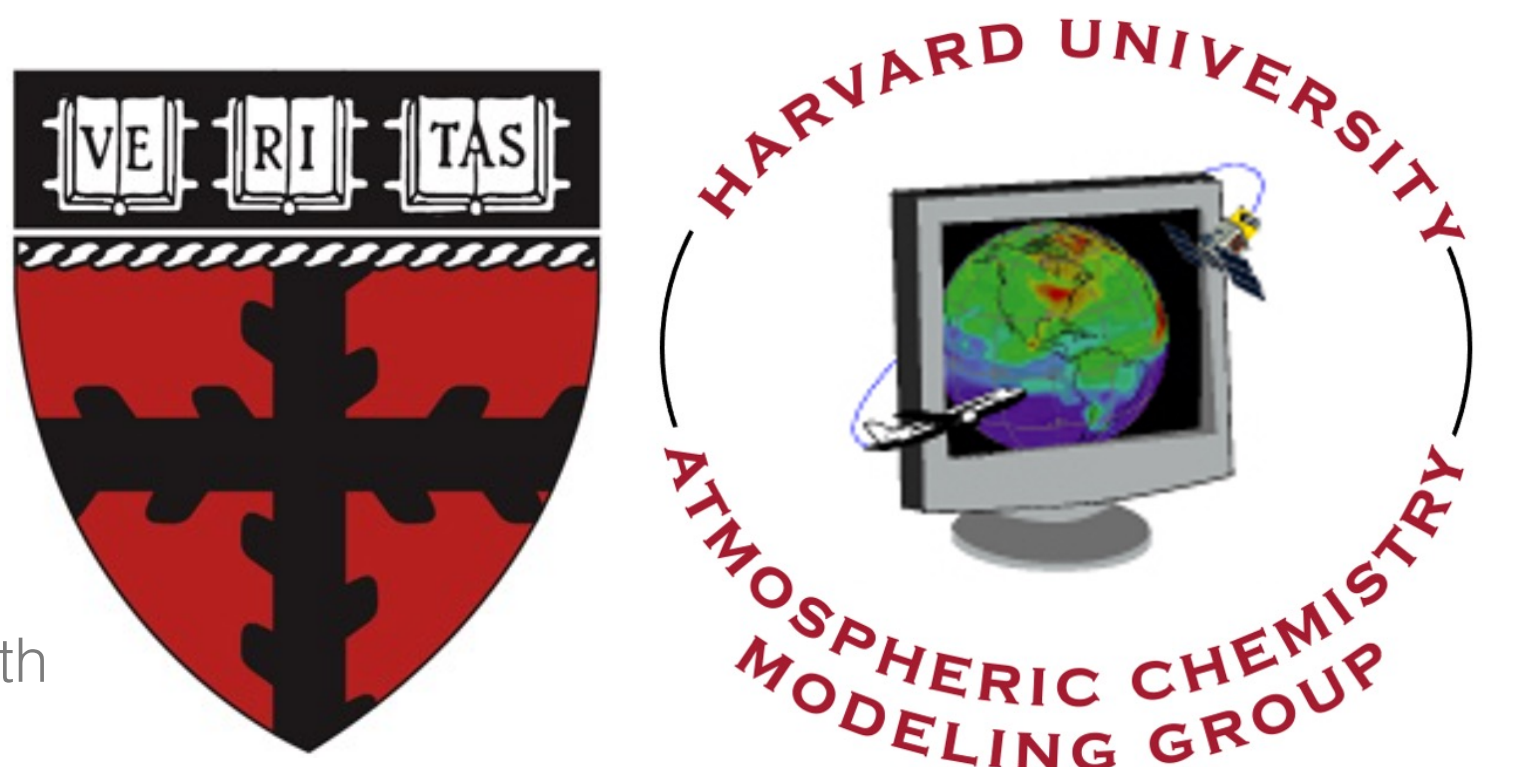


# A continuous 2011-2022 record of fine particulate matter (PM<sub>2.5</sub>) in East Asia at daily 2-km resolution from geostationary satellite observations: population exposure and long-term trends

