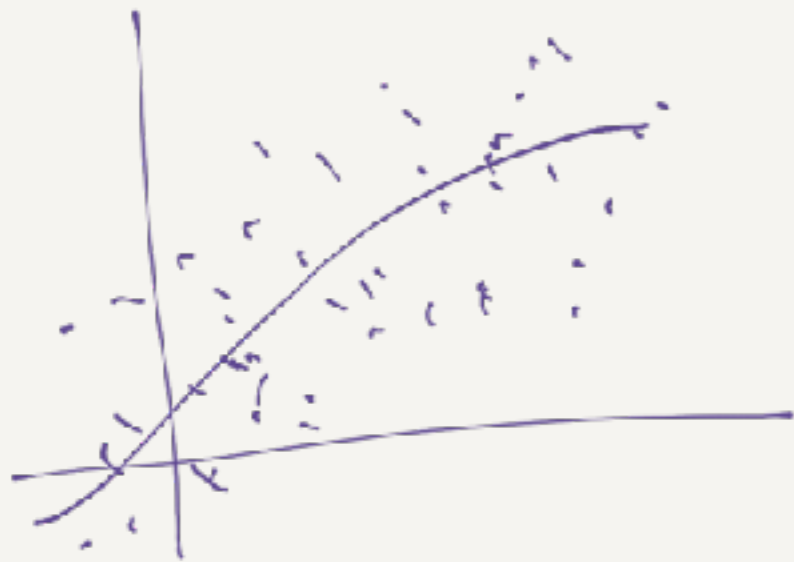


Data Science: The Good The Bad and The Future



Roger D. Peng
Department of **Biostatistics**
Johns Hopkins Bloomberg
School of Public Health
@rdpeng



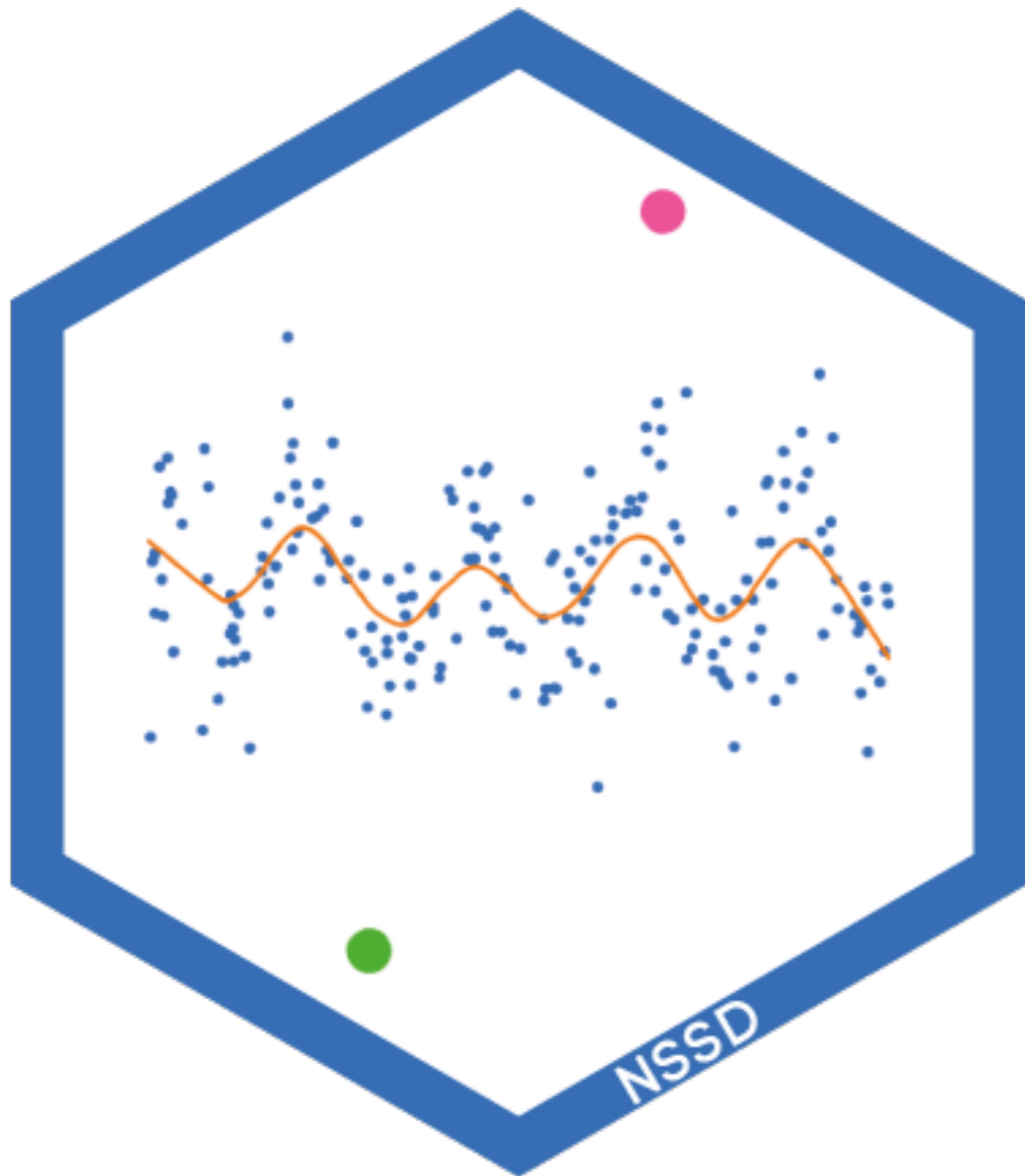
JOHNS HOPKINS

BLOOMBERG SCHOOL
of PUBLIC HEALTH



simplystatistics.org

@simplystats



Not So Standard Deviations

(with Hilary Parker of Stitch Fix)



