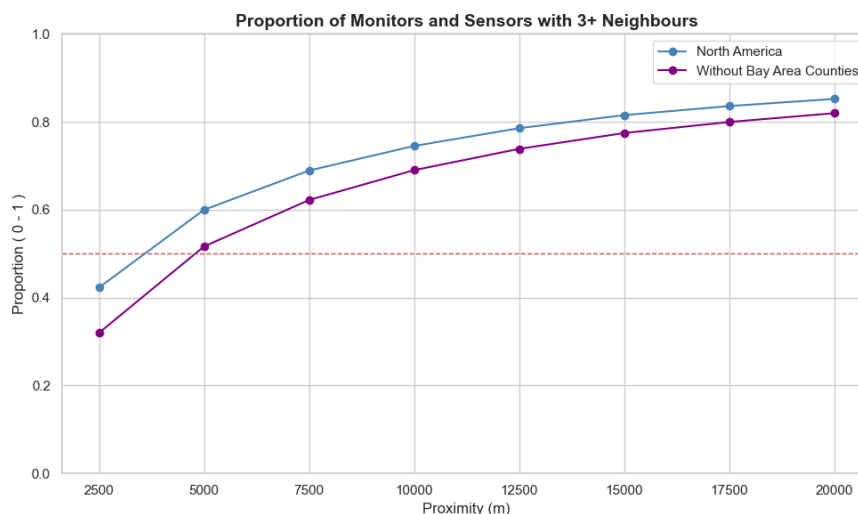
To construct local networks of monitors and sensors, one-to-many lookup tables were generated at multiple proximity thresholds, grouping any given PM_{2.5} monitor or sensor to others based at that proximity. First, each measurement location was assigned a unique identifier based on a combination of the devices’ given ID and latitude/longitude coordinate pair. This was necessary because many sensors and temporary monitors relocate or may report coordinate differences over time. Of the devices reporting to the Fire and Smoke Map, 984 sensors and 19 temporary monitors changed location (greater than 250



meters) at least once between April 2024 and February 2025. Device coordinates were reprojected into the North America Albers Equal Area Conic projection (ESRI:102008) to enable distance calculations in meters. Neighbouring units were then identified using a pairwise two-dimensional Euclidean distance threshold, implemented with the libpysal Python package (Rey and Anselin, 2007).

220 We considered network generation methods that account for terrain when linking sites, treating topographic features
as potential barriers that may separate distinct airsheds. In addition to standard two-dimensional distance thresholds, we
explored two other network generation methods that incorporated terrain complexity. First, we tried modifying the two-
dimensional distance metric by adding an absolute elevation gain or loss penalty, reducing the likelihood of linking devices
that are separated by significant elevation changes. Second, we implemented a least-cost path approach (Hart et al., 1968),
225 using slope from a digital elevation model (DEM; NASADEM 1 arc-second global DEM, NASA JPL, 2020) as a cost grid.
Connections between devices were determined by traversing a slope-derived cost surface, allowing paths to follow
topographically favourable routes rather than straight-line distances. These methods aimed to better capture connectivity in
topographically complex regions, to account for smoke transport along valleys and drainage features. Although least-cost path
networks appeared to yield slightly more cohesive groupings in mountainous areas, we ultimately selected the two-dimensional
230 distance method. This approach provided simpler implementation, broader generalizability, and removed the need to integrate
DEMs into the workflow. The implications of this decision will be explored in the discussion section, where we see networks
that include cross-drainage devices being falsely identified as outliers.

Initially, networks were generated at eight spatial thresholds, ranging from 2.5 to 20 km in 2.5 km increments. We
further reduced the number of network proximities for evaluation by considering the spatial distribution of devices and the
235 need to balance proximity with a sufficient number of neighbouring devices. Considering the first law of geography that
“everything is related to everything else, but near things are more related than distant things” (Tobler, 1970), smaller spatial
buffers may produce networks that capture shared conditions, however, identifying faulty behaviour within a network also
requires a minimum number of devices to establish consensus. This balance was complicated by the uneven distribution of
sensors across North America. The PurpleAir sensors are disproportionately concentrated in California, especially so in the 9
240 counties that comprise the Bay area – Alameda, Contra Costa, Marin, Napa, San Mateo, Santa Clara, San Francisco, Solano,
and Sonoma. At the time of analysis in early 2025, we found 18% of all PM_{2.5} monitors and sensors in North America were
located within these 9 counties, with the overwhelming majority being PurpleAir sensors. At 21,220 km², these counties cover
only 0.23% of the total United States land area, yet disproportionately skew our results for how many neighbouring devices a
sensor or monitor will have at a given proximity. To strike a balance between spatial proximity and sufficient network density,
245 we focused on thresholds where at least half of the devices located outside the Bay Area counties had three or more neighbours.
This ensured more consistent network structure across the broader sensor distribution. The impact of Bay Area sensor density
on neighbour count in networks is shown in **Fig. 2**. For our analysis, we elected to use networks 5 km and larger, ensuring that
at least half of devices had at least three neighbours.



250 **Figure 2:** Proportion of monitors and sensors with at least three neighbours across spatial thresholds from 2.5 km to 20 km. The blue line represents all devices across North America, while the purple line excludes devices located in the nine Bay Area counties. The high density of sensors in the Bay Area inflates network connectivity at smaller proximities.

An example of how connectivity within a network evolves over various spatial thresholds can be seen in **Fig. 3**, at
255 the Methow Valley, WA site. This example highlights how devices in separate drainages become associated with one another
at the higher thresholds, which has the potential to introduce network disagreement. This is particularly important in complex
terrain, where wildland fire smoke can become trapped in valleys due to radiative feedbacks that reinforce local stability and
suppress mixing (Kochanski et al., 2019). In very dense networks, such as Los Angeles, CA, these one-to-many networks can
become incredibly large and redundant, with devices having greater than 200 connections even at modest thresholds. To
260 maintain tractability and reduce computational requirements, we limited all networks to a maximum of 20 neighbouring
devices, keeping only those closest to our device of interest.



Methow Valley, WA, PM_{2.5} Networks

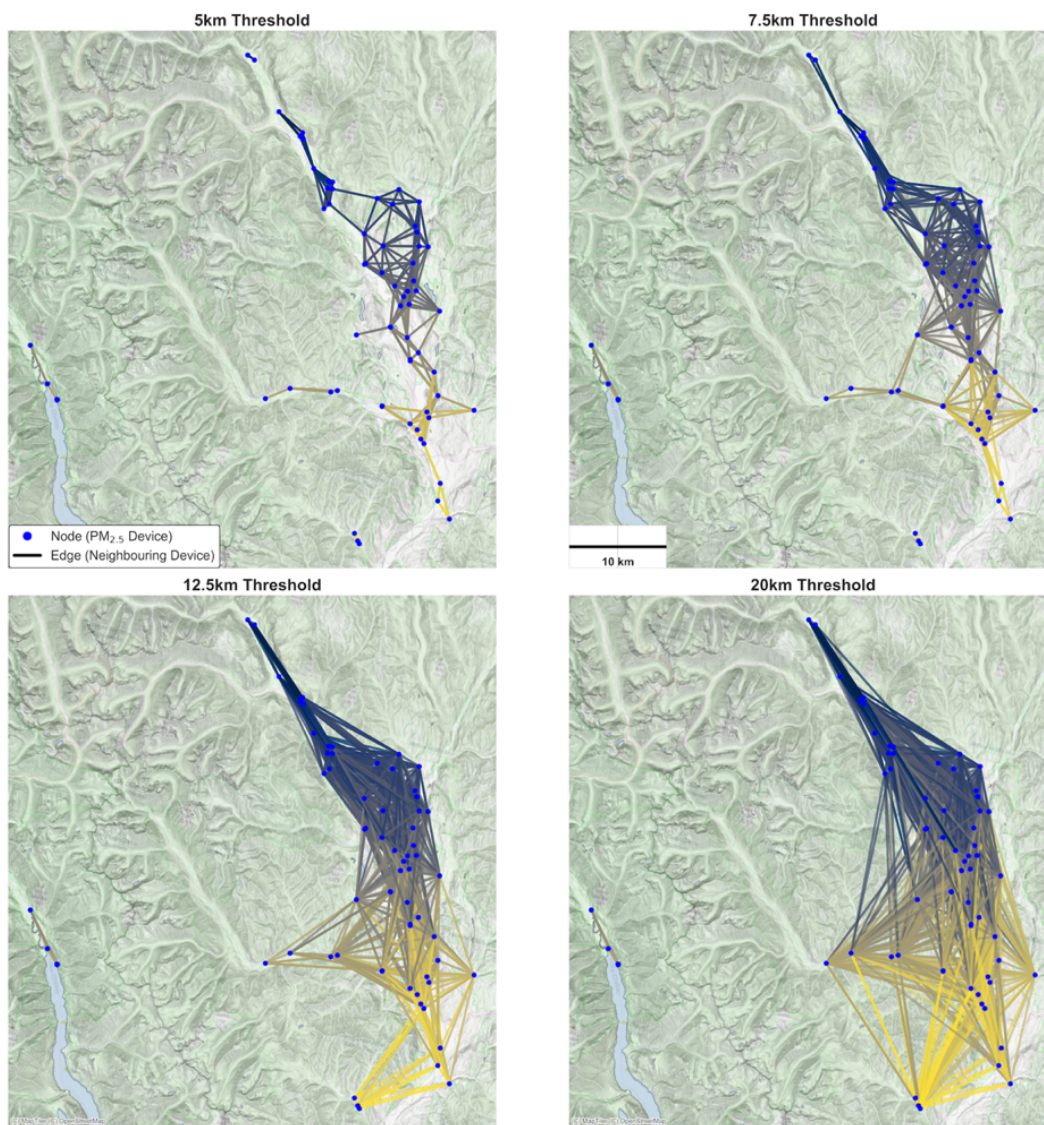


Figure 3: Example of PM_{2.5} sensor networks in the Methow Valley, WA, generated using four different spatial thresholds: 5 km, 7.5 km, 10 km, and 12.5 km. Each point represents a PM_{2.5} monitor or sensor, and lines indicate connections based on the selected proximity threshold. As the distance increases, overall network connectivity increases, linking devices across drainages. Line colour corresponds to device latitude, to help trace where connections originate. (Basemap: © MapTiler; map data © OpenStreetMap, distributed under the Open Data Commons Open Database License (ODbL) v1.0; for further information see <https://www.openstreetmap.org/copyright/en>)



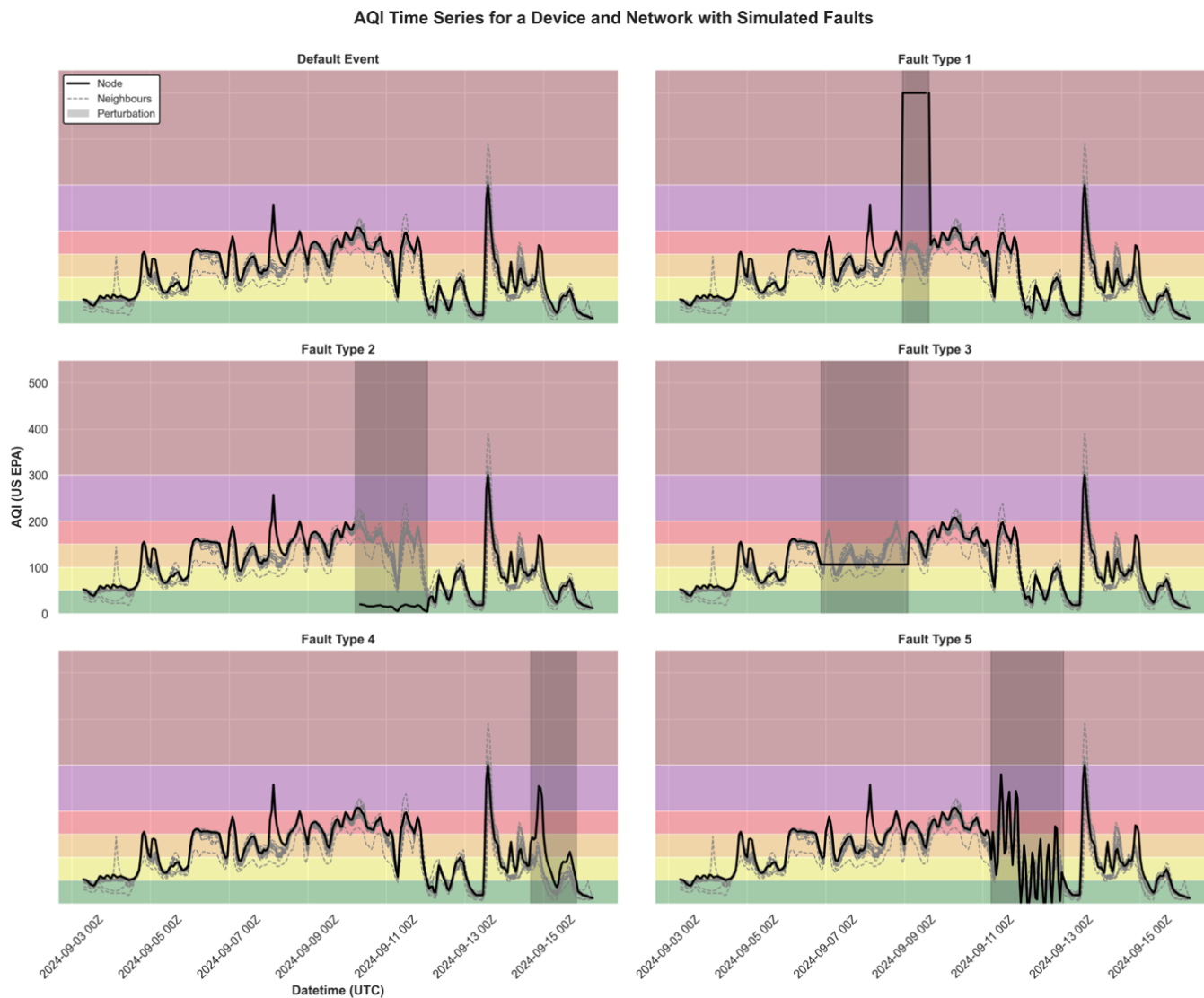
2.5 Simulated Fault Generation

270 A semi-synthetic dataset was constructed by injecting simulated faults into the real-world PM_{2.5} time series data from each of
the smoke events. For each event, we identified eligible devices (both monitors and sensors) based on a minimum number of
connected devices in their network at 5 km, selecting only those with at least three neighbours to ensure sufficient context. A
random subset of these eligible devices was then sampled, with the sampling fraction scaled to maintain coverage in both
sparse and dense networks. We sampled 40% of eligible devices if fewer than 25 devices were present within the site footprint,
275 10% if more than 100 were available, and 20% otherwise.

