

# London School of Hygiene and Tropical Medicine

Improving Health Worldwide

LONDON  
SCHOOL of  
HYGIENE  
& TROPICAL  
MEDICINE



# Big data in environmental epidemiology

Arturo de la Cruz Libardi

Environment and Health Modelling (EHM) Lab  
London School of Hygiene and Tropical Medicine

2024-03-13

# intended learning outcomes

by the end of this session (lecture-demonstration), you will be able to:

1. Critically define big data
2. Describe implications and applications of big data in public health and epidemiology
3. Evaluate sources of big health and environmental data
4. Think critically about data linkage in the context of exposure assessment

# lecture outline - [online slides link](#)

## 1. motivation

- brief history

## 2. big data

- definitions
- trends and implications
- epidemiology
- applications

## 3. health and environmental data

- source examples
- harmonization and modelling
- exposure assessment
- examples

big data → epidemiology → public health

# very brief history

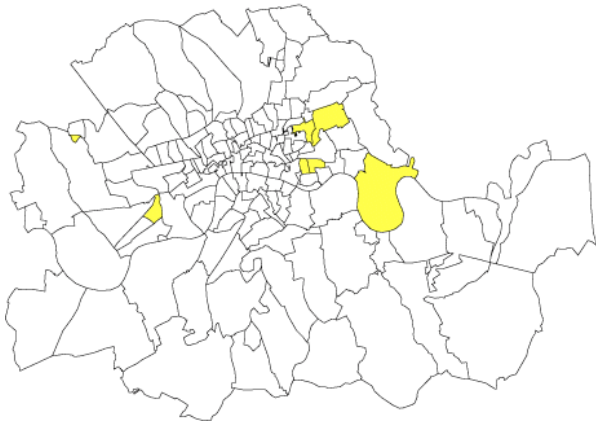
integration of subject matter knowledge, (large scale) data, and analysis

from weekly burial counts (1662) to maps (1854) and death certificates to 180k cohort (1952)

enabled by technology, **creativity**, individual and collective effort

ink and paper, punch-cards, telephone...

19/7 to 26/7



# the data line

## big data

Variety: many datasets merged

Volume: very large tables

Velocity: real-time updates

## not big data

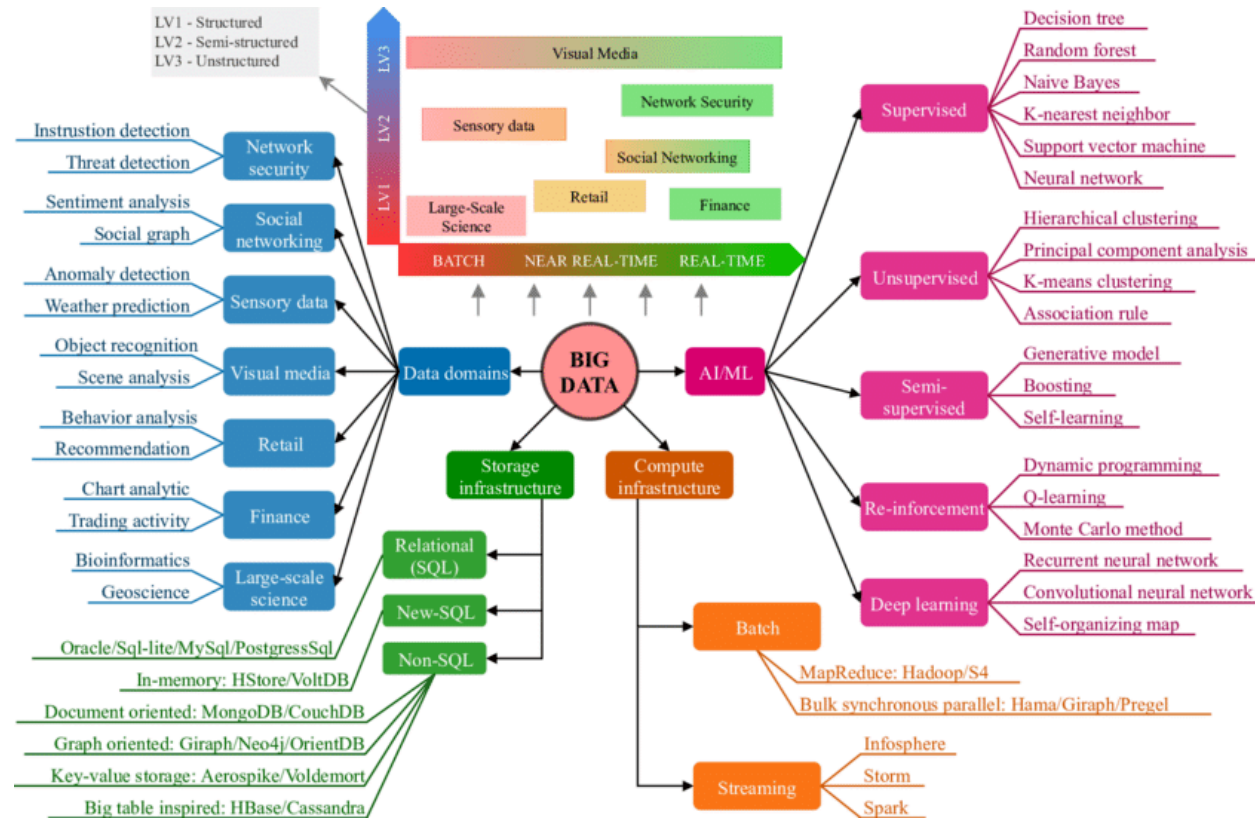
...