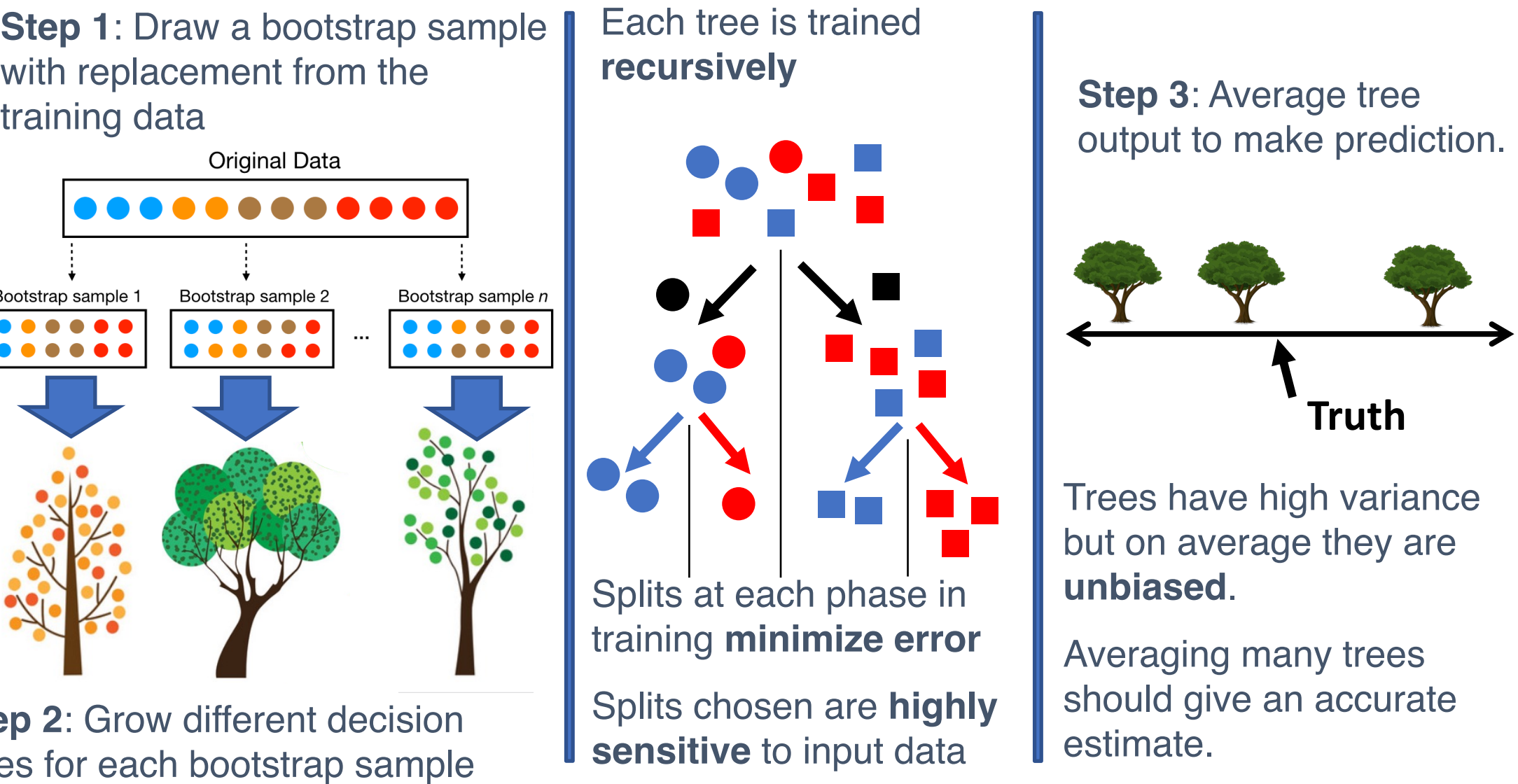
Drew C. Pendergrass<sup>1</sup>, Daniel J. Jacob<sup>1</sup>, Yujin Oak<sup>1</sup>, Jeewoo Lee<sup>2</sup>, Minseok Kim<sup>2</sup>, Jhoon Kim<sup>2,3</sup>, Seoyoung Lee<sup>4,5</sup>, Shixian Zhai<sup>6</sup>, Hitoshi Irie<sup>7</sup>, and Hong Liao<sup>8</sup>

<sup>1</sup>School of Engineering and Applied Sciences, Harvard University, Cambridge, Mass., USA; <sup>2</sup>Department of Atmospheric Sciences, Yonsei University, Seoul, South Korea; <sup>3</sup>Particulate Matter Research Institute, Samsung Advanced Institute of Technology (SAIT), Suwon, South Korea; <sup>4</sup>University of Maryland Baltimore County, Baltimore, Md., USA; <sup>5</sup>NASA Goddard Space Flight Center, Greenbelt, Md., USA; <sup>6</sup>Hong Kong University of Science and Technology, Hong Kong SAR, China; <sup>7</sup>Center for Environmental Remote Sensing (CEReS), Chiba University, Chiba, Japan; <sup>8</sup>Jiangsu Key Laboratory of Atmospheric Environment Monitoring and Pollution Control, Jiangsu Collaborative Innovation, Center of Atmospheric Environment and Equipment Technology, School of Environmental Science and Engineering, Nanjing University of Information Science and Technology, Nanjing, Jiangsu, China



**Abstract.** We applied a random forest (RF) algorithm to 2011-2022 aerosol optical depth (AOD) observations from the Geostationary Ocean Color Imager (GOCI) I and II instruments over East Asia to infer 24-h daily surface fine particulate matter (PM<sub>2.5</sub>) concentrations at continuous 2×2 km<sup>2</sup> resolution over eastern China, South Korea, and Japan. The RF uses PM<sub>2.5</sub> observations from the national surface networks as training data. Predictor variables along with AOD include meteorological data, land use indices, precursor emission inventories, and chemical transport model (CTM) output. Missing AOD data are gap-filled by a separate RF fit. For South Korea, PM<sub>2.5</sub> training data before 2015 (when the surface network began measuring PM<sub>2.5</sub>) are obtained with a separate RF trained on the available network data for other pollutants including coarse particulate matter (PM<sub>10</sub>). For China, PM<sub>2.5</sub> training data before 2014 are from the US embassy and consulates. The resulting dataset offers a unique continuous record of PM<sub>2.5</sub> over a period of rapid changes in the regulation of precursor emissions. The GOCI PM<sub>2.5</sub> data are successful in reproducing the surface network observations including extreme events. We find that after peaking in 2014 (China) or 2013 (South Korea, Japan), population-weighted PM<sub>2.5</sub> has been steadily decreasing in all three countries through 2022, and no region has been left behind. The Seoul region showed no decrease in winter PM<sub>2.5</sub> until 2019 but more recent years show a decrease. We evaluate our product over an extreme pollution event in Seoul and find that the predicted distribution is indistinguishable from observations, while our previous 6×6 km<sup>2</sup> product suffers from smoothing errors due to resolution. We find that in early years in Seoul and Shanghai weekdays are more polluted than weekends but the sign flips later in the study period, a result uniquely enabled by expanded temporal coverage.

