

# London School of Hygiene & Tropical Medicine

Improving Health Worldwide

LONDON  
SCHOOL of  
HYGIENE  
& TROPICAL  
MEDICINE



# Satellite-based machine learning models to estimate high-resolution environmental exposures across the UK

LONDON  
SCHOOL of  
HYGIENE  
& TROPICAL  
MEDICINE



**Rochelle Schneider**

Research Fellow in Geospatial Data Science

PhD Geospatial Analytics | MSc GIS | MRes Remote Sensing

Visiting Researcher at ECMWF

Centre on Climate Change and Planetary Health, LSHTM

**Antonio Gasparrini**

Professor of Biostatistics and Epidemiology

PhD Medical Statistics | Postgraduate Medical Statistics and Biometry | MSc Biostatistics

Centre for Statistical Methodology, LSHTM

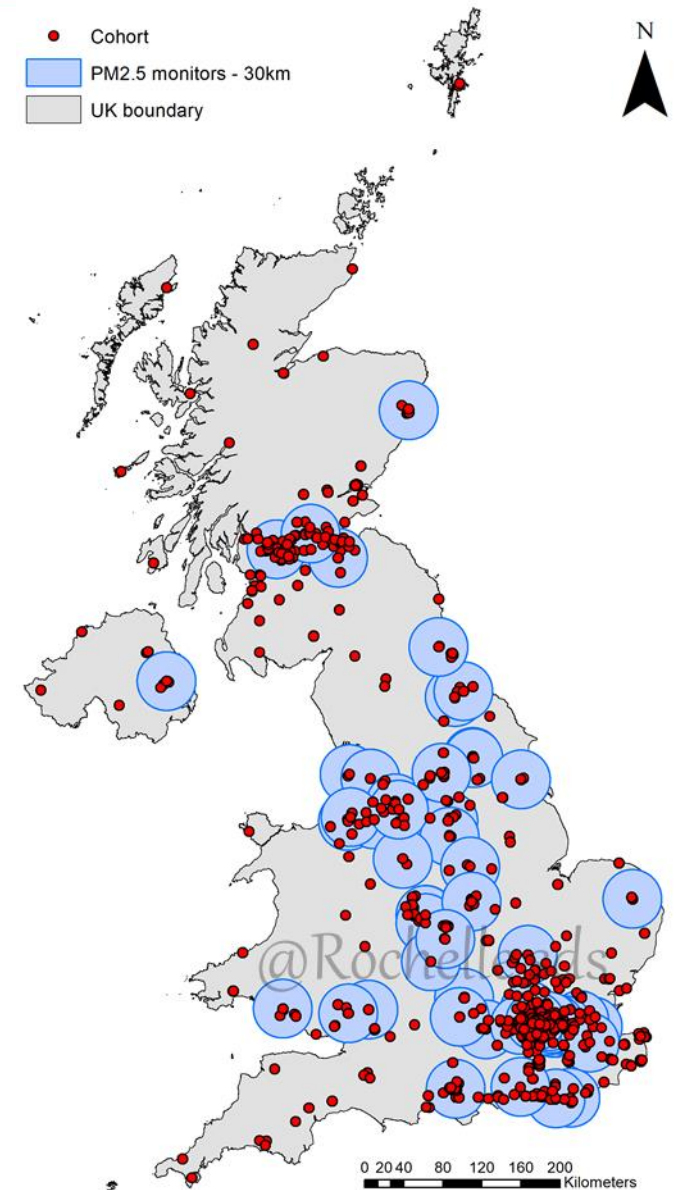
Centre on Climate Change and Planetary Health, LSHTM

## Specific aspects:

- **Widespread exposure** to environmental factors - e.g. air pollution - often affecting the whole population;
- Often **small risks** - e.g., a RR of 1.0051 (95%CI: 1.0007-1.0093) for an increase of  $10\mu\text{gr}/\text{m}^3$  of  $\text{PM}_{10}$  (Samet NEJM 2000);
- Need to perform epidemiological analyses on **large populations**.

## Traditional Limitations

- Outcomes/exposures **with low temporal and/or spatial resolution**, roughly aggregated, and lack of individual / small-area information;
- **Partial coverage**, especially in rural areas and in low/middle income countries, posing limitations in country-wide or global analyses.
- Linkage between various **temporal and spatial scales**





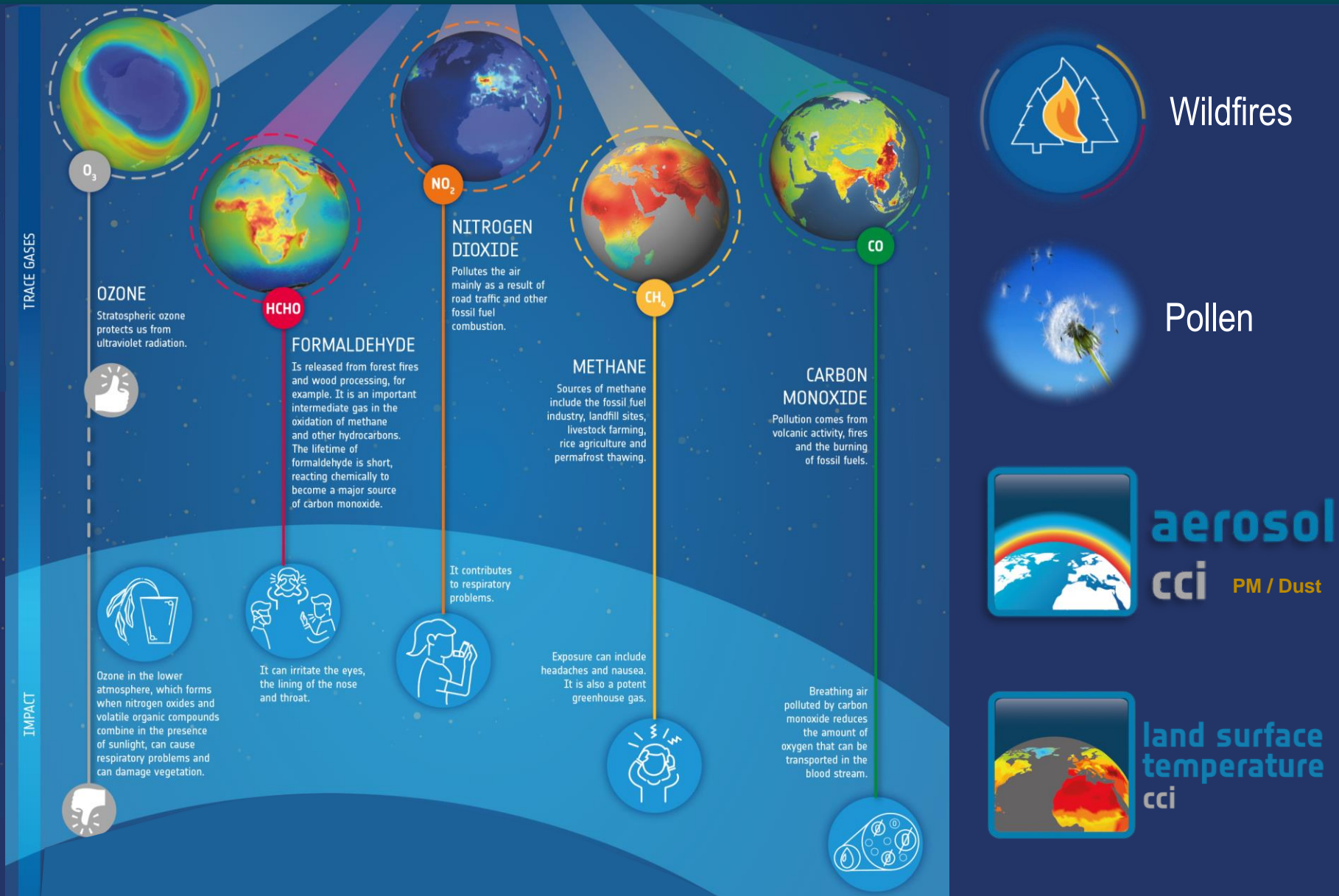
# Novel data sources for environmental epidemiological studies

## → THE AIR WE BREATHE

Air pollution is a major environmental health problem that affects millions of people around the world. Satellite data and computer models can show how pollution accumulates and how it is carried in the air. Mapping the global atmosphere every day, Sentinel-5P provides high-resolution data on a multitude of trace gases and information on aerosols that affect air quality and the climate. Offering advances in coverage and resolution, Sentinel-5P is set to take air-quality monitoring to a new level.



Sentinel-5P carries **TROPOMI**™ the most advanced multispectral imaging spectrometer to date



# Project ST- UK

(Spatio-temporal UK exposure modelling)

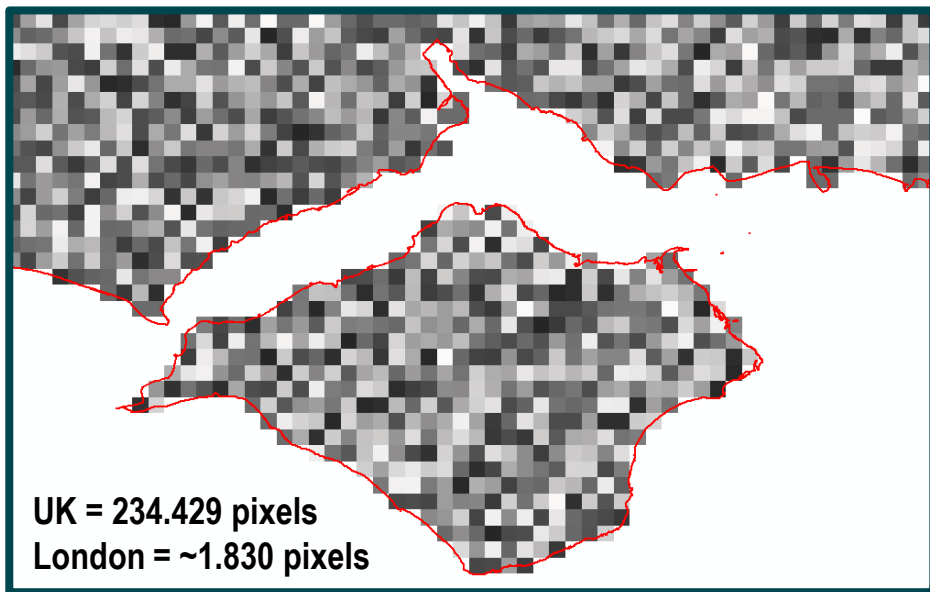
Estimation of daily  $PM_{10}$  and  $PM_{2.5}$  concentrations using  
satellite-based machine learning models

LONDON  
SCHOOL of  
HYGIENE  
& TROPICAL  
MEDICINE



# CASE STUDY

Case Study: Great Britain  
Time series: 2003-2018  
Temporal resolution: daily  
Spatial resolution: 1 x 1km<sup>2</sup>  
Variable estimated: PM<sub>10</sub> + PM<sub>2.5</sub>

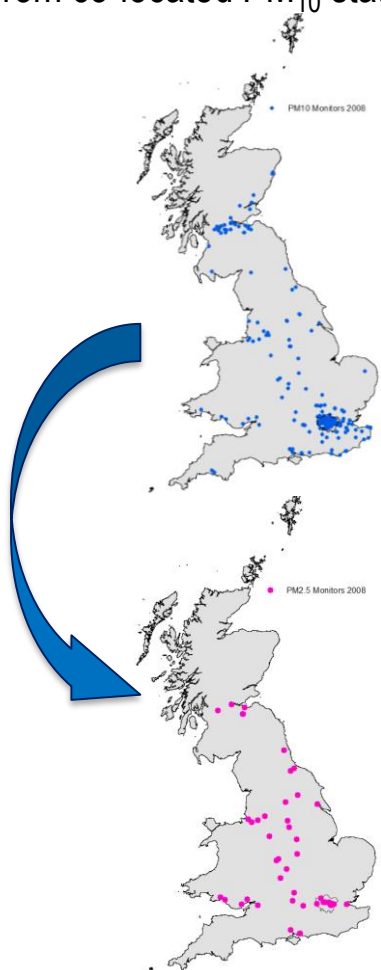




# Multi-stage machine learning spatio-temporal model for PM<sub>2.5</sub>

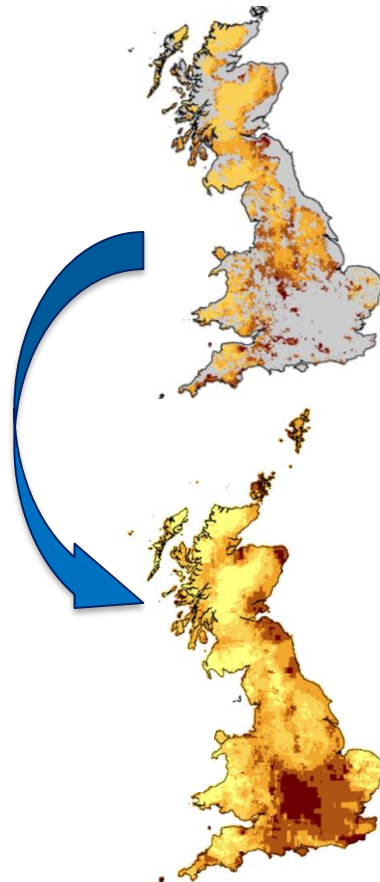
## Stage 1

Predict PM<sub>2.5</sub> concentrations from co-located PM<sub>10</sub> stations.