@NSSDeviations

<https://soundcloud.com/nssd-podcast>

Subscribe in iTunes: <https://goo.gl/ZhWYbd>

NSSD Episode 23



Not So Standard Deviations

Episode 23 - Special Guest Walt Hickey

1 month

Technology



<https://goo.gl/d8eszr>

Write a comment



Like



Share

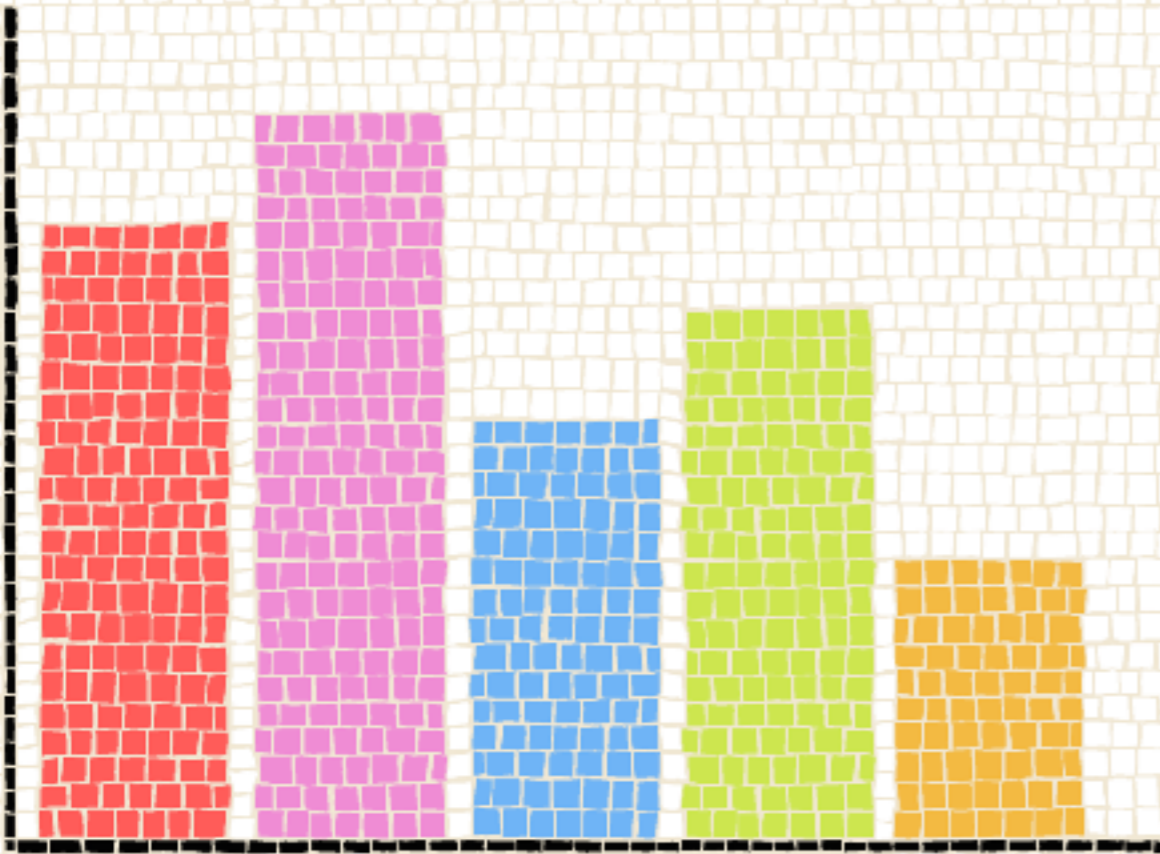


More

▶ 8,734 ♥ 6

The Art of Data Science

A Guide for Anyone Who Works with Data



Roger D. Peng & Elizabeth Matsui

CONVERSATIONS ON DATA SCIENCE



ROGER D. PENG
HILARY PARKER

leanpub.com/artofdatascience

leanpub.com/conversationsondatascience





Protecting Health, Saving Lives —



Protecting Health, Saving Lives —
Millions at a Time



Protecting Health, Saving Lives —

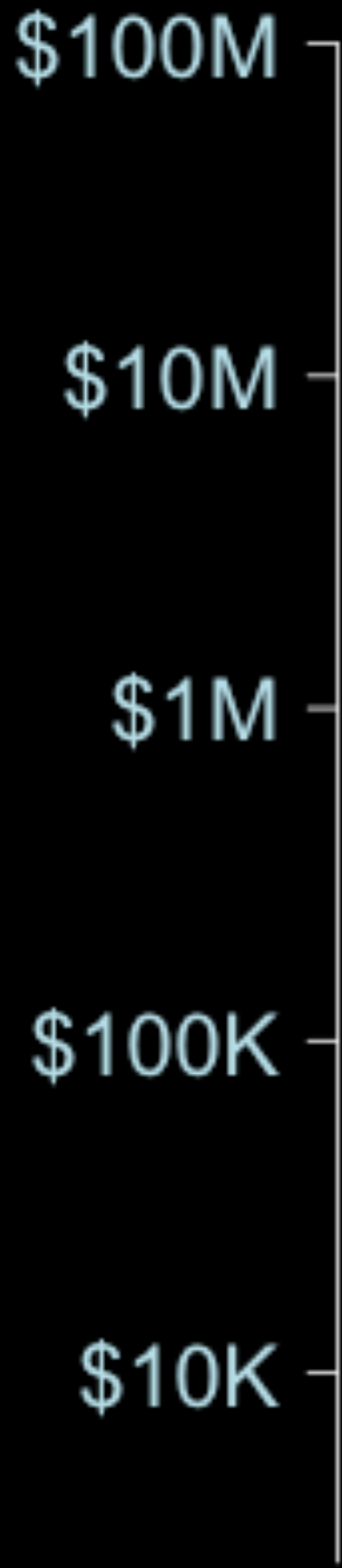
Millions at a Time



(of data points)

The Measurement Revolution...

\$ per (human) Genome



<http://www.genome.gov/sequencingcosts/>

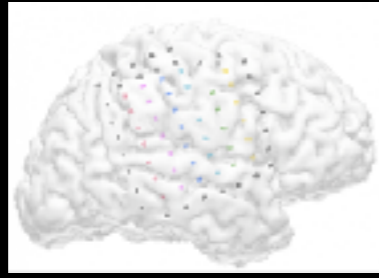
The Measurement Revolution...

Brian
Caffo



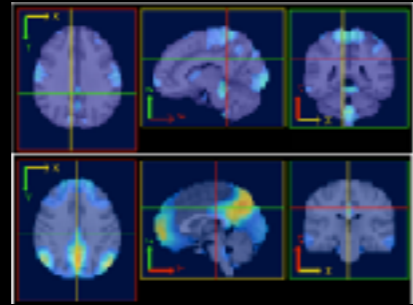
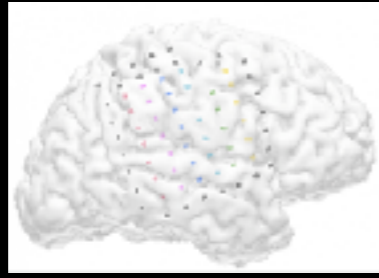
The Measurement Revolution...

Brian
Caffo



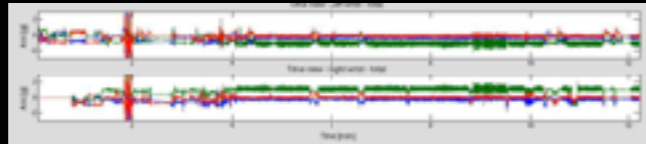
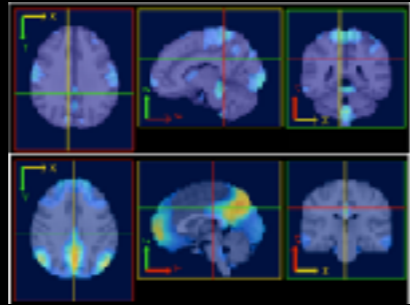
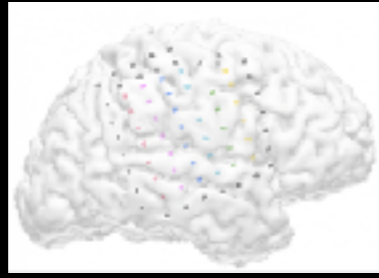
The Measurement Revolution...

Brian
Caffo



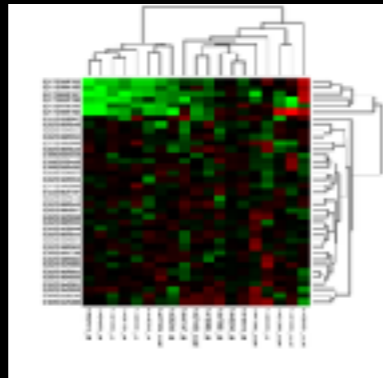
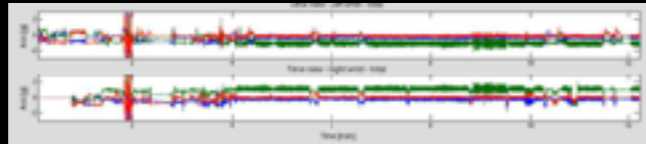
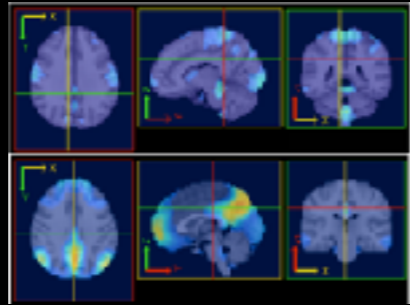
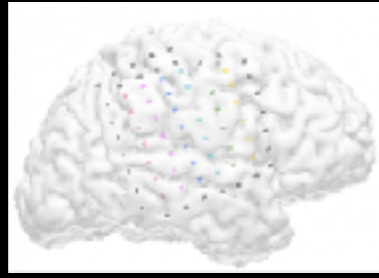
The Measurement Revolution...

Brian
Caffo



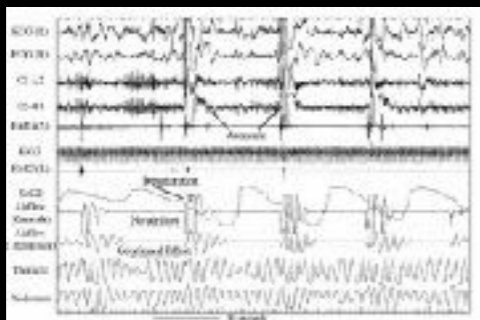
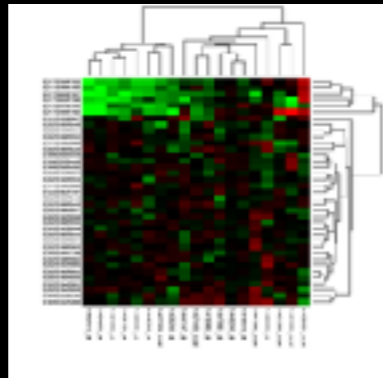
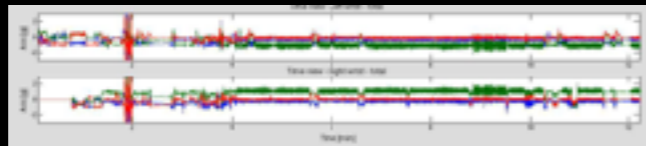
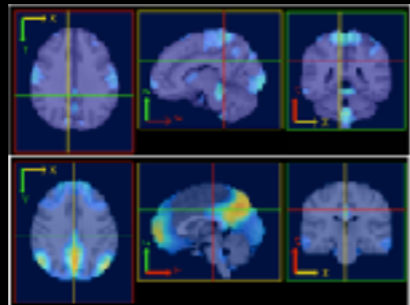
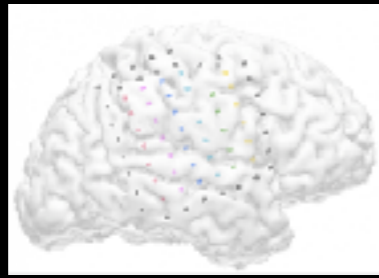
The Measurement Revolution...

