# Engaging Pre-College Students in Hypothesis Generation using a Citizen Scientist Network of Air Quality Sensors (Work in Progress)

**JAMES A MOORE**
**Matthew Dailey**

Matthew Dailey is a student at the University of Utah pursuing a B.S in Chemical Engineering. He is pursuing graduate school with a focus on biosensors.

**Mr. Zachary Wilhelm, AirU**

Zachary Wilhelm is pursing his undergrad in Chemical Engineering at the University of Utah and is aspiring to get his PhD to continue research to understand and address environmental challenges. For this project his primary focus was organizing visits to local schools and attending outreach events to engage citizen scientists across the Salt Lake City valley.

**Dr. Kerry Kelly, University of Utah**

Dr. Kerry Kelly is a professional engineer, an Assistant Professor of Chemical Engineering and Associate Director of the Program for Air Quality, Health, and Society at the University of Utah. She has a PhD in Environmental Engineering and a BS in Chemical Engineering, and she just completed 8 years of service on Utah's Air Quality Board. Her research focuses on air quality and the evaluation of emerging energy technologies including consideration of their associated health, environmental, policy and performance issues. Most recently she has been focusing on combustion particles, their associated health effects, low-cost air-quality sensing, and community engagement.

**Pascal Goffin, University of Utah**

Pascal Goffin received his PhD in Computer Science from Université Paris-Saclay in France in 2016. During his PhD he worked for the Aviz visualization group at Inria. He holds a Masters degree in Computer Science from ETH Zurich in Switzerland. His interest span information visualization, text visualization, and human com- puter interaction. His current research focuses on how to support the communication of air quality in urban environments to citizens. He also builds tools to assist the exploration of urban air quality data.

**Prof. Anthony Butterfield, University of Utah**

Anthony Butterfield is an Assistant Professor (Lecturing) in the Chemical Engineering Department of the University of Utah. He received his B. S. and Ph. D. from the University of Utah and a M. S. from the University of California, San Diego. His teaching responsibilities include the senior unit operations laboratory and freshman design laboratory. His research interests focus on undergraduate education, targeted drug delivery, photobioreactor design, and instrumentation.

**Prof. Jason Wiese,**

Jason Wiese is an Assistant Professor in the School of Computing at the University of Utah. His research takes a user-centric perspective of personal data, focusing on how that data is collected, interpreted, and used in applications. His work crosses the domains of machine learning, privacy, user-centered design, real-world data collection, and user study design. Dr. Wiese's research excellence has been recognized by awards including: recognition as a Yahoo Fellow in 2014, the Stu Card Fellowship in 2012, a Carnegie Mellon Usable Privacy and Security IGERT trainee, and the Yahoo! Key Scientific Challenges Award in 2011. He publishes work in top Computer Science and HCI venues including CHI, CSCW, and UbiComp. He received his Ph.D. in Human-Computer Interaction from Carnegie Mellon University in 2015.
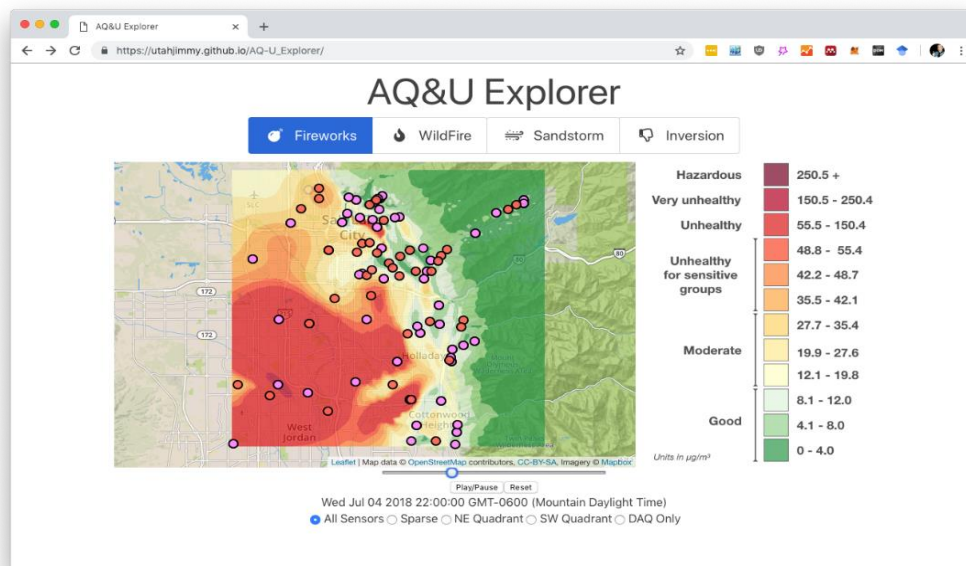
**Dr. WEI XING, UNIVERSITY OF UTHA**
**Katrina Myquyen Le, University of Utah**
**Mr. Thomas Becnel, University of Utah**