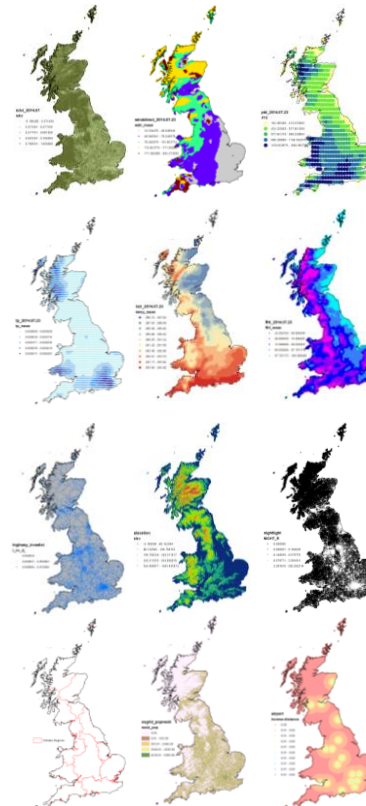
## Stage 2

Predict satellite-aerosol from Copernicus atmospheric reanalysis aerosol.



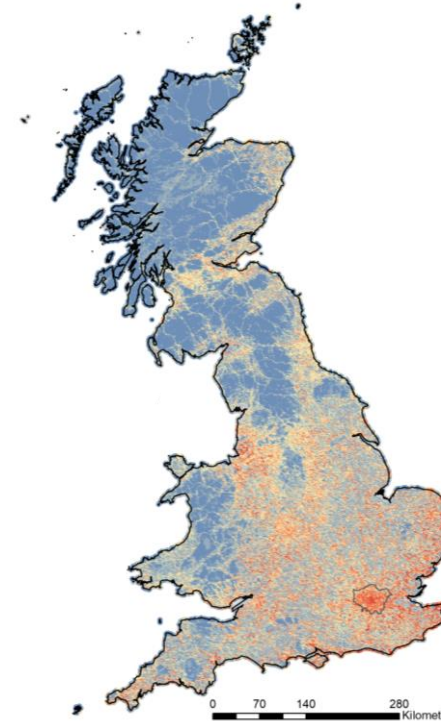
## Stage 3

Test different satellite-based machine learning models to predict PM<sub>2.5</sub> concentrations from the predictors.



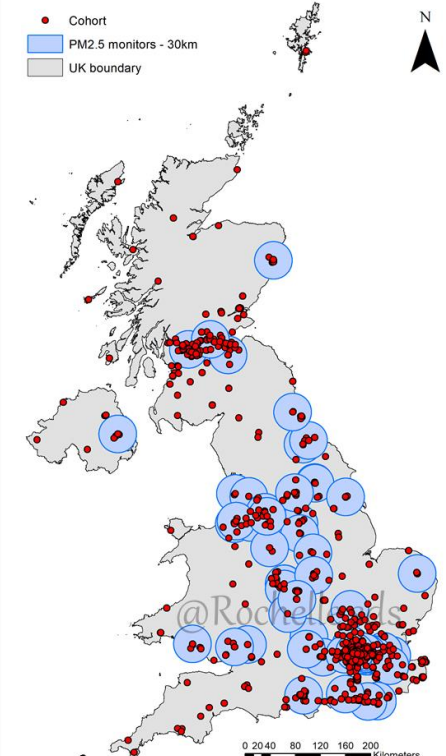
## Stage 4

Predict daily PM<sub>2.5</sub> concentrations at 1 km<sup>2</sup> using the parsimonious satellite-based machine learning model.



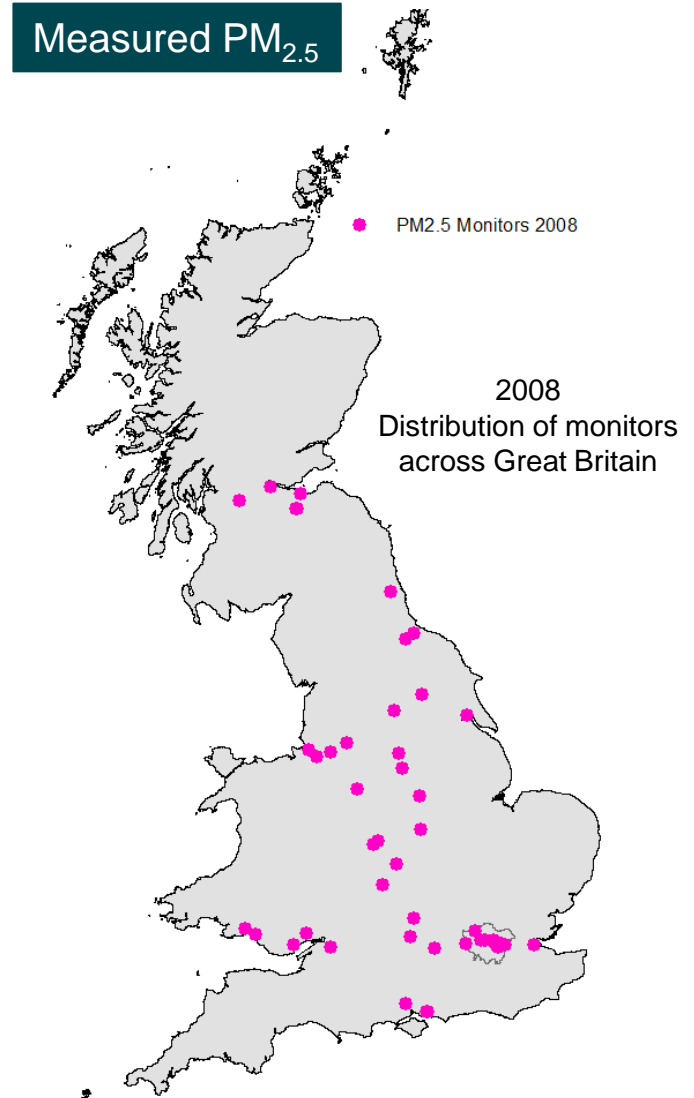
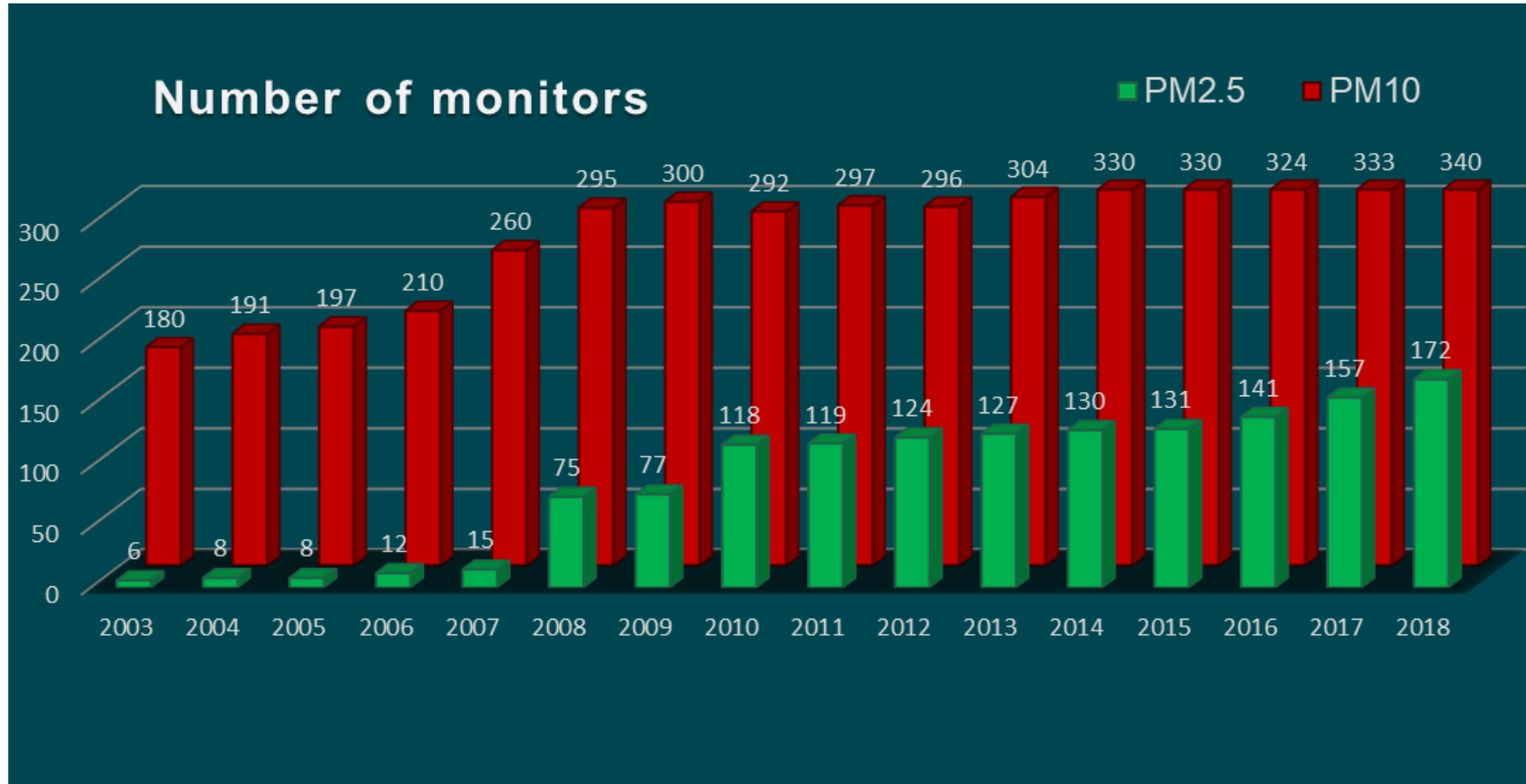
## Stage 5

Predict PM<sub>2.5</sub> concentrations at the monitor level using small-scale predictors