## The random forest algorithm



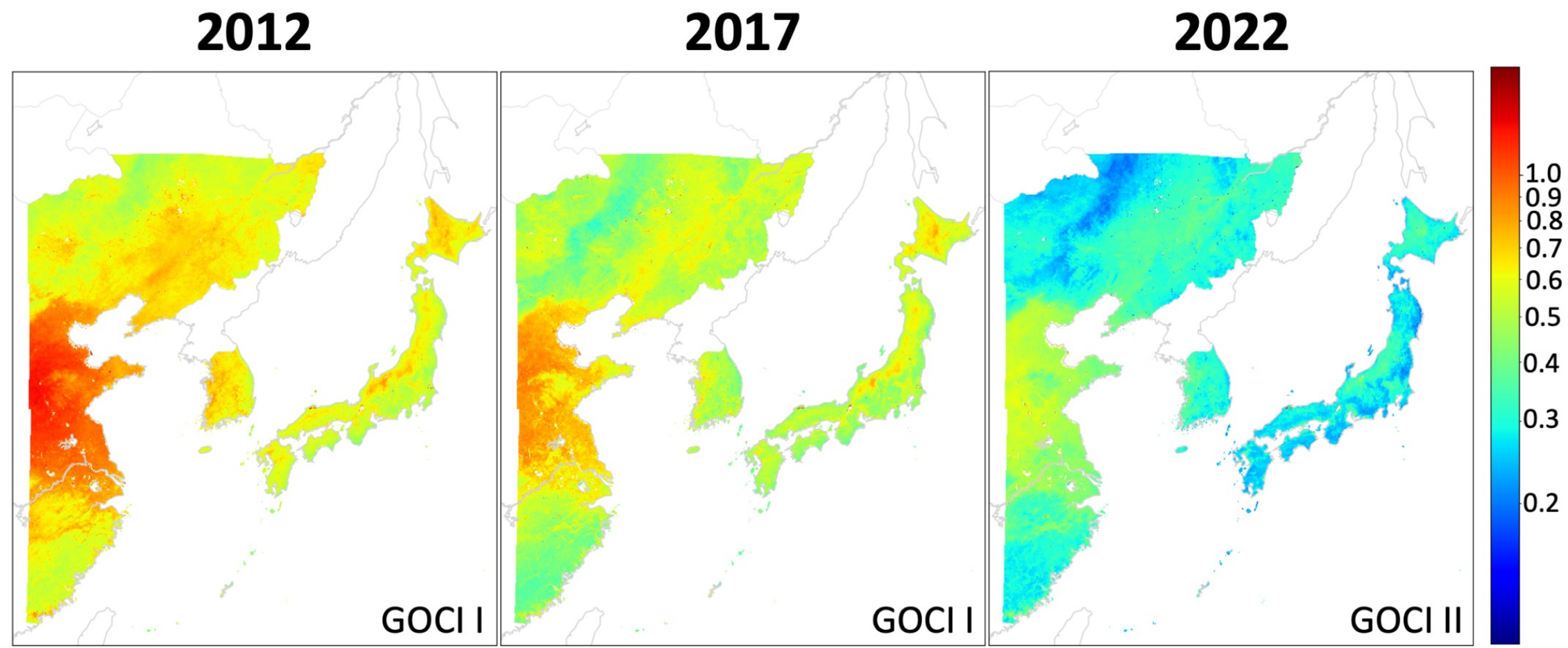
## Training data

We train our random forest (RF) machine learning algorithm to predict 24-hr surface PM<sub>2.5</sub> observed at sites in eastern China, South Korea, and Japan.

The RF algorithm predicts surface PM<sub>2.5</sub> using:

- GOCI I (2011-20) and GOCI II (2021-22) AOD, gap-filled with separate RF fit (daily crossvalidation R<sup>2</sup>: 0.91).
- ERA5 reanalysis (boundary layer height, sea level pressure) and ERA5-Land replay (2m temperature, relative humidity, 10m u/v winds)
- Emissions (KORUSv5 NO<sub>x</sub>, SO<sub>2</sub>, NH<sub>3</sub>) and bias-corrected GEOS-Chem CTM AOD monthly means
- Land use (land cover type, population density, elevation, NDVI) and time/country metadata