lll

llll

---

# Engaging Middle and High School Students in Hypothesis Generation using a Citizen Scientist Network of Air Quality Sensors



**Figure 1: Salt Lake City Valley-wide** air quality model of $PM_{2.5}$ concentrations, 7/4/18

## 1 Introduction

Polluted air afflicts 90% of the world's population and contributes to 7 million premature deaths every year [1]. Salt Lake City**,** Utah periodically experiences some of the worst air quality in the nation [2], yet is sparsely instrumented and subject to lengthy update intervals of one or more hours until this air-quality information is publicly available. To provide more finely resolved spatial and temporal air-quality data, we have deployed a low-cost sensor network and accompanying website [3] for improving public awareness. This network has generated a large corpus of fine particulate matter ($PM_{2.5}$) measurements that reveal how $PM_{2.5}$ concentrations evolve over time and space. Building on prior educational outreach and citizen science exercises [4], we explore an interactive, team-based teaching module using local real-world data. This teaching module's goal is to engage students in generating and testing hypotheses while also encouraging citizen use of real-time air quality data for their own interests, such as exploration, science fair projects, or environmental oversight.

We have piloted this module with over 500 students across 8 local high schools in various chemistry, engineering, environmental science, and physics classrooms. Structured around a data analysis exercise with local air quality data, the module helps guide students through creating and testing hypotheses about air quality under various conditions. The module also incorporates fundamental analysis tasks, such as loading and plotting data in a spreadsheet program to build students' familiarity with basic data analysis techniques. Using pre- and post-survey responses,

this work seeks to evaluate how a guided, team-based outreach module impacts students' ability to generate and test hypotheses, perceptions of outdoor air quality, sense of engagement with a data analysis exercise, and overall success of incorporating publicly available local data from distributed sensor networks into their curriculum. In addition, over half of the visited schools had an underrepresented enrollment exceeding 50% of the student body (according to state statistics), and all but two surpassed the state average of 25%. Consequently, this module helped introduce traditionally underrepresented students in STEM to distributed data collection, interpolation, modeling, and visualization concepts used for generating community-scale air quality models.