# STAGE 1: Gap-filling Model

Predict daily  $PM_{2.5}$  concentrations from co-located  $PM_{10}$  stations

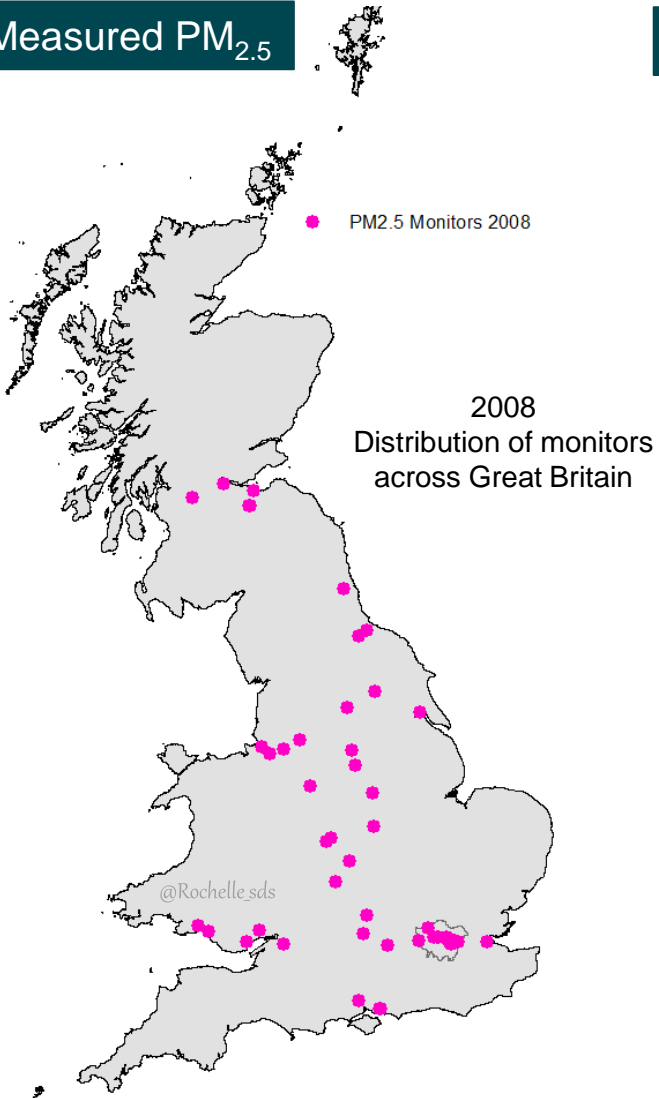




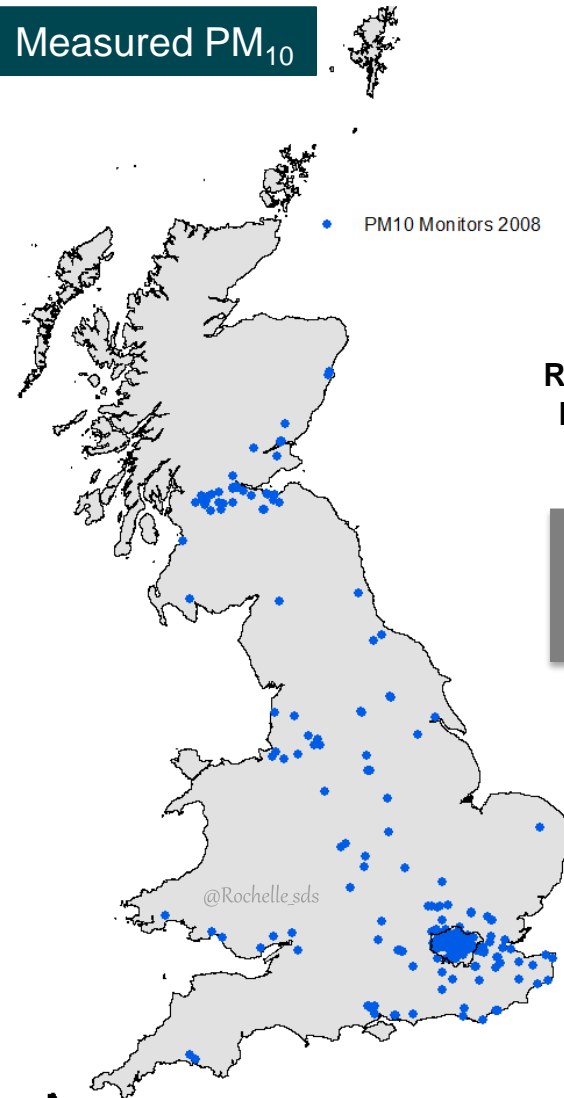
# STAGE 1: Gap-filling Model

Predict daily  $PM_{2.5}$  concentrations from co-located  $PM_{10}$  stations

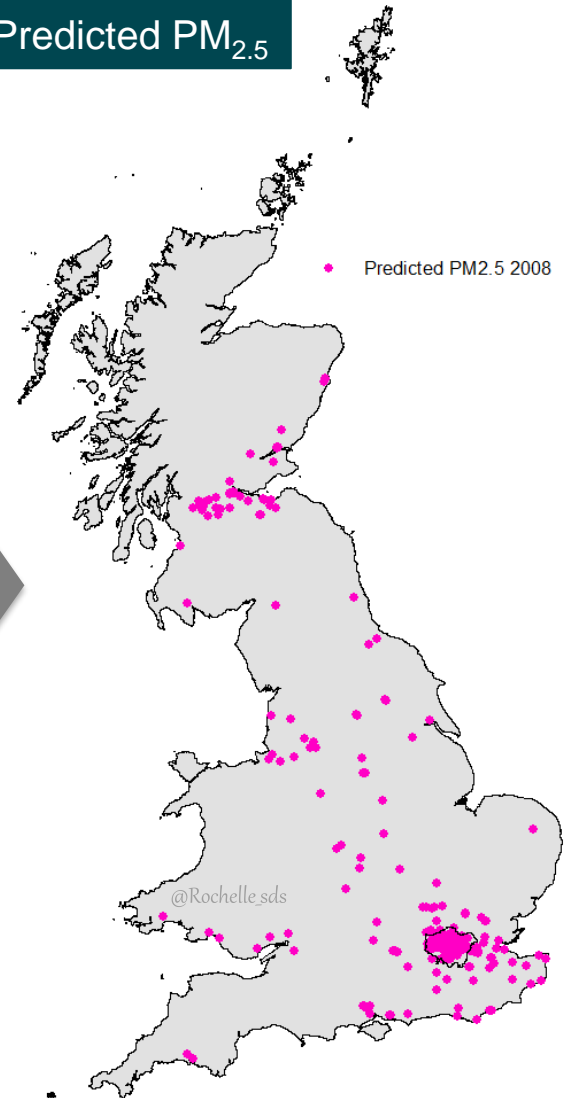
Measured  $PM_{2.5}$



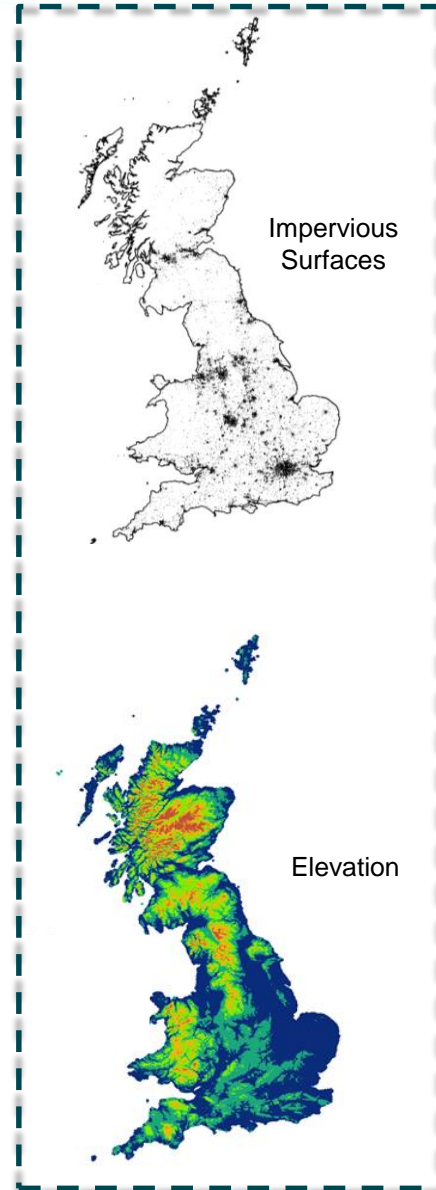
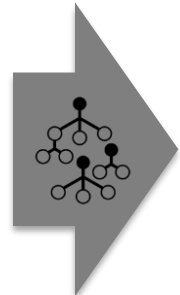
Measured  $PM_{10}$



Predicted  $PM_{2.5}$



Random Forest



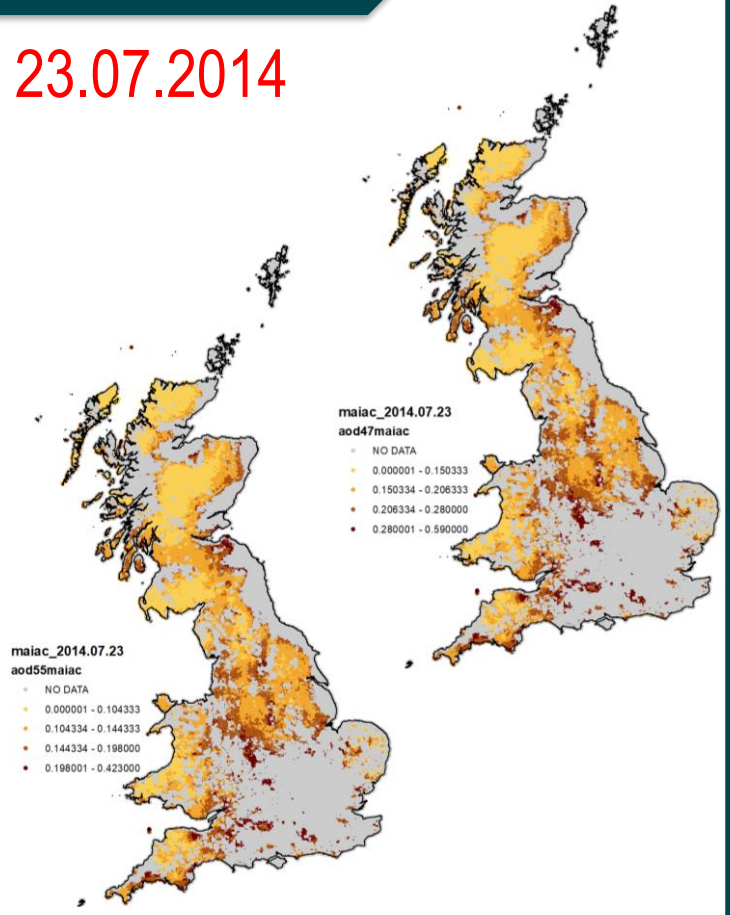
# STAGE 2: AOD GAP-FILLING MODEL

Predict AOD ( $2 \lambda$  / day) from MAIAC-MODIS using AOD ( $5\lambda, 7t$ ) from CAMS

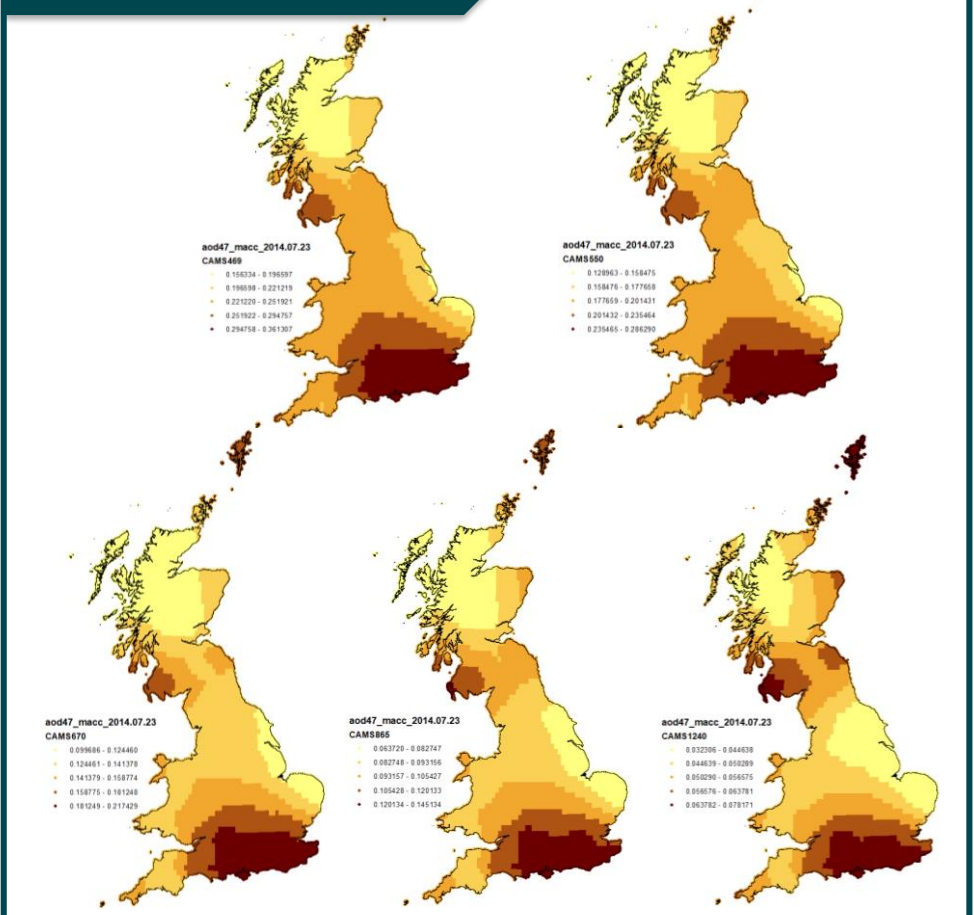


AOD MAIAC – 1km

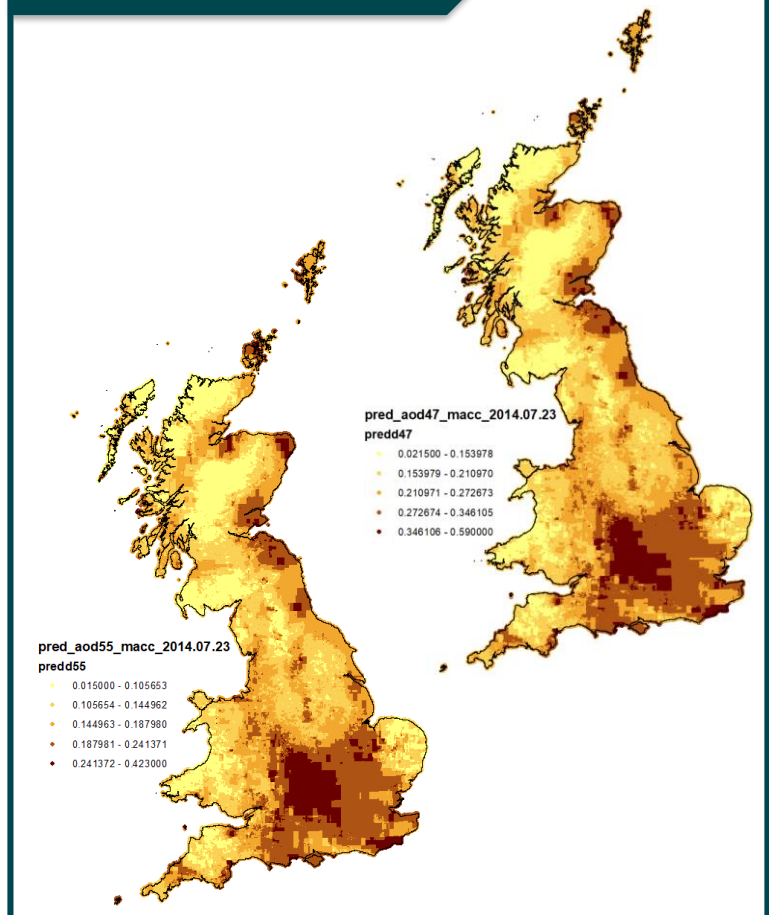
23.07.2014



AOD REANALYSIS – “10km”