For each selected unit, five unique fault types with different perturbations were applied, introducing only one into a
network within a smoke event at a time. The perturbations occurred for a continuous interval, randomly ranging from 8 to 48
hours in duration, keeping within a 3-hour buffer from the beginning or end of the time series. Fault types were designed to
emulate realistic sensor anomalies:

- 280 1. **Fault type 1:** The device becomes fixed at a high value, intermittently dropping out due to either failing
QA/QC or electronic malfunction.
2. **Fault type 2:** The device goes offline, then returns reporting a lower value – such as when an outdoor sensor
is moved indoors.
3. **Fault type 3:** The device becomes stuck and reports a single value without responding to ambient conditions.
- 285 4. **Fault type 4:** The device measurement is amplified by a random factor increasing it by 75-100% of the
original value.
5. **Fault type 5:** Erratic and high-frequency noise is applied to the value.

Figure 4 illustrates examples of each of these perturbations applied to a device. The code to generate similar synthetic faults
290 is provided in the supplemental material.



295 **Figure 4:** Time series showing how a randomized, synthetic fault is applied to a $PM_{2.5}$ device. Time series of hourly averaged $PM_{2.5}$ measurements converted into EPA AQI for a single sensor (solid black line) and devices within its network (dashed grey lines) during a smoke event. The top-left panel shows the unperturbed time series for reference. The remaining panels display the impact of each fault type, highlighted by the shaded period. Background is coloured by EPA AQI category ranges.

3 Information-theoretic detection

3.1 Applicability of Information Theory

300 For this work, we investigated whether the fundamental measures of information theory could validate a $PM_{2.5}$ sensor's measurements within the context of its local network. Information theory, first introduced by Claude Shannon in 1948



(Shannon, 1948), provides a mathematical framework for quantifying uncertainty and surprise in data. Originally developed to optimize the transmission of messages over communication channels, it has since been widely extended to fields as diverse as genetics, neuroscience, climate science, and machine learning. Today, information-theoretic concepts play a central role in how we model, compress, and extract meaning from complex, uncertain systems. The two principal concepts, entropy and information content, offer potential utility in interpreting sensor and network behaviour.

Entropy captures the degree of disorder, or unpredictability within a system. In the context of a $PM_{2.5}$ sensor network, it provides a measure of the variability of measurements within the network. High entropy reflects disorder or instability, such as conditions where smoke intrusion is causing nearby devices to report drastically different measurements. Conversely, low entropy indicates a broad agreement and stability across the network. The information content of a device within its network reflects the unexpectedness, or surprise, of a given value being seen relative to its network. This allows us to not only characterize the state of the network as a whole, but to identify individual devices that deviate from their local network. These information-theoretic concepts work together, allowing us to accommodate networks of varying size and device type, while providing measures of disorder that remain meaningful during smoke events.

It is important to note that Shannon entropy and Shannon information were formulated for discrete variables, which at first may seem at odds with real-valued $PM_{2.5}$ measurements. In both public health and regulatory contexts, $PM_{2.5}$ values are typically converted into discrete air quality categories, such as the EPA's Air Quality Index (AQI) (<https://www.airnow.gov/aqi/>) or the Canadian Air Quality Health Index (AQHI) (Stieb et al., 2008). These categorical bins are tractable, widely understood, and grounded in communicable risk levels, making them a natural starting point in determining whether two sensors are in agreement. This categorization avoids overemphasis on minor numerical differences, while capturing meaningful shifts in health-relevant air quality. We believe such categorization is sufficient for broadly detecting outliers, and supports practical communication of health protective actions.

Information theoretic methods appear particularly well suited to characterizing the state of sensor networks. These methods inherently scale with network density, adapting to situations with many or few neighbouring observations without the need for reparameterization. They can detect non-linear relationships in the data, making them robust to the complex, shifting patterns that occur during smoke events. One of the key advantages of applying information theory in this context is that it enables the detection framework to defer classification during periods of high uncertainty when the network becomes disordered, avoiding premature characterizations during rapidly evolving smoke events. Because the approach is contingent on network agreement, its sensitivity can be tuned to tolerate varying levels of disorder across spatial proximities or event types, from highly localized anomalies to region-wide smoke intrusions. This prevents the removal of potentially valuable data under unstable conditions, while allowing temporal averaging and threshold adjustments to modulate the aggressiveness of outlier detection as desired.



3.2 Rule-based Detection Equation

We define a rule-based outlier detection equation that uses a set of adjustable thresholds to flag measurements that match
335 specified criteria. The equation relies on conditional logic to assess a given device’s measurement relative to its local network.
This equation is applied to a single device at a time, and utilizes the current measurement of the device and the current
measurements within its network.

First, we compute the network entropy H_k for device k , which reflects the degree of disorder in the distribution of
observed categorical AQI bin values $x \in X$, reported by the device and its local network:

$$340 \quad H_k = - \sum_{x \in X} p(x) \log_2(p(x)) \quad (1)$$

where $p(x)$ is the proportion of neighbouring devices whose measurements fall into AQI bin x . A lower entropy value
indicates that neighbours are clustered in a small number of categorical bins (more ordered, less uncertain). Conversely, a
higher entropy value indicates a broader distribution of observed bins.

Second, we compute the information content I_k of the device’s bin category with respect to its network, representing
345 how uncommon that value is relative to observed bins in the network:

$$I_k = - \log_2(p_k) \quad (2)$$

Higher values of information content indicate observations that deviate more strongly from expectation, carrying greater
“surprise” in the information-theoretic sense. Such measurements are less probable within the local distribution. In this context,
they correspond to more surprising or infrequent bin classifications relative to those of neighbouring devices.

350 Third, we calculate the bin deviation D_k , defined as the absolute difference between a device’s bin index and the
mean bin index of all other devices in the same local network at the same time step. For a network of N devices, where B_k
denotes the bin index assigned to device k :

$$D_k = \left| B_k - \frac{1}{N-1} \sum_{\substack{j=1 \\ j \neq k}}^N B_j \right| \quad (3)$$

Here, the summation is taken over all devices in the network except device k , and the summation represents the mean
355 bin index of the remaining $N - 1$ devices. Note that B_k refers to the bin index (an integer from 1 to the number of bins),
distinct from the bin boundaries x used to define concentration intervals in Equation 1. Because information content is based
on the probability of discrete bin values, it reflects how rare a value is within the network but not how far that value is from
the norm. For example, an observation in a rarely occupied bin will yield the same information content whether it is close or
far in measurement space. A network with 49 observations in the lowest AQI bin and 1 in the highest produces the same
360 entropy as one with 49 in the lowest and 1 in the second-lowest bin. D_k captures this distinction, quantifying how far in
measurement space an observation of B_k may be from the network average.



Finally, the detection rule is implemented using a conditional logic statement, which flags a unit as an outlier only if all three conditions are met:

$$Outlier_k = \begin{cases} \text{True,} & \text{if } I_k^{(r)} \geq \theta \text{ and } H_k^{(r)} \leq S \text{ and } D_k^{(r)} \geq \left(\frac{B}{\beta}\right) \\ \text{False,} & \text{otherwise} \end{cases}$$

I_k = Information value of device k

H_k = Network entropy value for device k 's network

D_k = Absolute difference between device's bin and the mean network bin

r = spatial proximity threshold for defining the network around device k

θ = Information coefficient

S = Network entropy coefficient

β = Bin coefficient

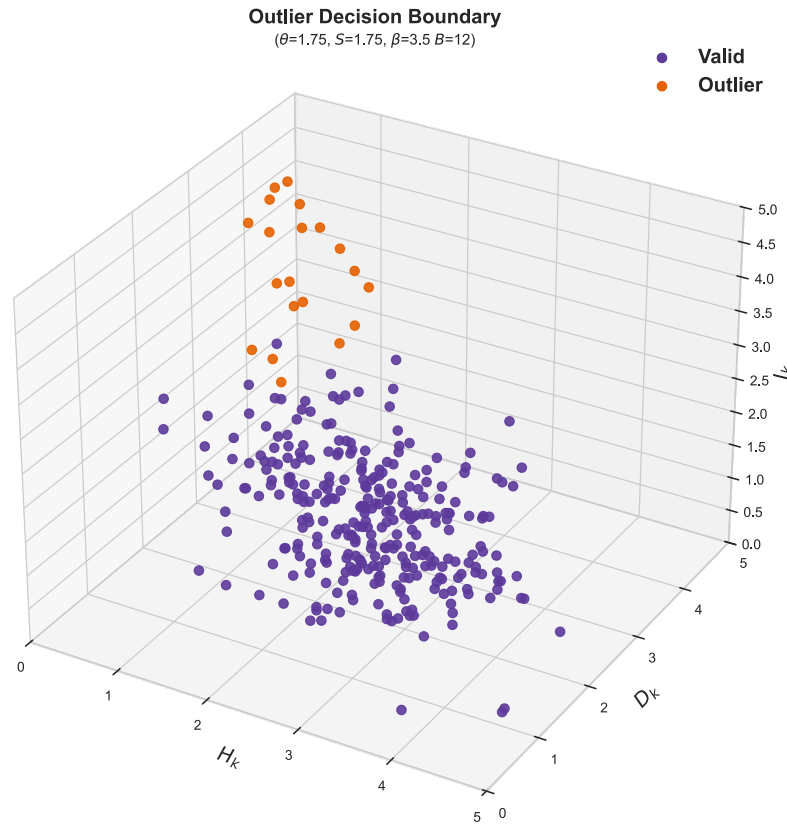
B = Number of possible bin values

(4)

365

This equation flags outliers by evaluating a device's measurement against its local network at a given spatial proximity r . We can adjust several inputs and thresholds, with the overall goal being to select instances where a device's measurement (i) exists in an ordered network (low entropy); (ii) has a high information content relative to the network; and (iii) is sufficiently distant from the average bin category of the network (see **Fig. 5**). We adjust the bin coefficient β by the number of possible bin values B . This is an attempt to help the equation remain transferrable across different binning schemes or air quality index classification systems.

370



375 **Figure 5:** Three-dimensional visualization of the decision boundary created by the outlier detection equation. The decision rule optimizes for outlier detection given the device measurements in the context of the device’s network. We tune the equation for measurements that (i) exist in an ordered network (low entropy); (ii) have a high information content relative to their network; and (iii) are sufficiently distant from the average bin category of the network.

3.3 Parameter Tuning

380 The goal of tuning our equation is to maximize the number of hours in this time series where we correctly identify the perturbed measurements, while not misclassifying any valid measurements. For our rule-based equation, there are a total of 5 configurable inputs, which may be tuned to optimize its performance. These are: 1) the local network at proximity r ; 2) the number of bins B and what concentration ranges they represent; 3) how much network entropy to allow S ; 4) how much information is required θ ; and 5) how many bin steps are sufficiently far away in measurement space D_k .

385 We selected both the F_1 score and F_β score ($\beta = 0.5$) as objective functions to tune our thresholds with, reflecting a trade-off between balanced performance and precision-heavy tuning. In a classification context, the F_1 score gives equal weight to precision and recall, $F_{\beta=0.5}$ increases the weight on precision relative to recall. Compared with accuracy, F_1 and F_β scores are



more appropriate for imbalanced datasets such as ours where outliers are relatively rare. Metrics such as area under the receiver operating characteristic curve (AUC-ROC) and precision-recall curves are less directly applicable in this context, as our decision rule outputs discrete class labels rather than continuous scores. The F_1 and F_β scores provide interpretable summaries of this trade-off of the classifier performance at a given parameter set.

To ensure meaningful optimization and keep computational runtime within reason, we constrained the parameter search space. For example, entropy thresholds had to allow some variability while excluding cases of complete disorder, and bin deviation thresholds needed to distinguish measurements that were meaningfully different in measurement space. Without such constraints, it is possible to arrive at configurations that minimize false positives by trivially classifying nearly all data as valid.

We conducted a grid search applying 8,400 unique parameter combinations to 880 time series containing a synthetic fault, totalling over 7.2 million evaluation runs. The parameter sweep utilized spatial proximity thresholds ranging from 5,000 to 20,000 m, S and θ limits from 1.4 to 2.3 (in 0.1 increments), D_k from 1.5 to 4.0 (0.5 increments), and three binning strategies: EPA AQI, the Canadian AQHI, and a modified EPA AQI where bins were added to overlap the breakpoints between categories (**Figure 6**). These breakpoints covered the lower 10% of the next highest category and the upper 10% of the next lower category, which result in a doubling of categorical bins, totalling 12 bins. This modified EPA AQI was done to test softening the transitions of the AQI categories, since for example a sensor could be only $1 \mu\text{g}/\text{m}^3$ higher than its neighbours and be in a different AQI category. This search enabled us to identify parameter configurations that achieved optimal trade-offs between detection sensitivity and false positive rates across fault types and network configurations.

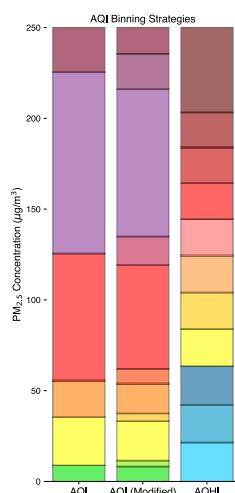


Figure 6: Three different binning strategies for $\text{PM}_{2.5}$ concentrations were used in this study. On the far left is the current EPA AQI, the middle bar is a modified EPA AQI, where additional categories were added to soften the breakpoints between categories, and the last bar is the Canadian AQHI scale.



410 3.4 Results on the Semi-Synthetic Dataset

We evaluate the performance of our equation with respect to two considerations: its ability to detect the perturbed device at least once during the event, and how well it classifies measurable outliers (perturbation induced outliers that meet the ruleset criteria for categorical deviation) across fault types and locations. **Figure 7** highlights the need for this measurable outlier definition and shows two cases where our equation is operating as designed, but not classifying all perturbed measurements. **Figure 7a** shows how the equation will not classify measurements during periods when the network entropy is high, and detection becomes uncertain. Similarly, **Fig. 7b** highlights a case where the synthetic fault never deviated from local conditions enough to qualify as an outlier, and therefore the perturbation is not a measurable outlier.