*more V's?*

*is it just about data?*

# specialized infrastructure, pipelines and jargon

Data - oceans, lakes, warehouses, bases



Gadekallu, Pham, Huynh-The et al. [1]





# concurrent global trends

- ageing population - urbanization (demographic change)
- environmentally complex climate change

# and technological developments

- powerful unknowable functions (machine-learning)
- smart(ish), cheap(er) and **pervasive monitoring**

# implications

- different training and emphasis
- widened research opportunities

# big data in epidemiology - challenges and opportunities

Variety - measurement error, confounding etc...

Volume - wide and tall datasets, methods, coverage, power relevant research questions ...

Velocity - highest impact potential, most dependent on infrastructure

*anything else?*

# applications in research and public health

- Research (genomics, electronic health records)
- Healthcare administration (logistics)
- COVID-19 (emergency response, tracking, data sharing)
- in references

# using EHRs to question maternal thyroid function and ASD link

- brain anatomy linked to autism is present at birth
- thyroid hormones play key role in brain development

**Q 1: is hypothyroidism associated with inc. risk of autism (430k births) ?**

**Q 2: does risk for medicated mothers differ from Q1 risk ?**

**Q 3: does risk for lab-tested medicated mothers differ from Q1 and Q2 risks ?**

**Q 4: are levels of TSH/FT4 associated with inc. risk of autism (50k births) ?**

Results indicate that maternal thyroid conditions are associated with increased ASD risk in progeny, but suggestively not due to direct effects of thyroid hormones. Instead, factors that influence maternal thyroid function could have etiologic roles in ASD through pathways independent of maternal gestational thyroid hormones and thus be unaffected by medication treatment. Factors known to disrupt thyroid function should be examined for possible involvement in ASD etiology.

# OpenSAFELY

## OpenSAFELY: the origin story

On 7th May 2020, the OpenSAFELY Collaborative pre-printed the world's largest study into factors associated with death from Covid-19, based on an analysis running across the full pseudonymised health records of 40% of the English population. This is an unprecedented scale of data.

... a huge collaboration including the Bennett Institute for Applied Data Science at the University of Oxford, the EHR research group at London School of Hygiene and Tropical Medicine, NHS England, and TPP. Over 42 days during the peak of the first wave of COVID-19 this team worked day and night to produce a fully open-source, privacy-preserving software platform, capable of running open and reproducible analytics across electronic health records, all held securely in situ. Since then the OpenSAFELY platform has expanded to a full scale analytic environment for secure data analysis, reproducible data curation, federated analysis, and code sharing, with every line of code for the platform, for data management, and for data analysis **all shared openly by default, in re-usable forms, automatically, and without exception.**

Couldn't load plug-in.