PREDICTED AOD – 1km





# STAGE 3: Test ≠ Satellite-based ML Models

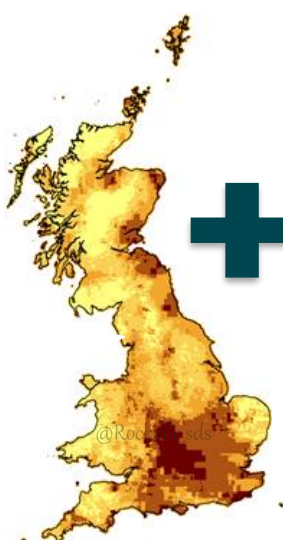
Predict  $PM_{2.5}$  at 1 km<sup>2</sup> using the parsimonious satellite-based ML model



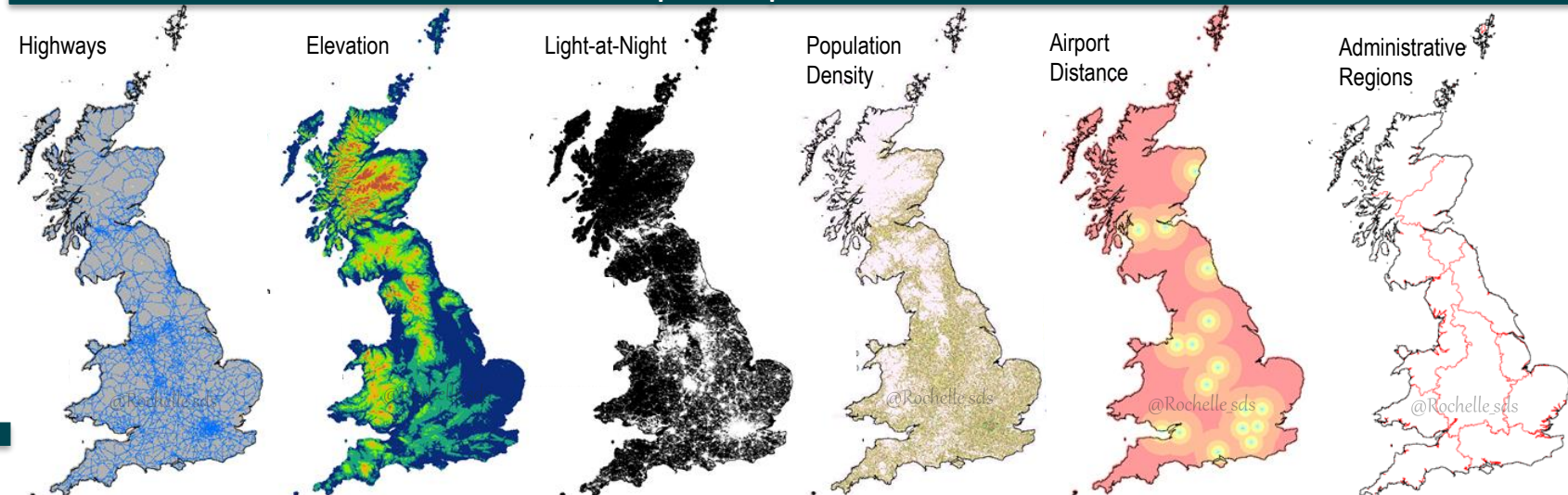
STAGE 1



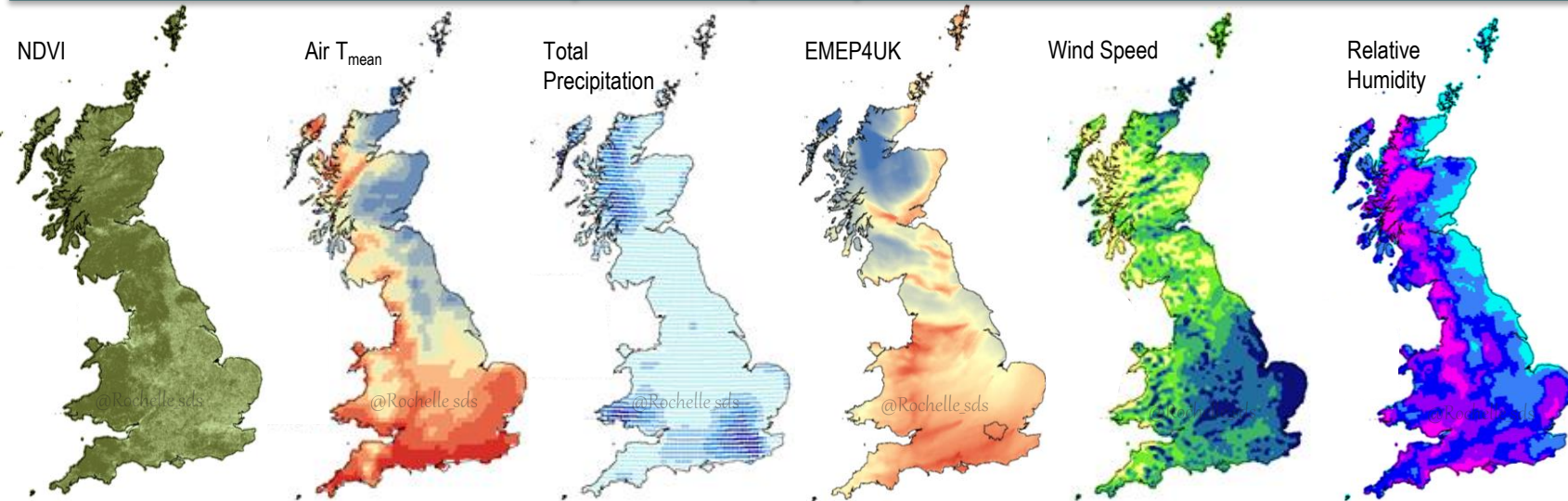
STAGE 2



## Spatial predictors

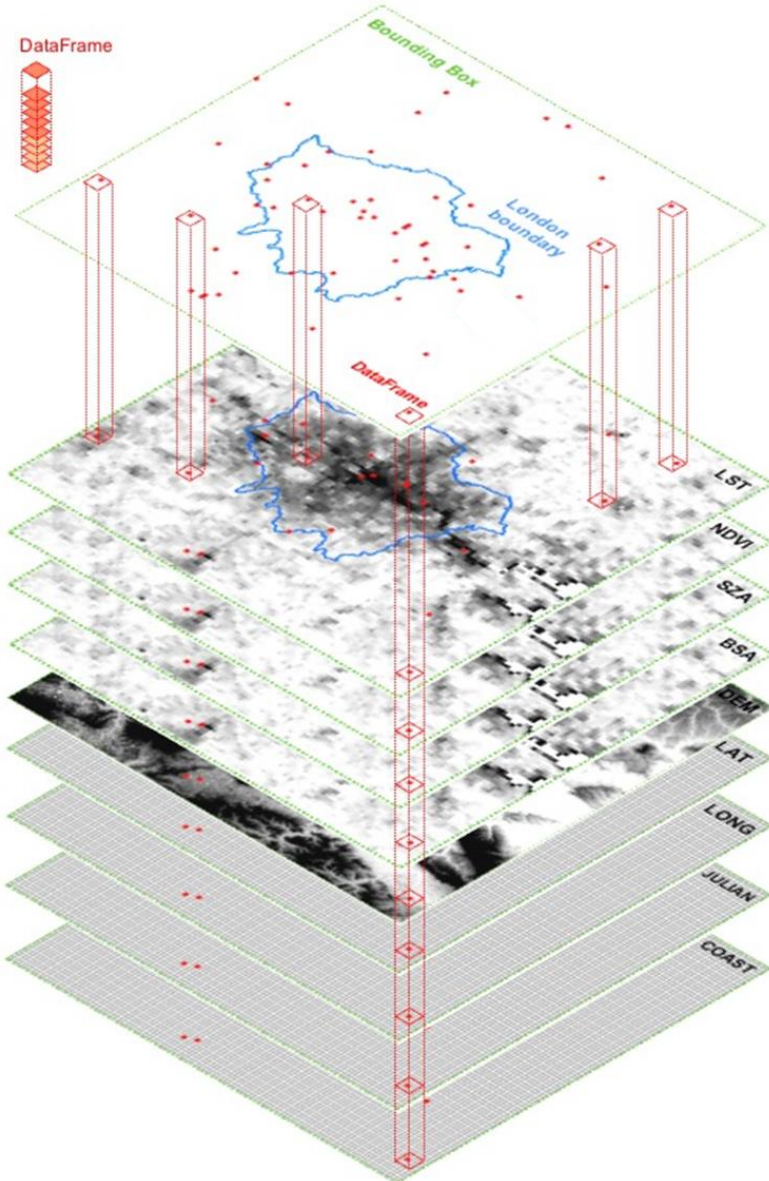


## Spatio-temporal predictors



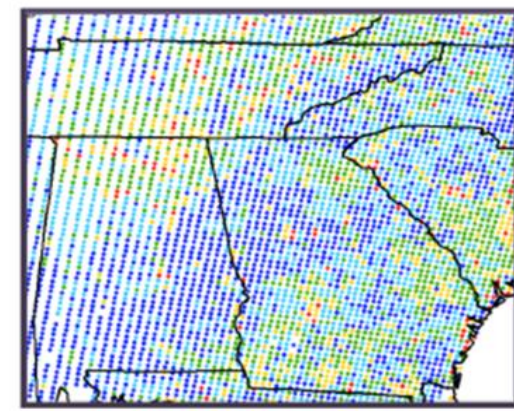



# Data Synchronisation + Gap-filling Model + Advanced Statistical Modelling

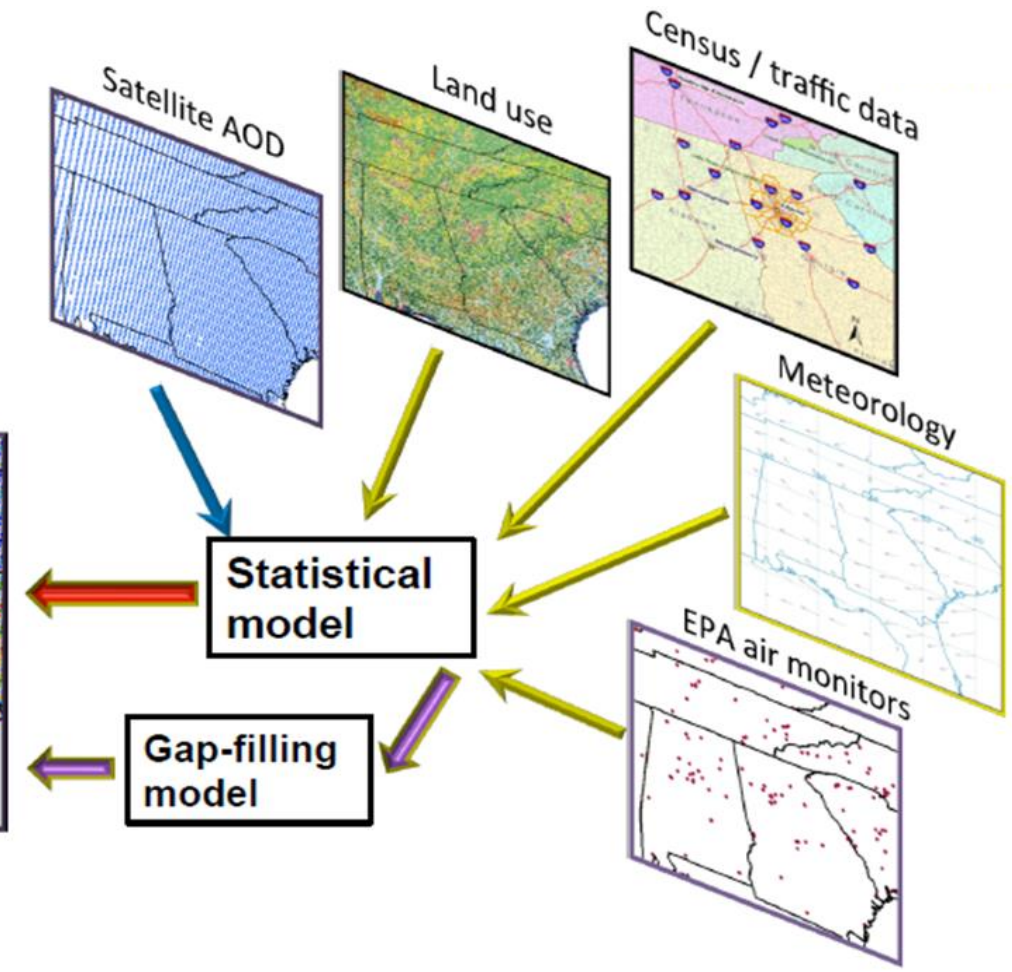


Particulate Matter Air Quality from Space – Advanced Statistical Modeling

Yang Liu, PhD



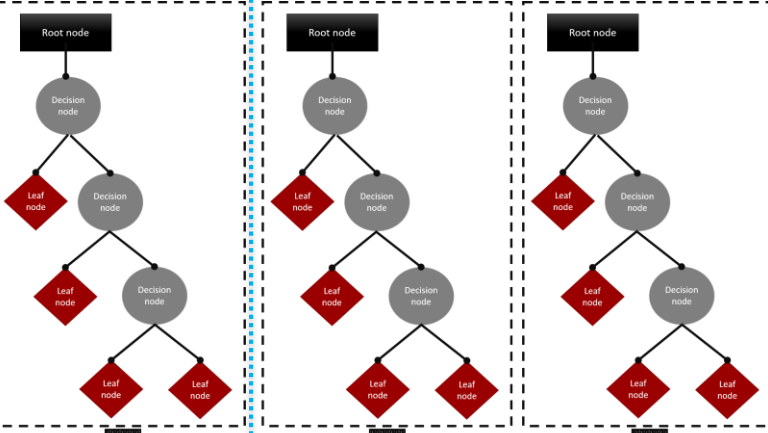
Predicted PM<sub>2.5</sub> surface





## Random Forest

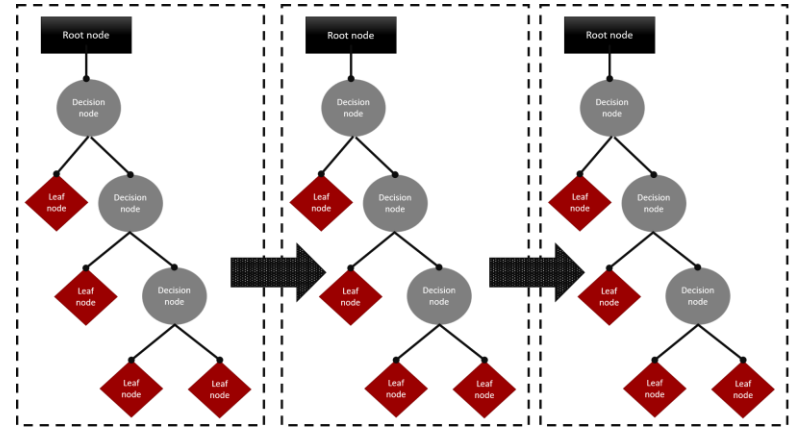
### Decision Tree



TREE 1 =  $R^2$ , MAE, RMSE      TREE 2 =  $R^2$ , MAE, RMSE      TREE N =  $R^2$ , MAE, RMSE

$$R^2_{mean} = \frac{1}{N} \sum_{i=1}^N R^2_i \quad MAE_{mean} = \frac{1}{N} \sum_{i=1}^N MAE_i \quad RMSE_{mean} = \frac{1}{N} \sum_{i=1}^N RMSE_i$$

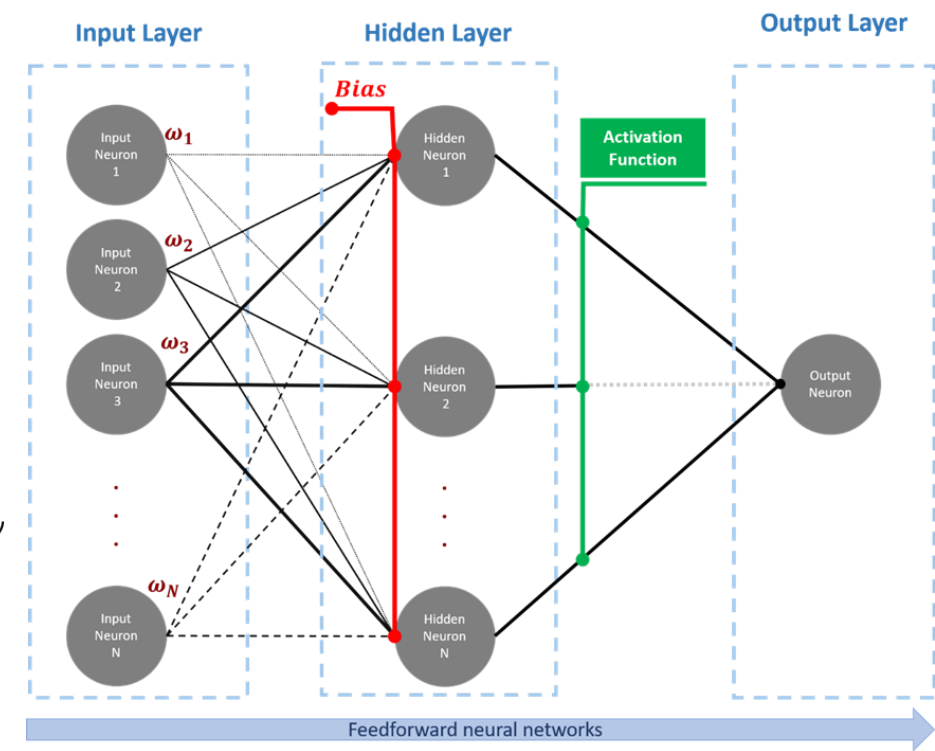
## Gradient Boosting



TREE 1 =  $R^2$ , MAE, RMSE      TREE 2 =  $R^2$ , MAE, RMSE      TREE N =  $R^2$ , MAE, RMSE

$$R^2_{mean} = \frac{1}{N} \sum_{i=1}^N w R^2_i \quad MAE_{mean} = \frac{1}{N} \sum_{i=1}^N w MAE_i \quad RMSE_{mean} = \frac{1}{N} \sum_{i=1}^N w RMSE_i$$

## Neural Network



Feedforward neural networks