### Evolution of gap-filled GOCI AOD

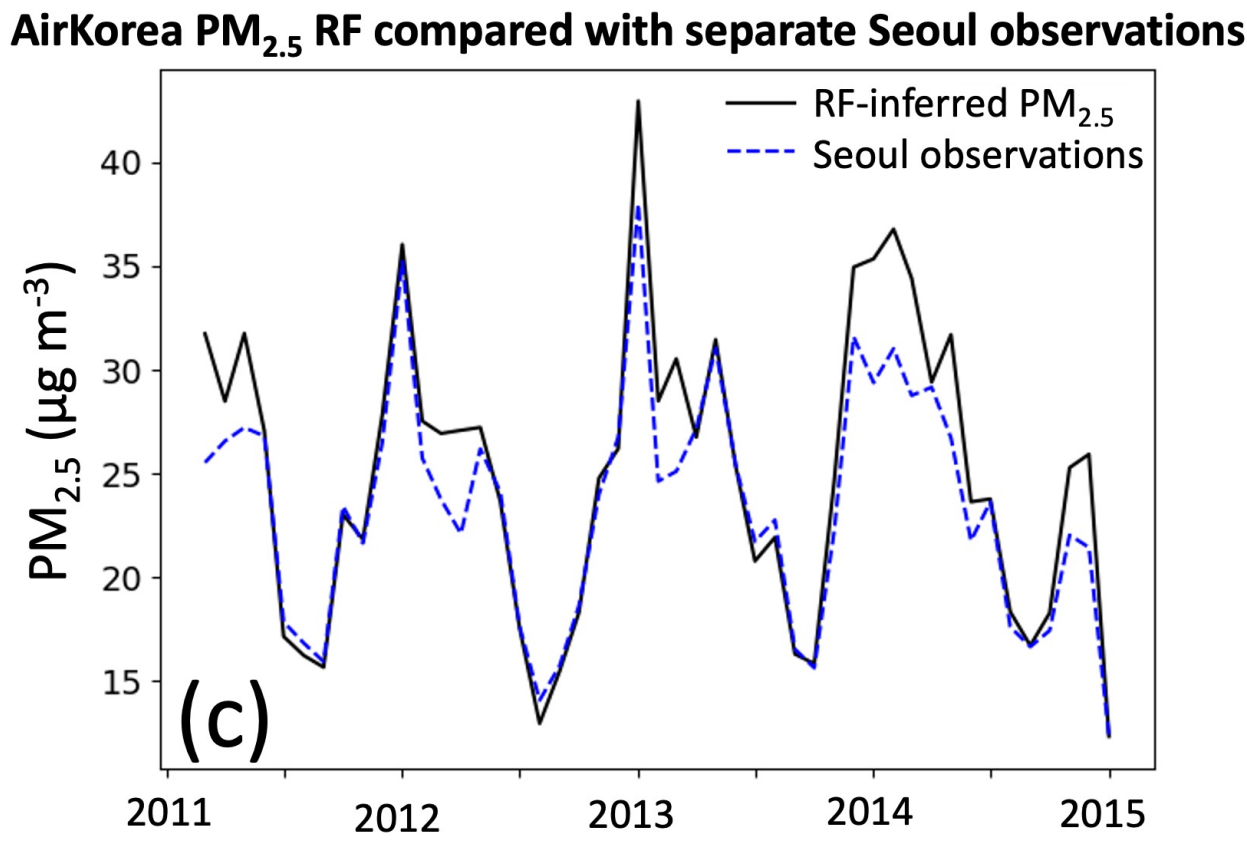


# We trained a machine learning algorithm on a fused 2011-2022 record of aerosol optical depth, air quality network data, meteorology, land use data, chemical transport model data, and more.

# We then produced a continuous 24-h PM<sub>2.5</sub> dataset for eastern China, South Korea, and Japan at 2x2 km<sup>2</sup> resolution.

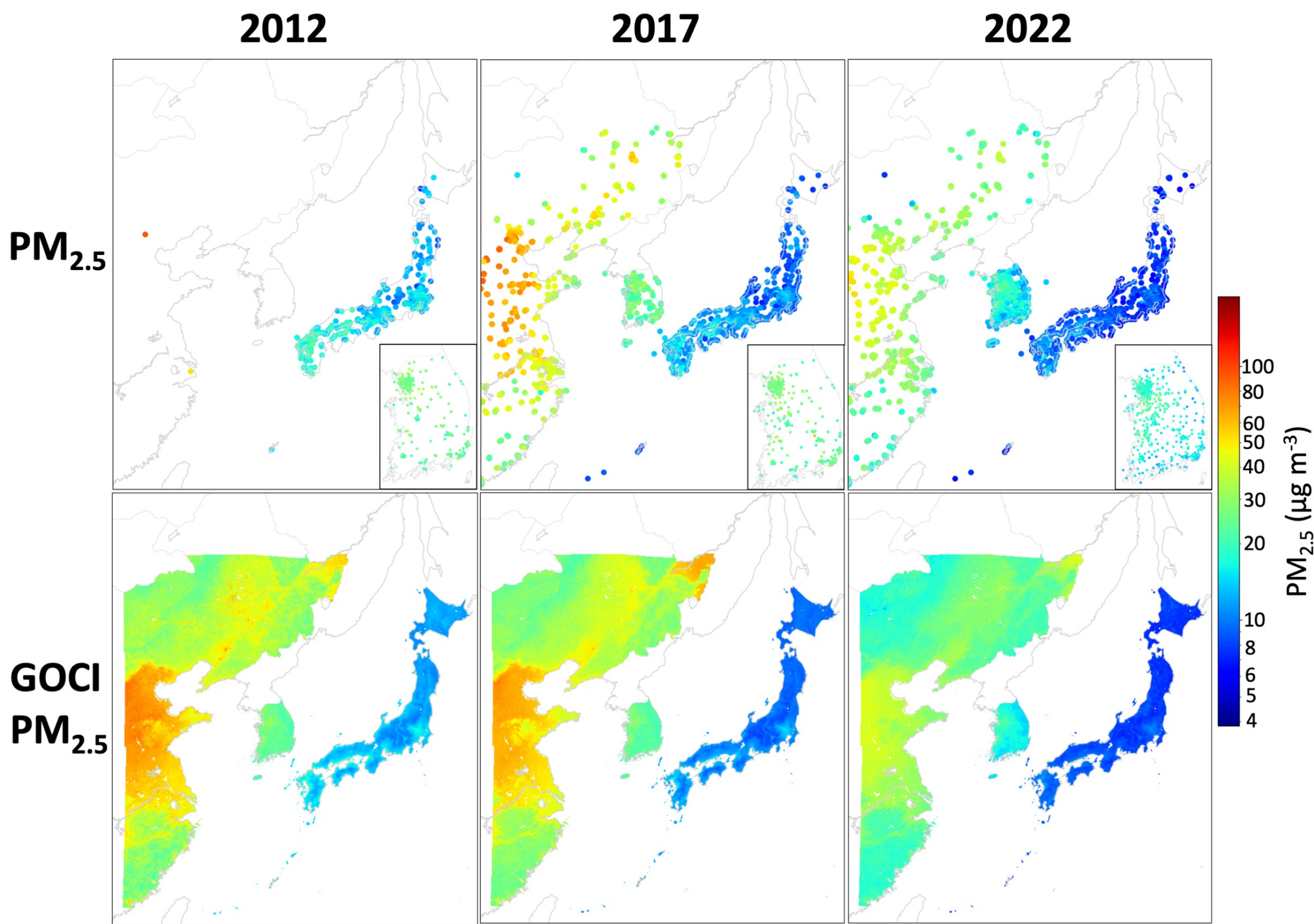
## Inferring S. Korean PM<sub>2.5</sub> prior to 2015