All LSHTM OpenSAFELY projects

# health data sources

- datasets: **ProjectTycho**
- cohorts: **BioBank** and **OurFutureHealth**
- platforms: **CPRD** and **OpenSAFELY**
- personal sensor data

**Table 1**

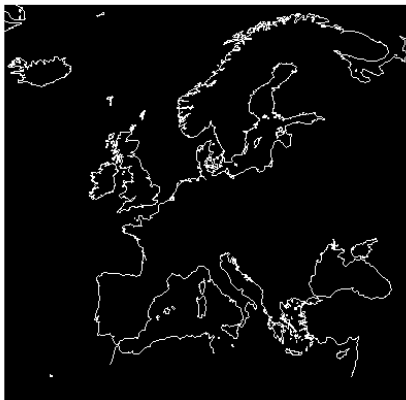
Examples of storage needs (per person).

Information	Size
Human genome	1 GB
+ structure determination of the proteins	Several PB <sup>a</sup>
Electronic health record	1 MB-5GB, expected to increase 50-fold from 2012 to 2020 <sup>b,c</sup>
Heart rate monitor (per month)	9 GB <sup>d</sup>
Continuous video life-logger (per month)	58GB <sup>e</sup>
Accelerometer (8-h a day, per month)	1 GB <sup>f</sup>
Medical image	MB to GB, up to 1 TB. e.g. 64/128-slice CT scan, 3.0 T MRI and PET often exceeding 100 MB <sup>b</sup> .

Tonne, Basagaña, Chaix et al. [3]

# environmental data

- modelled: atmospheric dispersion models, reanalysis, digital twins
- raw: ground monitors, mobile sensors, satellites
- raster vs vector



2018-01-01

\*\* figure from: <https://search.earthdata.nasa.gov>

\*\*

Search Results (73 Collections)

Sentinel-5P TROPOMI Radiance product  
band 3 (UVIS detector) L1B V1  
(S5P\_L1B\_RA\_BD3) at GES DISC

Showing 20 of 6,491 matching granules

<p>S5P_OFFL_L1B_RA_BD 3_20190806T003836_20190806T022006_09387_01_010000_20190806T040122.nc</p> <p>START 2019-08-06 01:00:11</p> <p>END 2019-08-06 01:58:34</p> <p>+ Download</p>	<p>S5P_OFFL_L1B_RA_BD 3_20190805T225707_20190806T003836_09386_01_010000_20190806T022059.nc</p> <p>START 2019-08-05 23:18:41</p> <p>END 2019-08-06 00:17:05</p> <p>+ Download</p>
<p>S5P_OFFL_L1B_RA_BD 3_20190805T211537_20190805T225707_09</p> <p>START 2019-08-05 21:15:37</p> <p>END 2019-08-05 22:57:07</p> <p>+ Download</p>	<p>S5P_OFFL_L1B_RA_BD 3_20190805T193408_20190805T225707_09</p> <p>START 2019-08-05 19:34:08</p> <p>END 2019-08-05 22:57:07</p> <p>+ Download</p>

Download All 6,491

# why use both environment and health data

- a part of disease etiology remains unexplained and is likely due to the environment
- big data processes offer great potential for environmental health research
- most of all data generated has a spatial and a temporal reference

## environment + health data synergy

1. research question
2. get health data
3. **get/harmonize/model environmental data**
4. **LINK**
5. analyse



# from data to exposure

## Env. data

(complexity)

(none) → continuous modelled output

(simple) → inverse distance weighted surface from point measurement

(complex) → multi-stage machine-learning models using harmonized features

## Linkage

(simple) → matching nearest

(simple/medium) → points on raster (bilinear interpolation) [4]

(medium) → aggregate over small area

(complex) → from a trajectory accounting for microenvironments [5]

Vanoli, Mistry, De La Cruz Libardi et al. [4] Smith, Mitsakou, Kitwiroon et al. [5]

