

## Cognitive development Respiratory Tract Illness and Effects of eXposure (CORTEX) project: Data processing challenges in combining high spatial resolution pollution level data with individual level health and education data

Lyons, J<sup>1\*</sup>, Mizen, A<sup>1</sup>, Rodgers, S<sup>1</sup>, Berridge, D<sup>1</sup>, Akbari, A<sup>1</sup>, Wilkinson, P<sup>2</sup>, Milojevic, A<sup>2</sup>, Doherty, R<sup>3</sup>, Dearden, L<sup>4</sup>, Lake, I<sup>5</sup>, Carruthers, D<sup>6</sup>, Strickland, S<sup>6</sup>, Mavrogianni, A<sup>7</sup>, and Davies, G<sup>1</sup>

<sup>1</sup>Swansea University

<sup>2</sup>The London School of Health and Tropical Medicine

<sup>3</sup>The University of Edinburgh

<sup>4</sup>Institute of Fiscal Studies

<sup>5</sup>University of East Anglia

<sup>6</sup>CERC

<sup>7</sup>University College London

### Background and Objectives

There is a lack of evidence of the adverse effects of air pollution and pollen on cognition for people with air quality-related health conditions. The CORTEX project combined routinely collected health and education data, high spatial resolution air pollution modelling, and daily pollen measurements for 18,241 pupils living in Cardiff, UK, between 2009 and 2015, to investigate the acute effects of air quality and respiratory conditions on education attainment.

### Datasets

Air pollutants PM<sub>2.5</sub>, PM<sub>10</sub>, NO<sub>2</sub>, and ozone levels were modelled for 157,361 home and school locations, anonymised into the Secure Anonymised Information Linkage (SAIL) Data-bank, and summarised into minimum, average and maximum readings for 4 daily time periods reflecting pupil home/school exposure. Adding a unique Residential Anonymised Linking

Field (RALF) allowed linkage of pollution estimates to individual level data. Annual pollution datasets contained 369 columns and 472,083-rows, with one column per location, pollutant, daily time-period and day of year. Dataset transformation produced a 5 column, 3,446,205,900-row matrix per year.

### Methods and Conclusions

An algorithm using Structured Query Language (SQL) to manage data held within a relational database management system, was designed to reduce dimensionality from 24 billion to 18,241 rows of data. The algorithm calculated average means for each pollutant (PM<sub>2.5</sub>, PM<sub>10</sub>, NO<sub>2</sub>, and ozone levels) over the revision and examination periods, and summarised data into one row per pupil. The algorithm adjusted for weekends, school, and bank holidays, it calculated daily pollutant exposure for each pupil, and successfully linked 95% of pupil pollution exposures to their health and education data.

\*Corresponding Author:

Email Address: [j.lyons@swansea.ac.uk](mailto:j.lyons@swansea.ac.uk) (L Lyons)