Table 2 summarizes the performance of each fault type. Across all induced faults where there was at least one instance of a synthetic measurement deviating from the original measurement's category, the equation successfully identified the outlier in 95.7% of time series events across all fault types and parameter configurations, demonstrating strong overall sensitivity to our outliers. When broken down by fault type, detection rates varied moderately. Fault type 1 achieved the highest success rate, with the outlier detected in 98.5% of perturbations. Fault types 2 and 3 followed, with detection rates of 90.7% and 90.0%, respectively. Fault type 4 was the most challenging, with a lower success rate of 82.8%, indicating that this class of fault may present more subtle or intermittent signatures. Fault type 5 matched the overall performance, with 95.7% of cases resulting in successful detection. These results suggest the general method is broadly effective across fault types as identifying an outlier from a composite time series, though performance may vary depending on the specific characteristics of the fault and parameter configurations used. These high detection rates were achieved while maintaining a very low rate of false positive detections. Across all configurations, events and fault types, the mean false positive rate was only 1.04%. This false positive rate remained consistently low across specific fault types: being the lowest at 0.91% for type 1, and the highest of 1.74% for type 5. When selecting the optimal parameter configuration for each fault type, in every test example the method demonstrated an ability to detect that given fault at least once during the smoke event.

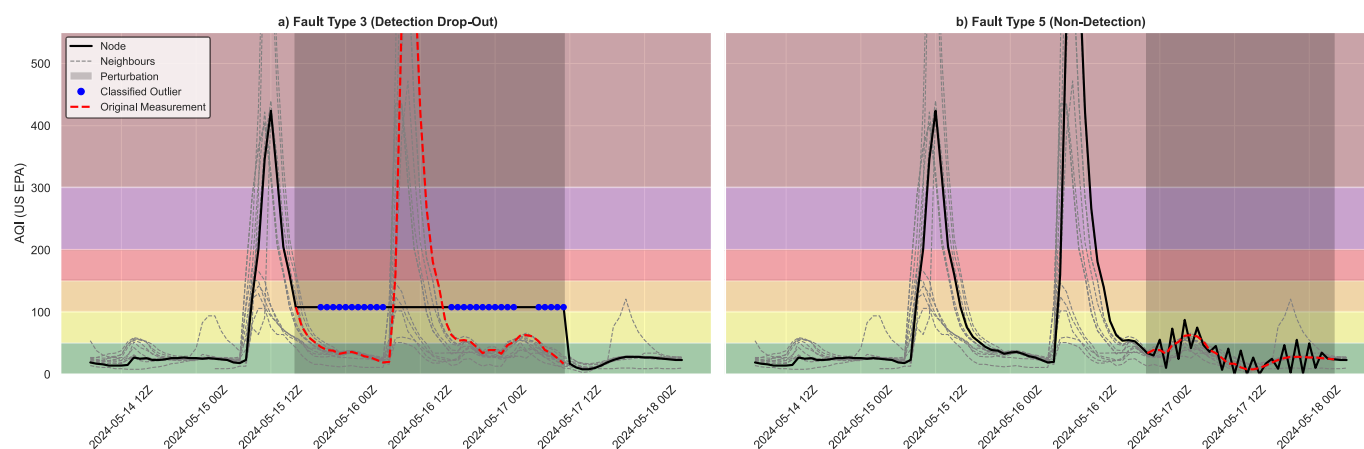




Figure 7: Time series of two different synthetic faults applied to a device (node), with its neighbouring network. The original device measurements shown as a black line. In **panel a**, we do not classify the perturbed measurements as an outlier while either entropy is high (disorder in network) or the local network measurement average is nearby in measurement space. In **panel b**, our equation skipped the perturbation entirely due to it being contained within the same categories as the network – we do not consider this event to contain measurable outliers.

The overall classification performance of the equation varied substantially in response to fault types and parameter configurations. **Figure 8** shows the full distribution of F_1 and $F_{\beta=0.5}$ scores across parameter configurations, plotted alongside the corresponding false positive rates. The optimal parameter configurations exhibit both high F_1 and $F_{\beta=0.5}$ scores, often exceeding 0.9, which is indicative of strong recall and precision. In its optimal parameter configuration, fault type 1 was the most reliably detected, achieving an F_1 score of 0.97. Fault types 2 and 3 followed with F_1 scores of 0.86, while fault type 4 was 0.83, and fault type 5 was 0.78 in its optimal configuration and when there were measurable outliers present (> 1). These results highlight that some fault types, such as type 1, are easier to separate using our equation. This variation is further illustrated in **Fig. 9**, which shows the distribution of F_1 and $F_{\beta=0.5}$ for each fault type. Fault type 1 not only achieves the highest median score but also exhibits the tightest distribution, indicating that many parameter configurations perform well across our events and different networks. In contrast, fault types 4 and 5 show wider distributions and lower central scores.

When examining the configurations that produced the highest F_1 scores for each fault type, no single parameter set was optimized for more than a single fault type. The information theoretic parameter values tended to be clustered in their test ranges in the top-performing configurations, suggesting partial generalizability. Notably, entropy limits were consistently in the upper end of the tested range, with all fault types preferring values ≥ 2.0 , indicating the importance of an ordered network. Conversely, information limits tended to cluster lower in the range suggesting that moderate to strict informational divergence thresholds are sufficient to distinguish perturbations from baseline variation. In terms of spatial configuration, the optimal proximity varied widely for each fault type, ranging from 5,000 to 20,000 m, implying that specific fault detectability may be at least partially scale-dependent. This diversity may reflect differences in how each fault type manifests relative to its neighbouring sensors. Finally, there was no clear optimal binning strategy across faults: fault types 1, 4, and 5 favoured the custom percentile binning, while fault types 2 and 3 favoured the standard EPA AQI categories. While the Canadian AQHI bins did not appear in the best parameter set, their results were closely aligned.

Table 2. Summary table for specific fault type detection performance where n measurable outliers > 1 . Shows the percentage of fault occurrences where our equation detected the outlier at least once under any parameter set tested, the optimal parameter set derived from our grid search, and the F_1 and false positive rate for the optimal parameter set.



Fault Type	Outlier Detected Once <i>(Any parameters)</i>	Optimal Parameter Set	F ₁ <i>(Optimal case)</i>	False Positive Rate <i>(Optimal case)</i>
1	98.5%	$\theta = 1.9 \ S = 2.3 \ B = 12 \ \beta = 3.5 \ r = 20km$	0.97	0.9%
2	90.7%	$\theta = 1.4 \ S = 2.0 \ B = 6 \ \beta = 2.5 \ r = 5km$	0.86	0.7%
3	90.0%	$\theta = 1.4 \ S = 2.1 \ B = 6 \ \beta = 2.5 \ r = 7.5km$	0.86	1.9%
4	82.8%	$\theta = 1.4 \ S = 2.0 \ B = 12 \ \beta = 3.5 \ r = 5km$	0.83	0%
5	95.7%	$\theta = 2.1 \ S = 2.3 \ B = 12 \ \beta = 3.5 \ r = 7.5km$	0.78	2.0%

465

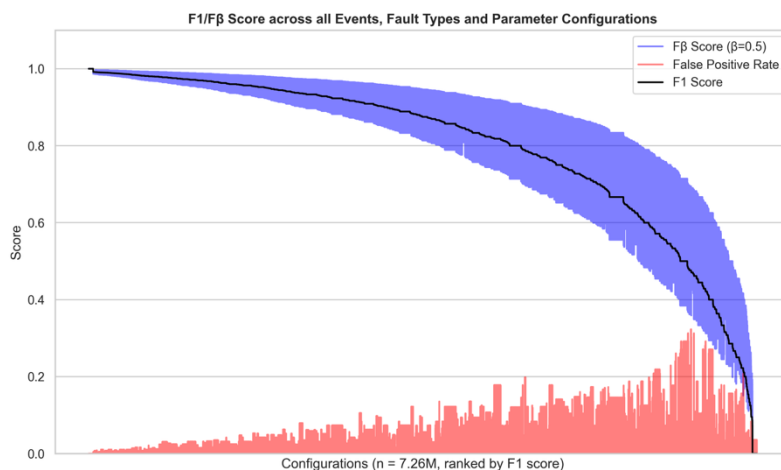


Figure 8: Curve of F₁ scores across all 7.2+ million event runs in the grid search, showing how the performance of our equation varied across all tested parameter configurations, synthetic fault types, network proximities, and events. Configurations are ranked by F₁ score from highest to lowest, showing how detection performance evolves across the parameter space.

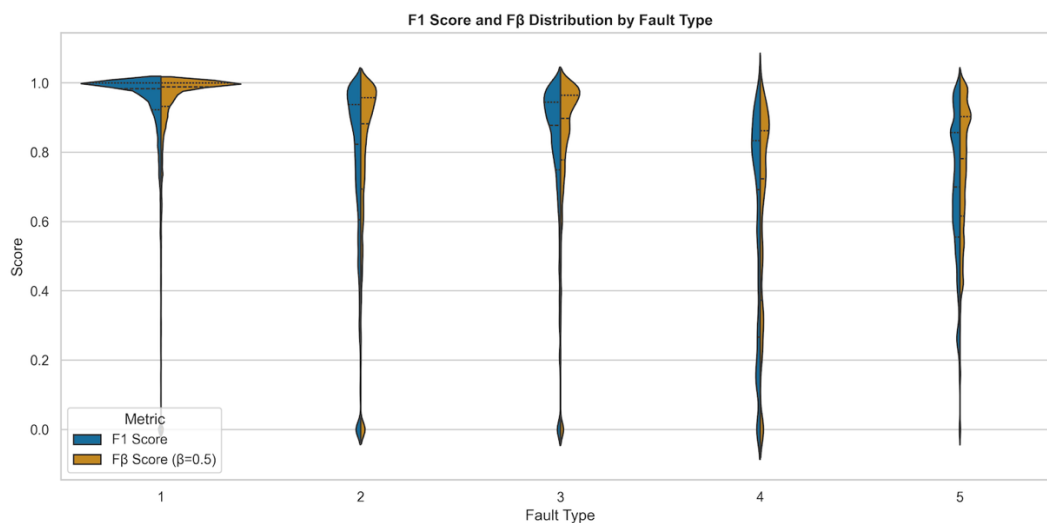


Figure 9: Distribution of F_1 scores and $F_{\beta=0.5}$ scores across each fault type under all parameter configurations, networks, and events where n measurable outliers > 1 . Fault type 1 was the most likely to be correctly identified using our equation, where fault types 4 and 5 were the most challenging.

Performance was poor for characterizing the portions of the perturbation window when there were no measurable outliers (see Fig 7b). When device measurements were perturbed but did not deviate categorically from the original valid measurement, they were commonly not detected. The F_1 score for fault type 1 remained the highest at 0.96, dropped to 0.33 for fault type 2, and then between 0.04 – 0.1 for the remaining faults. The higher F_1 score for fault type 1 is likely due to other devices within network not maintaining extreme $PM_{2.5}$ levels for extended periods of time.

4 Gradient-boosted decision tree detection

4.1 Overview

While the rule-based detection equation provided a straightforward and interpretable method for identifying outliers based on parameterized thresholds, its performance hinged on the perturbed device deviating substantially in measurement space, and needing to have a sufficient categorical distance between our bins. Additionally, performance varied substantially across fault types. The equation was most effective at capturing fault type 1, characterized by a persistent high measurement, but struggled with fault type 5, which produced possibly the most unrealistic pattern of $PM_{2.5}$ measurements -- that of a high frequency sinusoidal rhythm. This suggests that although the rule-based equation performs well under some common sensor fault scenarios, it may miss more nuanced anomalies or behaviours indicative of a device being faulty, or experiencing a hyper-local event.



495 In attempt to address these limitations, we sought to extend the detection framework by incorporating temporal
behaviour of the device and its local network. The objective is to capture not only measurement anomalies at the current-hour
snapshot, but to add antecedent context that may signal when a device's behaviour is diverging from its neighbours. This added
temporal context may improve detection performance throughout the perturbation window, particularly for faults that do not
deviate significantly in measurement space. For example, how can we identify a faulty device when both the $PM_{2.5}$
measurements of the device and its network are either very low or very high?

500 Machine learning is particularly well suited to this task, offering the ability to model complex, nonlinear relationships
from multivariate time series data. Gradient-boosted decision trees (GBDTs), in particular have stood out as a powerful and
commonly used class of model that excels in classification problems (Friedman, 2001; Chen et al., 2018; Bentéjac et al., 2021).
We chose to use the popular eXtreme Gradient Boosting (XGBoost) software package (Chen and Guestrin, 2016), due to its
ease of use and robust history of application to a wide range of problems in the environmental sciences and beyond.

4.2 Feature generation and Model Training

4.2.1 Features

505 Several model features were generated to characterize the recent-to-current behaviour of the measurement signal, with the idea
that naturally fluctuating $PM_{2.5}$ measurements during a smoke event might behave differently than those being recorded by a
faulty sensor. These include features on a per-device basis which measure the relative variability of $PM_{2.5}$ measurements, the
trending directionality of measurements over a time window, recent outages, the range of recent measurements within a time
window, the entropy of recent measurements within a time window, and how these values contrast with the behaviour of the
510 local network. Also included were the information theoretic measures of device entropy and information within its local
network, a true/false flag for whether our rule-based equation detected an outlier, and some network density data quantifying
how many devices existed within various proximities. The full set of features, along with the time windows used in their
generation, are listed in **Table 3**.

515

520



525

Table 3. Features generated for use in the XGBoost model, with the intention of characterizing recent measurement history in context to its neighbours, as well as inputs from the rule-based outlier detection equation described in Section 3. Time windows are referred to in hours below (3H = 3 hours, 6H = 6 hours), and a parameter set refers to the inputs specific to our information theoretic measures (AQI binning strategy, rule-based equation inputs).

530

Feature Type	Description	Thresholds/Variants
Log AQI	Log Transformed AQI	Current
Coefficient of Variation	Relative variability in AQI over rolling time windows	3H, 6H
Delta	Change from start to end of time window	3H
Range	Range (max - min) of values in a time window	3H
Directionality	Normalized AQI trend slope (Δ / mean) within the window	3H
Measurement Entropy	Rolling entropy of AQI values	3H / Specific to parameter set
Outages	Number of data dropouts within a time window	6H
Neighbour Count	Count of devices within a given proximity	7500–12500 m
Information Theoretic	Network Entropy and Information Value for device at given spatial proximity and binning type	Current / Specific to parameter set
Outlier Detection Flag	Output flag from the rule-based equation	Current / Specific to parameter set
Neighbour Deviation	Difference between device feature and median of neighbours in local network for: Log AQI / Delta / Range / Directionality / Measurement Entropy	3H rolling for all (current for Log AQI)

535

For our information theoretic features (measurement entropy, network entropy, device information, and outlier detection flag), we had to constrain our runs to use the parameter inputs discussed in **Section 3**. We trained the model using a total of 8 parameter sets, varying several factors. We used a central cluster of proximities (7500 – 12500 m), three binning strategies (EPA AQI, our custom percentile bins, and the Canadian AQHI), and both the central information theoretic threshold values as well as the mean value across best entropy/information thresholds.