# A Satellite-Based Spatio-Temporal Machine Learning Model to Reconstruct Daily PM<sub>2.5</sub> Concentrations across Great Britain [6]

- Ground observations of PM<sub>2.5</sub>
- A **lot** of environmental data
- Random forest (ML) algorithms

Schneider, Vicedo-Cabrera, Sera et al. [6]

# A Satellite-Based Spatio-Temporal Machine Learning Model to Reconstruct Daily PM<sub>2.5</sub> Concentrations across Great Britain [6]

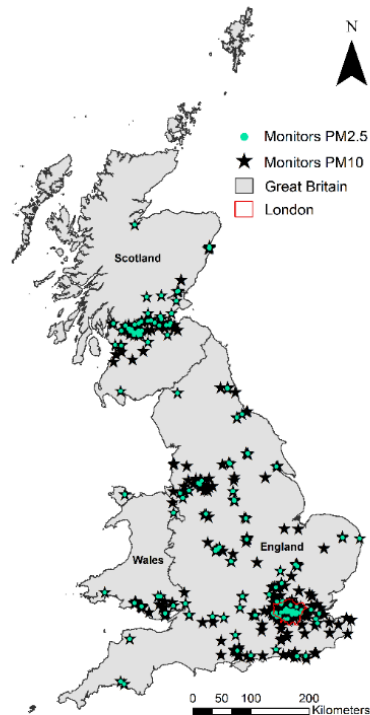


Figure 1. Spatial distribution of 581 PM<sub>10</sub> (black star) and 183 PM<sub>2.5</sub> (turquoise dots) monitors across Great Britain during the study period.

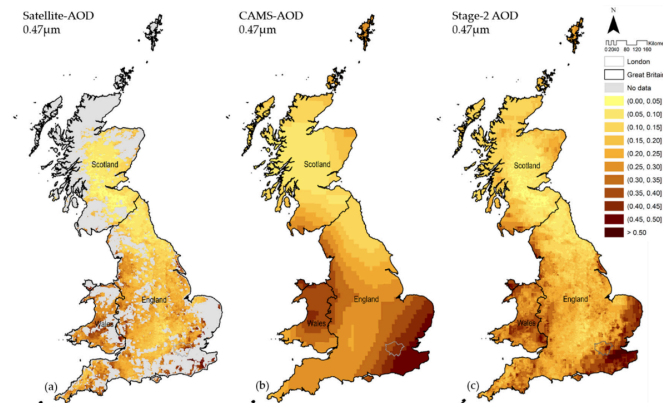


Figure 2. Satellite-AOD 0.47  $\mu\text{m}$  values are represented in Figure 2a (mean of all Terra- and

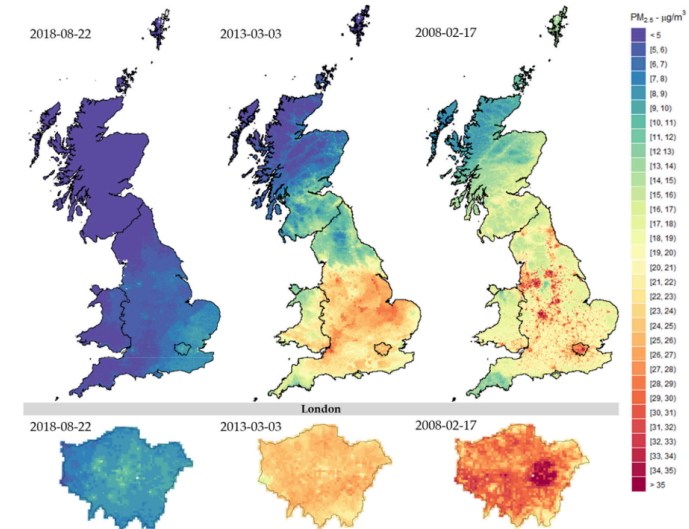


Figure 5. Stage-4 day-specific PM<sub>2.5</sub> estimations across Great Britain (Top) and London (Bottom).