# Results: STAGE 1 - Gap-filling Model

```
#####
### OPTIMISATION USING SPATIO-TEMPORAL CROSS-VALIDATION ###
#####

for (ntree in c(100, 500, 1000)){
  for (mtry in c(3,5,7)){
    cat("\n", "ntree=", ntree, "| mtry=", mtry, "\n")

    final = foreach(i=1:10, .packages=c("data.table", "ranger"))%dopar% {

      mod <- ranger(pm25 ~ pm10 + julian +
                    latitude + longitude +
                    SiteTypeCode + year +
                    month + dow,
                    data=subset(stg1, split!=i),
                    num.trees=ntree,
                    mtry=mtry,
                    respect.unordered.factors=TRUE,
                    verbose=TRUE )

      test <- subset(stg1, split==i)
      test$pred.rf <- predict(mod, test)$predictions
      test$itercv <- i
      test[, list(date, code.source, pm25, pred.rf, itercv)]

    }

    mod.cv = do.call(rbind, final)

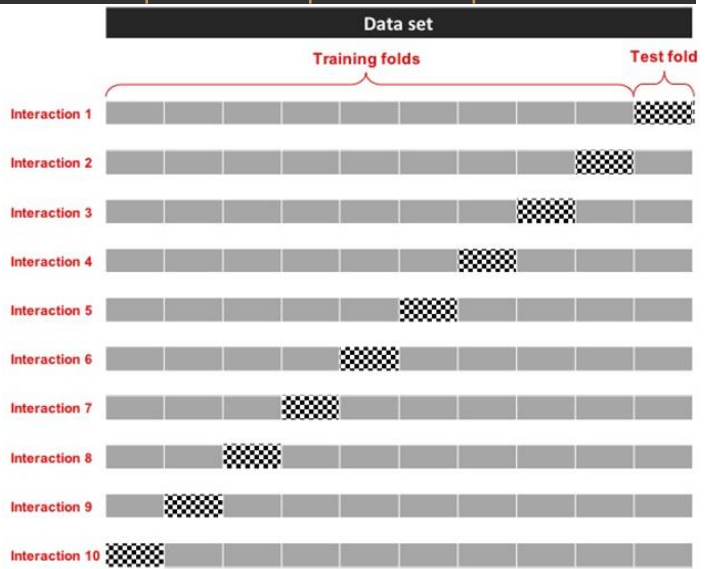
#####
### REGRESS PREDICTED AGAINST OBSERVED ###
#####

linear.pm25 <- lm(pm25~pred.rf, data=mod.cv)
r2.all <- summary(linear.pm25)$r.squared
```

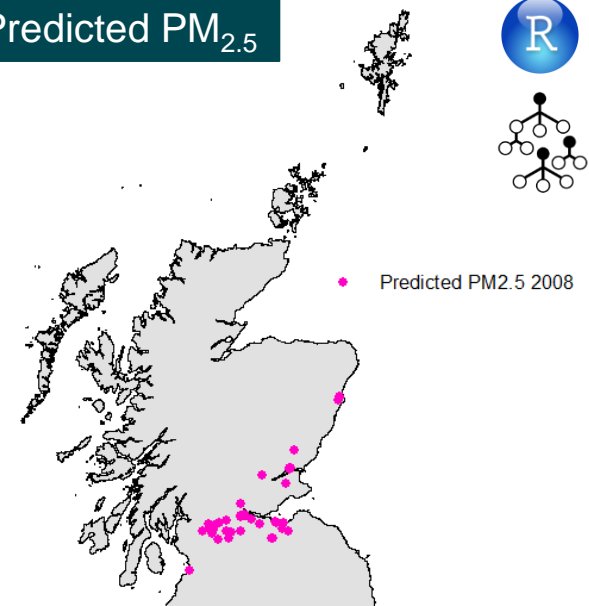
variable.importance

pm10	15351814.3
longitude	1158180.8
latitude	1066338.3
julian	875359.1
year	868112.3
month	505182.1
SiteTypeCode	292421.0
dow	266770.8

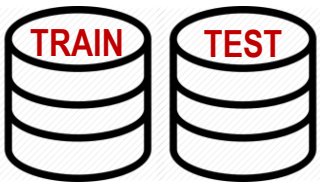
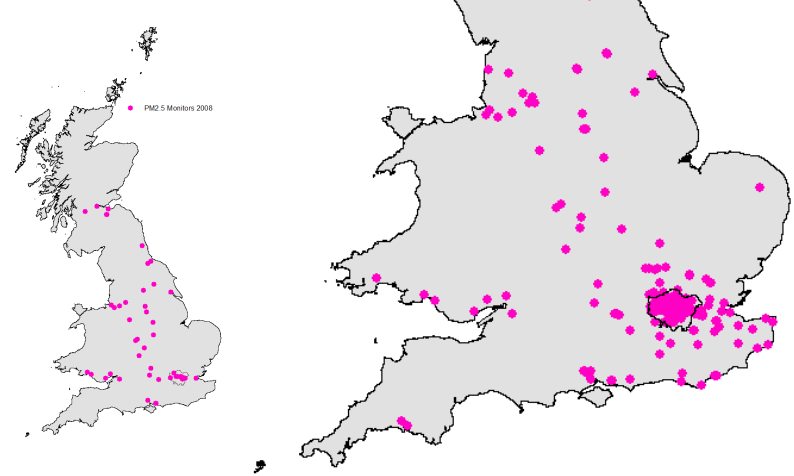
$R^2$ all	ntree=100	ntree=500	ntree=1000
mtry=3	0.860	0.861	0.861
mtry=5	0.849	0.850	0.850
mtry=7	0.842	0.843	0.843