4.2.2 Model Training

540

We trained several binary classification models using XGBoost to detect our synthetic faults, trying to identify measurements that had been perturbed by one of our fault methods. Each model was trained using differing combinations of features, parameterizations passed through from the rule-based equation methods used, and alternative PM_{2.5} binning strategies. All variations underwent identical training protocols, but do not need to be discussed at length. Performance of the top models will be described in the next section.

By design, the perturbed measurements (fault occurrence) in our training data were rare compared to true observations, and the dataset exhibited significant class imbalance, with only ~1% of the observations labelled as perturbations,



545 corresponding to a positive-to-negative ratio of approximately 1:90. To address this, the model was configured to apply increased weight to the minority class during training, directly proportional to this imbalance.

These models were trained and evaluated using five-fold cross-validation. In each fold, the data was split into training and validation subsets, allowing the model to be assessed on unseen data. For each fold, approximately 80% of the data were used for training and 20% for validation. We used a maximum tree depth of seven and applied conservative regularization, including subsampling of both rows and features, to reduce the risk of overfitting and encourage generalization. The models were trained using the binary logistic loss function and its performance was evaluated using multiple metrics, including the F_1 score, $F_{\beta=0.5}$, (AUC-ROC), overall classification accuracy, and the log loss of the predicted probabilities.

555 Given the class imbalance, we dynamically optimized the decision threshold for classification in each fold. This threshold was selected based on a precision-recall curve to maximize the F_1 score. To select the most representative and effective model during training, we retained the one that achieved the highest F_1 score across all folds. This approach, combined with regularization and threshold tuning, helped ensure that the final model maintained the best sensitivity to perturbations while limiting false positives.

After assessing model performance across parameter thresholds, binning strategies, and feature sets, configurations that consistently yielded the highest F_1 scores were advanced for evaluation on real-world data. The final feature sets evaluated used the features from **Table 3**, along with the two best parameter sets at spatial proximities of 7500m and 12500m.

4.3 Results on the Semi-Synthetic Dataset

We evaluated our XGBoost model performances relative to slightly different considerations than the information-theoretic equation. We considered its ability to detect the perturbed device at least once during the event, but also how well it classifies any perturbations (both measurable and non-measurable outliers) during the window of time where we altered a device's measurements. The goal is for the added contextual information from antecedent measurements and neighbouring devices to train and inform our ML model when devices are deviating in behaviour, independent of the categorical rulesets that constrain our rule-based, information-theoretic equation. We limit our evaluation of the training results to avoid broader claims about model performance across real-world conditions and reduce the risk of over-interpreting behaviour that may be biased toward the smoke events included in training.

570 After applying a classification threshold of 0.95 on the model's predicted probability (only flagging perturbations when the XGBoost model assigned $\geq 95\%$ probability to the positive class), we assess preliminary results. Across the induced faults and all trained models, the outlier was successfully identified in 96.4% of time series events. When broken down by fault type, detection rates varied moderately, and the performance, when segmented by fault type, aligned with the results from our information-theoretic equation. Fault type 1 was detected in 97.5% of events. Fault types 2 and 3 followed, with detection rates of 88.4% and 87.9%, respectively. Fault type 4 was the most challenging, with a lower success rate of 70.4%, and fault type 5 had 96.7% detection rate. Across all configurations, events and fault types, the mean false positive rate was far higher than for our information-theoretic flag (15.2%).



580 These results are indicative of a strong reliance on the information-theoretic outlier flag, and other information-theoretic inputs. In aggregate across all model variants, the most influential feature predictors were the rule-based equation outlier flag and the information value of the unit, followed by short-term variability measures (3-hour AQI range) and neighbour-based difference features. Feature importance is based on mean gain, which measures how much a feature improves the model's ability to correctly separate perturbed from non-perturbed measurements when it is used in a decision split.

585 Across all trained XGBoost models, we see significant advantages in ability to classify the perturbed measurements that did not meet the measurable outlier qualifications that constrained the rule-based equation. The XGBoost classifier significantly outperformed the rule-based detection equation in characterizing the perturbation window as a whole. Because the XGBoost learns patterns from a broader feature set, it could identify perturbed periods even when deviations were too subtle to qualify as measurable outliers under rule-based thresholds. Without the manual outlier constraint, the highest-performing XGBoost model achieved an F_1 score of 0.88 for accurately detecting a perturbation, compared to 0.42 for the rule-based method, while achieving a false positive rate below 1%. Interestingly, the difference between the methods flipped when 590 we limited our performance exclusively to the perturbation window where n measurable outliers > 1 (approximately 67% of the synthetic perturbations). In this case, our best XGBoost model had an F_1 score of 0.45, and our rule-based detection equation had an F_1 of 0.87. This discrepancy likely reflects differences in design intent: the XGBoost model was optimized to recover all labelled perturbations, including weak or ambiguous cases, whereas the manual method primarily flags perturbations with major categorical deviations. These findings confirm the semi-synthetic training regime can produce a generalizable signal, 595 but the operational utility of this model ultimately depends on its performance when applied to real sensor networks. Without a comprehensive manually validated set of known faulty sensors, synthetic perturbations can only approximate real-world conditions, and we remain cautious about overfitting to artifacts of the simulated faults.

5 Application to Continental-Scale Operational Data

5.1 Overview

600 To understand the broader utility of our methods, we expanded the scope beyond the 11 short-duration smoke events that were used in parameter tuning and model training. Here, we applied four parameterizations of the information-theoretic equation, and four parameterizations of the corresponding pre-trained XGBoost model (**Table 4**) to all hourly sensor and monitor data from the sources described in section 2, over the United States (including Alaska and Hawaii), and Canada. This data was collected from January 1 2025 through September 30 2025, totalling to 6,529 hours of air quality measurements, including 605 19,562 unique $PM_{2.5}$ devices (devices that recorded measurements at least once during this period). 74 known faulty sensors were removed from the dataset beforehand. For each device reporting a measurement at a given hour, we collected its network measurements and derived the appropriate statistics and features to include in the information-theoretic detection equation and the pretrained XGBoost model. If any method produced a positive outlier flag, that device-hour was retained as a candidate event for validation.



610 These 9 months of data span both network-scale baseline conditions and smoke events of varying magnitudes. This provided an opportunity to assess how our methods perform on real-world sensor networks, not only under conditions similar to those in training, but across diverse smoke events and sensor-network structures. This allows us to evaluate whether the tuning done with semi-synthetic perturbations generalizes to the variability present in an operational setting.

615 The information-theoretic parameters chosen represent a set that was stable and performed well across the various fault types and events for both our rule-based equation and XGBoost models. We used a high prediction-probability threshold of our XGBoost model to attempt limiting false positives.

Table 4. Outlier flag methods used in validation scheme.

Name	Flag Type	Info-Theory Parameters	Prediction Probability
info_7500_std	Information-Theoretic	$\theta = 1.75$ $S = 1.75$ $B = 6$ $\beta = 3.5$ $r = 7.5km$	N/A
info_7500_pct	Information-Theoretic	$\theta = 1.75$ $S = 1.75$ $B = 12$ $\beta = 3.5$ $r = 7.5km$	N/A
info_12500_std	Information-Theoretic	$\theta = 1.75$ $S = 1.75$ $B = 6$ $\beta = 3.5$ $r = 12.5km$	N/A
info_12500_pct	Information-Theoretic	$\theta = 1.75$ $S = 1.75$ $B = 12$ $\beta = 3.5$ $r = 12.5km$	N/A
xgb_7500_std	XGBoost	$\theta = 1.75$ $S = 1.75$ $B = 6$ $\beta = 3.5$ $r = 7.5km$	0.98
xgb_7500_pct	XGBoost	$\theta = 1.75$ $S = 1.75$ $B = 12$ $\beta = 3.5$ $r = 7.5km$	0.98
xgb_12500_std	XGBoost	$\theta = 1.75$ $S = 1.75$ $B = 6$ $\beta = 3.5$ $r = 12.5km$	0.98
xgb_12500_pct	XGBoost	$\theta = 1.75$ $S = 1.75$ $B = 12$ $\beta = 3.5$ $r = 12.5km$	0.98

620 To assess detection accuracy and characterize outlier types, we manually validated a stratified random sample of flagged events spanning consecutive-hour bins. These consecutive-hour bins, which we may refer to as persistence, are the duration of consecutive hours a device was flagged for by a given detection method (for example, info_7500_pct). Validation was conducted by hand using the best judgment possible, though we acknowledge this approach introduces inherent bias regarding what constitutes a malfunctioning monitor or sensor. Visual inspection utilized diagnostic plots showing (1) 24-hour AQI time series for the flagged device alongside its network neighbours, and (2) spatial context maps displaying network device locations coloured by contemporaneous AQI values. For each detection, we examined time series patterns, spatial coherence with neighbouring devices, nearby wildfire activity (using satellite fire detections, known fire locations, and smoke plume data) to assign two independent labels: an outlier classification (misclassified, hyper-local event, faulty sensor, or unknown) and a smoke status assessment (no smoke, local smoke, regional smoke, or unknown). In each of these events, we judged to the best of our ability whether there was smoke from a wildland fire source present.

630 We defined a hyper-local event, as a single devices' sharp deviation in behaviour from neighbouring devices, indicative of either a short-lived device malfunction, or PM_{2.5} conditions that would not be relevant to the population the device may represent, and perhaps solely of interest to the device owner (example shown in Fig. 12). These may include debris being



635 lodged in the sensor, a sensor being intentionally placed near an industrial emissions source, a bonfire, or even a sensor being moved indoors to monitor the operators air-purifier performance during a smoke event.

Our smoke presence assessment is defined as follows. No smoke event implies the absence of any smoke event, although ambient pollution sources, dust events, or local woodstove emissions may be present. Local smoke events were categorized as those where part of a network was being impacted by a smoke event, though some of the network was unaffected. Regional smoke events were defined as a smoke event that was impacting at least the entire network, although the impact does not need to be equal across the network.

645 Events were sampled to ensure representation across consecutive-hour bin durations. The validation process iteratively sampled unlabelled events while excluding previously classified detections, ultimately yielding 465 manually labelled episodes stratified across both detection methods and persistence categories. This retrospective analysis was conducted from remote, after-the-fact inspection. There is no way to definitively verify whether a given device was truly malfunctioning without physical inspection or understanding on-site conditions at the time of flagging. Despite these limitations, the standardized visual protocols and explicit labelling criteria help reduce uncertainty in our ground-truth determinations.

5.2 Detection Method and Parameter Performance

650 The manual validation of 465 detection events reveals distinct performance characteristics between the information-theoretic and XGBoost methods across different episode durations (**Table 5**). Unlike the model development phase, which employed F_1 scores, the validation analysis focuses on positive predictive value: the percentage of flagged events representing true outliers. Misclassification rates vary considerably with episode persistence: single-hour events show elevated false positive rates (31-51%), but accuracy improves rapidly with even modest persistence. By 3 consecutive hours, misclassification rates drop to 6% for the information-theoretic method and 15% for XGBoost, demonstrating that brief persistence filtering substantially improves precision. This relationship suggests two complementary approaches to reducing false positives: applying persistence thresholds that require multiple consecutive flagged hours, or extending the temporal averaging window beyond one hour (e.g., 3-hour rolling averages) to smooth transient fluctuations. For sustained episodes with persistence beyond 13 consecutive hours, both methods achieve performant accuracy, with 95-98% of detections representing genuine outliers (either faulty devices or hyper-local events). The information-theoretic method exhibits particularly high specificity for sensor malfunctions in longer episodes (98% faulty devices for 24+ hours), while XGBoost demonstrates more balanced detection across outlier categories.

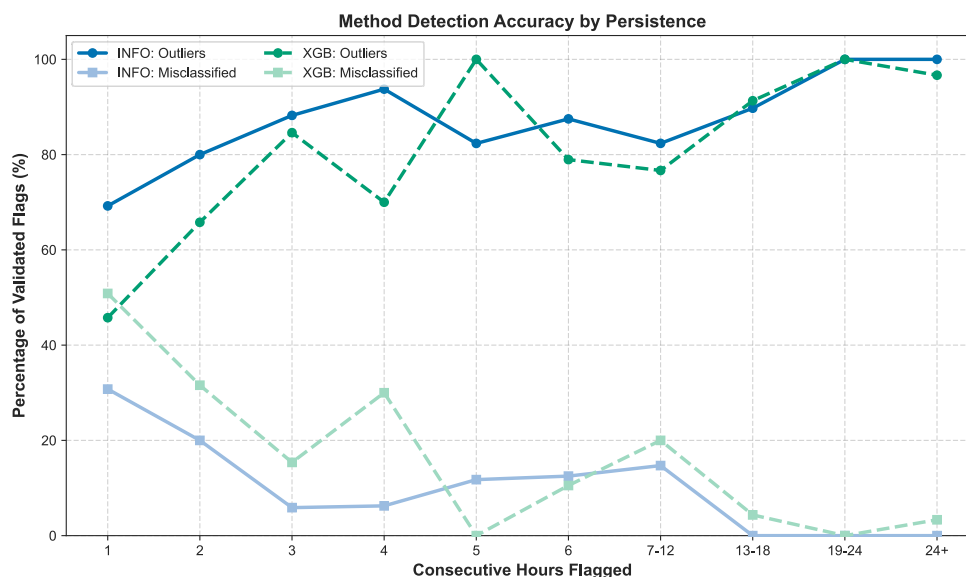
665 The duration-stratified validation results reflect fundamental characteristics of true sensor malfunctions and localized air quality phenomena. Short-duration events (2-6 hours) frequently capture legitimate hyper-local air quality variations, with 53-76% classified as hyper-local events in the 2-5 hour bins, suggesting these brief anomalies often represent real localized phenomena. Short episodes may also include intermittent hardware issues that self-resolve, such as temporary obstructions (debris, insects) or environmental interference. In contrast, genuine sensor malfunctions typically persist for extended periods, the predominance of broken sensors in episodes exceeding 13 hours (86-97%) reflects the sustained nature of hardware failures.



670 These findings indicate that even modest persistence thresholds (≥ 3 consecutive hours) effectively reduce misclassifications, while longer thresholds (≥ 12 -13 hours) reliably isolate sensor malfunctions and sustained outlier behaviour, with the caveat that brief flagging events may represent either transient real-world phenomena or self-resolving instrument issues rather than algorithmic error.

675 **Table 5.** Summary of validation performance for consecutive hour bins. Each consecutive hour category represents an event where a device was flagged for that duration continuously. Values show the total number of events flagged by each method, the subset that were manually validated, and the distribution of validation outcomes: hyper-local events (short-duration hardware faults or extreme environmental conditions), faulty devices, unknown, and misclassified detections (functional devices incorrectly flagged).

