# Using Crowd-Sourced Low-Cost Sensors in a Land Use Regression of PM$_{2.5}$ in 6 US Cities
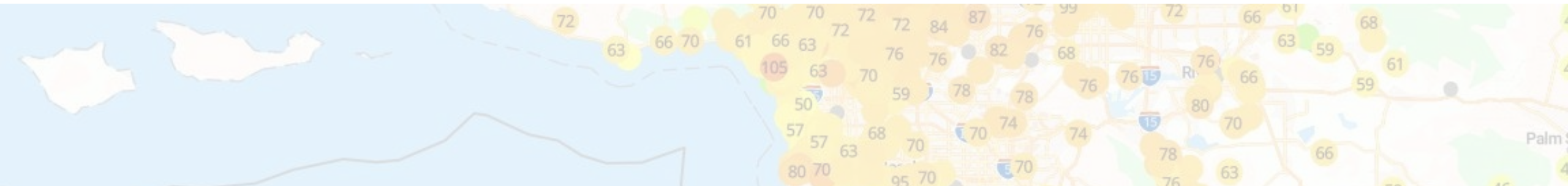
**Tianjun Lu[1]**, Matthew J Bechle[2], Albert A Presto[3], Steve Hankey[4]

[1]California State University, Dominguez Hills; [2]University of Washington
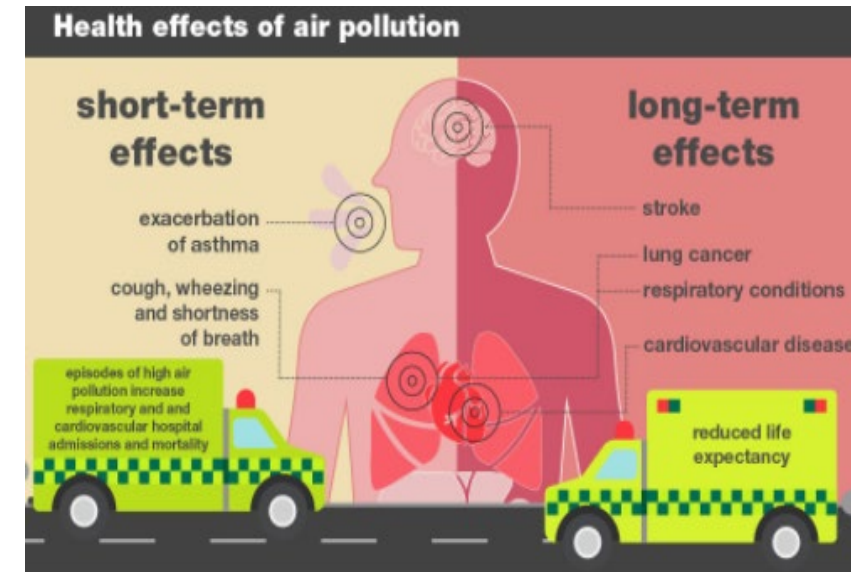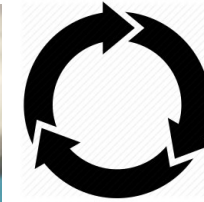
[3]Carnegie Mellon University; [4]Virginia Tech

ASIC 2022

May 11, 2022

**CSUDH**

CALIFORNIA STATE UNIVERSITY, DOMINGUEZ HILLS

# Background and Motivation

- **Health-promoting** cities and air quality.

- Health effects; policy; air quality **monitoring**.

- **Valuable** regulatory monitors.

- Growing global interest in **public data collection**.



National Ambient Air Quality Standards



The Clean Air Act

**EPA** United States Environmental Protection Agency



Health effects of air pollution

short-term effects

exacerbation of asthma

cough, wheezing and shortness of breath

episodes of high air pollution increase respiratory and and cardiovascular hospital admissions and mortality

long-term effects

stroke

lung cancer

respiratory conditions

cardiovascular disease

reduced life expectancy

**Crowd-sourced efforts in exposure assessment**

# Low-Cost Sensing

## Low-cost air quality sensing

- **Dense** fixed sensor network.
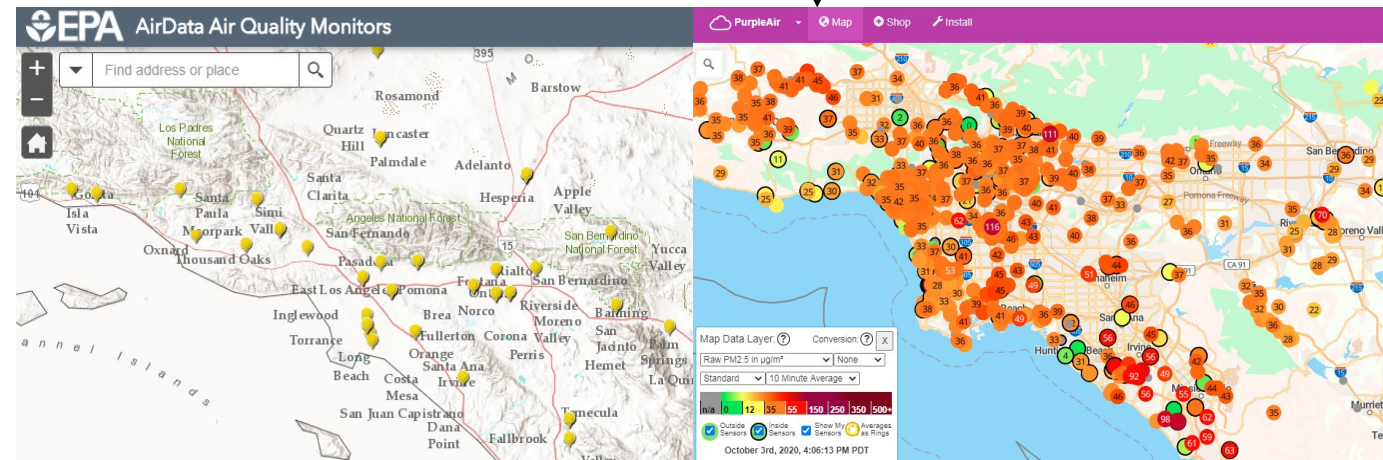- **Community** engagement.
- **"Open" data**.

## Data quality

- Relative humidity, temperature.
- Careful lab and in-field **calibrations**.
- **Well correlation** with reference measurements.
- **Emerging calibration** efforts.



Low-cost sensors

PurpleAir network



EPA air quality monitors vs. PurpleAir sensors

# How Low-cost Sensing Help?



**Regulatory monitors**



**Geographic variables**

**Traditional air quality models**

**+**



**Low-cost sensors**

**Crowd-sourced data**

- Little research assessed the **utility** of such growing network from **multiple** cities in land use regression (LUR).

- Possibility to **improve** the LUR model **to capture spatial variability**?