Brian
Caffo



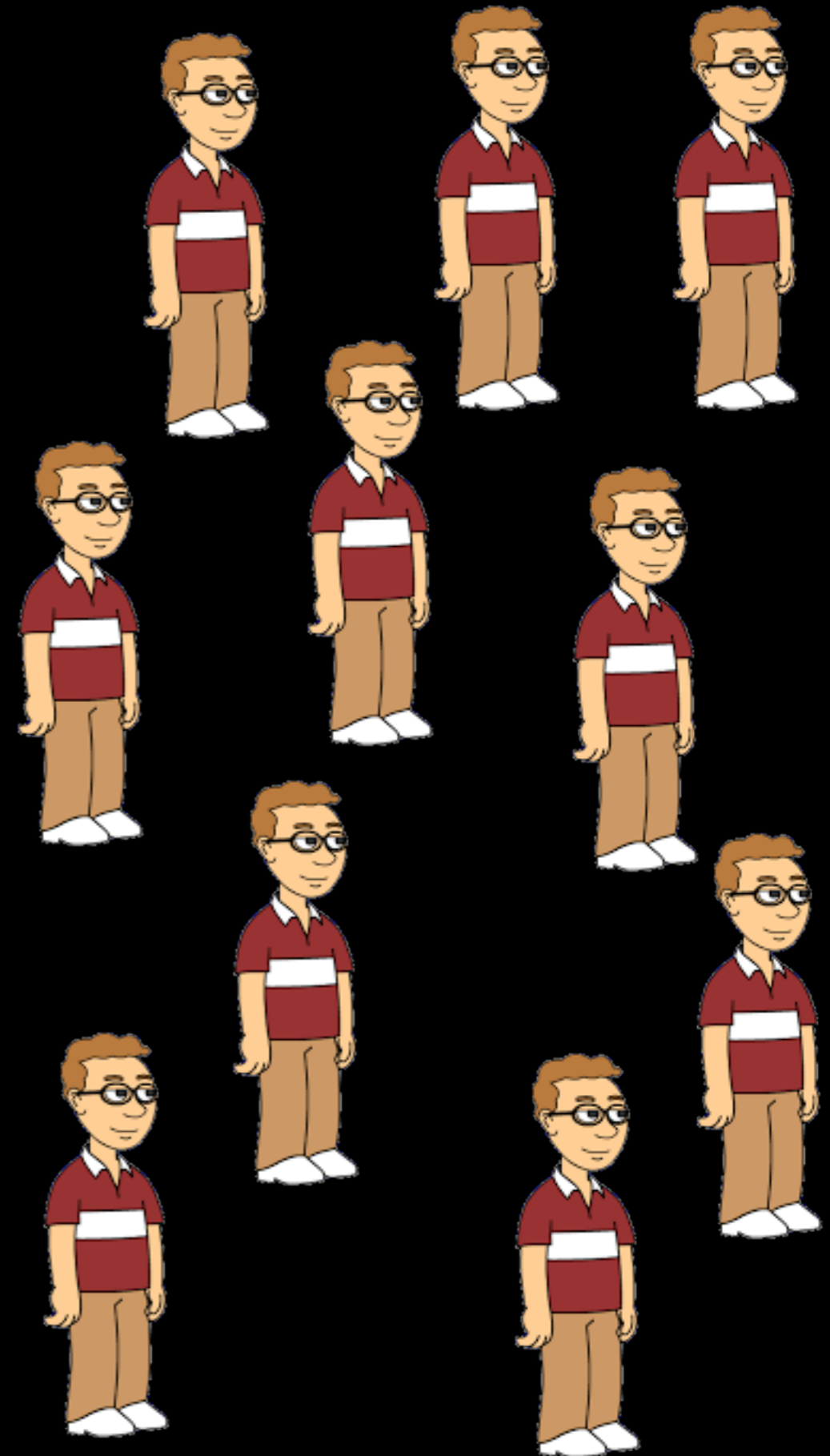
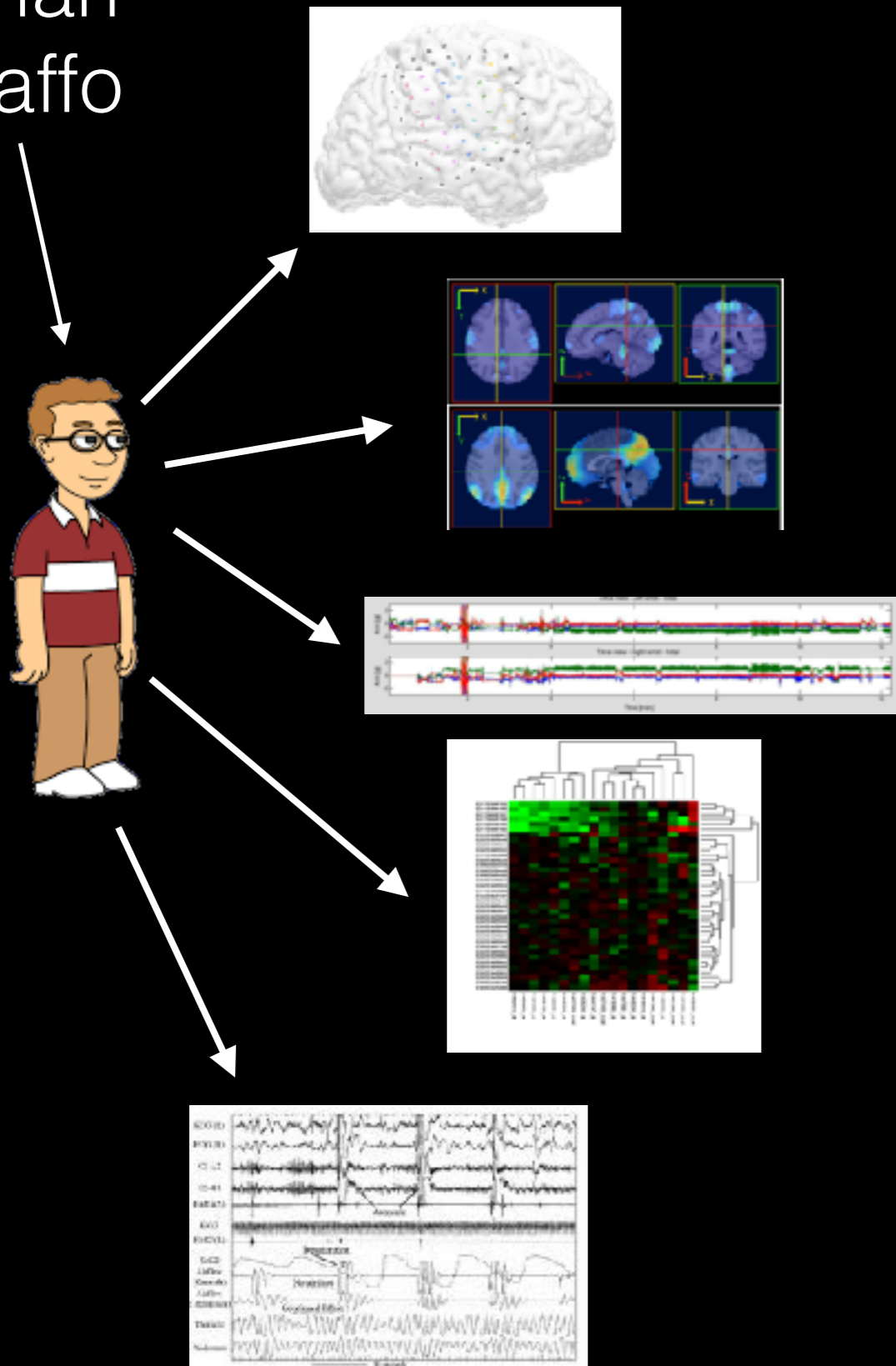
The Measurement Revolution...

Brian
Caffo



The Measurement Revolution... Multiplied!

Brian Caffo



A yellow Pac-Man character is positioned on the left side of the slide, facing right. Its mouth is open, and it appears to be eating the first letter of the text.

Data is Eating the World*

*see "Software is Eating the World" by Marc Andreessen



Data is Eating the World*

(but analysis isn't yet)

*see "Software is Eating the World" by Marc Andreessen

Not Enough Geeks

06/28
2011

Critical Shortage Of "Data Geek" Talent Predicted By 2018



McKinsey & Company

New research by the McKinsey Global Institute (MGI) forecasts a 50 to 60 percent gap between the supply and demand of people with deep analytical talent. These "data geeks" have advanced training in statistics, machine learning as well as the ability to analyze data sets. The study projects there will be approximately 140,000 to 190,000 unfilled positions for [data analytics](#) experts in the U.S. by 2018 and a

shortage of 1.5 million managers and analysts who have the ability to understand and make decisions using big data.

Not Enough Geeks

06/28
2011

Critical Shortage Of "Data Geek" Talent Predicted By 2018

McKinsey&Company

New research by the McKinsey Global Institute (MGI) forecasts a 50 to 60 percent gap between the supply and demand of people with deep analytical talent. These "data geeks" have advanced training in statistics, machine learning as well as the ability to analyze data sets. The study projects there will be approximately 140,000 to 190,000 unfilled positions for data analytics experts in the U.S. by 2018 and a

shortage of 1.5 million managers and analysts who have the ability to understand and make decisions using big data.

Epidemic of Bad Data Analysis

Opinion: Reproducible research can still be wrong: Adopting a prevention approach

Jeffrey T. Leek^{a,1} and Roger D. Peng^b

^aAssociate Professor of Biostatistics and Oncology and ^bAssociate Professor of Biostatistics, Johns Hopkins University, Baltimore, MD

Reproducibility—the ability to recompute results—and replicability—the chances other experimenters will achieve a consistent result—are two foundational characteristics of successful scientific research. Consistent findings from independent investigators are the primary means by which scientific evidence accumulates for or against a hypothesis. Yet, of late, there has been a crisis of confidence among researchers worried about the rate at which studies are either reproducible or replicable. To maintain the integrity of science research and the public's trust in science, the scientific community must ensure reproducibility and replicability by engaging in a more preventative approach that greatly expands data analysis education and routinely uses software tools.

We define reproducibility as the ability to

been some very public failings of reproducibility across a range of disciplines from cancer genomics (3) to economics (4), and the data for many publications have not been made publicly available, raising doubts about the quality of data analyses. Popular press articles have raised questions about the reproducibility of all scientific research (5), and the US Congress has convened hearings focused on the transparency of scientific research (6). The result is that much of the scientific enterprise has been called into question, putting funding and hard won scientific truths at risk.

From a computational perspective, there are three major components to a reproducible and replicable study: (i) the raw data from the experiment are available, (ii) the statistical code and documentation to reproduce the

computational tools such as knitr, iPython notebook, LONI, and Galaxy (8) have simplified the process of distributing reproducible data analyses.

Unfortunately, the mere reproducibility of computational results is insufficient to address the replication crisis because even a reproducible analysis can suffer from many problems—confounding from omitted variables, poor study design, missing data—that threaten the validity and useful interpretation of the results. Although improving the reproducibility of research may increase the rate at which flawed analyses are uncovered, as recent high-profile examples have demonstrated (4), it does not change the fact that problematic research is conducted in the first place.

The key question we want to answer when seeing the results of any scientific study is “Can I trust this data analysis?” If we think of problematic data analysis as a disease, reproducibility speeds diagnosis and treatment in

Epidemic of Bad Data Analysis

Opinion: Reproducible research can still be wrong: Adopting a prevention approach

Jeffrey T. Leek^{a,1} and Roger D. Peng^b

computational tools such as knitr, iPython

“The best way to prevent poor data analysis in the scientific literature is to (i) increase the number of trained data analysts in the scientific community and (ii) identify statistical software and tools that can be shown to improve reproducibility and replicability of studies.”

reproducible or replicable. To maintain the integrity of science research and the public's trust in science, the scientific community must ensure reproducibility and replicability by engaging in a more preventative approach that greatly expands data analysis education and routinely uses software tools.

We define reproducibility as the ability to

scientific enterprise has been called into question, putting funding and hard won scientific truths at risk.