Consecutive Hours flagged	Information-theoretic detection					
	Number of events flagged	Number of events manually validated	Distribution of validated events			
			Hyper-local event (%)	Faulty Device (%)	Unknown (%)	Misclassified (%)
1	2197	13	53.85	15.38	0.00	30.77
2	1628	15	66.67	13.33	0.00	20.00
3	1040	17	76.47	11.76	5.88	5.88
4	632	16	68.75	25.00	0.00	6.25
5	427	17	52.94	29.41	5.88	11.76
6	240	16	56.25	31.25	0.00	12.50
7-12	474	34	38.24	44.12	2.94	14.71
13-16	124	39	25.64	64.10	10.26	0.00
17-24	47	36	2.78	97.22	0.00	0.00
24+	91	41	2.44	97.56	0.00	0.00
Consecutive Hours flagged	XGBoost detection					
	Number of events flagged	Number of events manually validated	Distribution of validated events			
			Hyper-local event (%)	Faulty Device (%)	Unknown (%)	Misclassified (%)
1	8961	59	32.20	13.56	3.39	50.85
2	5283	38	47.37	18.42	2.63	31.58
3	2528	26	76.92	7.69	0.00	15.38
4	1102	20	45.00	25.00	0.00	30.00
5	578	10	50.00	50.00	0.00	0.00
6	331	19	47.37	31.58	10.53	10.53
7-12	443	30	20.00	56.67	3.33	20.00
13-16	81	23	4.35	86.96	4.35	4.35
17-24	23	14	7.14	92.86	0.00	0.00
24+	86	30	0.00	96.67	0.00	3.33



680 **Figure 10:** Accuracy of both information-theoretic detection and XGBoost detection across persistence bins. For each consecutive hour that a device is repeatedly flagged as an outlier by either method, accuracy increases that it is either a faulty device or experiencing a hyper-local event.

Standard and percentile AQI binning approaches demonstrated nearly identical performance when flagged
 685 independently (78.4% vs 76.7% correct classification). However, increased persistence substantially improved both approaches: filtering for events persisting beyond 3 consecutive hours reduced misclassification from 21-23% to 11.3% for either binning approach alone. Information-theoretic flags consistently outperformed XGBoost across all temporal scales, exhibiting 7.2% misclassification compared to 26.4% for XGBoost-only detections. This performance gap narrowed substantially when filtering for persistence > 3 hours (5.2% vs 12.6%). The elevated misclassification rate for XGBoost-only
 690 detections may reflect the method's sensitivity to patterns that are discussed in the next section. Information-theoretic flags maintained robust performance even for shorter-duration events.

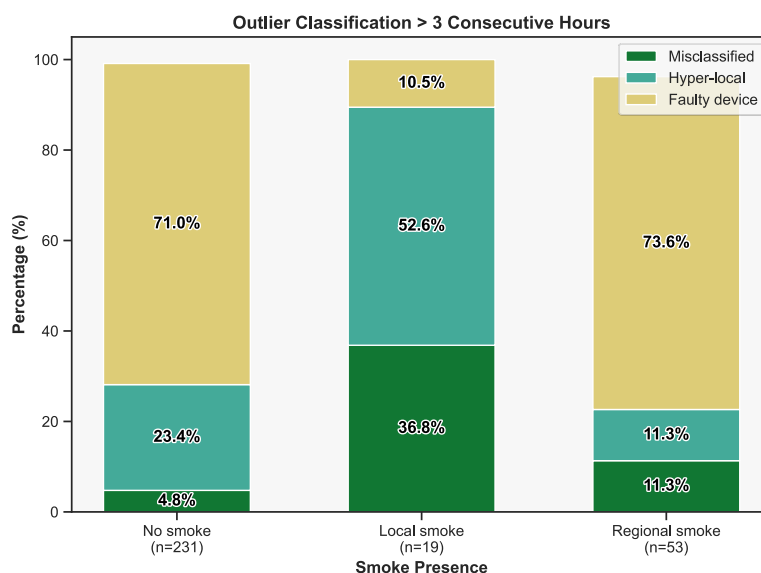
While individual parameter choices influenced performance, consensus between different parameter sets provided a strong validation signal. Outliers flagged by multiple methods, binning strategies, or network scales demonstrated substantially reduced misclassification rates across all configurations, suggesting that methodological agreement may offer robust
 695 validation. Consensus between AQI binning parameters demonstrated the strongest effect: events flagged by both standard and percentile binning approaches achieved 96.1% correct classification (n=152). AQI binning strategy consensus combined with increased persistence (> 3 hours) demonstrated exceptional performance (99.1% correct classification, n=106), since transient anomalies (< 3 hours) constitute the primary source of error in our classification.



5.3 Performance During Smoke Events

700 The following smoke-related findings based on the 465 manually validated detections had an imbalanced distribution across smoke conditions: no-smoke events comprise the majority of validations (76%, n=352), while regional smoke events represent more limited coverage (13%, n=60), and local smoke events are notably underrepresented (7%, n=32). This imbalance, particularly for local smoke, limits interpretation of observed patterns. This is further complicated by local smoke being the most ambiguous case to define from retrospective observation, and by its labelling criteria, it is characterized by
 705 network disagreement and transient conditions. Wood smoke sources were not included in the local smoke category, as we tried to identify only wildland fire sources.

While considering these limitations, we explore whether the presence of smoke may potentially influence both detection accuracy and the distribution of outlier types. Misclassification rates show variation by smoke context: local smoke events exhibit high false positive rates (50%, n=32), while regional smoke (15%, n=60) and conditions without smoke events
 710 (14%, n=352) display substantially lower misclassification. After increasing persistence to > 3 hours, misclassification drops for each smoke presence scenario, with local smoke still exhibiting the highest false positive rates (36.8%, n=19), and regional smoke (11.3%, n=53) and no smoke (4.8%, n=231) dropping substantially (**Fig. 11**).



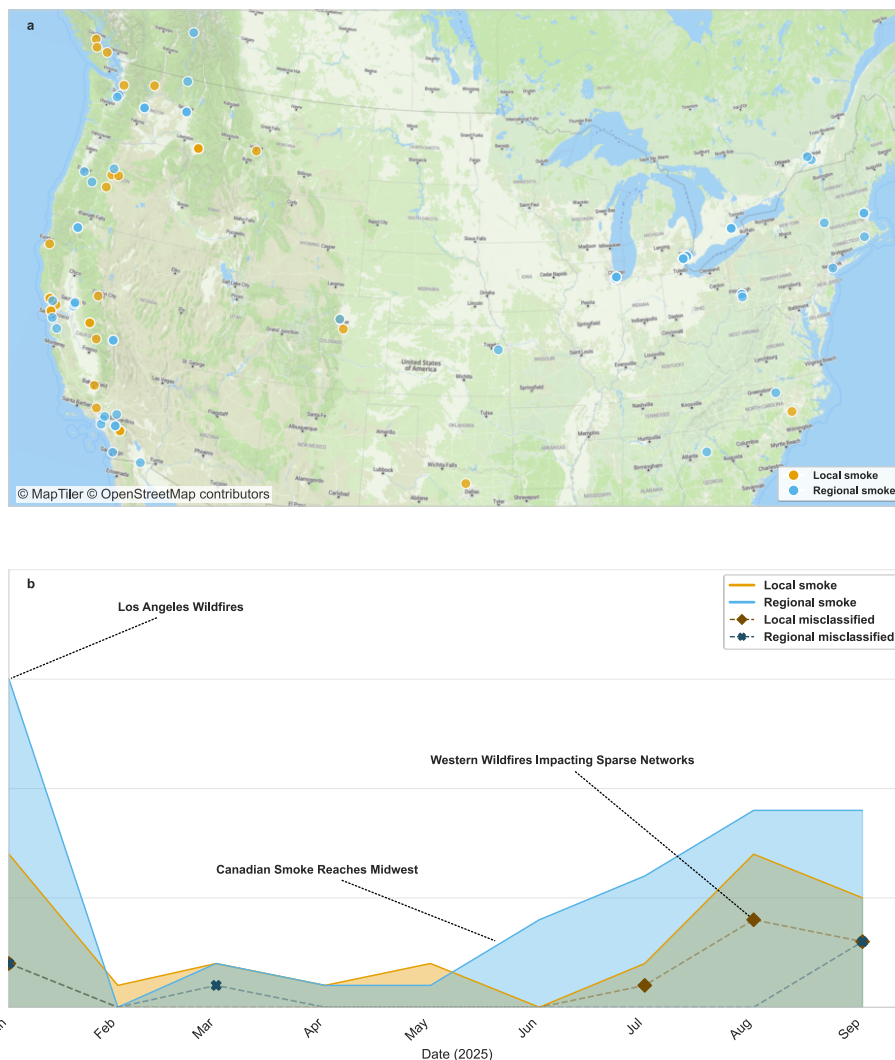
715 **Figure 11:** Distribution of outlier classifications that persist > 3 consecutive hours, based on whether there was suggestion of wildland fire smoke present. The no smoke class represents typical ambient conditions, local smoke is defined by a portion of the network being impacted from smoke, and regional smoke indicates that there is a smoke event that is larger than the network occurring.

720



In the selected cases with a persistence > 3 hours, local smoke was more often labelled in networks with approximately half the neighbour density of other smoke conditions (12.8% vs 23.6% average neighbours per network at 7500 m), suggesting these events are concentrated in more isolated or rural locations, closer to fire locations, or that these sparser networks did not fully capture the extent of the smoke. This sparsity is consistent across network scale, though less pronounced at 12500 m
725 (23.7% vs 34.0% neighbours), indicating that the error rate may be structural in part. This pattern suggests that smoke events that are not impacting this entire network may be more difficult to distinguish from sensor anomalies using these methods, whereas widespread regional smoke provides clearer spatial context for validation. When examining the network proximity used by our detection methods, events flagged exclusively by the 7500 m network showed 0% misclassification ($n=3$), while those flagged only by the 12500 m network exhibited 50% misclassification ($n=10$), and events flagged by both networks
730 showed intermediate rates (33%, $n=6$). This pattern suggests that local smoke phenomena may be more reliably validated when detected by denser networks, whereas farther devices introduce uncertainty.

When viewing the spatial and temporal distribution of the manually validated detections that had a smoke type association in Figure 12, these patterns become visible: local smoke events cluster in the mountainous West, while regional smoke events are more geographically dispersed. The temporal distribution shows misclassification rates (dashed lines)
735 tracking more closely with local smoke during the late-summer wildfire season, reinforcing that spatially limited smoke events in sparse networks pose the greatest challenge for accurate classification. However, given the limited sample size of local smoke events ($n=32$), this relationship warrants further investigation with a larger validation dataset.



740 **Figure 12:** Spatial and temporal distribution of manually validated detections that had an associated smoke event. (a) Geographic locations
of local smoke (orange) and regional smoke (blue) events. (b) Monthly counts of manually validated detections, constrained to be unique by
both device ID and smoke type to avoid duplicates. Dashed lines indicate misclassifications (false positives). Annotations highlight notable
smoke episodes during the study period (January–September 2025). (Basemap: © MapTiler; map data © OpenStreetMap, distributed under
the Open Data Commons Open Database License (ODbL) v1.0; for further information see <https://www.openstreetmap.org/copyright/en>)

745

Continuing with our cases with persistence > 3 hours, faulty sensors were more likely to be detected in either the no
smoke or regional smoke cases. Within our flagged events with no smoke present, 71% of flagged events represent genuine
sensor malfunctions, rising to 74% during regional smoke events but dropping to only 11% during local smoke conditions.
This pattern suggests that regional smoke may amplify the detectability of a malfunctioning device, while local smoke events
750 may mask sensor failures or are themselves frequently misclassified as sensor problems. Hyper-local classification shows



variation across smoke conditions (11-23% for regional smoke and no smoke; 53% for local smoke), indicating that truly localized phenomena might be identified regardless of local smoke conditions.

6 Discussion: Case Studies, Limitations, and Improvements

6.1 Separating Hyper-Local Events from Faulty Sensors

755 Both outlier detection methods were trained on synthetic measurement signals that arise not only from malfunctioning devices, but also from PM_{2.5} conditions that affect only the sensor in question. In reality, the hyper-local events validated in our 9 months of air quality data may be due to siting near a localized emission source (**Figure 13**), indoor placement during a smoke event (**Figure 14**), or other intermittent hardware issues that self-resolve. These events manifest in a single device and may not be representative of ambient air quality relevant to the local population. The distribution of classification type across persistence bins in **Table 5** demonstrates both methods' abilities to separate these hyper-local phenomena from truly faulty devices that provide no usable data.

760 For the information-theoretic method, the first persistence bins with acceptably low misclassification rates are those lasting 3 hours or more, where hyper-local events account for ~76% of validated detections and device faults remain comparatively uncommon. With more consecutive hours of flagging, the classification shifts decisively: once events persist beyond 13 hours, validated cases are overwhelmingly attributed to faulty devices, exceeding 90%. XGBoost shows the same transition. At 3 hours, hyper-local events represent ~77% of validated cases with misclassification well below 20%. But for longer-duration detections beginning around 13 hours, faulty devices dominate, rising to >85–90% of validated outcomes. Across both methods, these low-misclassification bins exhibit a consistent pattern: short-duration detections are primarily hyper-local phenomena, whereas long-duration detections correspond almost exclusively to a broken or faulty device. By stratifying flagged outliers by their persistence, the methodology provides a practical means of separating hyper-local events from faulty devices.

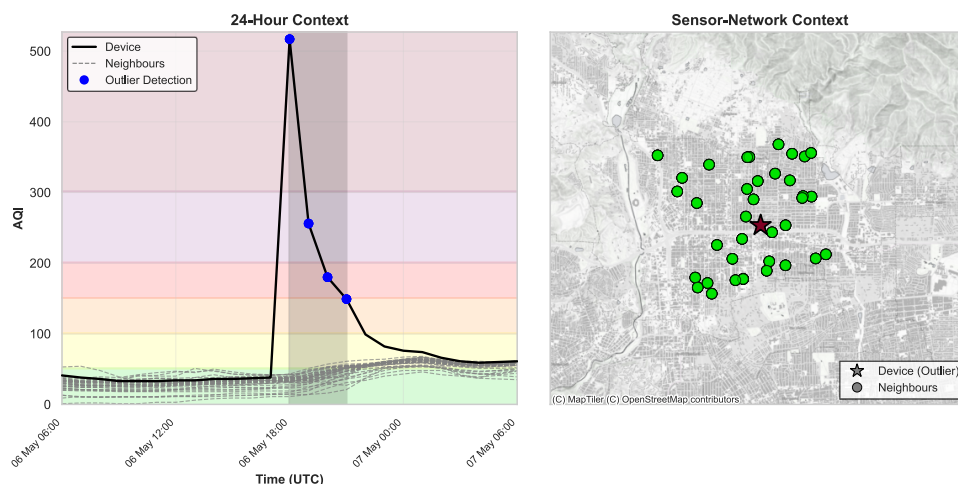
775 This stratification offers an operational lever for tailoring how detected outliers may be interpreted and used. For air-quality messaging purposes that prioritize a unified, population-level communication, a lower persistence threshold (e.g., > 3 hours) may be appropriate. At these durations, a large fraction of hyper-local events and almost all faulty devices may be filtered out while misclassification rates remain very low. This ensures that public messaging reflects regional air-quality conditions rather than transient, sensor-specific anomalies.