4

# Existing National LUR Model

## CACES LUR

PLS-UK partitions annual average concentrations into

- (1) a **variance** component that accounts for spatial and non-spatial variability.

- (2) a **mean** component based on a small number of reduced dimension variables from partial least squares of a large number of independent variables (Kim et al., 2020).

| Category | Measure | Note[a] |
|---|---|---|
| Traffic | Distance to the nearest road (0.05-15 km) | Any available road |
| Population | Sum (0.5-3 km) | Population in block groups |
| Land use/land cover (Urban) | Percent (0.05-15 km) | Urban or built-up land, etc. |
| Land use/land cover (Rural) | Percent (0.05-15 km) | Agriculture, forest, water, etc. |
| Position | Coordinates | Longitude, latitude |
| Source | Distance to the nearest source | Coastline, railroad, airport, etc. |
| Emission | Sum of cite-specific facility emissions (3-30 km) | $PM_{2.5}$ |
| Vegetation | Quantiles (0.5-10 km) | Normalized Difference Vegetation Index |
| Imperviousness | Percent (0.05-5 km) | Impervious surface value |
| Elevation | Counts of points above/below a threshold (1-5 km) | Elevation value |
| Satellite estimate | Grid-level estimates | $PM_{2.5}$ |

[a]Detailed information can be found from the CACES LUR modeling study (Kim et al., 2020).
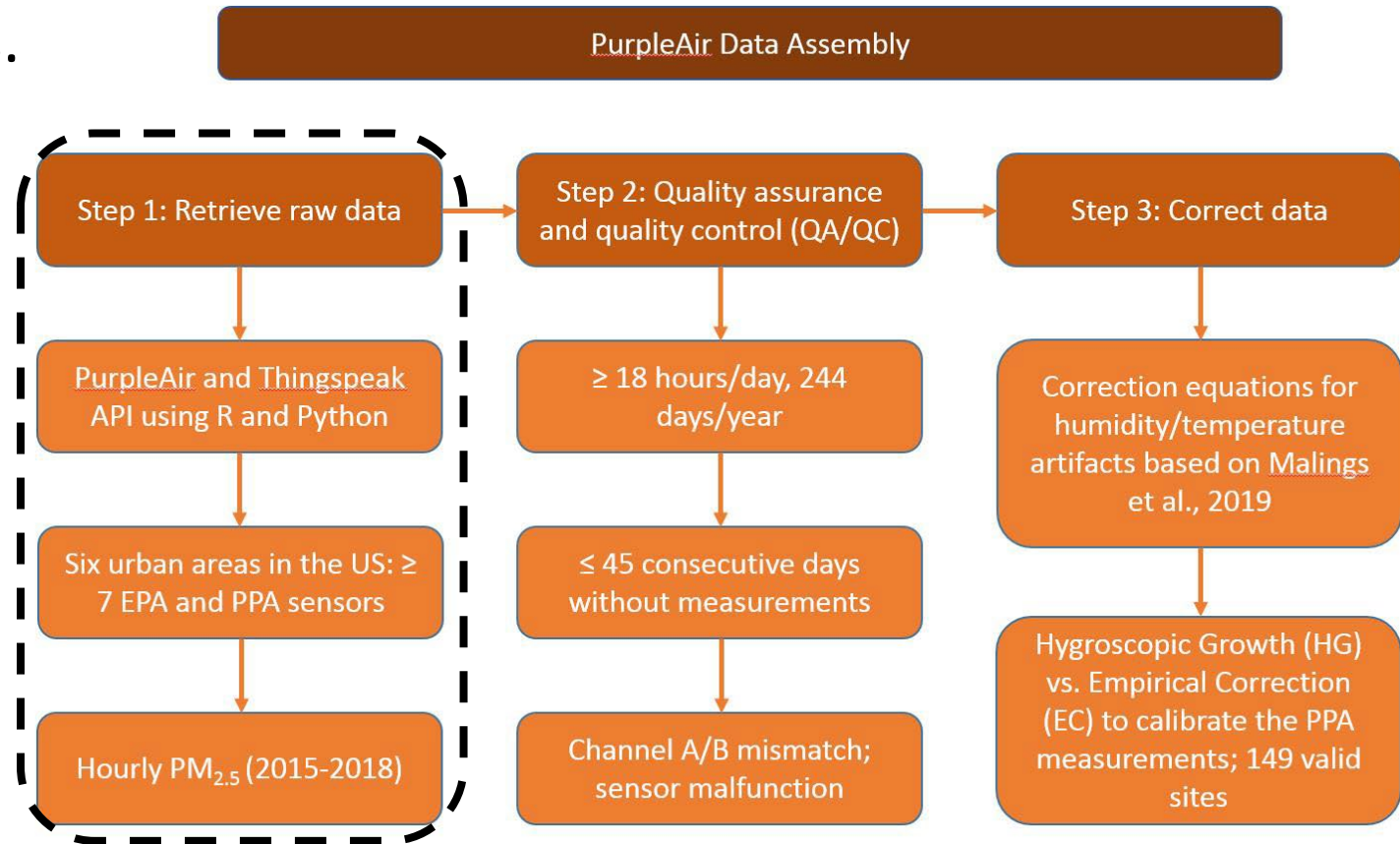
11 categories of geographic variables
339 independent variables
757 regulatory $PM_{2.5}$ monitoring sites

**CACES LUR (random 10-fold CV: $R^2$ = 0.83; standardized RMSE = 0.13)**

CACES: Center for Air, Climate, and Energy Solutions
PLS-UK: Partial Least Squares-Universal Kriging

# PurpleAir (PPA) Data Preparation

## PPA data assembly

- **Six cities**: ≥ 7 EPA and PPA sensors.



**PurpleAir Data Assembly**

| Step 1: Retrieve raw data | Step 2: Quality assurance and quality control (QA/QC) | Step 3: Correct data |
|---|---|---|
| PurpleAir and Thingspeak API using R and Python | ≥ 18 hours/day, 244 days/year | Correction equations for humidity/temperature artifacts based on Malings et al., 2019 |
| Six urban areas in the US: ≥ 7 EPA and PPA sensors | ≤ 45 consecutive days without measurements | Hygroscopic Growth (HG) vs. Empirical Correction (EC) to calibrate the PPA measurements; 149 valid sites |
| Hourly $PM_{2.5}$ (2015-2018) | Channel A/B mismatch; sensor malfunction | |