Prior to the January 2015 addition of PM<sub>2.5</sub> to the network, pollutants CO, O<sub>3</sub>, NO<sub>2</sub>, SO<sub>2</sub>, and PM<sub>10</sub> were measured at hourly resolution at sites across South Korea. We train a separate RF on the covariates at left plus non-PM<sub>2.5</sub> pollutants to predict daily 2011-2014 PM<sub>2.5</sub> at sites in Korea (daily crossvalidation R<sup>2</sup>: 0.88). We evaluate the results using a separate testing dataset which measures PM<sub>2.5</sub> in the Seoul Metropolitan Area (SMA) from 2006 (right). We find tight correlation (daily crossvalidation R<sup>2</sup>: 0.91) and a modest high bias of 1.6 µg m<sup>-3</sup>.



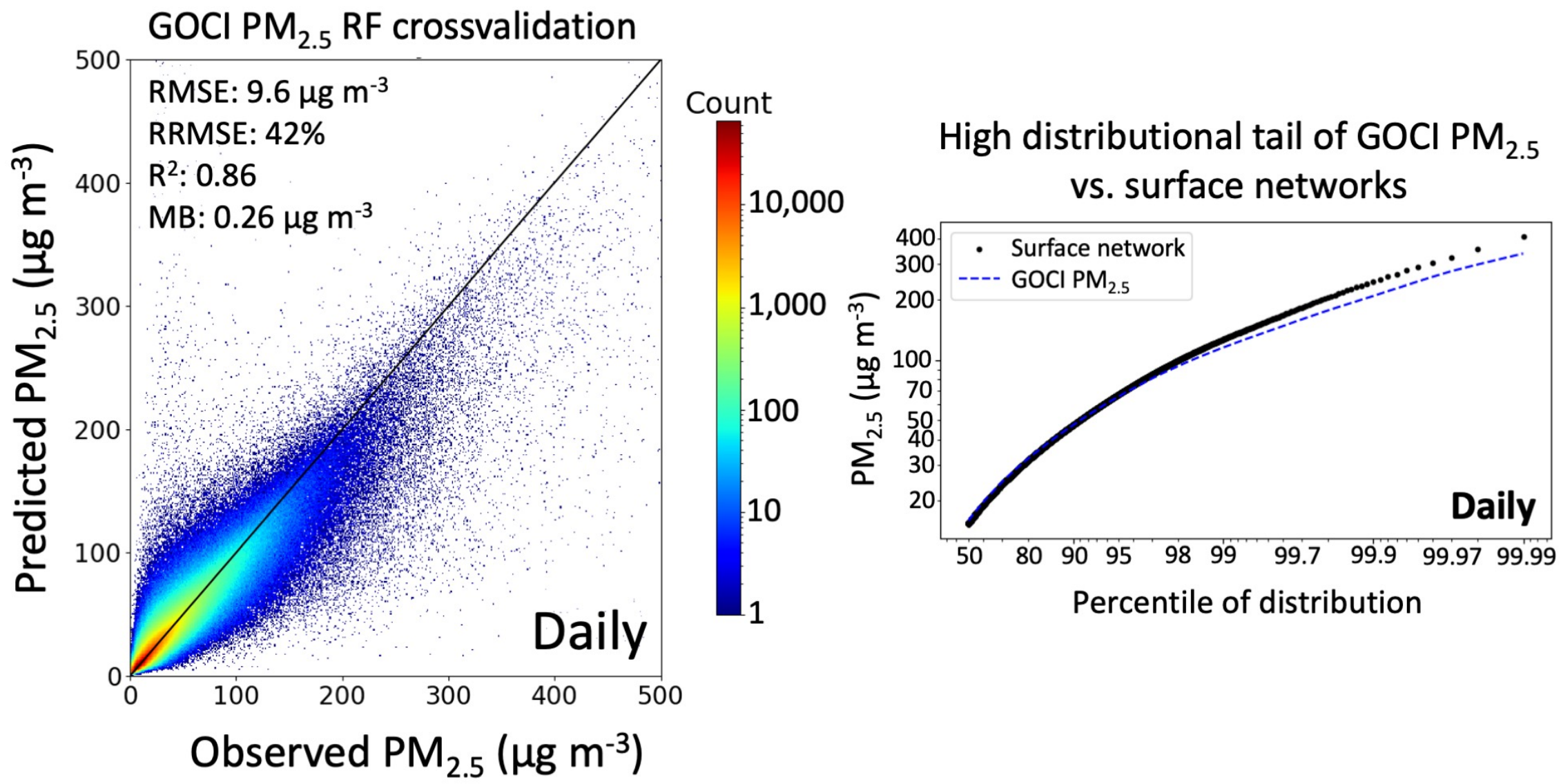
## GOCI PM<sub>2.5</sub> predictions

Over 2011-22, PM<sub>2.5</sub> networks increase dramatically in density; in China and Korea, records do not begin until 2014 and 2015 respectively. Our RF produces daily gap-filled coverage across this period of rapid change. The top row shows PM<sub>2.5</sub> as observed at surface network sites, with the inset showing PM<sub>2.5</sub> in Korea as supplemented by the AirKorea PM<sub>2.5</sub> RF detailed above. The bottom row shows annual mean GOCI PM<sub>2.5</sub> concentrations as predicted from our RF trained on data listed at left.