780 For applications that benefit from retaining fine-scale PM_{2.5} variability, such as neighbourhood level exposure assessments or high-resolution monitoring, a higher persistence threshold (e.g., >12 hours) is more suitable. At this threshold, misclassification is very low, and the vast majority of flagged devices are due to hardware or persistent faults. Configuring parameters and persistence thresholds thus allow this detection framework to be tailored to specific operational contexts, balancing trade-offs between public health messaging and fine-scale PM_{2.5} exposure.



785 Broken sensors exhibit significantly lower coefficient of variation ($CV = 30.5\%$, median = 11.5%) compared to hyper-local events ($CV = 67.6\%$, median = 52.4%), reflecting their tendency to report constant or near-constant values. This distinction in time series measurement variability may provide an additional signal for differentiating malfunctioning sensors from hyper-local events, complementing persistence thresholds.

Hyper-local event



790 **Figure 13:** Example plot showing what was classified as a hyper-local event. Colours in the right figure represent EPA AQI category. This device deviated from ambient conditions in an extreme way, and nearby devices remained unaffected. The event lasted roughly ~6 hours, with 4 hours being flagged by all the detection methods. (Basemap: © MapTiler; map data © OpenStreetMap, distributed under the Open Data Commons Open Database License (ODbL) v1.0; for further information see <https://www.openstreetmap.org/copyright/en>)

Sensor placed indoors

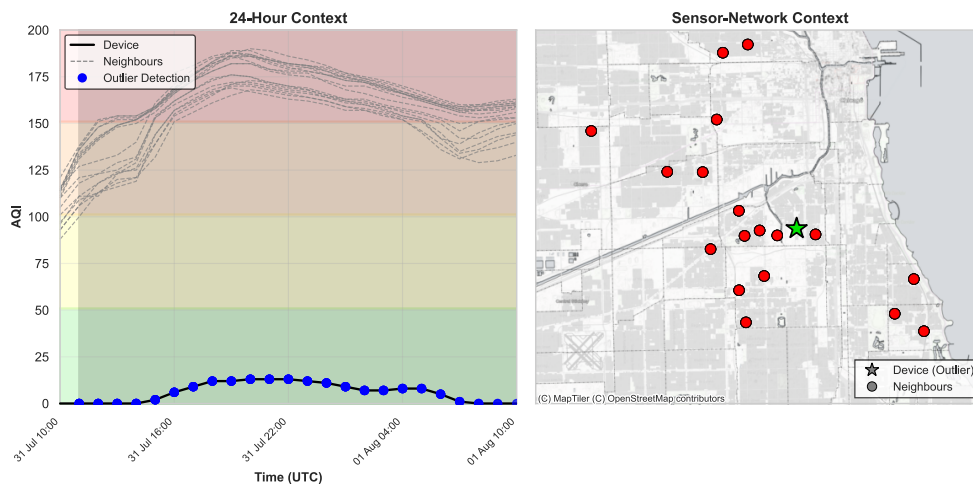




Figure 14: Example plot showing a classified hyper-local event of a sensor potentially moved indoors. Colours in the right figure represent EPA AQI category. During a regional smoke event, this device experienced significantly muted concentrations compared to other devices. As a publicly owned device, it was likely misclassified as an outdoor unit or brought inside. This outlier was only made apparent during elevated $PM_{2.5}$ concentrations and flagged as soon as neighbouring sensors deviated categorically. (*Basemap: © MapTiler; map data © OpenStreetMap, distributed under the Open Data Commons Open Database License (ODbL) v1.0; for further information see <https://www.openstreetmap.org/copyright/en>*)

6.2 Network Considerations

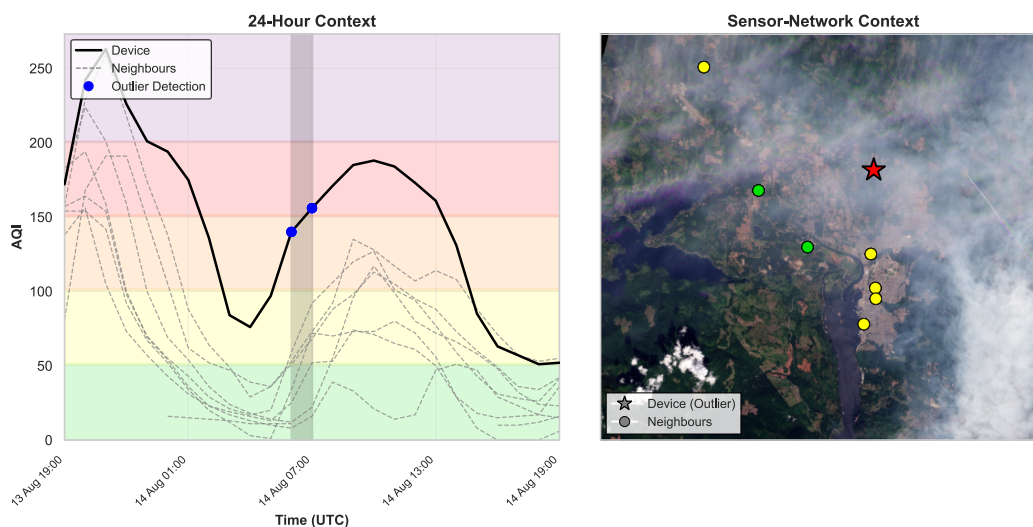
Defining sensor networks that are appropriately grouped, such that they behave predictably in response to a pollution source, is a significant challenge. Wildland fire events introduce unique smoke regimes that are tightly coupled to fire behaviour, local meteorology, and surrounding topography. As a result, each event produces novel $PM_{2.5}$ patterns across a network, and any previously established predictive relationships within networks may become irrelevant under new conditions.

To guard against unwanted outlier flagging, we must consider the sharp gradient of $PM_{2.5}$ concentrations may take across a network during a local smoke event, and tune our detection method parameters accordingly. In the local smoke case shown in **Figure 15**, the information-theoretic detections misclassified a valid sensor for 2 consecutive hours due to a nearby fire. To avoid potentially filtering out valid data, we might lower our network entropy threshold, increase our persistence requirements, or leverage the CV on recent measurements, using dynamic behaviour to verify.

Interestingly, the XGBoost method did not classify this event as an outlier; while speculative, this may be the use of antecedent measurement features used in the model, which assess whether a device's recent behaviour is both dynamic and consistent with that of neighbouring sensors. Future work could expand the machine learning training datasets to include either real or synthetic events of this nature that reflect transient smoke that may temporarily create a sharp gradient across a network. It might also be valuable to incorporate additional data sources profiling nearby fire and emissions sources.



Local-Smoke Misclassification



815 **Figure 15:** A misclassification due to local smoke from the nearby Mount Underwood fire near Port Alberni, British Columbia. Colours in the right figure represent EPA AQI category. This device was incorrectly classified as an outlier for 2 consecutive hours by information-theoretic detection methods. Nearby devices show similar trends in $PM_{2.5}$ concentrations, however this device was located where smoke was mixing down to the surface. Satellite imagery taken during the detection period shows the smoke plume. Image © 2025 Planet Labs PBC

820 Another complication arises when networks are grouped across drainages, which may function as separate airsheds. When the sensor distribution is sparse, imbalanced between drainages, or when the $PM_{2.5}$ source is localized, a Euclidean-based distance threshold, such as the one used in this study, may generate disorderly networks. Preliminary work explored the creation of terrain-aware networks using a least-cost path approach that used slope as a cost surface. Although this approach increases computational burden and complexity to network generation, early results indicate that these terrain-informed

825 networks exhibit stronger internal consistency among sensors, suggesting improved ability to reduce misclassifications in topographically complex regions.

Figure 16 illustrates this scenario, which shows a network of sensors generated by both methods, using a 7.5 km distance threshold. The primary node was frequently misclassified as an outlier in the winter months, when overnight woodsmoke pollution would settle in the western drainage and remain confined there. The lone sensor in the eastern drainage, left unimpacted, would be flagged using our detection methods if our network was comprised of cross-drainage devices. When terrain is incorporated into network generation, this device's neighbour count is reduced to the single station at the southern

830 end of its own drainage, which appropriately separates the distinct airsheds.

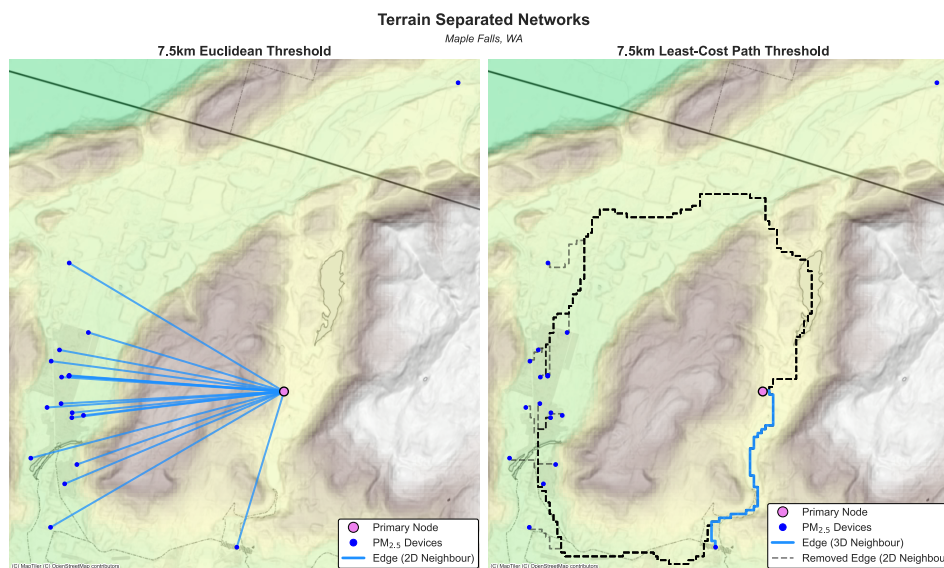


Figure 16: Example of a sensor network in which the primary node was frequently misclassified as an outlier. The left panel shows a network constructed using a Euclidean distance threshold, whereas the right panel shows a terrain-aware network generated using a least-cost path approach with a DEM as the cost surface. This difference in network construction reduces the sensor’s neighbour count from 16 to 1, preventing it from being flagged as an outlier in situations where air-quality impacts are confined to the neighbouring drainage. (Basemap: © MapTiler; map data © OpenStreetMap, distributed under the Open Data Commons Open Database License (ODbL) v1.0; for further information see <https://www.openstreetmap.org/copyright/en>. Terrain shading derived from NASADEM (NASA JPL, 2020))

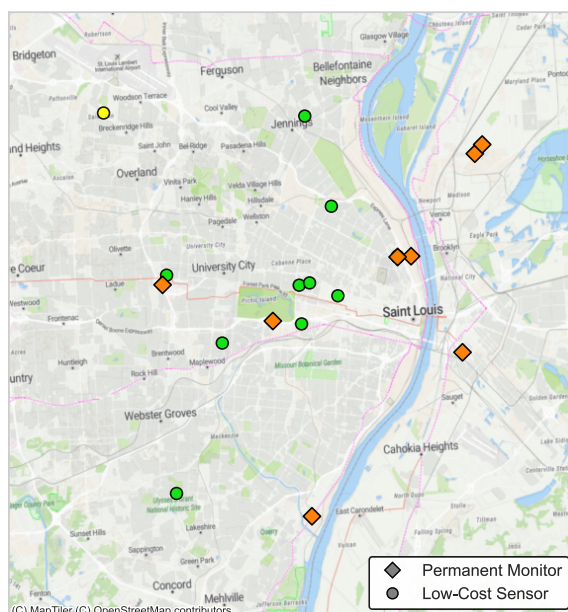
840 6.3 Sensitivity to Different PM_{2.5} Sources

The PM_{2.5} measurement technologies unique to each device type may respond differently to specific types of particulate sources, particularly when particle size distributions shift toward coarse diameters. Regulatory grade (FEM/FRM) monitors reliably capture PM_{2.5} concentrations during both dust storms and wildfire smoke, but many optical sensors common in networks such as PurpleAir (which use the Plantower PMS5003/6003) struggle to properly size and measure larger particles (roughly > 1µm) (Ouimette et al., 2024; Kuula et al., 2020), leading to an underestimation when coarse particles dominate (Jaffe et al., 2023). While steps are being taken towards dust-specific correction algorithms for PMS sensor data (Kaur et al., 2025), this remains a current limitation of our methodology, which assumes a relatively uniform response to particulate matter between network devices.

When a majority of sensors in a network share the same Plantower hardware, and a small number of nearby regulatory grade monitors respond correctly to a dust event, the “correct” devices can become outliers. This configuration can undermine methods that rely on neighbourhood agreement. As **Fig. 17** illustrates, a dust event near St. Louis produced elevated PM_{2.5} measurements at regulatory monitors while the surrounding PurpleAir sensors remained low, creating a pattern that our outlier detection methods could misinterpret as device malfunction. Although this mismatch is less problematic when dust and smoke



855 co-occur, it highlights a fundamental limitation: network-informed detection assumes that sensors and monitors respond similarly to the same pollutant source.



860 **Figure 17:** Dust event impacting St. Louis, Missouri. Different device types may respond to PM_{2.5} sources differently. In case where there are an overwhelming majority of lower-cost sensors compared to monitors, this would cause the valid measurements to be classified as outliers. (Basemap: © MapTiler; map data © OpenStreetMap, distributed under the Open Data Commons Open Database License (ODbL) v1.0; for further information see <https://www.openstreetmap.org/copyright/en>)

6.4 Sensor Type Considerations

865 This work evaluated outlier detection on a network of data from permanent monitors, temporary monitors, PurpleAir sensors, SensOR sensors, and SensWA sensors. All devices except for the PurpleAir sensors are solely operated and maintained by air quality experts. Sensors that are not set up by air quality experts may not be sited correctly leading to outlier results because they are located near a hyper-local emission source, the location may not be accurately reported, or the sensor may have been moved inside to understand indoor air quality without updating the location or sensor type. This may have led to many of the issues identified in the operational dataset. There are many other air sensors deployed across the United States with some having similar user deployment methods to PurpleAir, some deployed by users but applying further quality assurance (e.g., requiring a photo of the sensor siting), some deployed by state, local and tribal air quality agencies, and others deployed by community groups with expert support. Other networks may have much less frequent siting issues than PurpleAir depending on who and how they are set up.