# PurpleAir (PPA) Data Preparation
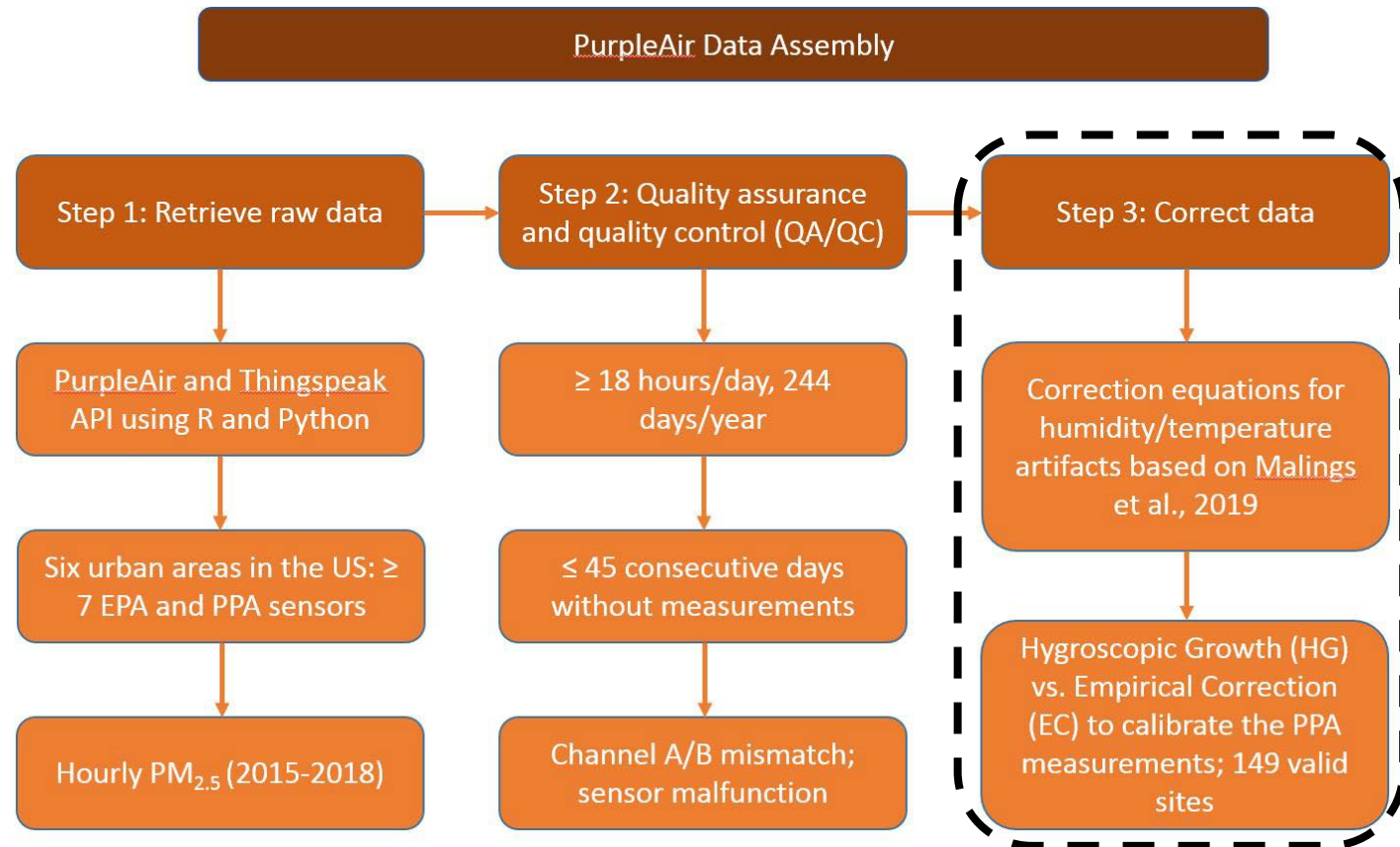
## PPA data assembly

- **Six cities**: ≥ 7 EPA and PPA sensors.

- **QA/QC**:
  - same criteria as the CACES LUR
  - channel mismatch (removing hours when the absolute difference was larger than **3 µg/m³** or **20%** of the maximum channel readings, whichever is greater (Malings et al., 2019).



PurpleAir Data Assembly

| Step 1: Retrieve raw data | Step 2: Quality assurance and quality control (QA/QC) | Step 3: Correct data |
|---|---|---|
| PurpleAir and Thingspeak API using R and Python | ≥ 18 hours/day, 244 days/year | Correction equations for humidity/temperature artifacts based on Malings et al., 2019 |
| Six urban areas in the US: ≥ 7 EPA and PPA sensors | ≤ 45 consecutive days without measurements | Hygroscopic Growth (HG) vs. Empirical Correction (EC) to calibrate the PPA measurements; 149 valid sites |
| Hourly PM$_{2.5}$ (2015-2018) | Channel A/B mismatch; sensor malfunction | |

# PurpleAir (PPA) Data Preparation

## PPA data assembly

- **Six cities**: ≥ 7 EPA and PPA sensors.

- **QA/QC**:

  - same criteria as the CACES LUR

  - channel mismatch (removing hours when the absolute difference was larger than **3 μg/m³** or **20%** of the maximum channel readings, whichever is greater (Malings et al., 2019).

- **Data correction**:

  - humidity and temperature artifacts;

  - colocation calibrations.

PurpleAir Data Assembly

| Step 1: Retrieve raw data | Step 2: Quality assurance and quality control (QA/QC) | Step 3: Correct data |
|---|---|---|
| PurpleAir and Thingspeak API using R and Python | ≥ 18 hours/day, 244 days/year | Correction equations for humidity/temperature artifacts based on Malings et al., 2019 |
| Six urban areas in the US: ≥ 7 EPA and PPA sensors | ≤ 45 consecutive days without measurements | Hygroscopic Growth (HG) vs. Empirical Correction (EC) to calibrate the PPA measurements; 149 valid sites |
| Hourly PM$_{2.5}$ (2015-2018) | Channel A/B mismatch; sensor malfunction | |