From a computational perspective, there are three major components to a reproducible and replicable study: (i) the raw data from the experiment are available, (ii) the statistical code and documentation to reproduce the

problematic research is conducted in the first place.

The key question we want to answer when seeing the results of any scientific study is “Can I trust this data analysis?” If we think of problematic data analysis as a disease, reproducibility speeds diagnosis and treatment in

Challenges:

- ① Measurement explosion!
- ② Not enough analysts!
- ③ Reproducibility crisis!

What should
we do ???

LEAD!



What is Data Science?

- Formulating a question that can be answered with data
- Assembling, cleaning, tidying data relevant to a question
- Exploring data, checking, eliminating hypotheses
- Developing a (statistical) model
- Making statistical inference
- Communicating findings

What does a data scientist do?

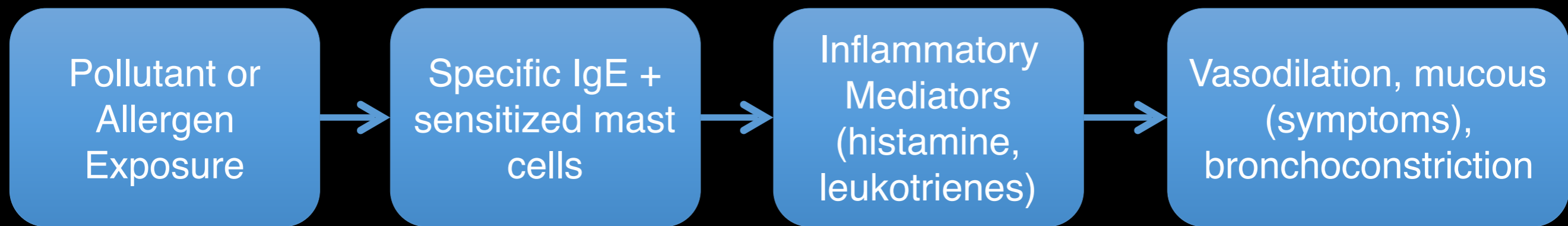
- ① Focus the Question
- ② Define the Measurement
- ③ Manage the process

Focus The Question

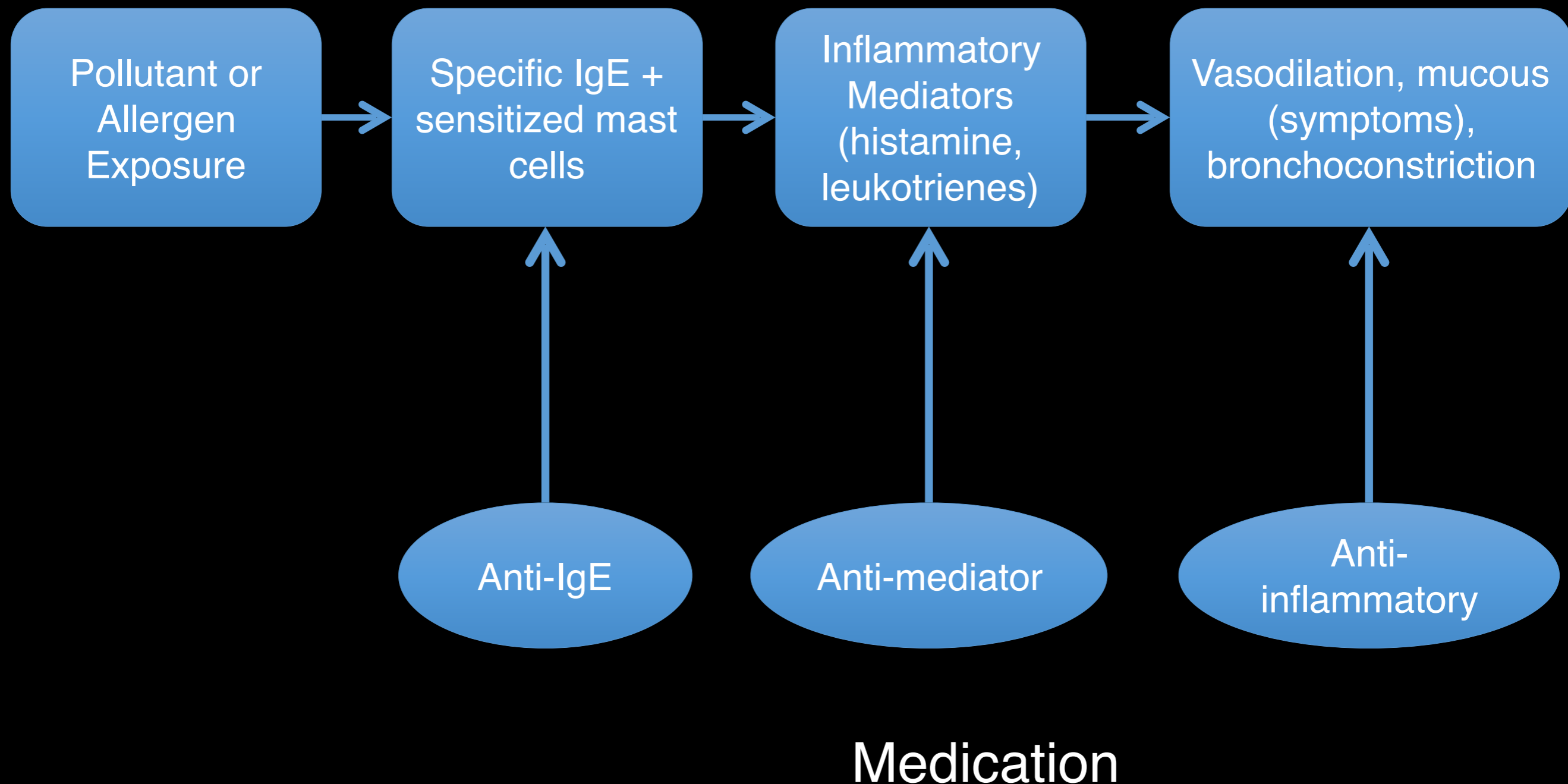
Inner-city Childhood Asthma

- Chronic inflammatory disorder of the airways
- Inflammation associated with (1) airways hyper-responsiveness; (2) airflow limitation (at least partially reversible); (3) respiratory symptoms (wheeze, cough)
- Airway inflammation can be present even in mild disease
- Racial/ethnic minorities often comprise majority of residents in inner-cities
- Asthma prevalence rates 25-28% in inner-cities

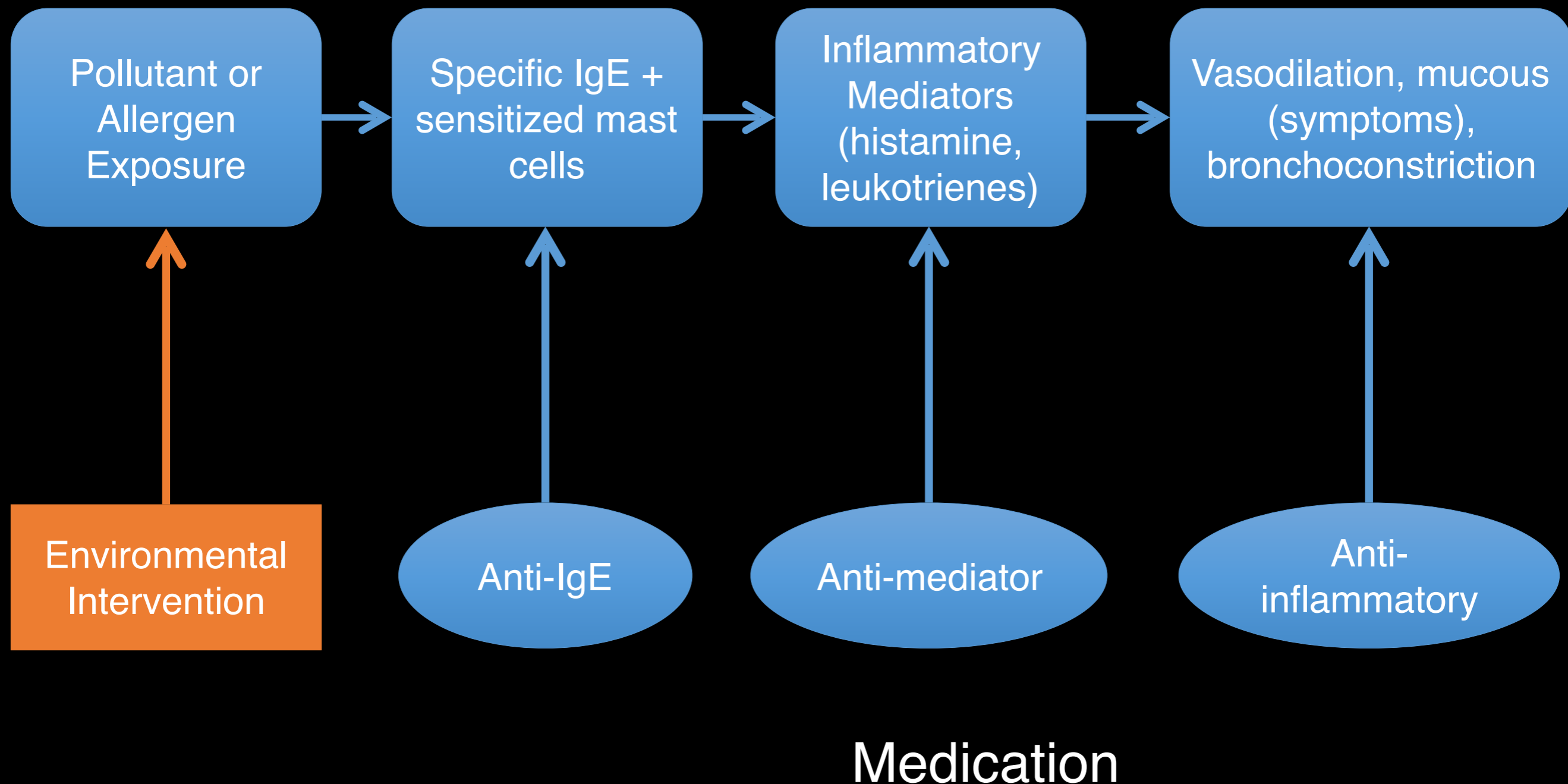
Environmental Intervention and (Allergic) Asthma Pathogenesis