Classroom activities were designed around helping improve students' analytical competency through interpreting local PM$_{2.5}$ measurements. Preliminary results showed this module to be highly engaging, and effective for improving students' awareness of air quality's geospatial and temporal variations during a variety of pollution episodes. Survey results also showed this module was effective at introducing hypothesis generation and testing techniques. All classrooms reported wanting to host the module again and having plans to incorporate their local air quality data into future activities and assignments.
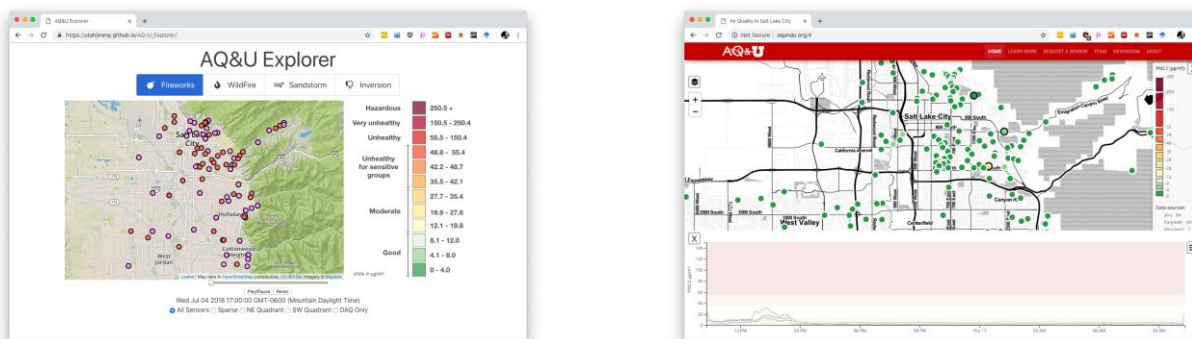
## 2 Background

Citizen science enlists the public to participate in data collection across an array of locations and time spans, and has contributed to scientific discoveries in a variety of fields from bird watching to environmental oversight [5]. Research involving citizen science has also grown within scientific literature, as evidenced by an increased discussion of citizen science in peer reviewed articles [6]. Despite its growth, relatively few citizen science projects have focused on engineering disciplines [6, 7]. Prior citizen science efforts have developed curricula for low-cost, air-quality sensors in schools [8] and a recent study enabled citizen scientists to monitor and report unlawful air quality emissions from local industry [9]. One challenge integrating air quality measurement with citizen science initiatives is over sensors' perceived "black box" operation, with citizen scientists having little understanding of how these sensors function [10]. While prior outreach has helped expose the inner workings of sensor hardware, specifically highlighting key operational principles and design trade-offs [4], this teaching module helps students reason about the larger sensing infrastructure. Through interactive presentations and guided activities, students get a glimpse of how raw, multi-sensor air-quality data are incorporated within simulation models to provide data in locations without sensor coverage. Though there are analytical limits to what the general public can do on their own, statistical and computational tools are being developed to assist citizens in analyzing these complex data sets [11].

This work attempts to involve students in hypothesis generation and testing, in order to engage the higher cognitive domains of Bloom's taxonomy. Through the outreach activity, students are taught how to interact with vast amounts of data of concern to their communities and use it to support their conclusion about various air quality events.

## 3 Study Design

This teaching module pairs an interactive classroom presentation with guided, hands-on activities. Our aims with this module are to (a) generate interest in STEM fields using an environmental problem that students experience in their everyday lives; (b) help students use real-world data for developing and testing hypotheses about a pressing local and national challenge: air quality; (c) promote citizen use of real-time air quality data for their own interests, such as data exploration, science fair projects, and environmental oversight; and (d) introduce students to new STEM concepts such as distributed data collection, interpolation, modeling, and visualization.



**Figure 2:** Websites used in our classroom activities. (a) **AQ&U Explorer:** View air quality model simulations. (b) **AQ&U website**: Public-facing interface for viewing real-time air quality measurements.

### 3.1 AQ&U Infrastructure

The AQ&U infrastructure provides public access to data from its air quality sensing network and supports citizen scientist participation. A primary goal of this infrastructure is to provide dense, spatiotemporal estimates of air quality to researchers and the general public. Our teaching module utilizes this infrastructure to help students gain a better understanding of distributed collection of real-world data for developing and testing hypotheses. The **AQ&U** infrastructure integrates measurements from the Utah Division of Air Quality (DAQ) (2 gold-standard measurements) and over 100 citizen- and school-hosted $PM_{2.5}$ sensors. Dynamic data-fusion algorithms and visualization techniques process these data streams to provide highly resolved $PM_{2.5}$ concentration information to the public facing **AQ&U** website (Figure 2(b)).