# LUR Model Development

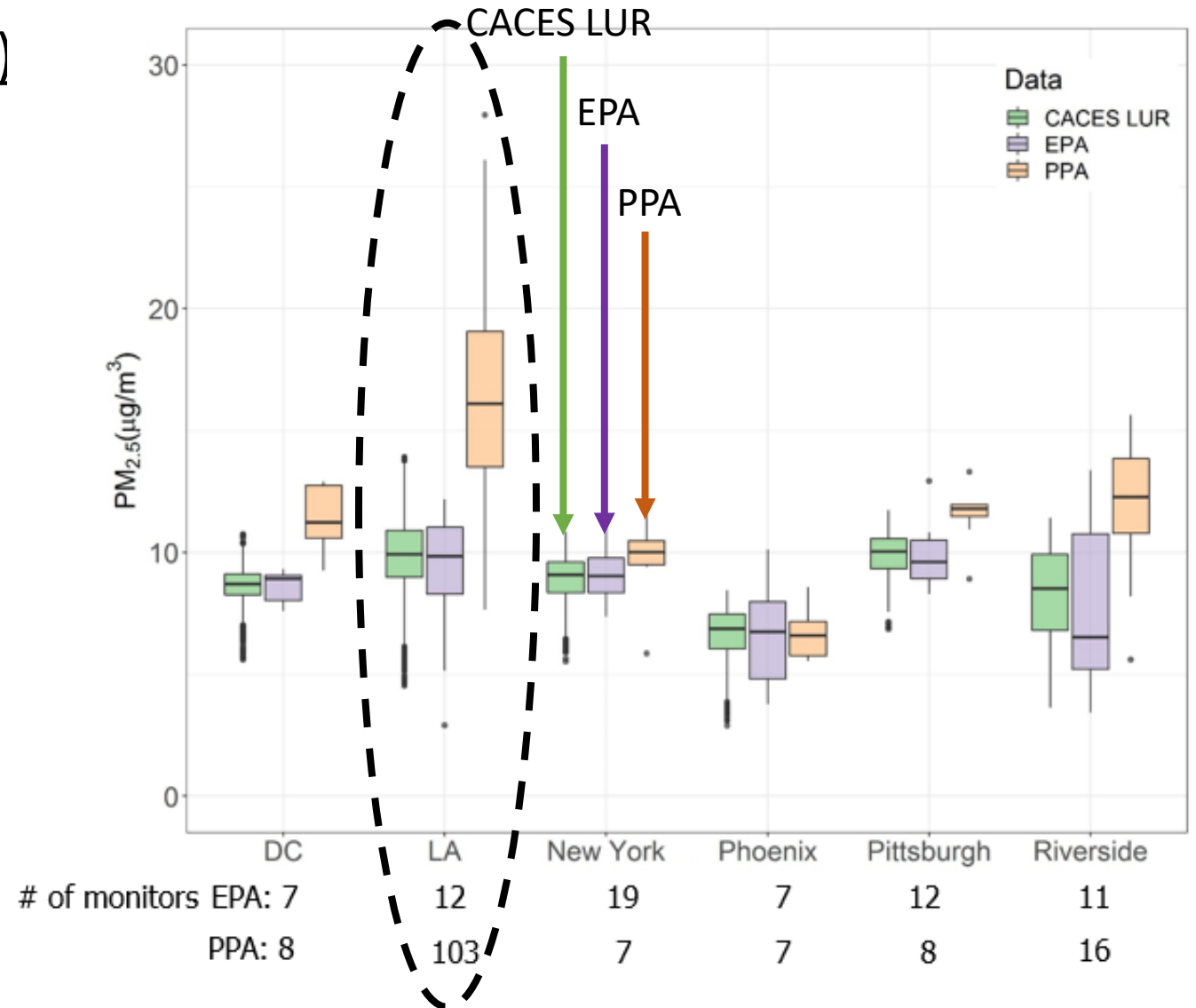**Dependent variables (annual averages)**

- EPA data (national and 6 cities).

- PPA data (6 cities).
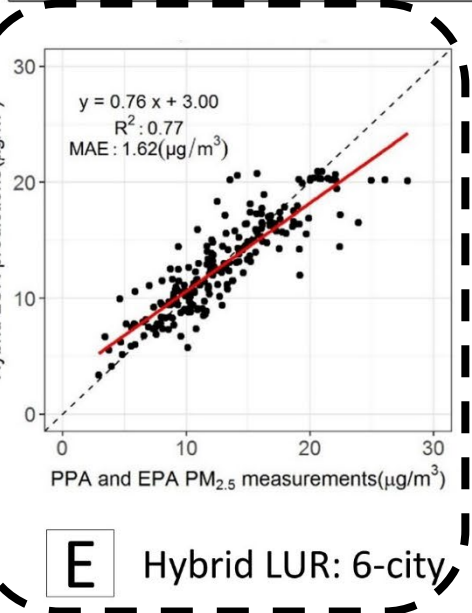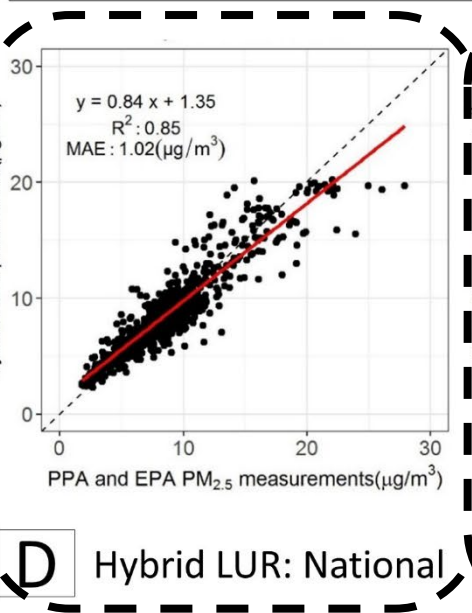
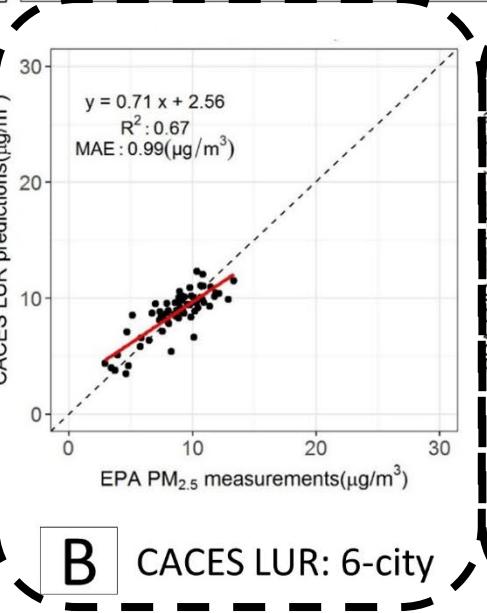- Hybrid (EPA + PPA data).

**Independent variables**

- 11 categories (e.g., traffic, population, land use).

**Modeling approach**

- PLS-UK.

# LUR Model Comparison (Pop-weighted)



Population-weighted PM$_{2.5}$ concentration maps
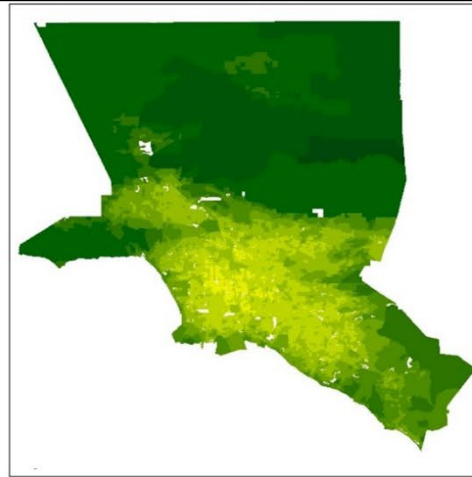
A    CACES LUR: National

B    CACES LUR: 6-city

C    PPA LUR

D    Hybrid LUR: National

E    Hybrid LUR: 6-city

# LUR Model Comparison (Transect Plots)

- Transect plot of the five LUR predictions.