PurpleAir removes problematic data based on duplicate measurements. SensOR and SensWA sensors also have duplicate
875 sensors to more easily identify issues. Many other sensors on the market do not have duplicate measurements (Barkjohn et al.,
2025a) potentially making these outlier detection methods useful for other air sensor types. In addition, different types of
sensors may have different common failure patterns than PurpleAir (e.g., repeat zeros, outliers, baseline shift, influences
leading to variable relationships with air monitors) (Barkjohn et al., 2025a). They may also have different outlier patterns due
to different sampling or averaging intervals (e.g., < 1 hour sampled for solar devices). These outlier removal methods may
880 perform better or worse for other sensor types depending on their common failure patterns.

Many sensors and monitors have varying uncertainty depending on the measurement type and the concentration range
(Agrawal et al., 2025; Khan et al., 2024; Hagler et al., 2022). Sensor response often decreases at high concentration due to
saturation or particle coincidence (Agrawal et al., 2025; Kim et al., 2025) and uncertainty can be further influenced by the
correction applied to sensor data. The correction applied to PurpleAir sensors in this paper was developed using collocation
885 PM_{2.5} data up to 1500 µg/m³ and it accounts for nonlinearity above 300 µg/m³ (Barkjohn et al., 2022). However, uncertainty
may still be higher at these elevated concentration levels compared to lower concentrations, where more training data is
available and monitor measurements are more reliable. It is unknown what concentrations range was used in developing the
SensWA corrections and it will vary since many different corrections were used. Although a median regression is used for
PM_{2.5} data > 100 µg/m³ it is unknown if it will extrapolate accurately to any data outside the concentration range used for
890 correction development and/or if higher concentration data may be underestimated. SensOR use a linear fit so higher
concentration data may be underestimated. If we are most interested in AQI category estimation, the highest break point (Very
Unhealthy to Hazardous) is 225.5 µg/m³. So even if some sensor types underestimate at concentrations > 225.5 µg/m³, it will
not matter for AQI categorization. However, variations in correction and uncertainty across concentration ranges may result
in differing outlier detection performance among sensor types.

895 The demonstrated outlier detection framework offers some resilience to these differences through the bin deviation
measure D_k shown in Eq 3 and the flexibility of binning strategy (e.g., standard AQI, our modified percentile AQI, custom).
These can be tailored to the variability of a given sensor type or network, particularly at concentration extremes where
correction differences are largest, but bin categorization or AQI classification differences are minimal.

900 6.5 Expanding Machine Learning Applications

Although our machine learning method produced substantial “chatter,” with a high volume of detections overall, and
approximately half of the 1-hour detections ultimately validated as misclassifications, this behaviour likely reflects limitations
in the chosen parameter sets and the perturbation types used in the semi-synthetic training dataset. Retraining the model on a
broader context of event types, including human validated events would likely improve classification performance and reduce
905 false positives. Despite these limitations, results in the higher persistence range demonstrate that machine learning models
trained on semi-synthetic data can achieve strong performance.



Importantly, the XGBoost methods identified several types of faulty behaviour that the information-theoretic approach could not capture. The models were able to (a) anticipate emerging outliers based on antecedent device behaviour before categorical divergence from neighbouring sensors occurred, and (b) detect problematic signals associated with device malfunction, such as unchanging measurements, implausible variability, or frequent outages. These capabilities suggest that machine learning methods may serve as an advanced mode of detection, particularly when trained on a more comprehensive set of real and synthetic fault scenarios.

The capacity of machine learning to leverage a broader and more diverse set of features than rule-based methods is a significant advantage. With quality labelled training data (which requires considerable curation effort), machine learning models could also incorporate contextual features, such as nearby fire activity or seasonal patterns, to assess whether flagged readings reflect genuine smoke influence rather than sensor malfunction, or distinguish between hyper-local events and true device faults. However, these capabilities come at the cost of interpretability. Rule-based equations offer auditable decision logic, whereas machine learning models are less transparent, making them better suited to contexts where explainability is not imperative. Future work should focus on improved training datasets and targeted feature engineering to develop a more robust and generalizable framework than rule-based detection alone can provide.

7 Conclusions

This work presents a generalizable framework for identifying outlier behaviour in PM_{2.5} sensor networks during wildfire smoke events by integrating information-theoretic metrics into both a rule-based equation, as well as a gradient-boosted decision-tree classifier. We applied synthetic perturbations designed to mimic the sensor failures commonly observed in the field to real-world air quality data during 11 different smoke events, demonstrating that an information-theoretic rule-based method can reliably detect statistical outliers with low false-positive rates, while remaining tolerant of the network disorder that occurs during smoke intrusions. Extending this approach with machine learning features that incorporate short-term dynamics and network-aware behaviour allowed us to capture a broader range of anomalous signatures, including subtle precursor cues and faulty sensor behaviour that may have otherwise gone undetected. When these methods were applied to 9 months of operational data from more than 19,000 devices across the United States and Canada, both approaches showed strong performance in detecting outliers, with accuracy increasing substantially when devices were flagged between consecutive hours. Stratifying flagged events by their persistence duration provides an effective method of separating outlier types between hyper-local events and sustained sensor malfunctions. Persistence filtering provides a practical mechanism for tailoring detection aggressiveness to specific operational contexts: lower thresholds (3+ hours) effectively balance precision and recall for applications prioritizing population-level air quality communication, while higher thresholds (12+ hours) reliably isolate hardware malfunctions. Extending temporal averaging on measurements, with windows beyond one hour (e.g., 3-hour rolling averages), may similarly focus detection on sustained anomalies.



While overall performance was encouraging, this application was intentionally non-prescriptive. The flexibility of both detection approaches, (particularly the tunability of parameter thresholds, persistence criteria, or additional rulesets), means that these tools can be readily adapted to specific operational priorities. For example, the framework could be configured conservatively for regional public-health messaging, where false positives must be minimized, or tuned to retain maximum sensitivity for fine-scale PM_{2.5} monitoring in research contexts. Both methods are computationally efficient, requiring minimal processing time per device-hour, and consensus between different parameter sets or binning strategies substantially reduced misclassification rates while adding negligible computational overhead. Future work could expand the representation of real-world fault scenarios, refine network generation methods in complex terrain, and potentially incorporate contextual information such as nearby wildland fire events. These developments could support more reliable automated quality assurance for PM_{2.5} sensors, aiding in exposure assessment and supporting public health communication during wildfire smoke and other pollution events.

Code availability

The core outlier detection methods using the information-theoretic framework, generating sensor networks, and deriving the features used in the XGBoost models, are implemented in the elwood-spatial Python package (Illson, 2026a; <https://doi.org/10.5281/zenodo.18856271>). A companion repository provides a reproducible example notebook, implementation of the perturbations used, pre-trained XGBoost models and a small data subset from one of the events to allow users to recreate and extend the analysis presented here (Illson, 2026b; <https://doi.org/10.5281/zenodo.18897333>). Further code can be made available upon reasonable request.

Data availability

The underlying air quality data were obtained from sources with varying data use policies and cannot be redistributed in full. Data from public sources (e.g., AirNow) can be re-downloaded directly by users. Privately owned sensor data may not be available for redistribution. The spatial extent of the event sites, a site and event characterization document, and a sample subset of data from one smoke event are included in the companion repository (Illson, 2026b; <https://doi.org/10.5281/zenodo.18897333>). Additional data and methods can be made available upon reasonable request to the corresponding author.



Author Contributions

965 SJI conceptualised the study; SJI designed the framework and analysis methodology with feedback from KKB; SJI wrote the software; SJI performed the analysis and validation; Both authors reviewed and interpreted the results. SJI prepared the manuscript, with KKB contributing writing and providing feedback and contributions throughout.

Competing Interests

The authors declare that they have no conflict of interest.

970 Disclaimer

The views expressed in this paper are those of the author(s) and do not necessarily represent the views or policies of the U.S. Environmental Protection Agency, the U.S. Forest Service, or the University of Washington. Any mention of trade names, products, or services does not imply an endorsement by the U.S. Government, the U.S. Environmental Protection Agency, the U.S. Forest Service, or the University of Washington. None of these institutions endorse any commercial products, services,
975 or enterprises.

Some authors are affiliated with, or contribute to, the AirFire research program, Mazama Science software, and the AirNow Fire and Smoke Map. References to data products, tools, or visualizations associated with these efforts are included for transparency and continuity with ongoing research activities and do not imply preferential treatment or endorsement. All
980 analyses and interpretations presented here are intended to be scientifically objective.

Acknowledgements

We thank PurpleAir for providing access to PurpleAir data (MTA #1261-19), for their collaboration on the AirNow Fire and Smoke Map, and for creating and maintaining this incredible dataset. We are grateful to the joint U.S. Environmental Protection Agency and U.S. Forest Service Fire and Smoke Map team for their continued commitment to maintaining an invaluable public
985 resource for wildfire smoke information. We thank Andrea Clements, Rachelle Duvall, and Marc Houyoux from the U.S. EPA for their reviews of this work, adding to its rigor. We thank Dr. Sim Larkin for leadership and insight, the AirFire team, and the Interagency Wildland Fire Air Quality Response Program for their support. Dr. Brian Potter for an early review of this paper, and Dr. Jonathan Callahan for developing some of the software used to collect and manage air quality data. At the University of Washington, we thank Dr. Larry Pierce for his mathematical expertise and Dr. Ernesto Alvarado for his
990 mentorship.



References

- 40 CFR Part 58, Appendix E: Probe and Monitoring Path Siting Criteria for Ambient Air Quality Monitoring, Code of Federal Regulations, Title 40, Chapter I, Subchapter C, U.S. Environmental Protection Agency, available at: <https://www.ecfr.gov/current/title-40/chapter-I/subchapter-C/part-58/appendix-Appendix%20E%20to%20Part%2058>,
995 accessed: [December 5th, 2025].
- Agrawal, D., Saini, A. K., Rai, A. C., and Kala, P.: Study on calibration of low-cost particulate matter sensors for hydrophilic and hydrophobic particles under varying relative humidity, *Air Quality, Atmosphere & Health*, 10.1007/s11869-025-01852-y, 2025.
- AirNow Fire and Smoke Map: <https://fire.airnow.gov/>, last accessed: 7 April 2025.
- 1000 Barkjohn, K. K., Gantt, B., and 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*, 14, 4617–4637, <https://doi.org/10.5194/amt-14-4617-2021>, 2021.
- Barkjohn, K. K., Holder, A. L., Frederick, S. G., and Clements, A. L.: Correction and accuracy of PurpleAir PM_{2.5} measurements for extreme wildfire smoke, *Sensors*, 22, 9669, <https://doi.org/10.3390/s22249669>, 2022.
- 1005 Barkjohn, K. K., Plessel, T., Yang, J., Pandey, G., Xu, Y., Krabbe, S., Seppanen, C., Bichler, R., Tran, H. N. Q., Arunachalam, S., Clements, A. L.: Air Sensor Network Analysis Tool: R-Shiny Application, *Atmosphere*, 16, <https://doi.org/10.3390/atmos16111270>, 2025b.
- Barkjohn, K. K., Yaga, R., Thomas, B., Schoppman, W., Docherty, K. S., and Clements, A. L.: Evaluation of Long-Term Performance of Six PM_{2.5} Sensor Types, *Sensors*, 25, 1265, <https://doi.org/10.3390/s25041265>, 2025a.
- 1010 Bentéjac, C., Csörgő, A., and Martínez-Muñoz, G.: A comparative analysis of gradient boosting algorithms, *Artif. Intell. Rev.*, 54, 1937–1967, <https://doi.org/10.1007/s10462-020-09896-5>, 2021.
- Bi, J., Wildani, A., Chang, H. H., and Liu, Y.: Incorporating Low-Cost Sensor Measurements into High-Resolution PM_{2.5} Modeling at a Large Spatial Scale, *Environ. Sci. Technol.*, 54, 2152–2162, <https://doi.org/10.1021/acs.est.9b06046>, 2020.
- 1015 Brook, R. D., Rajagopalan, S., Pope, C. A., Brook, J. R., Bhatnagar, A., Diez-Roux, A. V., Holguin, F., Hong, Y., Luepker, R. V., Mittleman, M. A., Peters, A., Siscovick, D., Smith, S. C., Whitsel, L., and Kaufman, J. D.: Particulate Matter Air Pollution and Cardiovascular Disease, *Circulation*, 121, 2331–2378, <https://doi.org/10.1161/CIR.0b013e3181dbee1>, 2010.
- Burke, M., Driscoll, A., Heft-Neal, S., Xue, J., Burney, J., and Wara, M.: The changing risk and burden of wildfire in the United States, *Proc. Natl. Acad. Sci. U.S.A.*, 118, e2011048118, <https://doi.org/10.1073/pnas.2011048118>, 2021.
- 1020 Burke, M., Childs, M. L., de la Cuesta, B., Qiu, M., Li, J., Gould, C. F., Heft-Neal, S., and Wara, M.: The contribution of wildfire to PM_{2.5} trends in the USA, *Nature*, 622, 761–766, <https://doi.org/10.1038/s41586-023-06522-6>, 2023.