Individual air-quality events possess different measurement signal characteristics. For instance, $PM_{2.5}$ levels during winter cold air pools tend to exhibit consistency in time, whereas summer events (fires, fireworks) change more rapidly. We therefore utilize a moving window including several days' worth of data to generate estimates of $PM_{2.5}$ within the modeling domain. The low-cost $PM_{2.5}$ sensor [1] measurements are also seasonally corrected using factors derived from co-

---

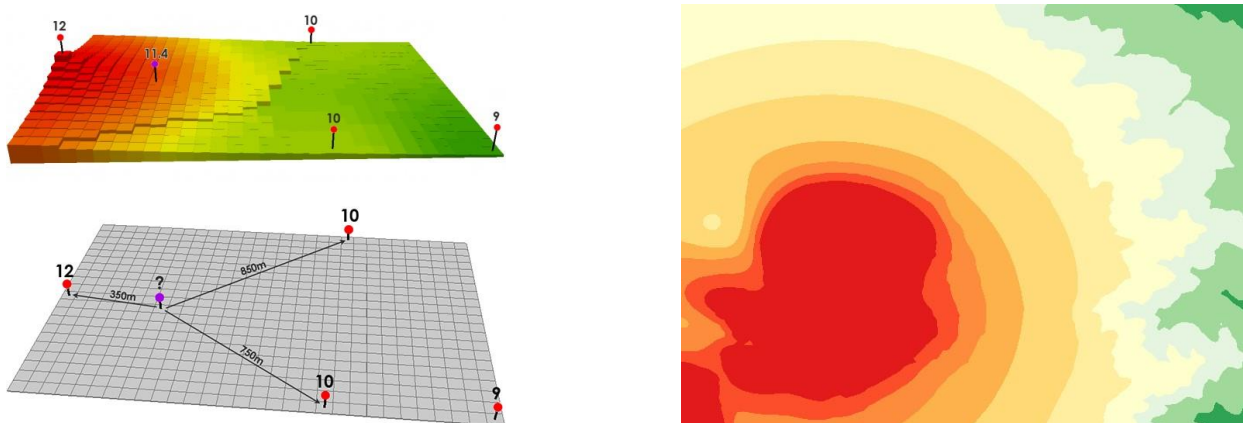[1] Plantower PMS sensor: http://www.plantower.com/en/list/?118 1.html

located reference measurements and applied to a classical Gaussian process model for computing spatiotemporal regression. [12, 13].

These results are translated into contours (Figure 3(b)) using standard plotting techniques from Matplotlib [14]. The visualization encodes the contours using a color map based on the EPA air quality index color scheme, subdividing each of the EPA categories into three ranges in order to provide more resolved concentration information. These divisions differentiate it slightly from EPA's health-related color scheme, which is based on 24-hr average pollutant concentrations.

## 3.2 Outreach

The air quality teaching module adopts a student-led approach [15] and engages chemical engineering undergraduates who give classroom presentations that introduce core concepts including $PM_{2.5}$ health impacts, sources, and geographical variation. Students learn about point- and area-sources of $PM_{2.5}$ pollution, motivating the concept of spatiotemporal variations in air quality over a region. Using an analogy to image resolution, undergraduate facilitators discuss how more sensors can provide a "clearer image" of air quality variability, and underscores the importance of sensor distribution to accurately capture this information.

Students are then shown the AQ&U website illustrating how this information can be used in practice (Figure 2(b)). Students use this site to explore individual air quality measurements from specific locations before being asked how they might characterize the air quality of the entire instrumented region.



**Figure 3:** Interpolating sensor measurements and simulation output contour plot. (a) Interpolation in two dimensions [16]. (b) Still frame from air quality model output.