Advantages

- Models with the PPA data were **more spatially variable** than models without.

- Models with the PPA data alone is **not recommended**.

# Variable Importance



**A**  CACES LUR: National

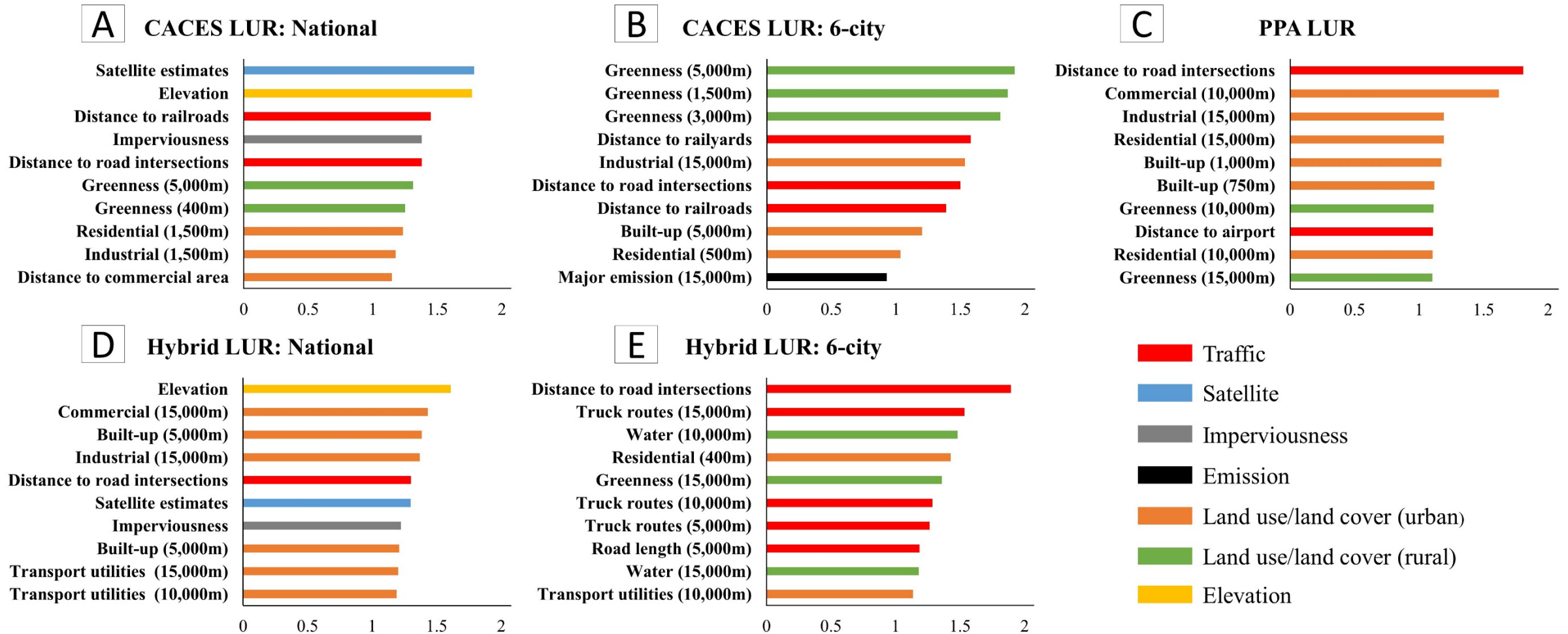| Variable | |
|---|---|
| Satellite estimates | (Satellite) |
| Elevation | (Elevation) |
| Distance to railroads | (Traffic) |
| Imperviousness | (Imperviousness) |
| Distance to road intersections | (Traffic) |
| Greenness (5,000m) | (LULC rural) |
| Greenness (400m) | (LULC rural) |
| Residential (1,500m) | (LULC urban) |
| Industrial (1,500m) | (LULC urban) |
| Distance to commercial area | (LULC urban) |

**B**  CACES LUR: 6-city

- Greenness (5,000m)
- Greenness (1,500m)
- Greenness (3,000m)
- Distance to railyards
- Industrial (15,000m)
- Distance to road intersections
- Distance to railroads
- Built-up (5,000m)
- Residential (500m)
- Major emission (15,000m)

**C**  PPA LUR

- Distance to road intersections
- Commercial (10,000m)
- Industrial (15,000m)
- Residential (15,000m)
- Built-up (1,000m)
- Built-up (750m)
- Greenness (10,000m)
- Distance to airport
- Residential (10,000m)
- Greenness (15,000m)

**D**  Hybrid LUR: National

- Elevation
- Commercial (15,000m)
- Built-up (5,000m)
- Industrial (15,000m)
- Distance to road intersections
- Satellite estimates
- Imperviousness
- Built-up (5,000m)
- Transport utilities (15,000m)
- Transport utilities (10,000m)

**E**  Hybrid LUR: 6-city

- Distance to road intersections
- Truck routes (15,000m)
- Water (10,000m)
- Residential (400m)
- Greenness (15,000m)
- Truck routes (10,000m)
- Truck routes (5,000m)
- Road length (5,000m)
- Water (15,000m)
- Transport utilities (10,000m)

Legend:
- Traffic (red)
- Satellite (blue)
- Imperviousness (gray)
- Emission (black)
- Land use/land cover (urban) (orange)
- Land use/land cover (rural) (green)
- Elevation (yellow)

- Traffic and land use variables were important variables for models with the PPA data; strength of capturing "hotspots".

# Summary and Implications

- Hybrid models may capture **small-scale variations** that may be **missed** by the regulatory-based models

- Valuable dataset for LUR if data is **carefully** cleaned and calibrated.


- With available national correction approaches (Barkjohn et al., 2021), additional cities would help assess tradeoffs in **national vs. local corrections**.

- Calibrations by **co-locating** PPA sensors with regulatory-grade monitors in additional cities may help reduce bias.

- Further empirical investigation is warranted in hybrid models with **additional sensors** from **larger areas and multiple cities**.

- Neighborhood **planning and design**; clean streets; guidance on **outdoor** activities; interventions.