## Predicted PM<sub>2.5</sub>



## Measured PM<sub>2.5</sub>



**Predictors:** 8 variables  
**Train/Test data:** Air Quality monitors  
**Total Size:** ~ 260.000  
**Model Design:** 1 model for 16 years

# Results: STAGE 2 – AOD Gap-filling Model

```
#####
# ||||| CROSS-VALIDATION ||||| #
#####

ranger.model <- caret::train (formula,
                              data      = stage2.train,
                              method    = "ranger",
                              metric     = "RMSE",
                              num.trees = 50,
                              tuneGrid  = expand.grid(.mtry = 20),
                              trControl = trainControl(
                                  method="cv",
                                  number=10,
                                  allowParallel=T,
                                  verboseIter=T)
                              )

#####
### REGRESS THE PREDICTED AGAINST OBSERVED ###
#####

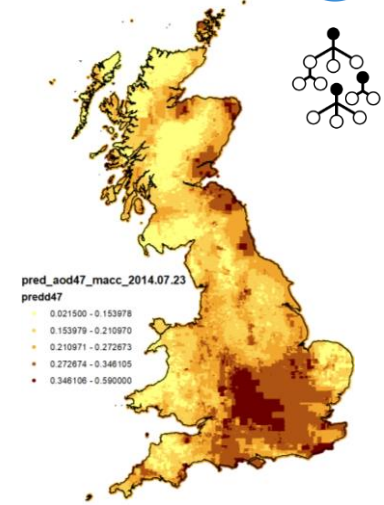
stage2.valid$pred <- predict(ranger.model, stage2.valid)
linear.model     <- lm(formula.linear, data=stage2.valid)
r2.stage2       <- summary(linear.model)$r.squared
```

```
#####
### RESULTS ###
#####

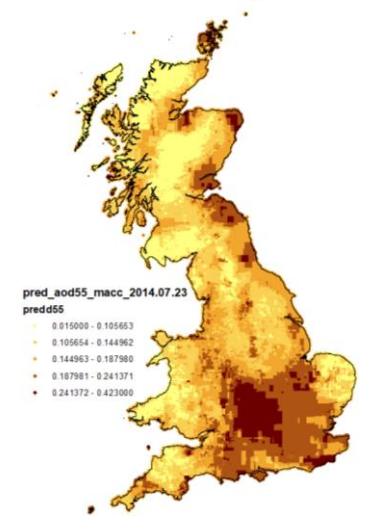
variable.importance
Julian          9926.3528
CAMS_469_12    4722.1932
CAMS_550_12    4340.4057
y              2176.8599
x              1818.0451
CAMS_AOD670_12 874.2431
CAMS_AOD865_12 572.2148
CAMS_AOD469_15 531.0689
CAMS_AOD1240_12 513.8520
CAMS_AOD1240_18 479.6304
```

	R <sup>2</sup> .stage 2	AOD 470	AOD 550
2003	0.960	0.960	0.960
2004	0.926	0.926	0.926
2005	0.935	0.934	0.934
2006	0.955	0.955	0.955
2007	0.960	0.961	0.961
2008	0.932	0.931	0.931
2009	0.935	0.935	0.935
2010	0.919	0.918	0.918
2011	0.957	0.970	0.970
2012	0.939	0.939	0.939
2013	0.942	0.942	0.942
2014	0.921	0.920	0.920
2015	0.914	0.914	0.914
2016	0.923	0.923	0.923
2017	0.911	0.910	0.910
2018	0.921	0.921	0.921

## Predicted AOD 470

## Predicted AOD 550



- Optimised RF : ntree=50 | mtry=20



**Predictors:** 43 variables  
**Train/Test data:** Satellite 1km<sup>2</sup> grid  
**Total Size:** ~ 7 Million by year and AOD type (112 Million for 16 years)  
**Model Design:** 1 model by year and AOD type (i.e. 16 models for AOD 470 and 16 for AOD550)



# Results: STAGE 3 – Test ≠ Satellite-based ML Models

```
#####
## OPTIMISATION USING SPATIO-TEMPORAL CROSS-VALIDATION ##
#####

for (ntree in c(100, 500,1000)){
  for (mtry in c(7,70,20)){
    ## DEFAULT: sqrt(42) = 7

    cat("\n", "ntree=", ntree, "| mtry=", mtry, "\n")

    final=foreach (i=1:10, .packages=c("data.table","ranger")) %dopar% {
      print(i)

      mod <- ranger(formula.pm25,
                    data=subset(stg3, split!=i),
                    num.trees=ntree,
                    mtry=mtry,
                    respect.unordered.factors=TRUE,
                    verbose=TRUE)

      test <- subset(stg3, split=i)
      test$logpred.rf <- predict(mod, test)$predictions
      test$exppred.rf <- exp(test$logpred.rf)
      test$itercv <- i
      test[, list(day, osgrid.id, code.source, pm25pred.meas,
                  exppred.rf, logpm25, logpred.rf, itercv)]
    }

    stage3.cv = do.call(rbind, final)

    #####
    ### REGRESS THE PREDICTED AGAINST OBSERVED ###
    #####

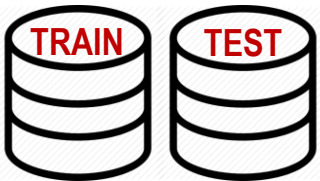
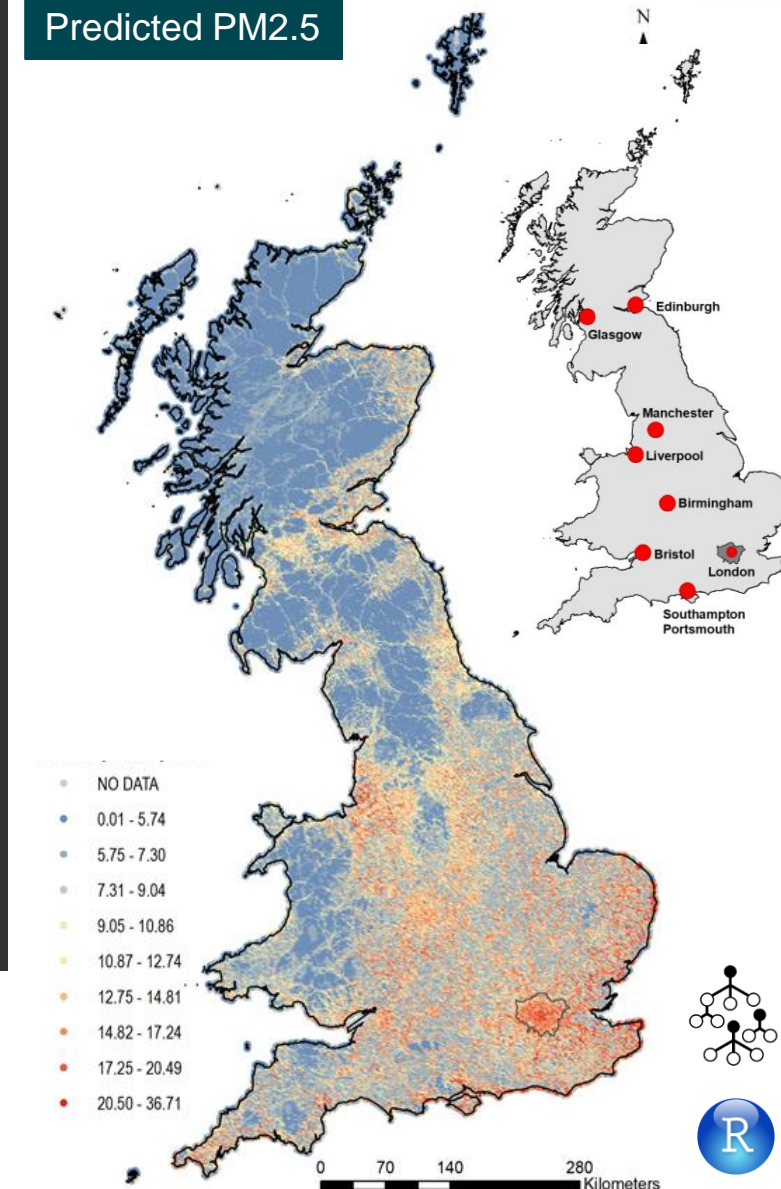
    linear.all <- lm(pm25pred.meas~exppred.rf, data=stage3.cv)
    r2.all <- summary(linear.all)$r.squared
  }
}
```

## ##### ### RESULTS ### #####

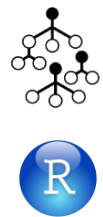
### Variable Importance

Wind Direction	3774.673120
Distance from Sea	2921.701911
Wind Speed	1803.228131
Day of year	1780.716491
Mean Precipitation	1725.783224
Month(factor)	1700.169657
Planetary Boundary Layer_12h	1471.014341
Mean Temperature	1408.845640
Planetary Boundary Layer_0h	1291.013732
Mean Sea Pressure	1103.636722
Elevation	873.211102
Stage2-AOD550	833.542341
Stage2-AOD470	822.952578

	$R^2.all$	ntree=100	ntree=500	ntree=1000
mtry=7	0.670	0.670	0.675	0.675
mtry=10	0.687	0.687	0.691	0.691
mtry=20	0.700	0.700	0.702	0.702



**Predictors:** 42 variables  
**Train/Test data:** Air Quality monitors  
**Total Size:** ~ 100.000 by year  
**Model Design:** 1 model by year