Environmental Intervention and (Allergic) Asthma Pathogenesis



Environmental Intervention and (Allergic) Asthma Pathogenesis

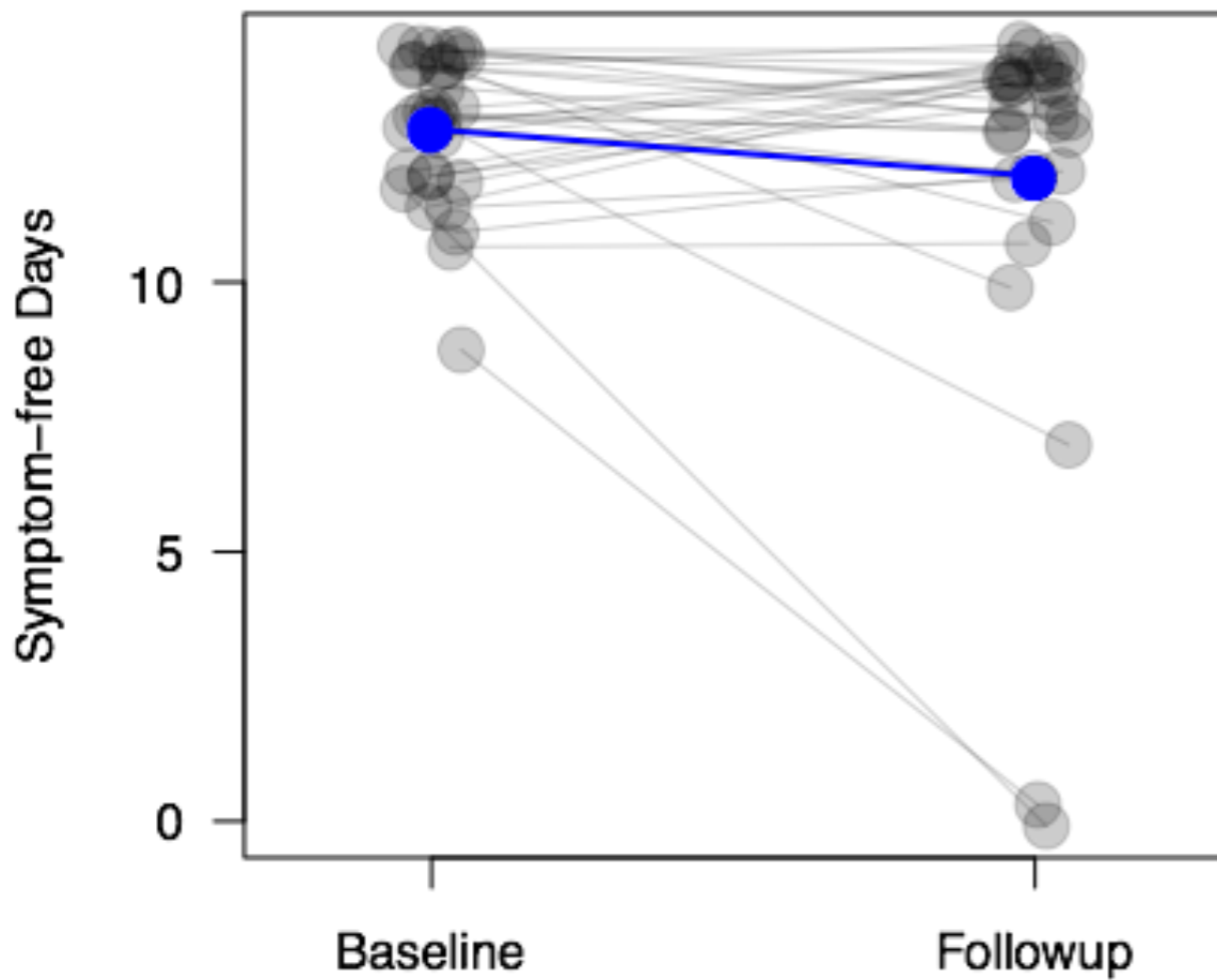


PREACH Study

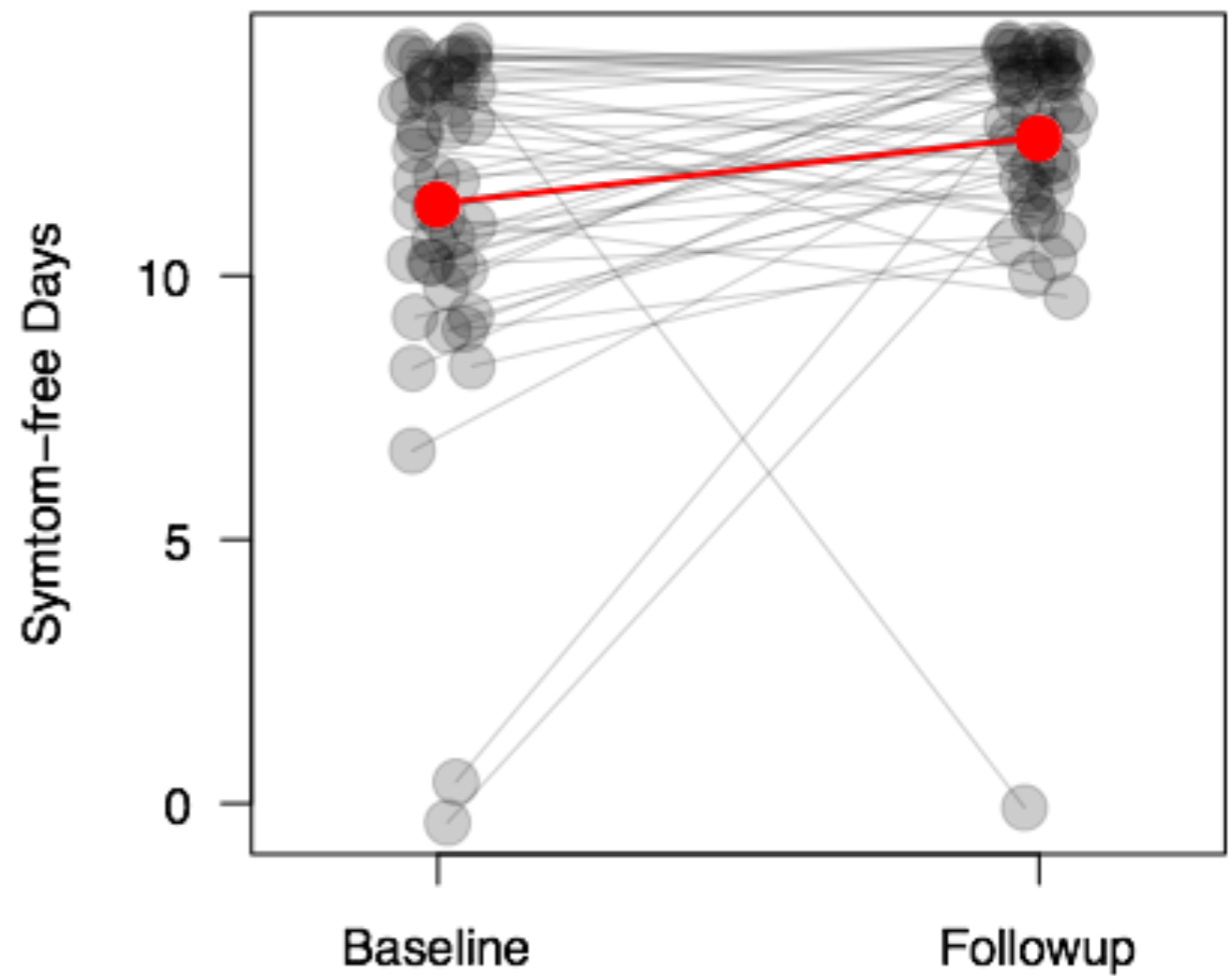
- Randomized intervention in homes in East Baltimore to lower indoor PM levels
- Two groups: **Control, Air Cleaner**
- Baseline and 6-month clinic and home visit
- 126 children 6-12 yrs old with asthma enrolled
- Homes had to have a smoker (> 5 cigs/day) living there at least 4 days/week
- **Goal:** Decrease PM_{2.5} and increase symptom-free days

PREACH Results (outcome)