# Acknowledgement and Contact

**Using crowd-sourced low-cost sensors in a land use regression of $PM_{2.5}$ in 6 US cities**

Tianjun Lu[1] · Matthew J. Bechle[2] · Yanyu Wan[3] · Albert A. Presto[3] · Steve Hankey[4]

2  Department of Civil & Environmental Engineering, University of Washington, 201 More Hall, Seattle, WA 98195, USA

3  Department of Mechanical Engineering, Carnegie Mellon University, 2115 Doherty Hall, Pittsburgh, PA 15213, USA

4  School of Public and International Affairs, Virginia Tech, 140 Otey Street, Blacksburg, VA 24061, USA

**Tianjun (Luke) Lu**

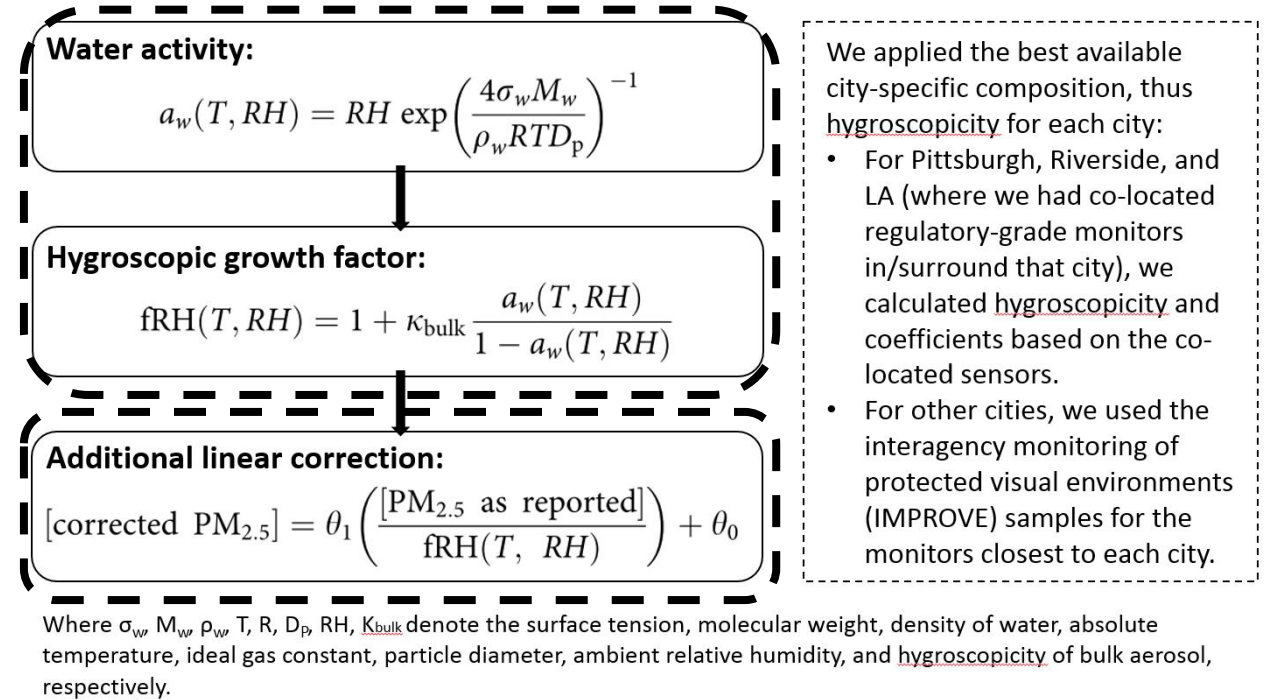tilu@csudh.edu

https://tianjunlu.weebly.com/

14

# Supplemental Material

# Hygroscopic Growth (HG) Correction

## HG correction

- Adjusted to be "Beta Attenuation Monitors (BAM) equivalent".

- **Over** prediction at high RH and **under** prediction of particles < 300 nm.

- Cities **with/without** co-located PPA sensors.

- Either the **Pittsburgh** (New York, DC) or the **Riverside** regression (LA, Phoenix) based on similarities in climate and $PM_{2.5}$ composition.
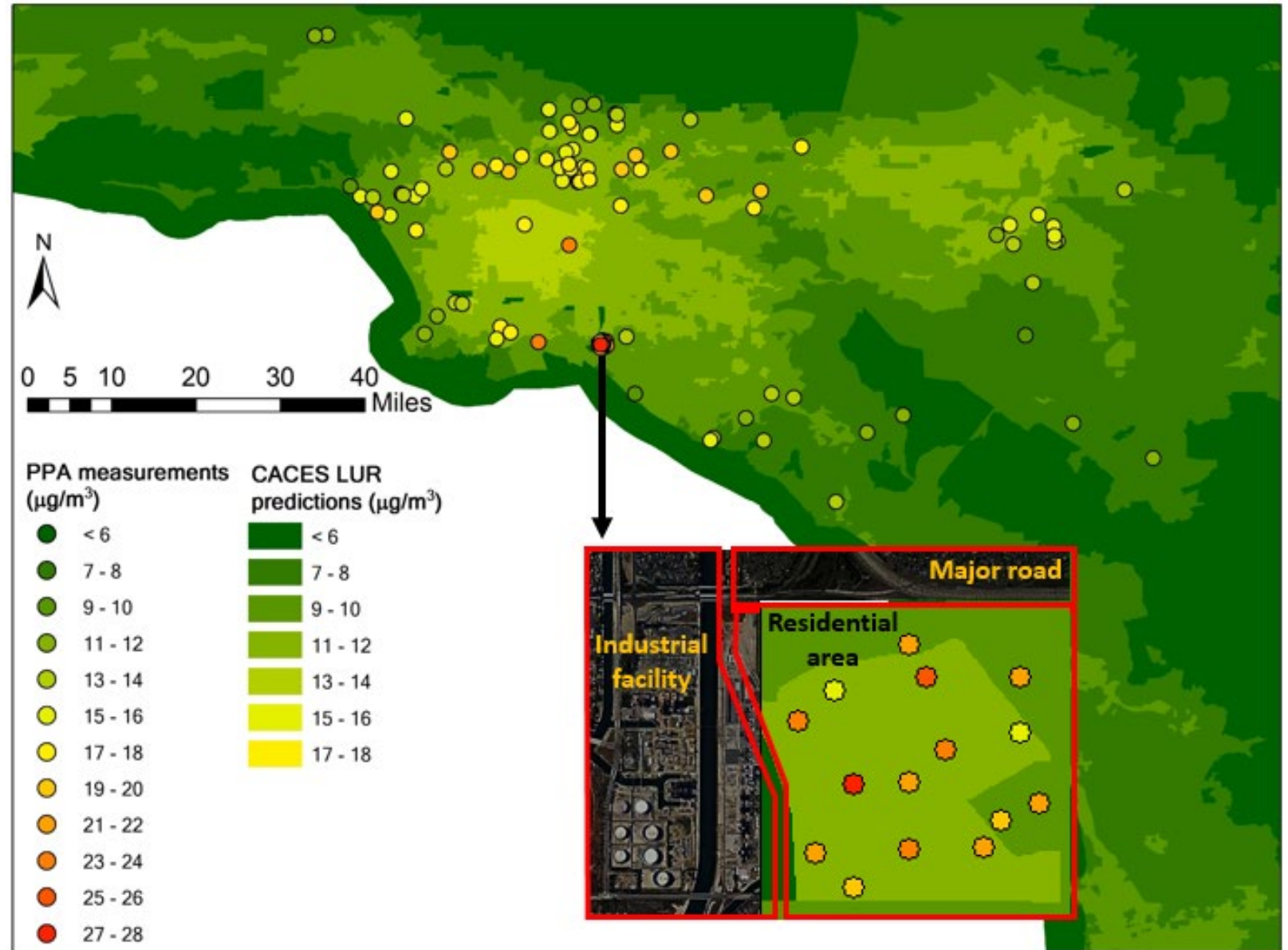
**Water activity:**

$$a_w(T, RH) = RH \exp\left(\frac{4\sigma_w M_w}{\rho_w RTD_{\mathrm{p}}}\right)^{-1}$$