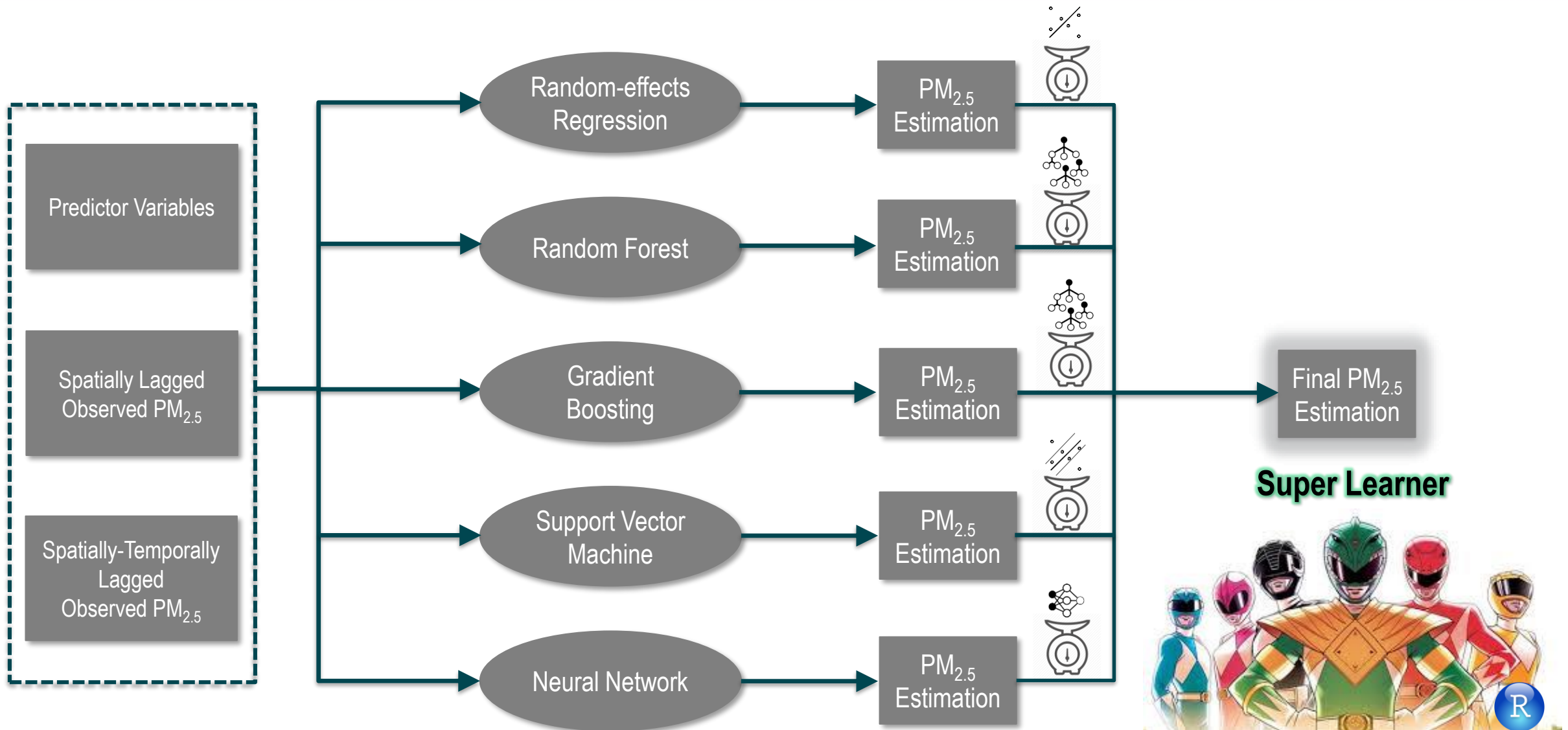


# Next Steps

LONDON  
SCHOOL of  
HYGIENE  
& TROPICAL  
MEDICINE



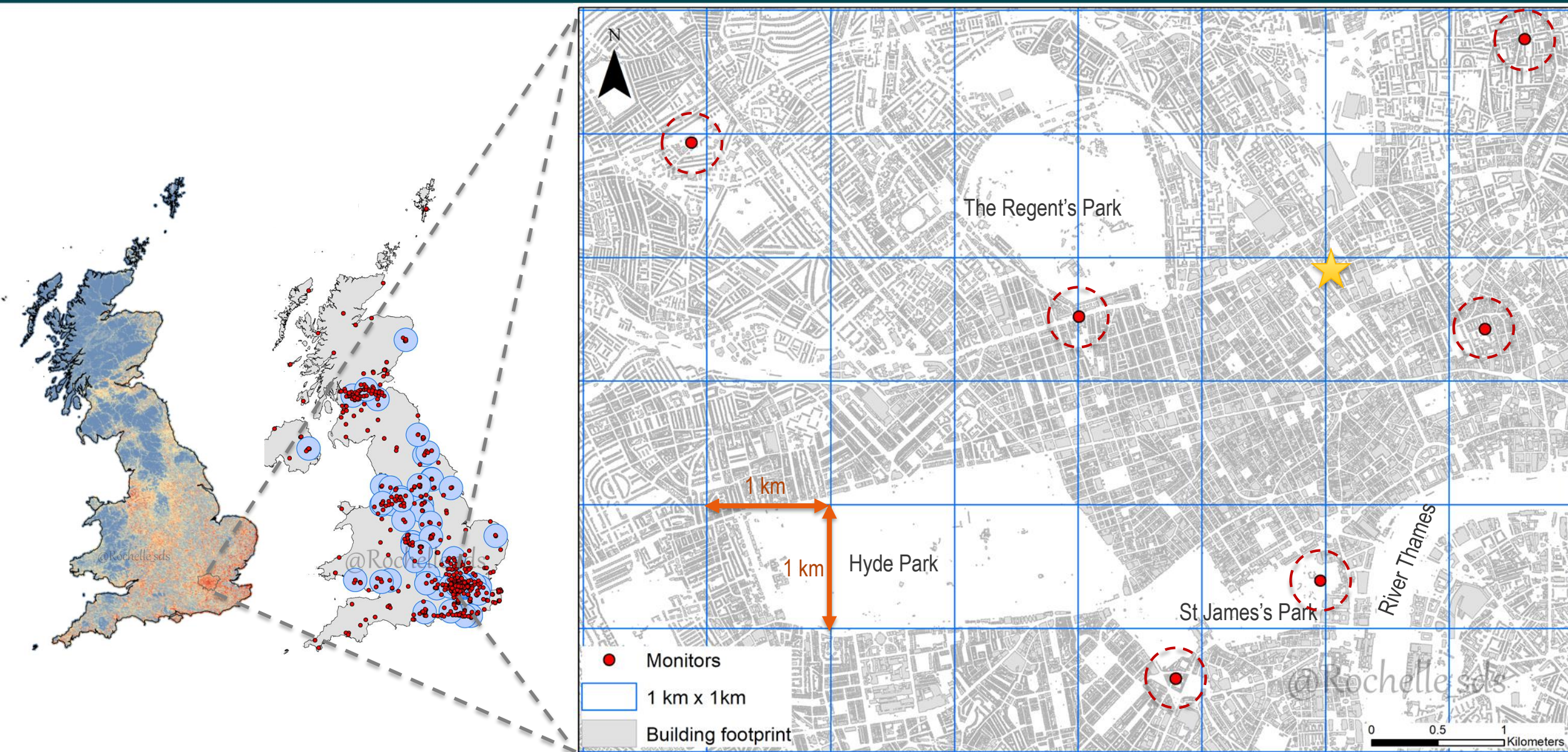
# Explore the Ensemble ML Models





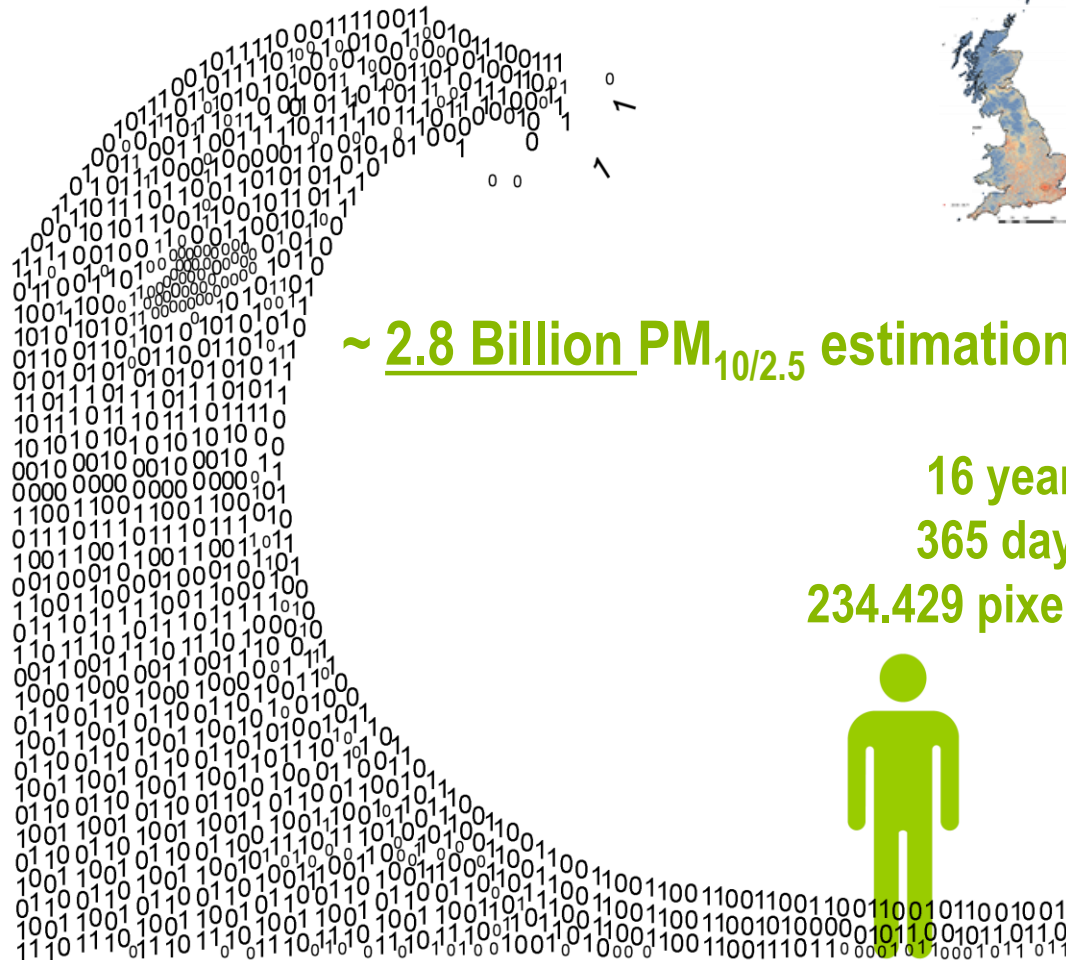
# STAGE 5: Downscaling the PM<sub>2.5</sub> Estimations

Predict PM<sub>2.5</sub> concentrations at the monitor level using small-scale predictors





# Connect with Health Data



Longitudinal studies (cohorts)



# Connect with **Worldwide** Health Data

MCC Collaborative Research Network  
An international research program on the associations  
between environmental stressors, climate, and health