Measurement interpolation and contour plots (Figure 3) are motivated by way of analogy to topographic hiking maps and weather data. Students are then shown a second interactive tool, the AQ&U Explorer (Figure 2(a)), and use this to play air quality model output generated from the distributed sensor network measurements. The AQ&U Explorer lets users select between four air quality events effecting the entire Salt Lake City Valley: 4th of July fireworks,

regional wildfires, a persistent cold air pool (also known as an "inversion"), and a dust storm. This tool incorporates valley-wide sensor output to generate and display an air quality model as a color-coded contour plot. Simulation output illustrates how the underlying measurements change over space and time to highlight air quality's variability. For each air quality event, students can choose from five underlying sensor distributions: All deployed monitors, a 50% deployment ("sparse"), the northeast or southwest quadrants, and "DAQ-only" sensors, which are official government measurements from the Division of Air Quality. Viewing model differences from different sensor distributions helps illustrate the importance of sensor density and placement.

Students use this interface to complete a guided exercise for analyzing separate air quality events to consider strengths and weaknesses of the air-quality model estimates related to sensor technology, sensor location, and sensor density.

### 3.3 Guided group activities

After familiarizing themselves with the visualization interface, students form groups of three to five to analyze separate types of air quality events, focusing on how air quality measurements vary over space and time for each event. Figure 4 gives snapshots of the animations that student teams may view on the AQ&U Explorer interface to develop expertise on different air quality events. As can be seen, each type of event has different characteristic behaviors in the valley regarding timing, and location (particularly elevation). For instance, fireworks tend to affect areas without fireworks restrictions and last briefly, whereas inversions tend to pool and slosh in



**Figure 4:** Example events which student teams may observe in animated form to develop expertise on the characteristics of each AQ event type.

the valley and last significant time. Furthermore, students can detect both the benefits and limitations of such a sensor network when compared to the existing DAQ stations. For instance, in Figure 4 it can be seen that the low-cost sensors, while generating more resolution during inversions and wildfire, have a difficult time detecting the large particles common in dust storms.

**Table 1:** Classroom Visit Timeline. If a 50 min class period is not available, the items in orange may be removed to fit into a smaller portion of a class period.

| Activity | Time (min) |
|---|---|
| Introduction to air quality science and sensor networks | 5 |
| Familiarization with visualization interface | 5 |
| Group air quality event analysis | 5 |
| Group presentation on unique event types | 10 |
| Shuffle groups | - |
| Familiarization with network API for data download | 5 |
| Group mystery data set analysis | 10 |
| Group presentation and discussion of hypothesis & evidence | 10 |

Individual teams then take turns presenting their analysis to the rest of the class, which then collectively consider how the data from low-cost sensors compares to the more accurate (but more sparse) government monitoring stations. Students also discuss socioeconomic factors, such as locations of communities that are predominantly affected by poor air quality, the underlying sensor distribution, and how they may be related.

After this exercise, teams shuffle and reform so that there is an 'expert' from each prior air quality event in the new group. Together this team tries to classify a mystery data set as one of the previous four air quality scenarios. Student groups receive three unique sensor data streams from locations in the valley with the goal of hypothesizing which air quality event was captured. Students may also be asked to determine the possible locations of their sensor (e.g. in the lower elevation suburbs, downtown, on the foothills, etc.).

Students download air quality datasets through the AQ&U API and plot these measurements in a spreadsheet program. Data properties such as measurement values and time stamps are compared to the simulation output to help discern the captured events, while individual sensor responses can help localize the sensor's position. Figure 5 shows a replicated example of typical student work from this section of the module. Analysis of the downloaded csv file also allows students to get some practice using Google Sheets or MS Excel. In this example, both the 4th of July occurs with fireworks and it is followed immediately the next day by a wild fire. Through observations of the plots they generate, student teams are meant to come to conclusions about event type and sensor locations.

In addition to engaging students' hypothesis generation and evaluation process, this activity is also meant to familiarize students with the AQ&U API, with the aim being to allow them to use the interface and historical data for individual projects.