# London Hybrid Exposure Model: Improving Human Exposure Estimates to NO<sub>2</sub> and PM<sub>2.5</sub> in an Urban Setting [5]

...the London Hybrid Exposure Model (LHEM), (...) calculates exposure of the Greater London population to outdoor air pollution sources, in-buildings, in-vehicles, and outdoors, using survey data of when and where people spend their time.

- London Travel Demand Survey, trip route simulation

Exposure to outdoor air pollution was provided by CMAQ-urban, which couples the Weather Research and Forecasting (WRF) meteorological model, the Community Multiscale Air Quality (CMAQ) regional scale model, and the Atmospheric Dispersion Modeling System (ADMS) roads model

- I/O ratio for indoor air levels

- for in-vehicle levels:  $\frac{dC_{in}}{dt} = \lambda_{in}(C_{out} - C_{in}) - n\lambda_{HVAV} \cdot C_{in} - V_g\left(\frac{A^*}{V}\right) \cdot C_{in} + \frac{Q}{V}$

- constant value for the underground

- **"microenvironments"**

Smith, Mitsakou, Kitwiroon et al. [5]

# London Hybrid Exposure Model: Improving Human Exposure Estimates to NO<sub>2</sub> and PM<sub>2.5</sub> in an Urban Setting [5]

## residential vs modelled exposure

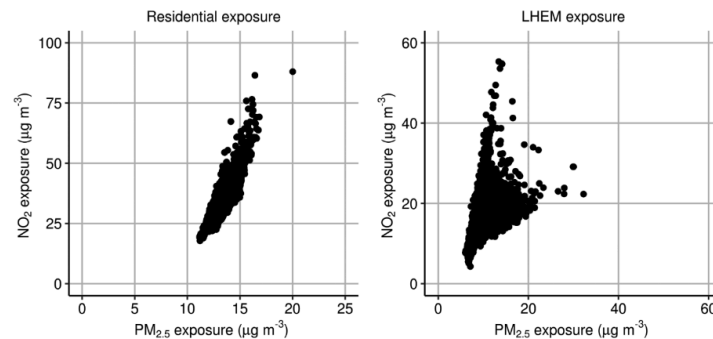


Figure 4. Scatter plots of NO<sub>2</sub> and PM<sub>2.5</sub> exposure outdoors at the residential address (left) ( $R = 0.90$ ) and using the LHEM (right).

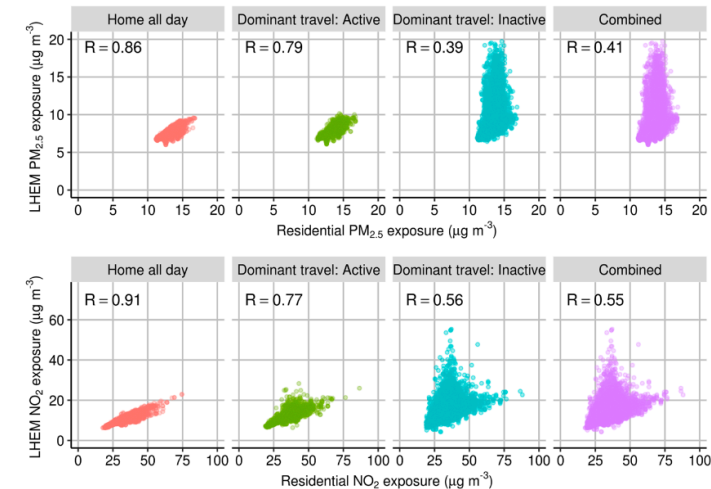


Figure 3. Scatter plots of NO<sub>2</sub> (bottom) and PM<sub>2.5</sub> (top) LHEM exposure versus exposure at the residential address - demonstrating the relative strength of the relationship between those who undertake active travel (cycle and walk), those that stay at home, and those who undertake inactive travel (car, motorcycle, bus, train, and tube).