Control



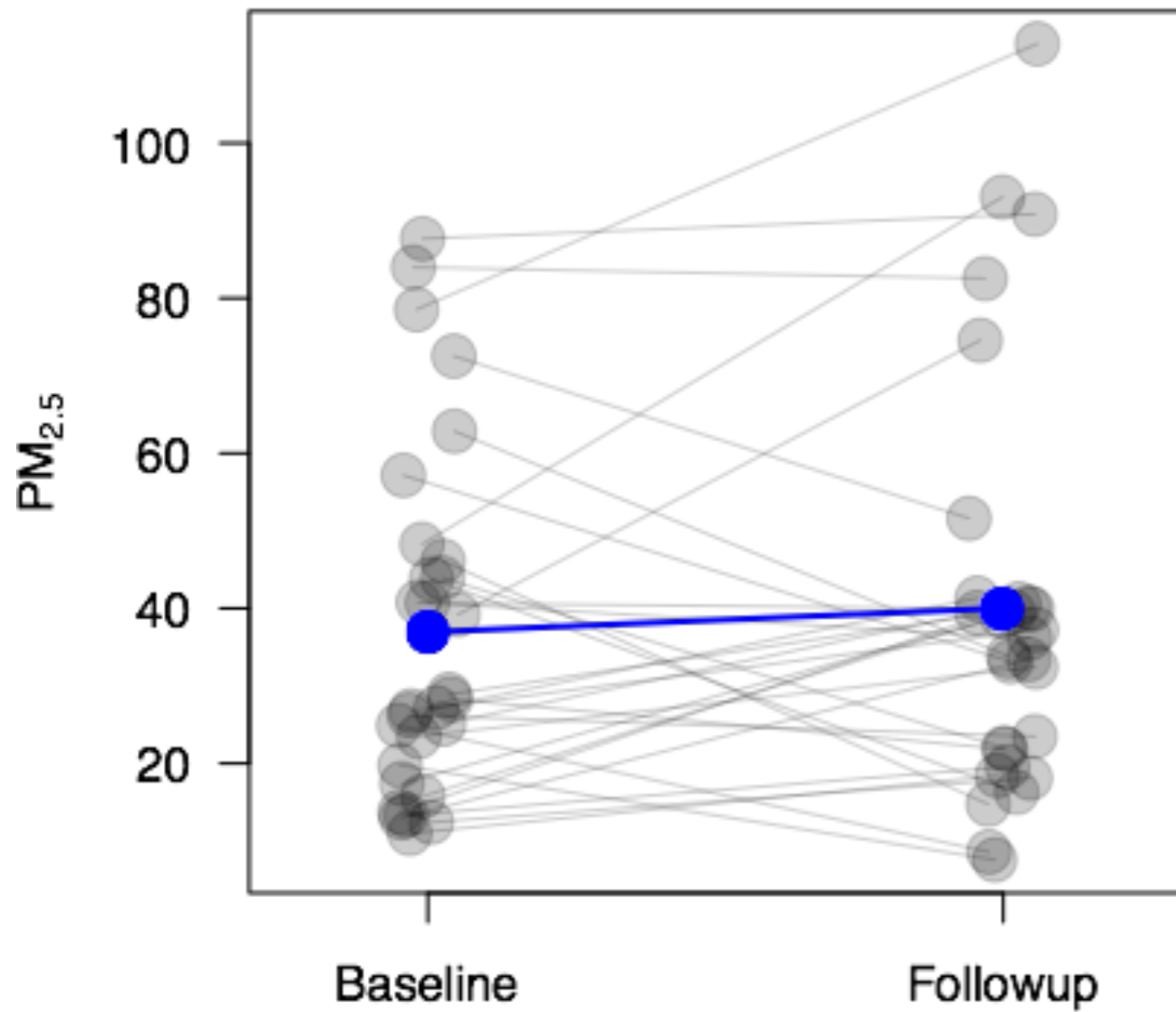
Air Cleaner



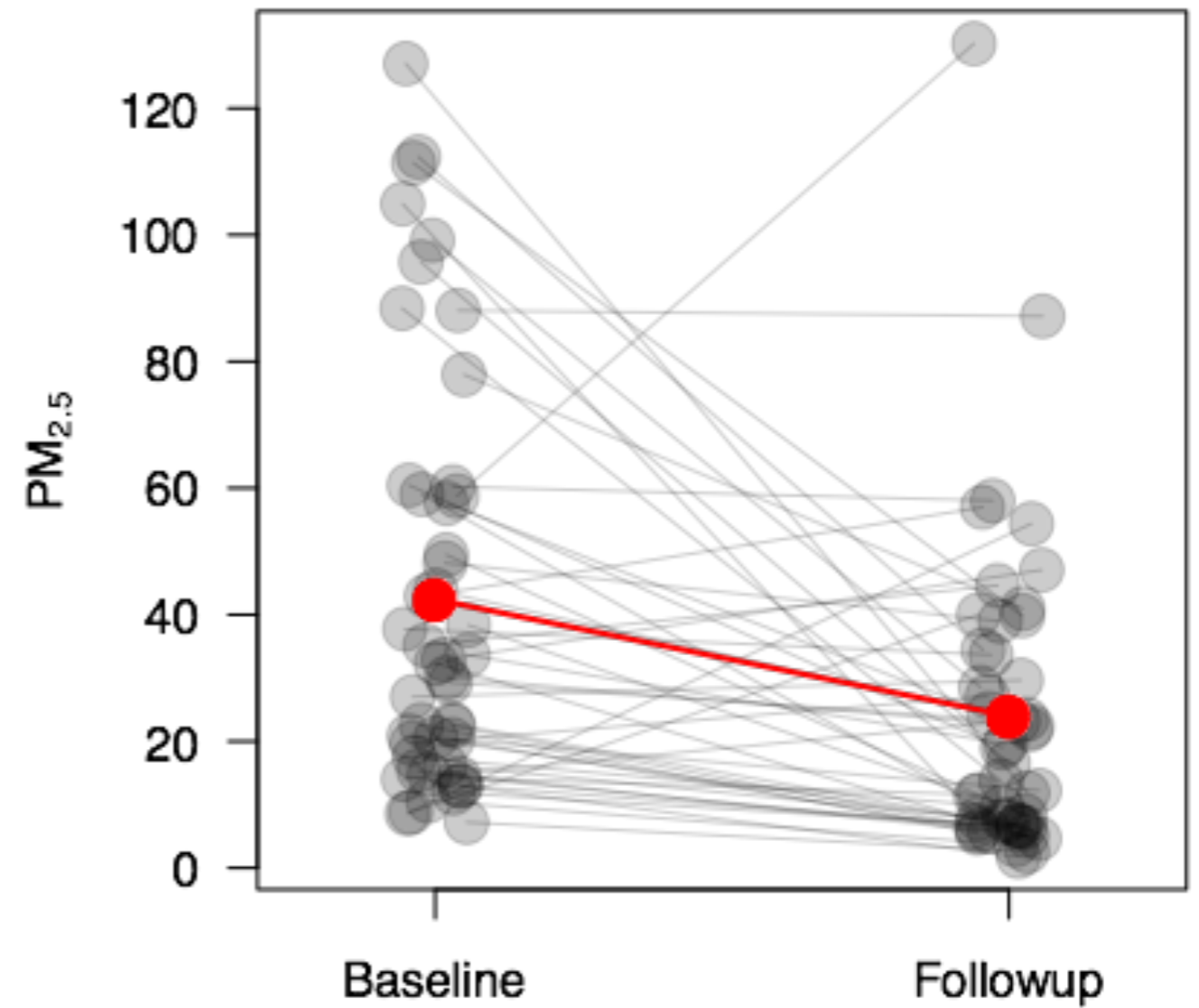
Improvement of 1.9
symptom-free days

PREACH Results (PM_{2.5})

Control



Air Cleaner



Bayesian Mixture Models

$$f(y | S = 1, D = 1) = \frac{\pi_a}{\pi_a + \pi_c} \phi(y | \mu_{a1}, \sigma_a) + \frac{\pi_c}{\pi_a + \pi_c} \phi(y | \mu_{c1}, \sigma_c)$$

$$f(y | S = 1, D = 0) = \frac{\pi_n}{\pi_n + \pi_d} \phi(y | \mu_{n1}, \sigma_n) + \frac{\pi_d}{\pi_n + \pi_d} \phi(y | \mu_{d1}, \sigma_d)$$

$$f(y | S = 0, D = 1) = \frac{\pi_a}{\pi_a + \pi_d} \phi(y | \mu_{a0}, \sigma_a) + \frac{\pi_d}{\pi_a + \pi_d} \phi(y | \mu_{d0}, \sigma_d)$$

$$f(y | S = 0, D = 0) = \frac{\pi_n}{\pi_n + \pi_c} \phi(y | \mu_{n0}, \sigma_n) + \frac{\pi_c}{\pi_n + \pi_c} \phi(y | \mu_{c0}, \sigma_c)$$

Comparison of Model Estimates: Change in Symptom-Free Days

Model	Always-taker	Never-taker	Complier
1			5.2 (-0.1, 11.8)
2	-0.3 (-1.4, 0.9)		5.5 (0.4, 13.3)
3		3.0 (-2.5, 10.2)	4.1 (0.1, 10.8)
Original		1.9 (0.2, 3.6)	

Define the Measurement

Thing I
Want to
Measure

?