A timeline of an example classroom visit is shown in Table 1. If less than a 50-minute class time is available, this module may be shortened by only going up to the first student team presentation.
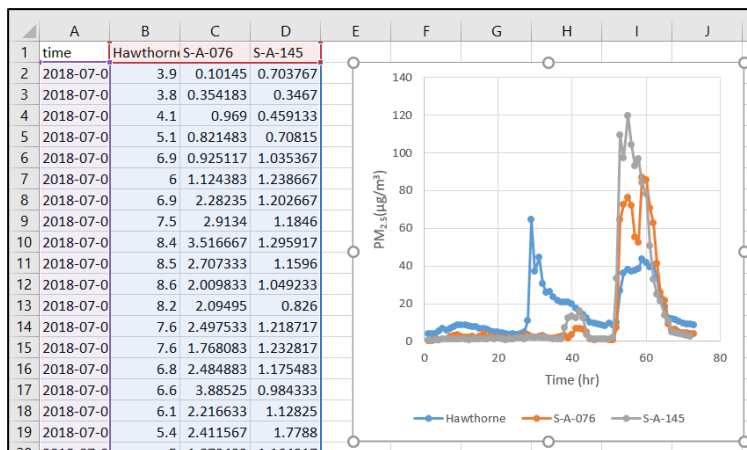
### 3.4 Teaching module surveys

Pre- and post-module surveys are distributed to the classroom teachers before and after the air quality teaching module,



**Figure 5:** Replicated example of student analysis of mystery air quality event. There are two main peaks: the first due to fireworks and the second due to wildfires. The data shows how the relative location to the pollution event causes different concentration readings. The students would use this data and knowledge gained from previous parts of the module to determine what pollution events are happening and where the sensors are likely located.

respectively. These surveys collect self-reported data on teachers' expectations for hosting the module, along with their appraisal of student's analytical abilities, experience with data analysis, hypothesis generation, visualization, and overall air quality awareness. Post-surveys are issued 1 to 2 weeks after a site visit to gauge the teaching module's effectiveness relative to the teacher's expectations. These surveys also assess the instructor's perception of students' understanding and engagement of the exercise (Figure 6).

### 4 Results and discussion

To date, our outreach team has visited 8 schools and 22 classes, and has reached over 540 students. Table 2 outlines each of the school and classroom demographics. Nine teachers to our pre- and post-surveys.

Pre-survey results (not plotted) indicate that both AP and non-AP teachers are strongly interested in hosting the teaching module. However, AP teachers rank the following motivational goals as higher importance than non-AP teachers: gaining a better understanding of how location affects air quality, learning how to use spreadsheets, and satisfying learning objectives of the class. Post-survey results begin to highlight the module's most effective parts. Based on the survey, 100% of teachers strongly agreed that they would use the module again in their classroom. This feedback in particular suggests that the module was of value to the teachers and the knowledge gained by the students was worth the time dedicated to the module. Furthermore, all teachers

**Table 2:** Teaching module site visit statistics.

| School ID | Course Curriculum | # of Students | Minority Enrollment (%) | Student Age |
|---|---|---|---|---|
| HS1 | Intro to Engineering | 12 | 51 | 14-18 |
| HS2A | Env. Science (AP) | 25 | 66 | 16 |
| HS2B | Env. Science (AP) | 25 | 66 | 16 |
| HS2C | Env. Science (AP) | 25 | 66 | 16 |
| HS3 | Env. Science (AP) | 5 | 15 | 17 |
| HS4 | Physics | 173 | 9 | 16 |
| E1 | Gifted and Talented | 7 | 7 | 11-12 |
| HS5 | Physics | 236 | - | 16 |
| HS6 | Chemistry | 27 | 46 | 15-17 |

either agreed or strongly agreed that this module was effective at integrating with their class objectives.