# we have learned to

1. Critically define big data as **big data processes**
2. Describe implications and applications of big data in public health and epidemiology
  - classical (**measurement error, confounding**) challenges
  - new (**comprehensive health data, real-time action**) opportunities
3. Evaluate sources of big health and environmental data
  - health **genetic data, EHRs, wearable sensors**
  - environment **reanalyses, satellites, ground sensors**
4. Think critically about data linkage in the context of exposure assessment

# references

- [1] T. R. Gadekallu, Q. Pham, T. Huynh-The, et al. *Federated Learning for Big Data: A Survey on Opportunities, Applications, and Future Directions*. En. 2021.
- [2] R. S. Rotem, G. Chodick, V. Shalev, et al. "Maternal Thyroid Disorders and Risk of Autism Spectrum Disorder in Progeny". En-US. In: *Epidemiology* 31.3 (May. 2020), p. 409. ISSN: 1044-3983. DOI: [10.1097/EDE.0000000000001174](https://doi.org/10.1097/EDE.0000000000001174). URL: [https://journals.lww.com/epidem/fulltext/2020/05000/maternal\\_thyroid\\_disorders\\_and\\_risk\\_of\\_autism.15.aspx](https://journals.lww.com/epidem/fulltext/2020/05000/maternal_thyroid_disorders_and_risk_of_autism.15.aspx) (visited on 03/13/2024).
- [3] C. Tonne, X. Basagaña, B. Chaix, et al. "New frontiers for environmental epidemiology in a changing world". In: *Environment International* 104 (Jul. 2017), pp. 155-162. ISSN: 0160-4120. DOI: [10.1016/j.envint.2017.04.003](https://doi.org/10.1016/j.envint.2017.04.003). URL: <https://www.sciencedirect.com/science/article/pii/S0160412017301459> (visited on 02/15/2024).
- [4] J. Vanoli, M. N. Mistry, A. De La Cruz Libardi, et al. "Reconstructing individual-level exposures in cohort analyses of environmental risks: an example with the UK Biobank". En. In: *Journal of Exposure Science & Environmental Epidemiology* (Jan. 2024). ISSN: 1559-0631, 1559-064X. DOI: [10.1038/s41370-023-00635-w](https://doi.org/10.1038/s41370-023-00635-w). URL: <https://www.nature.com/articles/s41370-023-00635-w> (visited on 03/10/2024).

# references

- [5] J. D. Smith, C. Mitsakou, N. Kitwiroon, et al. "London Hybrid Exposure Model: Improving Human Exposure Estimates to NO<sub>2</sub> and PM<sub>2.5</sub> in an Urban Setting". En. In: *Environmental Science & Technology* 50.21 (Nov. 2016), pp. 11760-11768. ISSN: 0013-936X, 1520-5851. DOI: [10.1021/acs.est.6b01817](https://doi.org/10.1021/acs.est.6b01817). URL: <https://pubs.acs.org/doi/10.1021/acs.est.6b01817> (visited on 02/02/2023).
- [6] R. Schneider, A. Vicedo-Cabrera, F. Sera, et al. "A Satellite-Based Spatio-Temporal Machine Learning Model to Reconstruct Daily PM<sub>2.5</sub> Concentrations across Great Britain". En. In: *Remote Sensing* 12.22 (Nov. 2020), p. 3803. ISSN: 2072-4292. DOI: [10.3390/rs12223803](https://doi.org/10.3390/rs12223803). URL: <https://www.mdpi.com/2072-4292/12/22/3803> (visited on 02/03/2022).
- [7] D. Cox, C. Kartsonaki, and R. H. Keogh. "Big data: Some statistical issues". In: *Statistics & Probability Letters* 136 (May. 2018), pp. 111-115. ISSN: 0167-7152. DOI: [10.1016/j.spl.2018.02.015](https://doi.org/10.1016/j.spl.2018.02.015). URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5992743/> (visited on 03/11/2024).
- [8] E. J. Williamson, A. J. Walker, K. Bhaskaran, et al. "Factors associated with COVID-19-related death using OpenSAFELY". En. In: *Nature* 584.7821 (Aug. 2020). Publisher: Nature Publishing Group, pp. 430-436. ISSN: 1476-4687. DOI: [10.1038/s41586-020-2521-4](https://doi.org/10.1038/s41586-020-2521-4). URL: <https://www.nature.com/articles/s41586-020-2521-4> (visited on 03/12/2024).