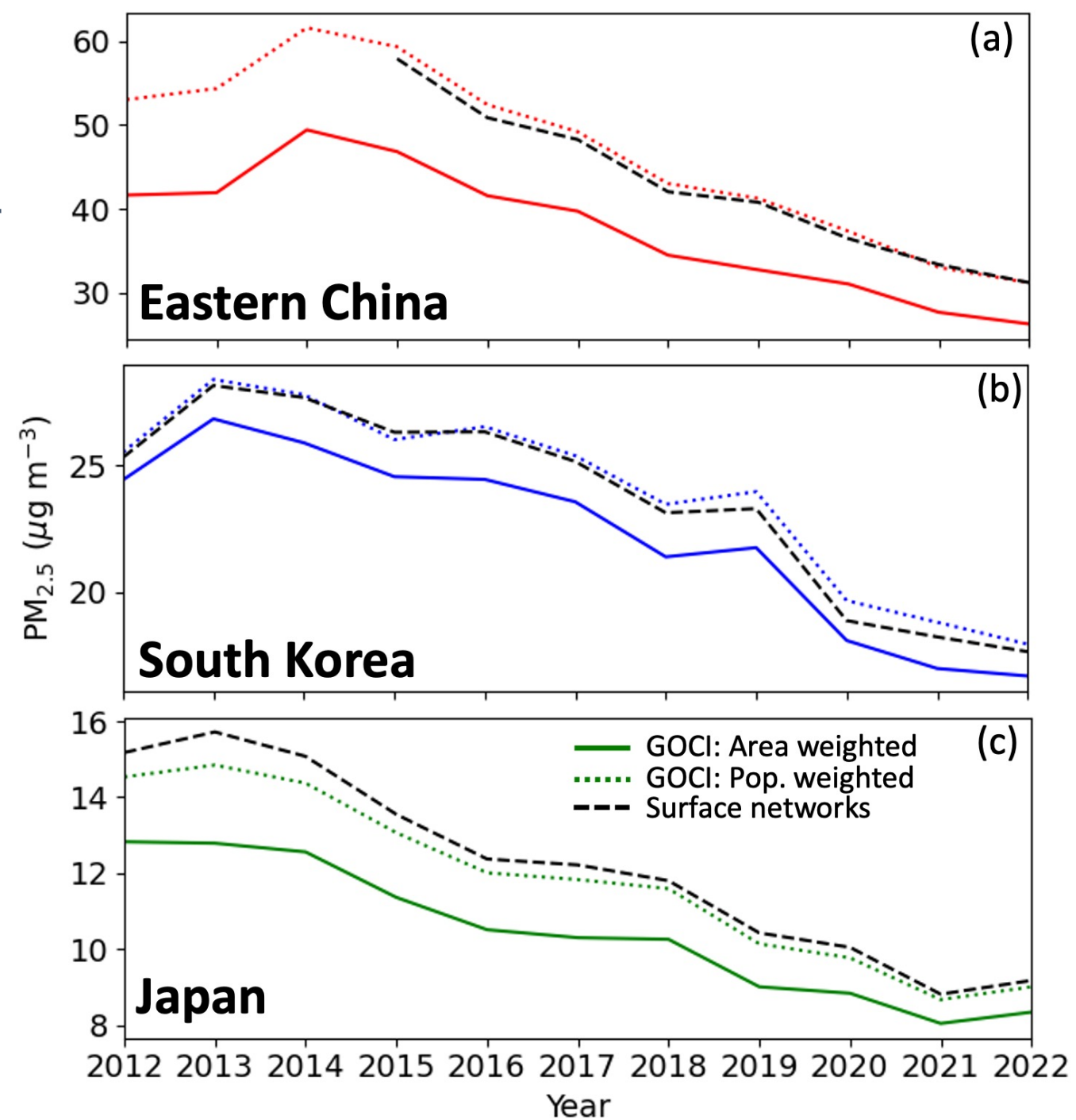
## GOCI PM<sub>2.5</sub> performance

The left panel below shows the crossvalidation performance of GOCI PM<sub>2.5</sub> against surface network observations. Root mean square error (RMSE) between observed and predicted 24-h PM<sub>2.5</sub> is 9.6 µg m<sup>-3</sup> (annual, 3.6 µg m<sup>-3</sup>) corresponding to a relative RMSE (RRMSE) of 42% (annual, 16%). The prediction captures 86% of the observed 24-h variance (R<sup>2</sup> = 0.86) and 95% of annual (R<sup>2</sup> = 0.95). Air quality managers are interested in high pollution events because they can create public health emergencies, but because ensemble learning algorithms like the RF involve averaging many simpler predictors they tend to struggle to predict extremes. However, due to the very large training set ingested in the training process, we find that, although tail bias does exist and worsen in the extreme high tail (>99.9 percentile, right panel), the RF captures a broad distributional range.