**Hygroscopic growth factor:**

$$\mathrm{fRH}(T, RH) = 1 + \kappa_{\mathrm{bulk}} \frac{a_w(T, RH)}{1 - a_w(T, RH)}$$

**Additional linear correction:**

$$[\text{corrected } PM_{2.5}] = \theta_1\left(\frac{[PM_{2.5} \text{ as reported}]}{\mathrm{fRH}(T, \ RH)}\right) + \theta_0$$

We applied the best available city-specific composition, thus hygroscopicity for each city:
- For Pittsburgh, Riverside, and LA (where we had co-located regulatory-grade monitors in/surround that city), we calculated hygroscopicity and coefficients based on the co-located sensors.
- For other cities, we used the interagency monitoring of protected visual environments (IMPROVE) samples for the monitors closest to each city.

Where $\sigma_w$, $M_w$, $\rho_w$, T, R, $D_P$, RH, $K_{bulk}$ denote the surface tension, molecular weight, density of water, absolute temperature, ideal gas constant, particle diameter, ambient relative humidity, and hygroscopicity of bulk aerosol, respectively.

**Hygroscopic Growth (HG) Correction Method**

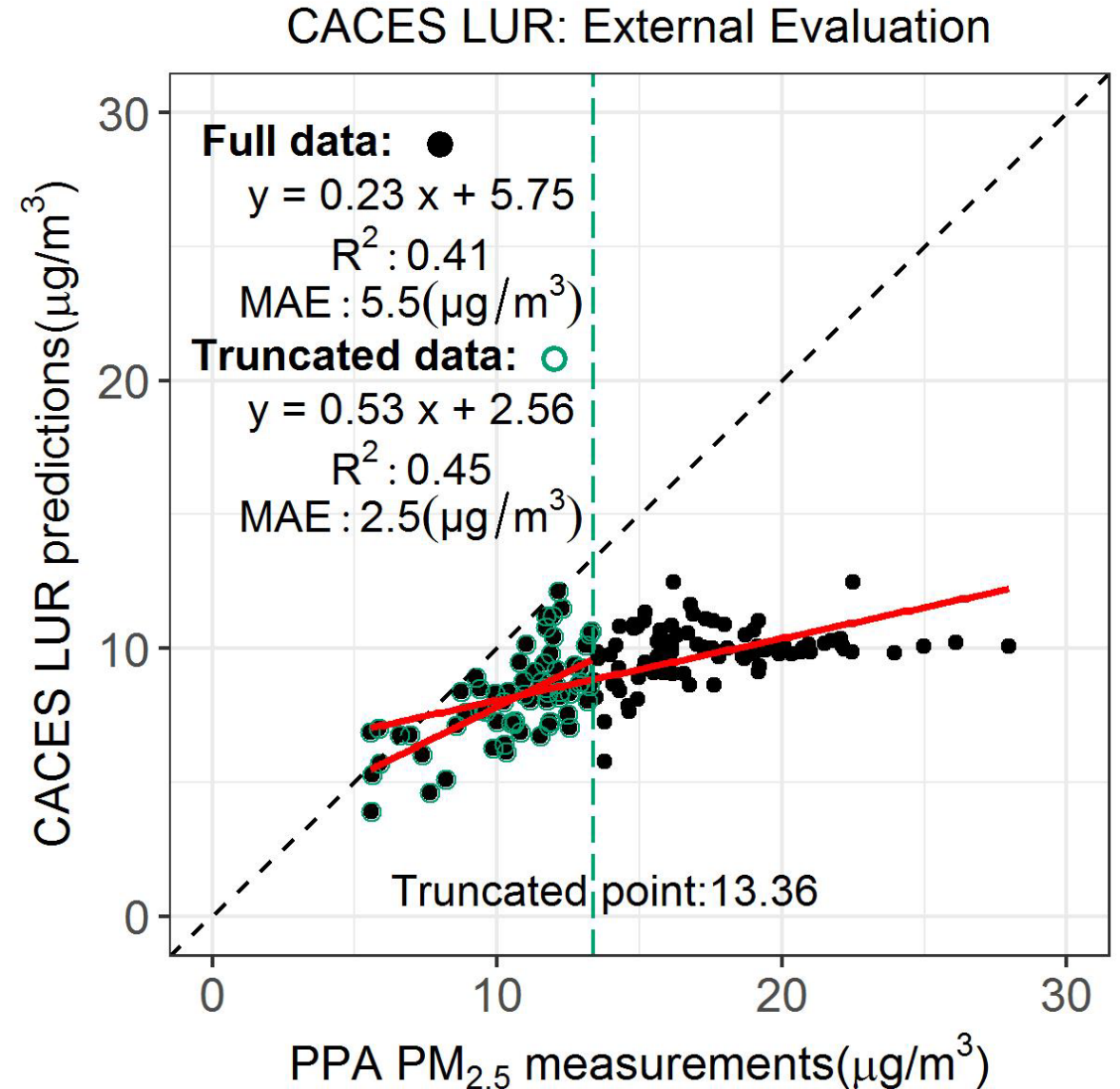# CACES LUR Estimates vs. PPA Measurements

## Differences

- Spatial **mismatch**.