# references

- [9] M. J. Khoury and J. P. A. Ioannidis. "Big data meets public health". En. In: *Science* 346.6213 (Nov. 2014), pp. 1054-1055. ISSN: 0036-8075, 1095-9203. DOI: [10.1126/science.aaa2709](https://doi.org/10.1126/science.aaa2709). URL: <https://www.science.org/doi/10.1126/science.aaa2709> (visited on 03/08/2024).
- [10] S. J. Mooney and V. Pejaver. "Big Data in Public Health: Terminology, Machine Learning, and Privacy". In: *Annual Review of Public Health* 39.1 (2018). \_ eprint: <https://doi.org/10.1146/annurev-publhealth-040617-014208>, pp. 95-112. DOI: [10.1146/annurev-publhealth-040617-014208](https://doi.org/10.1146/annurev-publhealth-040617-014208). URL: <https://doi.org/10.1146/annurev-publhealth-040617-014208> (visited on 03/05/2024).

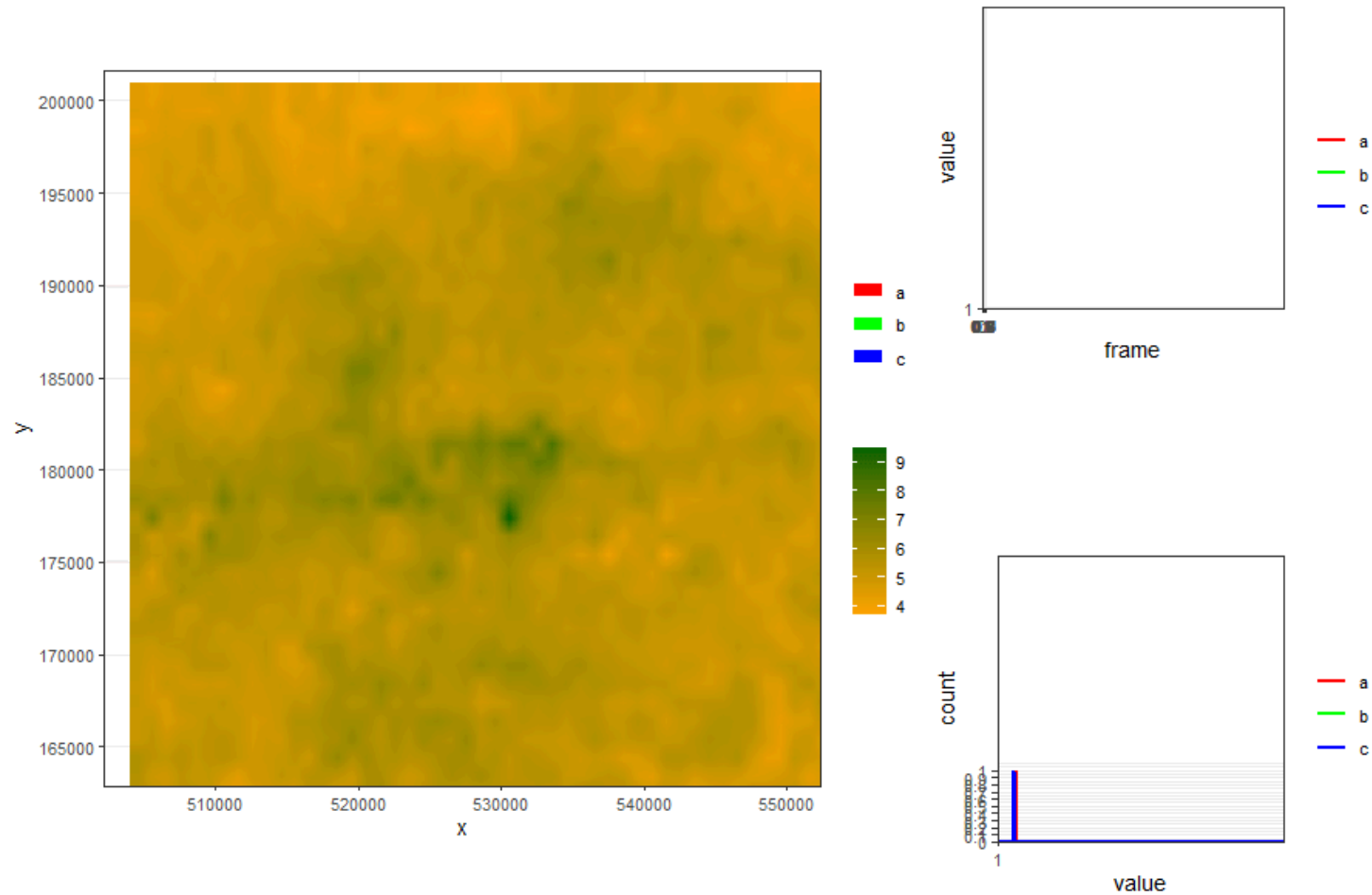
## other info

Presentation made with [xaringan](#) in RStudio.

Contact: [Arturo.de-la-Cruz-Libardi@lshtm.ac.uk](mailto:Arturo.de-la-Cruz-Libardi@lshtm.ac.uk)