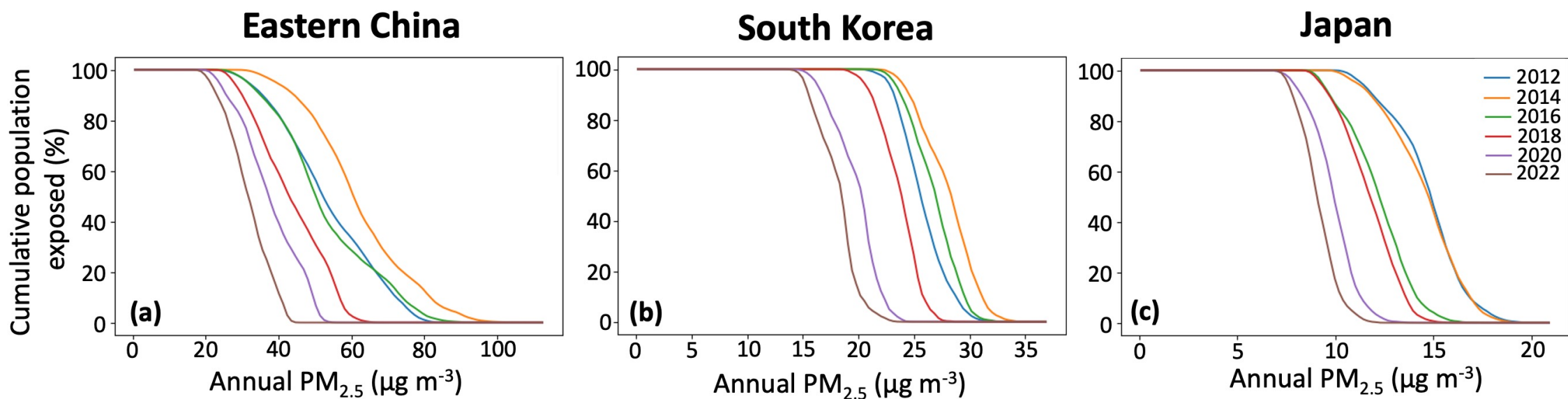
## National trends in PM<sub>2.5</sub>

The PM<sub>2.5</sub> networks show decreasing trends in all three countries and these trends are consistent with the GOCI PM<sub>2.5</sub> for both areal and population-weighted means. We find the decrease in China is consistent through 2022; indeed, economic shutdowns due to COVID-19 are difficult to discern at the annual temporal scale, possibly because emissions decreases are offset by an increase in oxidants producing secondary aerosol. Population-weighted means are more reflective of surface networks and are systematically higher than areal averages, reflecting the tendency of networks to sample urban conditions.



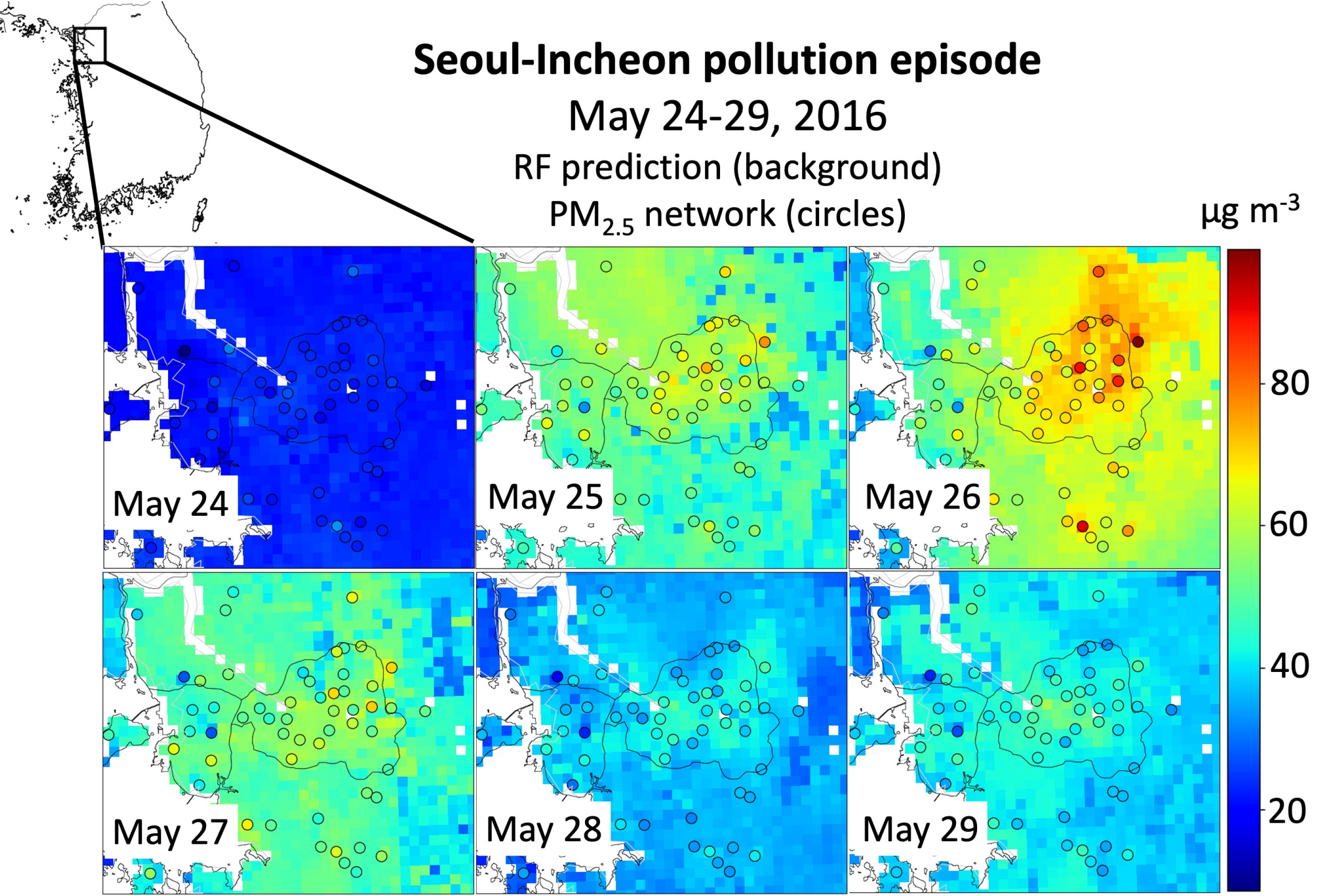
## Fall in population exposure everywhere