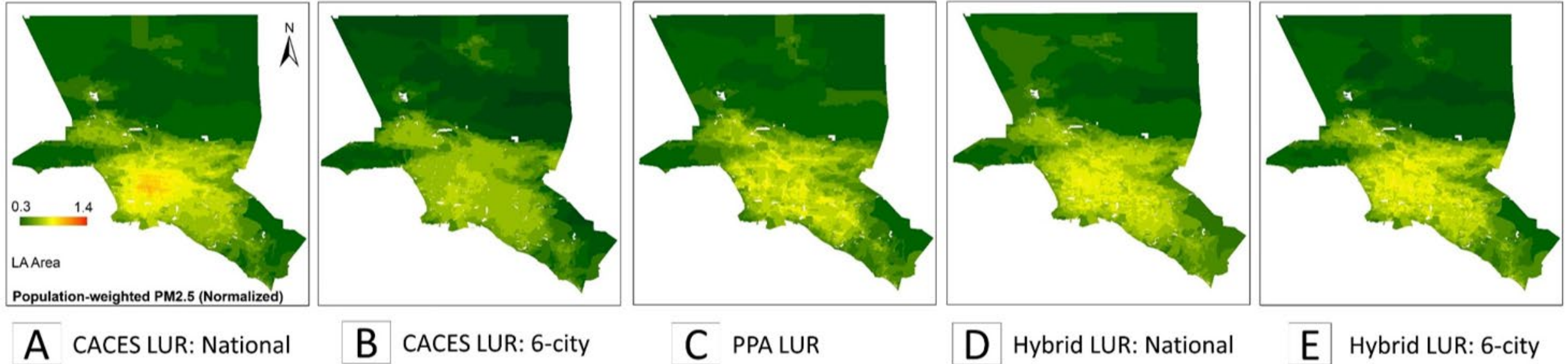- Uncaptured "**hotspots**": industrial facilities and highway.

# External Evaluation of CACES LUR Estimates

<u>Uncaptured</u>

- Miss "**hotspots**".



CACES LUR: External Evaluation

**Full data:** ●
$$y = 0.23\,x + 5.75$$
$$R^2 : 0.41$$
$$MAE : 5.5\,(\mu g/m^3)$$
**Truncated data:** ○
$$y = 0.53\,x + 2.56$$
$$R^2 : 0.45$$
$$MAE : 2.5\,(\mu g/m^3)$$

Truncated point: 13.36

CACES LUR predictions $(\mu g/m^3)$

PPA PM$_{2.5}$ measurements $(\mu g/m^3)$

# LUR Model Comparison (Normalized Pop-weighted)



| A | CACES LUR: National | B | CACES LUR: 6-city | C | PPA LUR | D | Hybrid LUR: National | E | Hybrid LUR: 6-city |

**Hybrid models not only benefit from capturing "hotspots" but are also consistent with the regional spatial trends in the CACES LUR models.**
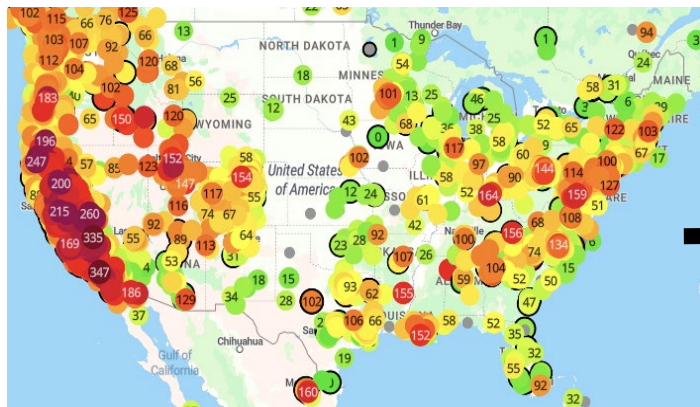
Normalized Population-weighted PM$_{2.5}$ concentration maps

# Hybrid LUR: Mitigating Uncertainty

- LUR using only the PPA data may be reasonable; however, **consistently higher** predictions.

- Hybrid models suggest the **value** of combinations.

- Future LUR models: investigating **factors behind** model improvement.
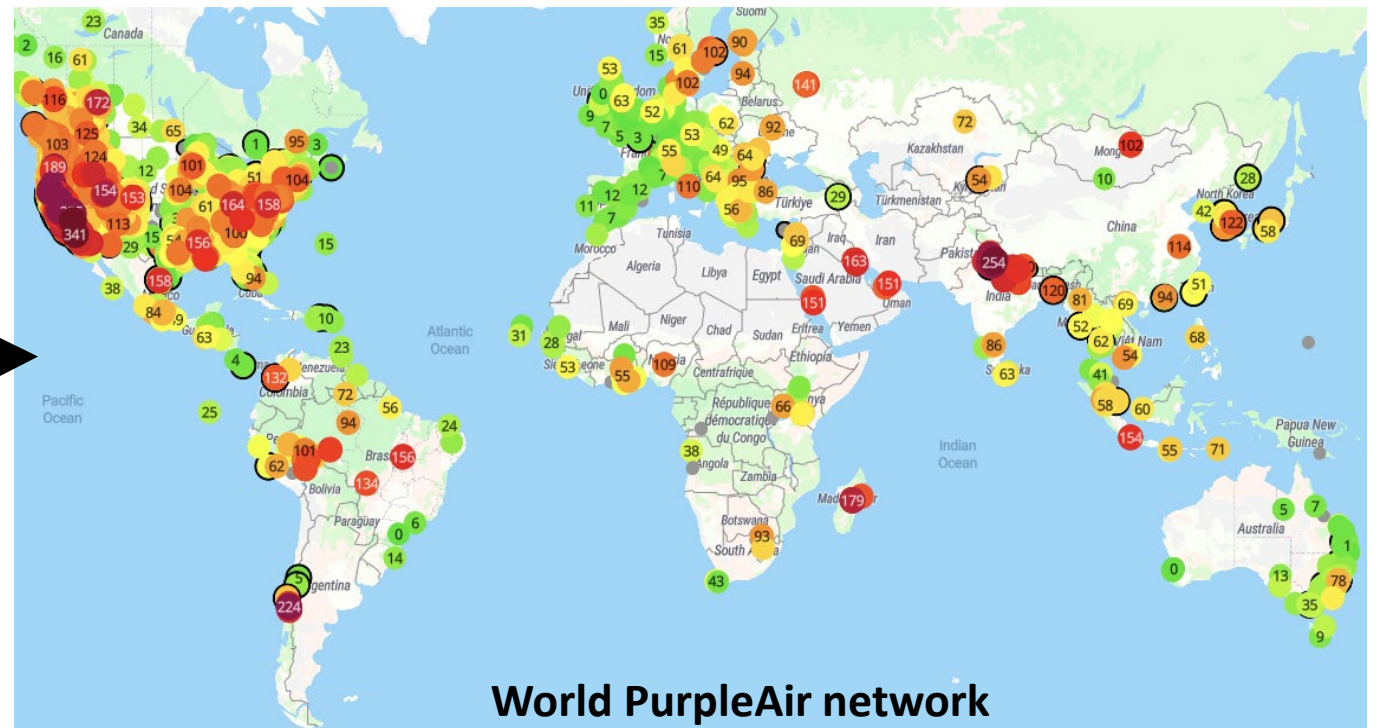
# Low-Cost Sensing in Air Quality Models

- **Representative** samples.
- **Fast-growing** network.
- **Rural areas and low- and middle-income countries** (sparse regulatory monitors).

- Neighborhood **planning and design**; clean streets; guidance on **outdoor** activities; interventions



US PurpleAir network



World PurpleAir network