Slides: [https://adlcruz.github.io/linked\\_content/pres\\_bigdataenvepi\\_2024/bdee\\_slides.html](https://adlcruz.github.io/linked_content/pres_bigdataenvepi_2024/bdee_slides.html)

## points trajectories on dynamic map and corresponding exposure



## finnish meteorological institute reanalysis

## SILAM v.5.7

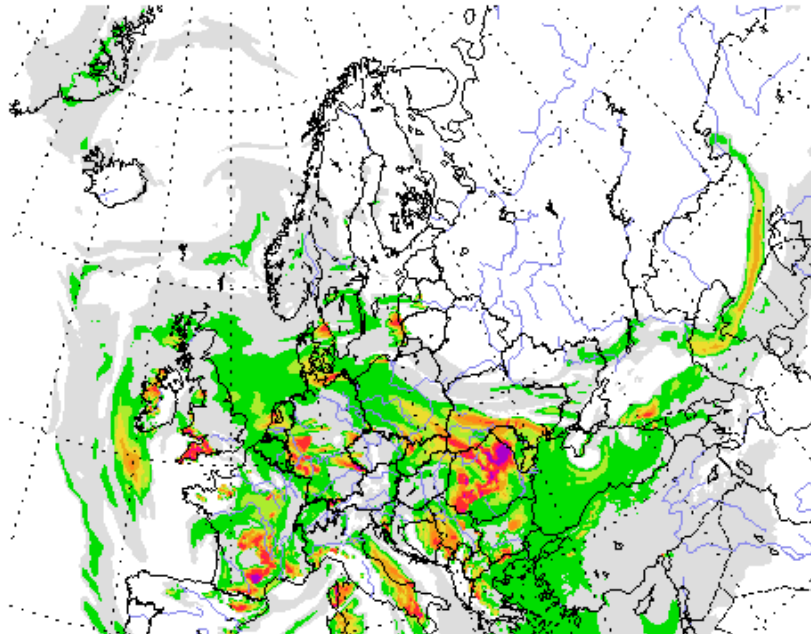
System for Integrated modelLing of Atmospheric  
coMposition



FINNISH METEOROLOGICAL INSTITUTE (<http://en.ilmatieteenlaitos.fi/>)

☁ Region: europe

SILAM model forecast: alder pollen  
(#/m3) 15Z10MAR2024



☰ Control panel

Pollen taxon:

Alder

Birch

Grass

Olive

Ragweed

Mugwort

Pollen Index

Allergy Risk

Region:

Europe

Hour: -57



how is OpenSAFELY testing their new features? with chatGPT of course.

- OpenSAFELY query (ehrQL) reliability **testing** using **generative artificial intelligence!**

# suggestions?

- DASH 26th March opening event
- hundreds of hours of free and open resources
- a lot of local and global circumstances to improve