Below, we show trends in cumulative population exposure disaggregated by population. Panels show on the y axis the cumulative populations exposed to at least the annual PM<sub>2.5</sub> level given on the x axis, with year indicated by color. In China and Japan, we find the greatest improvements in PM<sub>2.5</sub> are achieved by populations exposed to the highest pollution, leading to a narrowing spread of annual exposures, as shown by the sharpening slope of the cumulative distribution. Distributions across the domain shift left over time, with the maximum to which any population is exposed decreasing everywhere.



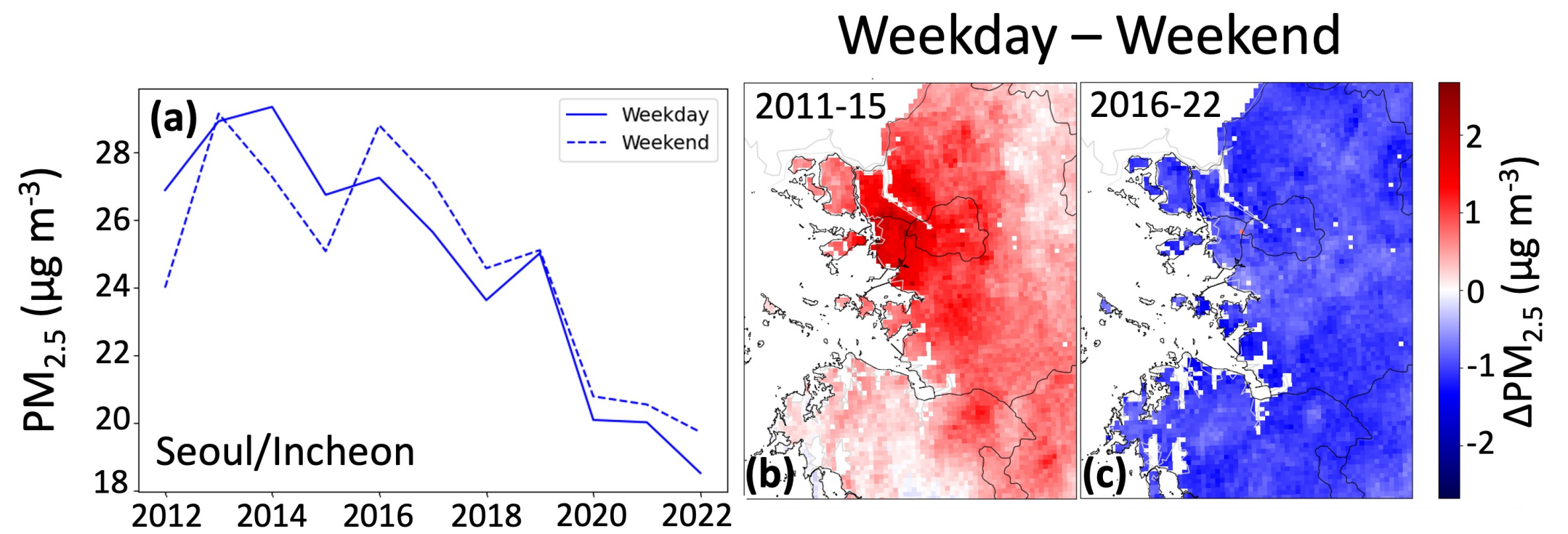
## An extreme urban PM<sub>2.5</sub> event

Here we analyze severe pollution event over Seoul on 24–29 May 2016 which corresponds with the one sampled during the KORUS-AQ field campaign. For this period, our 2×2 km<sup>2</sup> product better captures daily changes in neighborhood-level spatial gradients (R<sup>2</sup> of 0.97 with observations) as compared with our previous generation 6×6 km<sup>2</sup> product (R<sup>2</sup> of 0.77). Observational distributions are also better captured. While a two-sample Kolmogorov-Smirnov test suggests the 6×6 km<sup>2</sup> product has a statistically significantly different distribution than the observations over Seoul (p < 0.001), the 2×2 km<sup>2</sup> product is indistinguishable (p = 0.52).



## Reversing weekend effect

Previous work has shown a reversed weekend effect in Seoul and in some cities in China, where PM<sub>2.5</sub> levels are more elevated on weekends than on weekdays. However, the GOCI PM<sub>2.5</sub> product suggests this pattern was reversed in the period prior to surface network observations, with weekdays more polluted by weekends (region surrounding Seoul/Incheon shown on maps). A CTM analysis would be necessary to separate the role of emissions from synoptic meteorology, but controls on industrial emissions might have reduced the difference between weekday and weekend pollution levels.



## Acknowledgements

This work was funded by the Samsung PM<sub>2.5</sub> Strategic Research Program and the Harvard-NUIST Joint Laboratory for Air Quality and Climate (JLAQC). GOCI data was provided by Korea Institute of Ocean Science and Technology (KIOST). DCP was funded in part by a US National Science Foundation Graduate Fellowship.

## Contact information

Contact Drew Pendergrass at [pendergrass@g.harvard.edu](mailto:pendergrass@g.harvard.edu)

We plan to release the data in NetCDF form later this year for free use in public health, air quality, and related studies. To be notified on release, email the above address.