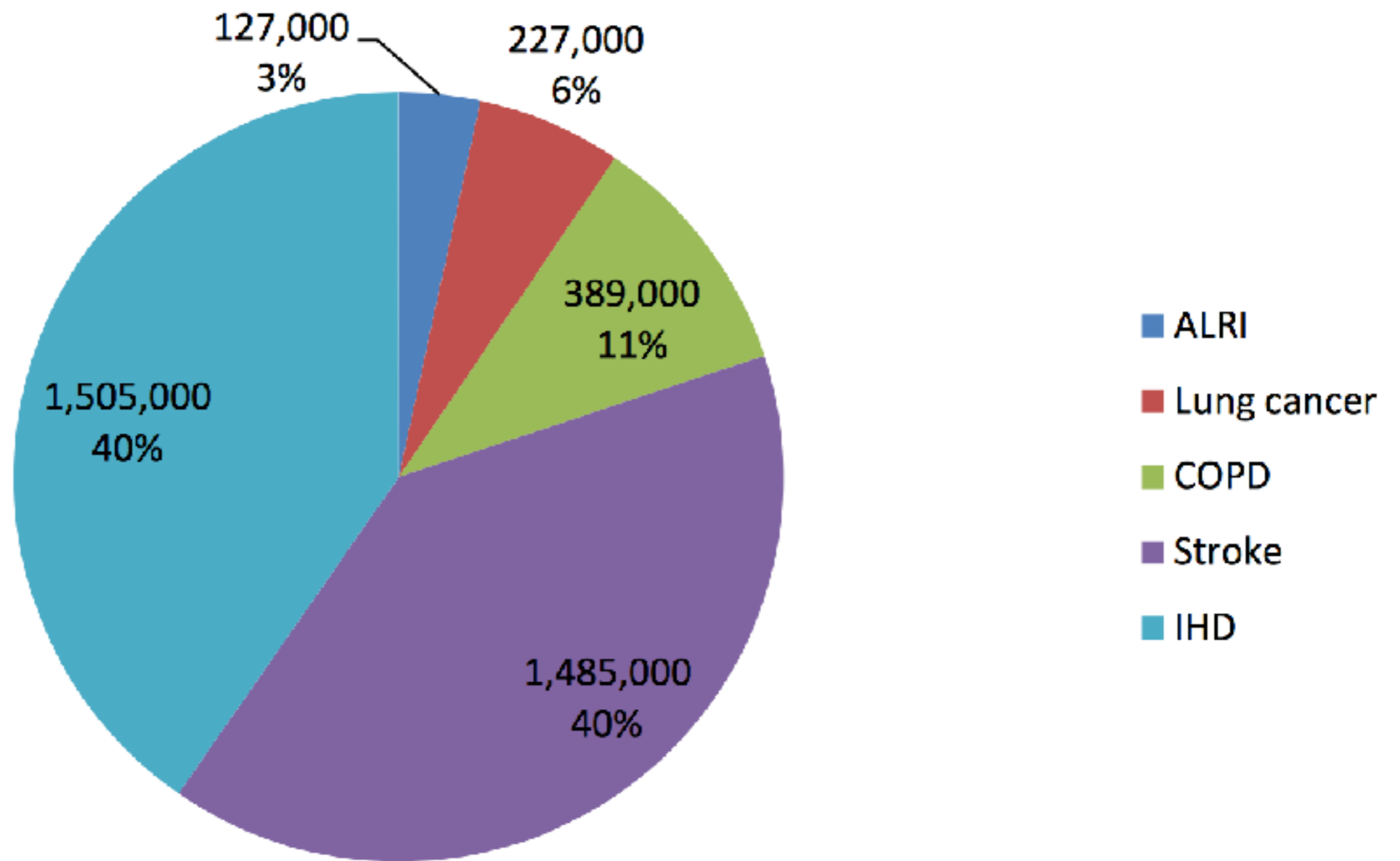
Thing I
Measured





Air Pollution

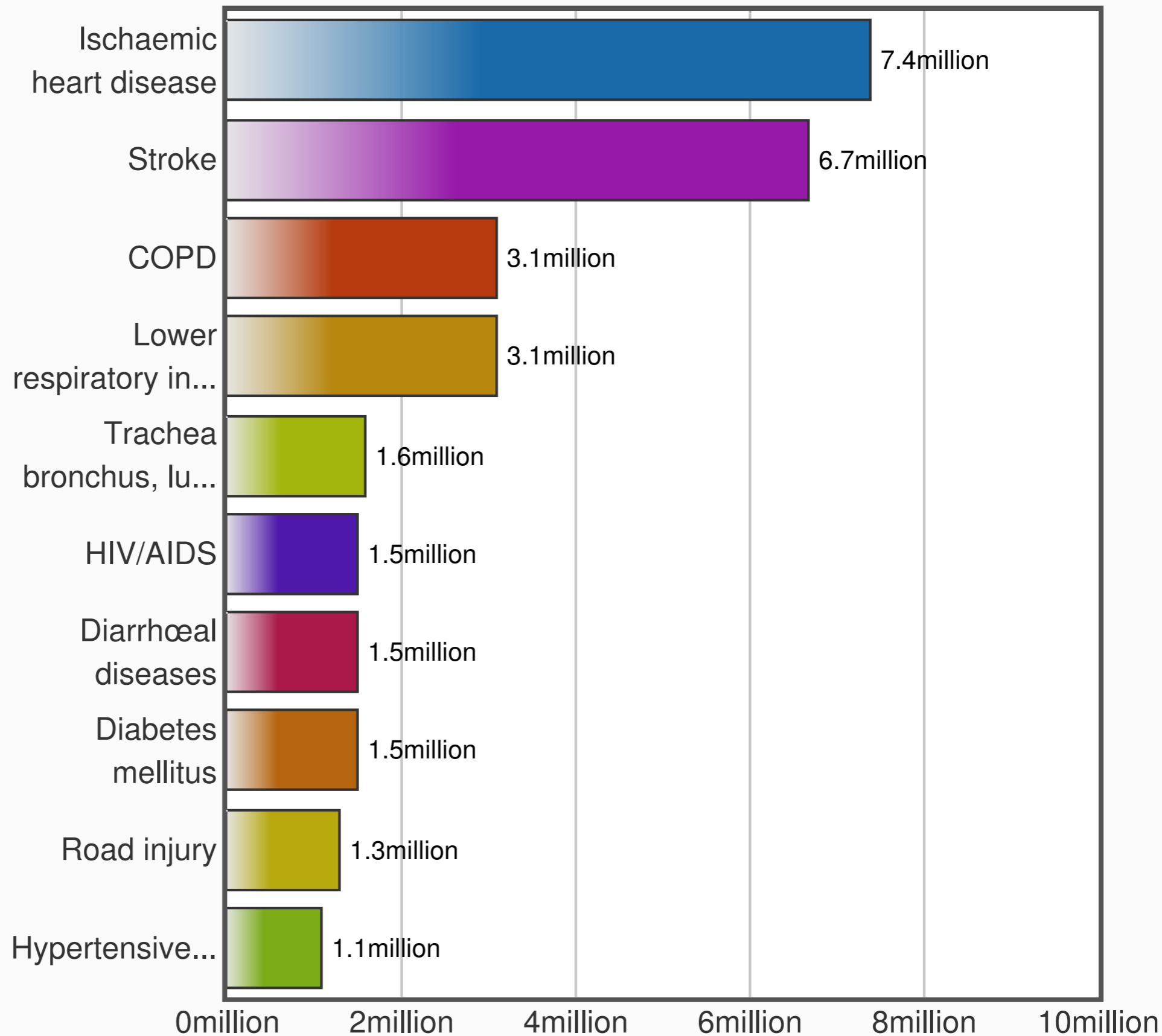
The Landscape



3.7 million deaths due to ambient air pollution (2012)