- Chen, T. and Guestrin, C.: XGBoost: A Scalable Tree Boosting System, in: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, arXiv:1603.02754 [cs], 785–794, <https://doi.org/10.1145/2939672.2939785>, 2016.
- 1025 Chen, L.-J., Ho, Y.-H., Hsieh, H.-H., Huang, S.-T., Lee, H.-C., and Mahajan, S.: ADF: An Anomaly Detection Framework for Large-Scale PM_{2.5} Sensing Systems, *IEEE Internet of Things Journal*, 5, 559–570, <https://doi.org/10.1109/JIOT.2017.2766085>, 2018.
- Clements, A. L., Griswold, W. G., Rs, A., Johnston, J. E., Herting, M. M., Thorson, J., Collier-Oxandale, A., and Hannigan, M.: Low-Cost Air Quality Monitoring Tools: From Research to Practice (A Workshop Summary), *Sensors*, 17, 2478, <https://doi.org/10.3390/s17112478>, 2017.
- 1030 deSouza, P., Kahn, R. A., Limbacher, J. A., Marais, E. A., Duarte, F., and Ratti, C.: Combining low-cost, surface-based aerosol monitors with size-resolved satellite data for air quality applications, *Atmospheric Measurement Techniques*, 13, 5319–5334, <https://doi.org/10.5194/amt-13-5319-2020>, 2020.
- deSouza, P. and Kinney, P. L.: On the distribution of low-cost PM_{2.5} sensors in the US: demographic and air quality associations, *J Expo Sci Environ Epidemiol*, 31, 514–524, <https://doi.org/10.1038/s41370-021-00328-2>, 2021.
- 1035 Ferrer-Cid, P., Barceló-Ordinas, J. M., and García-Vidal, J.: Volterra Graph-Based Outlier Detection for Air Pollution Sensor Networks, *Environmental Science & Pollution Research*, 29, 37614–37629, <https://doi.org/10.1007/s11356-022-19452-z>, 2022.
- Friedman, J. H.: Greedy Function Approximation: A Gradient Boosting Machine, *The Annals of Statistics*, 29, 1189–1232, 2001.
- 1040 Gould, C. F., Heft-Neal, S., Johnson, M., Aguilera, J., Burke, M., and Nadeau, K.: Health Effects of Wildfire Smoke Exposure, *Annu. Rev. Med.*, 75, 277–292, <https://doi.org/10.1146/annurev-med-052422-020909>, 2024.
- Hagler, G., Hanley, T., Hassett-Sipple, B., Vanderpool, R., Smith, M., Wilbur, J., Wilbur, T., Oliver, T., Shand, D., Vidacek, V., Johnson, C., Allen, R., and D'Angelo, C.: Evaluation of two collocated federal equivalent method PM_{2.5} instruments over a wide range of concentrations in Sarajevo, Bosnia and Herzegovina, *Atmospheric Pollution Research*, 13, 101374, <https://doi.org/10.1016/j.apr.2022.101374>, 2022.
- 1045 Hart, P. E., Nilsson, N. J., and Raphael, B.: A Formal Basis for the Heuristic Determination of Minimum Cost Paths, *IEEE Transactions on Systems Science and Cybernetics*, 4, 100–107, <https://doi.org/10.1109/TSSC.1968.300136>, 1968.
- Illson, S.: elwood-spatial: Information-Theoretic Outlier Detection for Spatial Networks, Zenodo [code], <https://doi.org/10.5281/zenodo.18856271>, 2026a.
- 1050 Illson, S.: Companion Materials: Outlier Detection in PM_{2.5} Air Sensor Networks, Zenodo [code/data], <https://doi.org/10.5281/zenodo.18897333>, 2026b.



- Jaffe, D. A., O'Neill, S. M., Larkin, N. K., Holder, A. L., Peterson, D. L., Halofsky, J. E., and Rappold, A. G.: Wildfire and prescribed burning impacts on air quality in the United States, *Journal of the Air & Waste Management Association*, 70, 583–615, <https://doi.org/10.1080/10962247.2020.1749731>, 2020.
- 1055 Jaffe, D. A., Miller, C., Thompson, K., Finley, B., Nelson, M., Ouimette, J., and Andrews, E.: An evaluation of the U.S. EPA's correction equation for PurpleAir sensor data in smoke, dust, and wintertime urban pollution events, *Atmos. Meas. Tech.*, 16, 1311–1322, <https://doi.org/10.5194/amt-16-1311-2023>, 2023.
- Jayaratne, R., Liu, X., Thai, P., Dunbabin, M., and Morawska, L.: The influence of humidity on the performance of a low-cost air particle mass sensor and the effect of atmospheric fog, *Atmos. Meas. Tech.*, 11, 4883–4890, <https://doi.org/10.5194/amt-11-4883-2018>, 2018.
- 1060 Johnson, D., "Quality Assurance in an Optical Sensor Network for Public Health Information". National Air Quality Conference. 2024.
- Kaur, K., Mangin, T., and Kelly, K.: Identification of Dust-Dominated Periods and a PM_{2.5} Correction Based Solely on Plantower PMS Sensor Observations, *EGUsphere* [preprint], <https://doi.org/10.5194/egusphere-2025-5063>, 2025.
- 1065 Keyes, T., Domingo, R., Dynowski, S., Graves, R., Klein, M., Leonard, M., Pilgrim, J., Sanchirico, A., and Trinkaus, K.: Low-cost PM_{2.5} sensors can help identify driving factors of poor air quality and benefit communities, *Heliyon*, 9, e19876, <https://doi.org/10.1016/j.heliyon.2023.e19876>, 2023.
- Khan, T. R., Emerson, Z. I., and Mentz, K. H.: Evaluation of Fine Particulate Matter (PM_{2.5}) Concentrations Measured by Collocated Federal Reference Method and Federal Equivalent Method Monitors in the U.S, *Atmosphere*, 15, 978, 2024.
- 1070 Kim, K. T., Kim, H., Jeong, S., Lee, Y. S., Zhao, X., and Kim, J. Y.: Coincidence effect of a low-cost particulate matter sensor: Observations from environmental chamber tests at diverse particle concentrations, *Atmospheric Pollution Research*, 16, 102581, <https://doi.org/10.1016/j.apr.2025.102581>, 2025.
- Kochanski, A. K., Mallia, D. V., Fearon, M. G., Mandel, J., Souri, A. H., and Brown, T.: Modeling Wildfire Smoke Feedback Mechanisms Using a Coupled Fire-Atmosphere Model With a Radiatively Active Aerosol Scheme, *Journal of Geophysical Research: Atmospheres*, 124, 9099–9116, <https://doi.org/10.1029/2019JD030558>, 2019.
- 1075 Kuula, J., Mäkelä, T., Aurela, M., Teinilä, K., Varjonen, S., González, Ó., and Timonen, H.: Laboratory evaluation of particle-size selectivity of optical low-cost particulate matter sensors, *Atmos. Meas. Tech.*, 13, 2413–2423, <https://doi.org/10.5194/amt-13-2413-2020>, 2020.
- Larkin, N. K.: Modeling, monitoring, and messaging wildfire smoke for air quality and public health, Health Effects Institute 2019 Annual Conference, Seattle, WA, 6 May 2019.
- 1080 Madhwal, S., Tripathi, S. N., Bergin, M. H., Bhave, P., de Foy, B., Reddy, T. V. R., Chaudhry, S. K., Jain, V., Garg, N., and Lalwani, P.: Evaluation of PM_{2.5} spatio-temporal variability and hotspot formation using low-cost sensors across urban-rural landscape in lucknow, India, *Atmospheric Environment*, 319, 120302, <https://doi.org/10.1016/j.atmosenv.2023.120302>, 2024.



- 1085 Malings, C., Tanzer, R., Hauryliuk, Aliaksei, Saha, Provat K., Robinson, Allen L., Presto, Albert A., and Subramanian, R.: Fine particle mass monitoring with low-cost sensors: Corrections and long-term performance evaluation, *Aerosol Science and Technology*, 54, 160–174, <https://doi.org/10.1080/02786826.2019.1623863>, 2020.
- Manikonda, A., Zíková, N., Hopke, P. K., and Ferro, A. R.: Laboratory assessment of low-cost PM monitors, *Journal of Aerosol Science*, 102, 29–40, <https://doi.org/10.1016/j.jaerosci.2016.08.010>, 2016.
- 1090 Mehadi, A., Moosmüller, H., Campbell, D. E., Ham, W., Schweizer, D., Tarnay, L., and Hunter, J.: Laboratory and field evaluation of real-time and near real-time PM_{2.5} smoke monitors, *Journal of the Air & Waste Management Association*, 70, 158–179, <https://doi.org/10.1080/10962247.2019.1654036>, 2020.
- Molina Rueda, E., Carter, E., L'Orange, C., Quinn, C., and Volckens, J.: Size-resolved field performance of low-cost sensors for particulate matter air pollution, *Environ. Sci. Technol. Lett.*, 10, 247–253, <https://doi.org/10.1021/acs.estlett.3c00030>,
1095 2023.
- Morawska, L., Thai, P. K., Liu, X., Asumadu-Sakyi, A., Ayoko, G., Bartonova, A., Bedini, A., Chai, F., Christensen, B., Dunbabin, M., Gao, J., Hagler, G. S. W., Jayaratne, R., Kumar, P., Lau, A. K. H., Louie, P. K. K., Mazaheri, M., Ning, Z., Motta, N., Mullins, B., Rahman, M. M., Ristovski, Z., Shafiei, M., Tjondronegoro, D., Westerdahl, D., and Williams, R.: Applications of low-cost sensing technologies for air quality monitoring and exposure assessment: How far have they gone?,
1100 *Environ. Int.*, 116, 286–299, <https://doi.org/10.1016/j.envint.2018.04.018>, 2018.
- NASA JPL: NASADEM Merged DEM Global 1 arc second V001, NASA EOSDIS Land Processes DAAC [data set], https://doi.org/10.5067/MEASURES/NASADEM/NASADEM_HGT.001, 2020.
- Ouimette, J., Arnott, W. P., Laven, P., Whitwell, R., Radhakrishnan, N., Dhaniyala, S., Sandink, M., Tryner, J., and Volckens, J.: Fundamentals of low-cost aerosol sensor design and operation, *Aerosol Sci. Technol.*, 58, 1–15, <https://doi.org/10.1080/02786826.2023.2285935>, 2024.
- Pope III, C. A., Burnett, R. T., Thun, M. J., Calle, E. E., Krewski, D., Ito, K., and Thurston, G. D.: Lung cancer, cardiopulmonary mortality, and long-term exposure to fine particulate air pollution, *JAMA*, 287, 1132–1141, <https://doi.org/10.1001/jama.287.9.1132>, 2002.
- 1110 Qiu, M., Li, J., Gould, C. F., Jing, R., Kelp, M., Childs, M. L., Wen, J., Xie, Y., Lin, M., Kiang, M. V., Heft-Neal, S., Diffenbaugh, N. S., and Burke, M.: Wildfire smoke exposure and mortality burden in the US under climate change, *Nature*, <https://doi.org/10.1038/s41586-025-09611-w>, 2025.
- Raheja, G., Nimo, J., Appoh, E. K.-E., Essien, B., Sunu, M., Nyante, J., Amegah, M., Quansah, R., Arku, R. E., Penn, S. L., Giordano, M. R., Zheng, Z., Jack, D., Chillrud, S., Amegah, K., Subramanian, R., Pinder, R., Appah-Sampong, E., Tetteh, E. N., Borketey, M. A., Hughes, A. F., and Westervelt, D. M.: Low-Cost Sensor Performance Intercomparison, Correction Factor Development, and 2+ Years of Ambient PM_{2.5} Monitoring in Accra, Ghana, *Environ. Sci. Technol.*, 57, 10708–10720, <https://doi.org/10.1021/acs.est.2c09264>, 2023.
- Raysoni, A. U., Pinakana, S. D., Mendez, E., Wladyka, D., Sepielak, K., and Temby, O.: A Review of Literature on the Usage of Low-Cost Sensors to Measure Particulate Matter, *Earth*, 4, 168–186, <https://doi.org/10.3390/earth4010009>, 2023.



- 1120 Reid, C. E., Brauer, M., Johnston, F. H., Jerrett, M., Balmes, J. R., and Elliott, C. T.: Critical review of health impacts of wildfire smoke exposure, *Environ. Health Perspect.*, 124, 1334–1343, <https://doi.org/10.1289/ehp.1409277>, 2016.
- Rey, S. J. and Anselin, L.: PySAL: A Python Library of Spatial Analytical Methods, *The Review of Regional Studies*, 37, 2007.
- 1125 Sablan, O., Ford, B., Gargulinski, E., Hammer, M. S., Henery, G., Kondragunta, S., Martin, R. V., Rosen, Z., Slater, K., van Donkelaar, A., Zhang, H., Soja, A. J., Magzamen, S., Pierce, J. R., and Fischer, E. V.: Quantifying Prescribed-Fire Smoke Exposure Using Low-Cost Sensors and Satellites: Springtime Burning in Eastern Kansas, *Geohealth*, 8, e2023GH000982, <https://doi.org/10.1029/2023GH000982>, 2024.
- Schwartz, J., Dockery, D. W., and Neas, L. M.: Is Daily Mortality Associated Specifically with Fine Particles?, *Journal of the Air & Waste Management Association*, 46, 927–939, <https://doi.org/10.1080/10473289.1996.10467528>, 1996.
- 1130 Schweizer, D., Cisneros, R., and Shaw, G.: A comparative analysis of temporary and permanent beta attenuation monitors: The importance of understanding data and equipment limitations when creating PM_{2.5} air quality health advisories, *Atmospheric Pollution Research*, 7, 865–875, <https://doi.org/10.1016/j.apr.2016.02.003>, 2016.
- Shannon, C.E.: A Mathematical Theory of Communication. *Bell System Technical Journal*, 27, 379-423. <http://dx.doi.org/10.1002/j.1538-7305.1948.tb01338.x>, 1948.
- 1135 State of Oregon Department of Environmental Quality: Fact Sheet DEQ PM_{2.5} Sensors, <https://www.oregon.gov/deq/FilterDocs/aqwsensors.pdf>, n.d.
- State of Oregon Department of Environmental Quality: Sampling and Analysis Plan SensOR Site Selection and Installation, <https://www.oregon.gov/deq/aq/Documents/labSensorSAP.pdf>, 2023.
- 1140 Stieb, D. M., Burnett, R. T., Smith-Doiron, M., Brion, O., Shin, H. H., and Economou, V.: A new multipollutant, no-threshold air quality health index based on short-term associations observed in daily time-series analyses, *J. Air Waste Manage. Assoc.*, 58, 435–450, <https://doi.org/10.3155/1047-3289.58.3.435>, 2008.
- Tobler, W. R.: A Computer Movie Simulating Urban Growth in the Detroit Region, *Economic Geography*, 46, 234, <https://doi.org/10.2307/143141>, 1970.
- 1145 Washington State Department of Ecology: SensWA Correlations for PM_{2.5} Reporting, Publication No. 24-02-052, Washington State Department of Ecology, Olympia, WA, <https://apps.ecology.wa.gov/publications/documents/2402052.pdf>, 2024a.
- Washington State Department of Ecology: SensWA Quality Assurance Project Plan, Publication No. 24-02-024, Washington State Department of Ecology, Olympia, WA, <https://apps.ecology.wa.gov/publications/documents/2402024.pdf>, 2024b.
- 1150 Zamora, M., Xiong, F., Gentner, D., Kerkez, B., Kohrman-Glaser, J., and Koehler, K.: Field and Laboratory Evaluations of the Low-Cost Plantower Particulate Matter Sensor, *Environ. Sci. Technol.*, 53, 838–849, <https://doi.org/10.1021/acs.est.8b05174>, 2019.

<https://doi.org/10.5194/egusphere-2026-1273>

Preprint. Discussion started: 3 June 2026

© Author(s) 2026. CC BY 4.0 License.



Zimmerman, N., Presto, A. A., Kumar, S. P. N., Gu, J., Hauryliuk, A., Robinson, E. S., Robinson, A. L., and R. Subramanian: A machine learning calibration model using random forests to improve sensor performance for lower-cost air quality monitoring, *Atmospheric Measurement Techniques*, 11, 291–313, <https://doi.org/10.5194/amt-11-291-2018>, 2018.