The post-survey not only highlighted the effective parts of the module, it also exposed the current flaws. One of the larger gaps in the module was the effectiveness of explaining how pollution varies by elevation, location, and over the course of the day. Questions 5 and 6 (Figure 6) suggest that students understood that pollution varies over the course of the day although they were less clear about the effects of location and elevation.

The open-ended response portion of the post-survey provided additional insight into these trends. Feedback on the module's most effective aspects included its ability to pair data analysis tools with real-world data. All post-survey responses indicated that some aspect of the visualizations were the most engaging component of the module, whether exploring how fireworks impact local air quality or observing the variation of pollution levels in different neighborhoods. One teacher claimed that the most effective aspect of the module was "Connecting engineering with real world problems students can relate to." Feedback like this shows the module is effectively conveying engineering aspects to the students.

Several teachers and our student outreach team suggested the following improvements: allowing 75 minutes per class rather than 50 minutes, organizing the high school students into teams of two rather than four so that each student has a chance at analyzing data, and pairing in-class activities with an accompanying take-home assignment for reinforcing key concepts. Lastly, teachers and team members reported the highest student engagement with the fireworks and inversion datasets, which captured the largest air quality impacts over Salt Lake City **Salt Lake City**. Students' interest in this module underscores how incorporating local and personally meaningful data in outreach programs can help foster and maintain student engagement.

**Figure 6:** Module feedback from post-survey results.

## Conclusion and future work

This project develops, pilot tests, and performs a preliminary evaluation of an interactive, team-based teaching module incorporating locally significant, real-world data collected though a citizen-hosted network of air-quality sensors. The goal of this teaching module is to engage students in hypothesis generation and evaluation, and to encourage citizen use of real-time air quality data. We have piloted this module to over 500 students at 9 local high school classrooms. Preliminary survey results suggest that the most effective aspects of this module are the engaging visualizations and the use of real-world data to generate and test hypotheses.

The feedback from teachers and outreach students suggest several improvements.

- Encourage greater student engagement through smaller group sizes.
- Developing a follow-up assignment to reinforce core ideas behind the teaching module.
- Have students download $PM_{2.5}$ concentrations from the DAQ and our co-located sensor through the API, and compare the $PM_{2.5}$ concentrations by calculating accuracy, error, and noise.

A follow-up assignment is being developed that will assess the students' knowledge of the core ideas behind the teaching module. The assignment will contain questions that focus on the hypotheses generated, the impacts elevation has on air quality, how contour plots are generated and so on. As the module is refined, we will continue to visit classrooms and perform evaluations with an expectation of developing a more rigorously reviewed teaching module.

Lastly, in response to the COVID-19 virus, most of the country is currently under quarantine and many universities and K-12 schools have turned to teaching exclusively online. This has, of course, had greatly limited outreach efforts. Our outreach team has had to turn to conducting outreach through teleconferencing, and is in the midst of organizing outreach "visits" to virtual classrooms. In organizing these efforts, we have, of course, found many of our hands-on teaching modules would be inappropriate. However, teaching modules such as the module described in this work are ideal for virtual outreach visits. In future work, this module will be used as one of the few outreach activities that may be executed under quarantine, while at the same time engaging more sophisticated cognitive domains, such as those involved in hypothesis generation and testing. Even once this current crisis passes, we envision using such outreach modules to conduct virtual outreach "visits" to rural parts of the country.

## 5 Acknowledgements

References

[1] World Health Organization, "How air pollution is destroying our health," https://www.who.int/air-pollution/news-and-events/how-air-pollution-is-destroying-our-health, January 2019, online; accessed 17 January 2019.

[2] American Lung Association, "State of the Air 2018 - Most Polluted Cities," https://www.lung.org/our-initiatives/healthy-air/sota/city-rankings/most-polluted-cities.html, 2019, online; accessed 17 January 2019.

[3] "AQ & U," http://aqandu.org/, January 2018, online; accessed 24 January 2019.