## United Kingdom



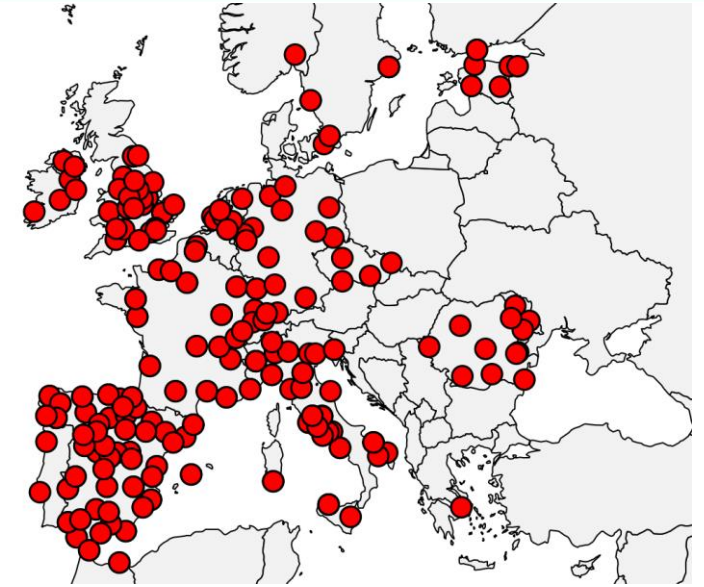
Digital

BCS70

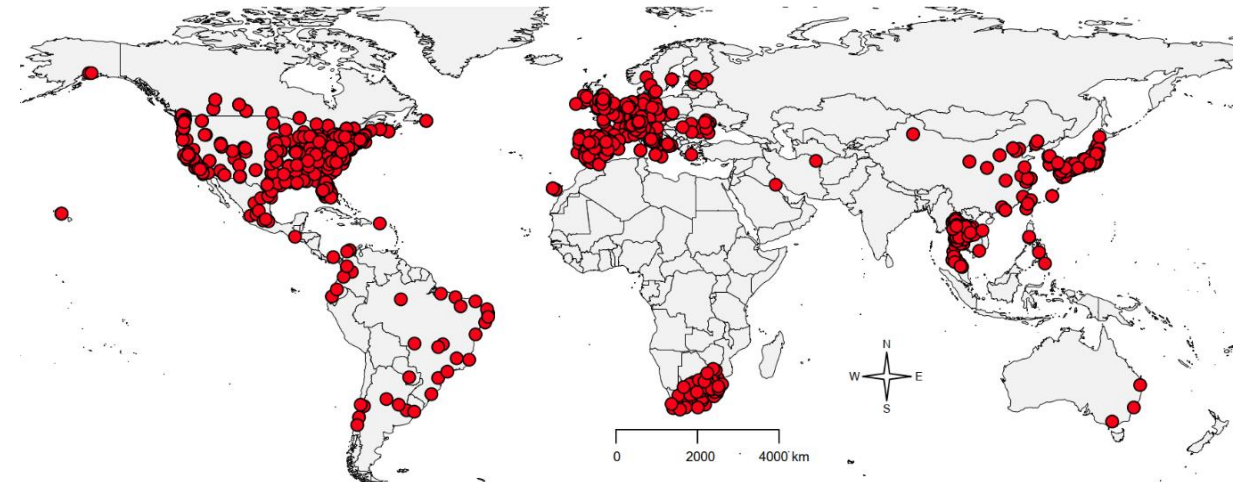
1970 British Cohort Study



## European MCC Locations



## Global MCC Locations

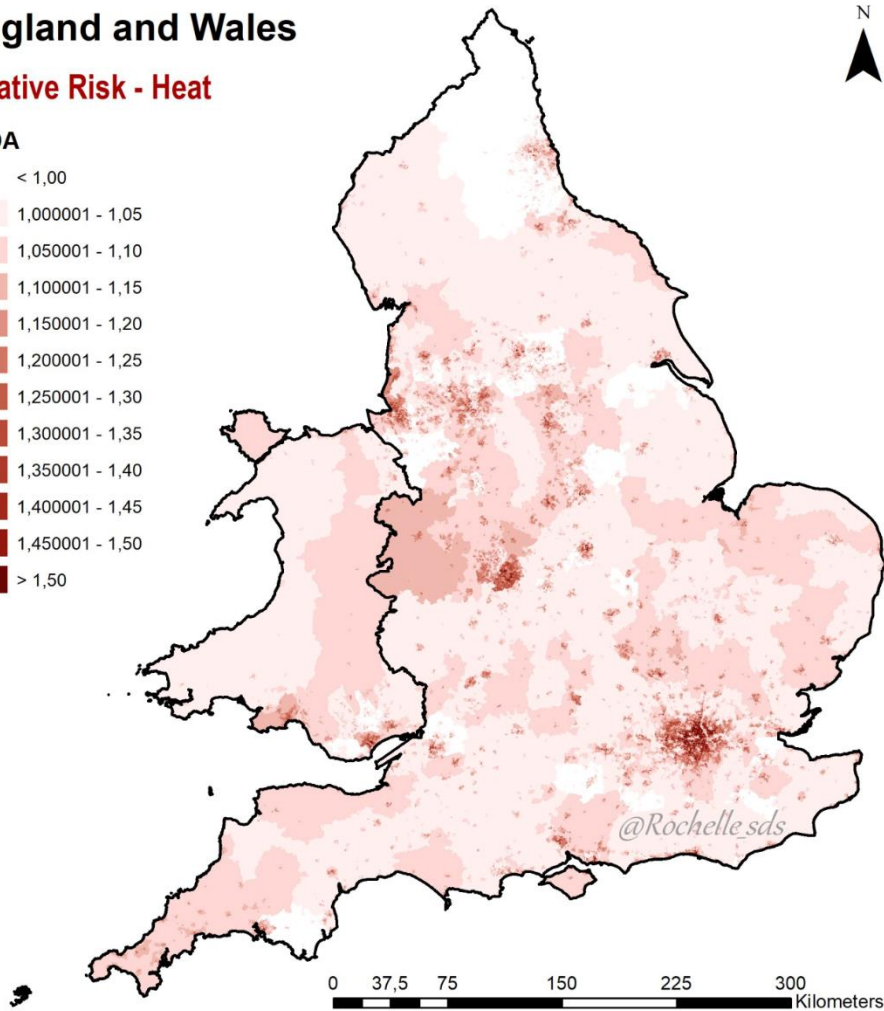
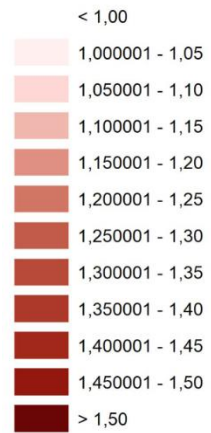


# Modelling Area-specific Environmental Risks

## England and Wales

### Relative Risk - Heat

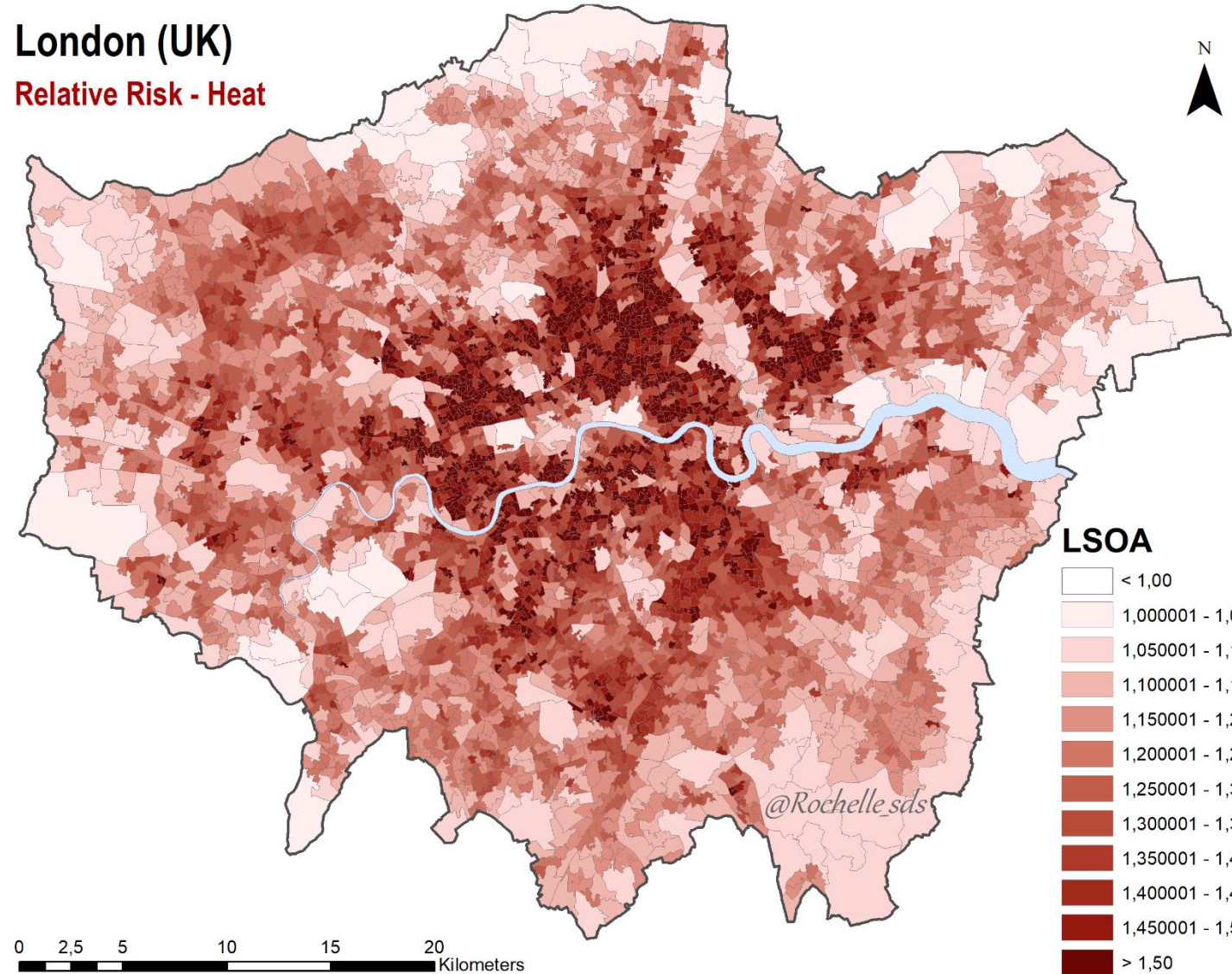
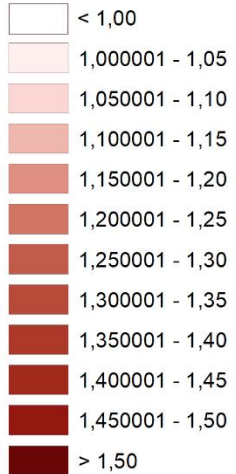
#### LSOA



## London (UK)

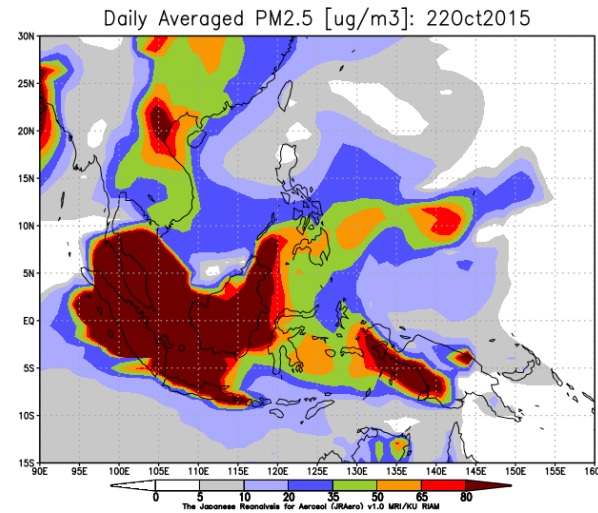
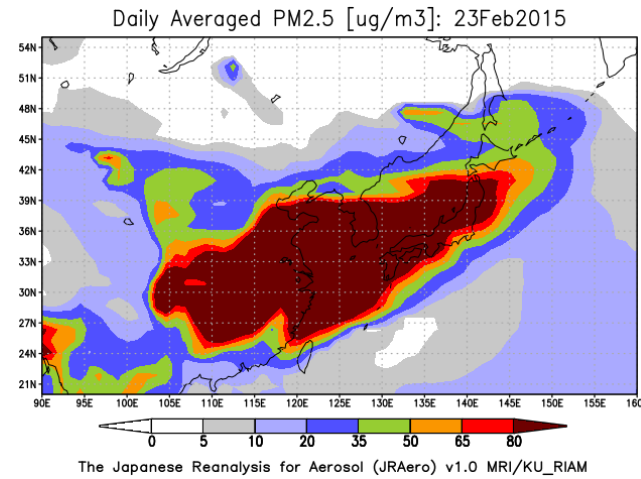
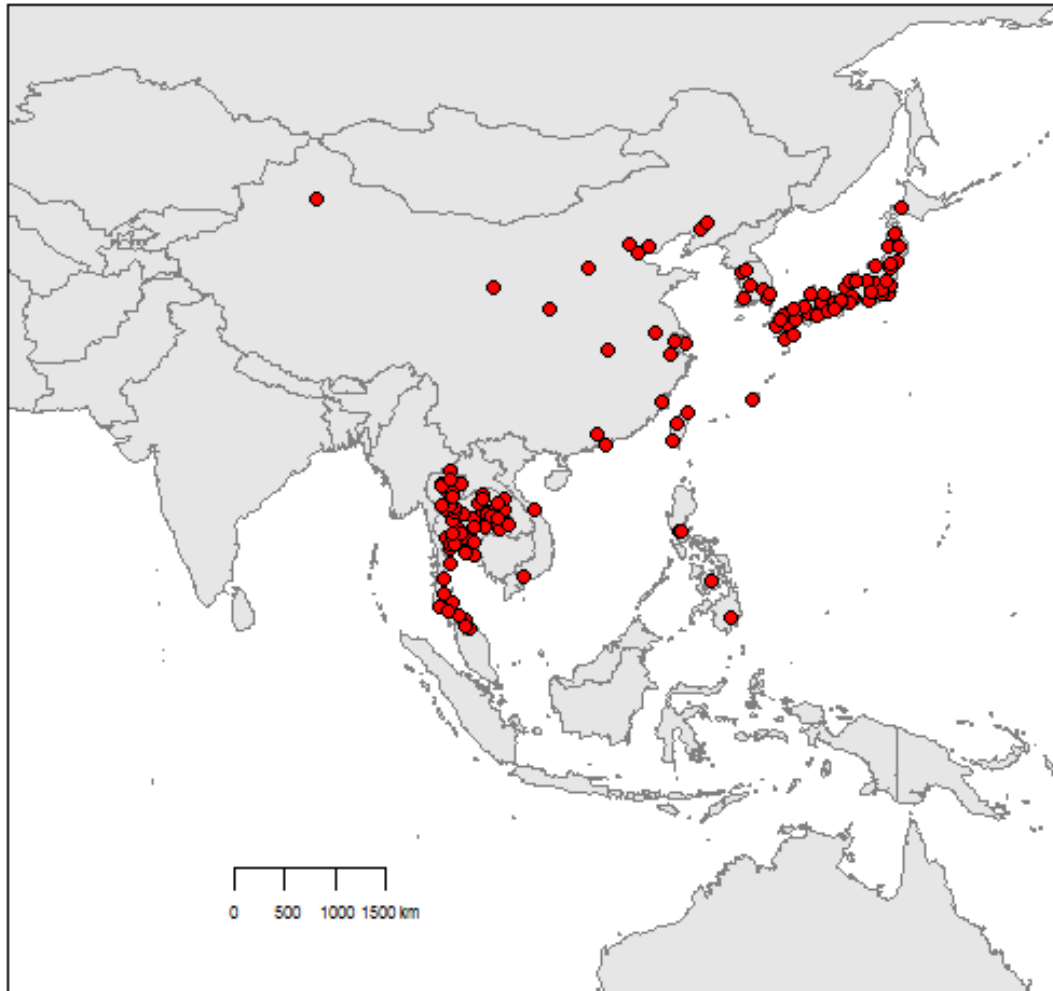
### Relative Risk - Heat

#### LSOA





# Analysis of Health Effects of Transboundary Pollution





Thank you

**Rochelle Schneider**

Senior Research Fellow in Geospatial Data Science

**Antonio Gasparrini**

Professor of Biostatistics and Epidemiology

LONDON  
SCHOOL of  
HYGIENE  
& TROPICAL  
MEDICINE