[4] A. Butterfield, K. My Quyen Le, K. Kelley, P. Goffin, T. Becnel, and P.-E. Gaillardon, "Citizen scientists engagement in air quality measurements," in *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, ser. ASEE '18. ASEE, 2018. [Online]. Available: https://www.asee.org/public/conferences/106/papers/22785/view

[5] R. Bonney, C. B. Cooper, J. Dickinson, S. Kelling, T. Phillips, K. V. Rosenberg, and J. Shirk, "Citizen Science: A Developing Tool for Expanding Science Knowledge and Scientific Literacy," *BioScience*, vol. 59, no. 11, pp. 977–984, dec 2009. [Online]. Available: https://academic.oup.com/bioscience/article-lookup/doi/10.1525/bio.2009.59.11.9

[6] R. Follett and V. Strezov, "An analysis of citizen science based research: Usage and publication patterns," *PLOS ONE*, vol. 10, no. 11, pp. 1–14, 11 2015. [Online]. Available: https://doi.org/10.1371/journal.pone.0143687

[7] C. Kullenberg and D. Kasperowski, "What is citizen science? a scientometric meta-analysis," *PLOS ONE*, vol. 11, no. 1, pp. 1–16, 01 2016. [Online]. Available: https://doi.org/10.1371/journal.pone.0147152

[8] E. Adams, G. Smith, M. Henthorn, T. J. Ward, D. Vanek, N. Marra, D. Jones, and J. Striebel, "Air toxics under the big sky: A real-world investigation to engage high school science students," *Journal of Chemical Education*, vol. 85, no. 2, p. 221, 2008. [Online]. Available: https://doi.org/10.1021/ed085p221

[9] Y.-C. Hsu, P. Dille, J. Cross, B. Dias, R. Sargent, and I. Nourbakhsh, "Community empowered air quality monitoring system," in *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, ser. CHI '17. New York, NY, USA: ACM,2017, pp. 1607–1619. [Online]. Available: http://doi.acm.org/10.1145/3025453.3025853

[10] "Use of low-cost sensor technology to monitor air quality & engage citizens," in *SECURE Workshop*, S. R. partnership for Air Pollution health Effects (SHAPE), Ed., COSLA Edinburgh, Scotland, Mar 2016.

[11] R. Bonney, J. L. Shirk, T. B. Phillips, A. Wiggins, H. L. Ballard, A. J. Miller-Rushing, and J. K. Parrish, "Next Steps for Citizen Science," *Science*, vol. 343, no. 6178, pp. 1436 LP – 1437, mar 2014. [Online]. Available: http://science.sciencemag.org/content/343/6178/1436.abstract

[12] K. Kelly, J. Whitaker, A. Petty, C. Widmer, A. Dybwad, D. Sleeth, R. Martin, and A. Butterfield, "Ambient and laboratory evaluation of a low-cost particulate matter sensor," *Environmental Pollution*, vol. 221, pp. 491–500, 2017.

[13] B. A. Sayahi, T. and K. Kelly, "Long-term field evaluation of the Plantower PMS lowcost particulate matter sensors," *Environmental Pollution*, vol. 245, pp. 932–940, 2019.

[14] J. D. Hunter, "Matplotlib: A 2d graphics environment," *Computing in science & engineering*, vol. 9, no. 3, pp. 90–95, 2007.

[15] S. Kassab, M. F. Abu-Hijleh, Q. Al-Shboul, and H. Hamdy, "Student-led tutorials in problem-based learning: educational outcomes and students' perceptions," *Medical Teacher*, vol. 27, no. 6, pp. 521–526, 2005, pMID: 16199359. [Online]. Available: https://doi.org/10.1080/01421590500156186

[16] J. Martino, D. Rodriguez, Orr, Don, and Ruru, "Inverse Distance Weighting (IDW) Interpolation," *GIS Geography*, 21-Feb-2018. [Online]. Available: https://gisgeography.com/inverse-distance-weighting-idw-interpolation