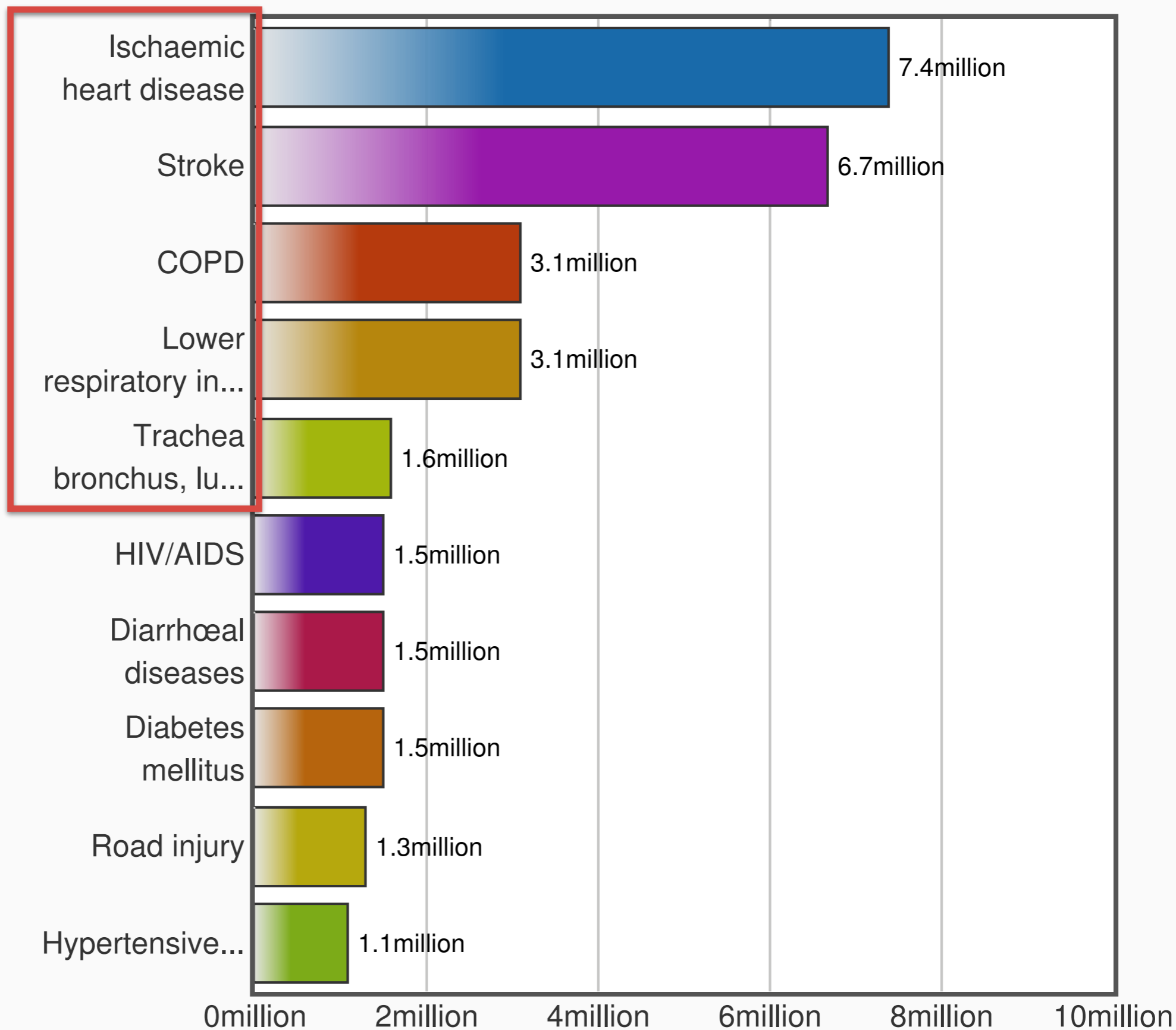
The Landscape

The 10 leading causes of death in the world 2012



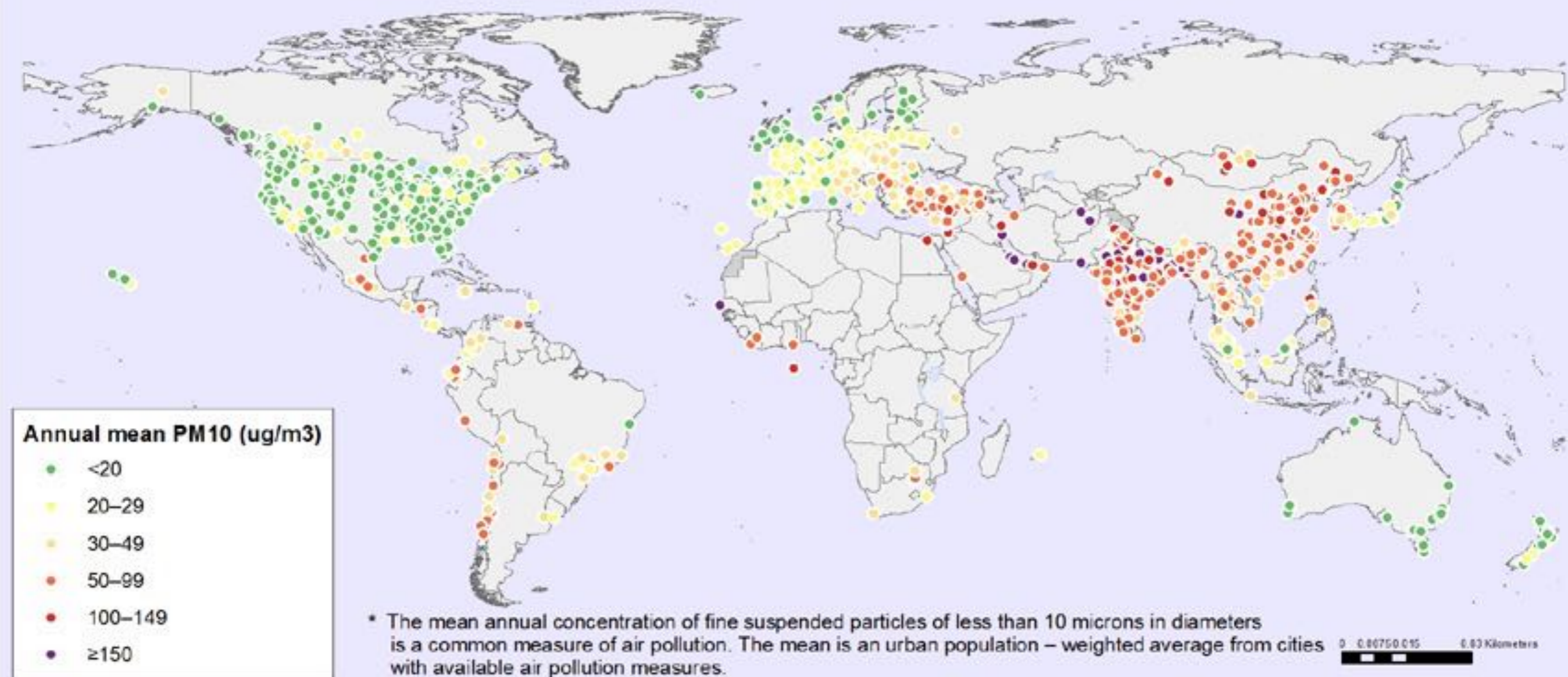
The Landscape

The 10 leading causes of death in the world 2012



The Global Landscape

Exposure to particulate matter with an aerodynamic diameter of 10 μm or less (PM10) in 1600 urban areas*, 2008–2013



The boundaries and names shown and the designations used on this map do not imply the expression of any opinion whatsoever on the part of the World Health Organization concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. Dotted and dashed lines on maps represent approximate border lines for which there may not yet be full agreement.

Data Source: World Health Organization
Map Production: Health Statistics and Information Systems (HSI)
World Health Organization



© WHO 2014. All rights reserved

What is an air pollution study?

Beijing, August 18, 2011



What is an air pollution study?

Beijing, December 5, 2011



Beijing, August 18, 2011



What is an air pollution study?

Beijing



What is an air pollution study?

Shanghai



Beijing





Air Pollution Monitoring



Air Pollution Monitoring

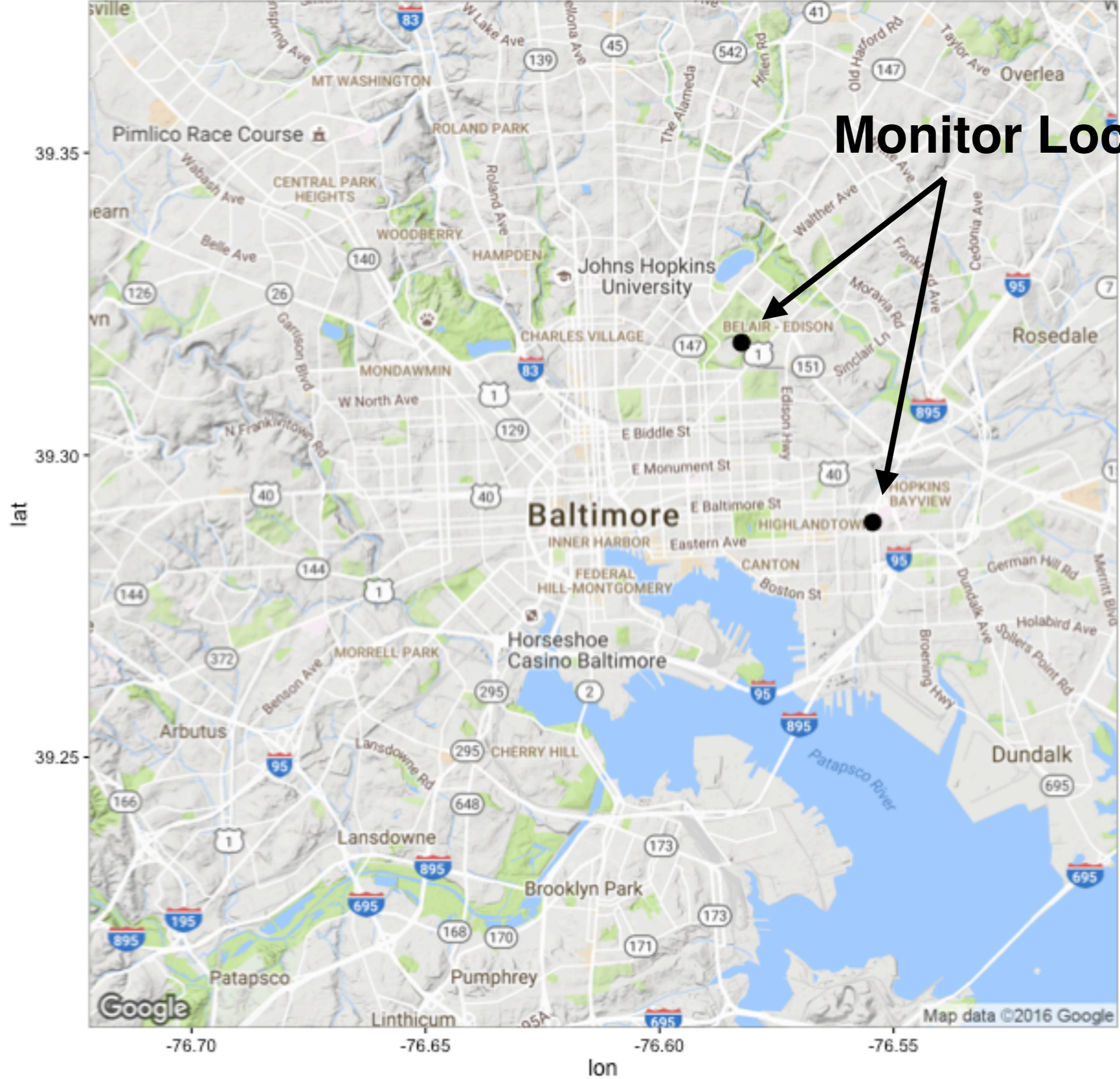


BeijingAir ✓

@BeijingAir

MetOne BAM 1020 and Ecotech EC9810 monitors, reporting PM2.5 and ozone readings. Format for each: pollutant type; concentration; AQI; definition.

📍 Chaoyang District, BEIJING



Monitor Locations

lat

lon

Google

Map data ©2016 Google

Spatial—Temporal Model

$$w(s, t) = \mu(s, t) + z(s, t)$$

Pollutant level at
location s and time t

Fixed effects

Gaussian process with
correlation function

$$\rho(\cdot, \cdot | \phi, \kappa)$$

County-wide block average
pollutant level

$$x_t = \frac{1}{\|A\|} \int_A w(s, t) ds$$

Predictive Distribution for \mathbf{x}_t

Joint distribution of monitor values and block average is Normal

$$\begin{pmatrix} w(\mathbf{v}_1, t) \\ \vdots \\ w(\mathbf{v}_m, t) \\ x_t \end{pmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_t \\ \mu_{x,t} \end{bmatrix}, \sigma^2 \begin{bmatrix} H_{11} & H_{12} \\ H_{21} & H_{22} \end{bmatrix} \right)$$

$$[H_{11}]_{ij} = \rho(\|\mathbf{v}_i - \mathbf{v}_j\|; \phi, \kappa)$$

$$[H_{12}]_i = \frac{1}{\|A\|} \int \rho(\|\mathbf{v}_i - \mathbf{s}\|; \phi, \kappa) d\mathbf{s}$$

$$H_{21} = H'_{12}$$

$$H_{22} = \frac{1}{\|A\|^2} \iint \rho(\|\mathbf{s} - \mathbf{s}'\|; \phi, \kappa) d\mathbf{s} d\mathbf{s}'$$

Bayesian and Plug-in Approaches

- Predictive distribution of block average given monitor values

$$x_t | \mathbf{w}_t \sim \mathcal{N}(\mu_{x,t} + H'_{12} H_{11}^{-1} (\mathbf{w}_t - \boldsymbol{\mu}_t), \sigma^2 (H_{22} - H'_{12} H_{11}^{-1} H_{12})).$$

- *Bayesian Model*: We can use MCMC to sample from

$$p(\theta, \mathbf{x} | \mathbf{y}, \mathbf{w}) \propto p(\mathbf{y} | \theta, \mathbf{x}, \mathbf{w}) r(\mathbf{x} | \mathbf{w}) \pi(\theta)$$

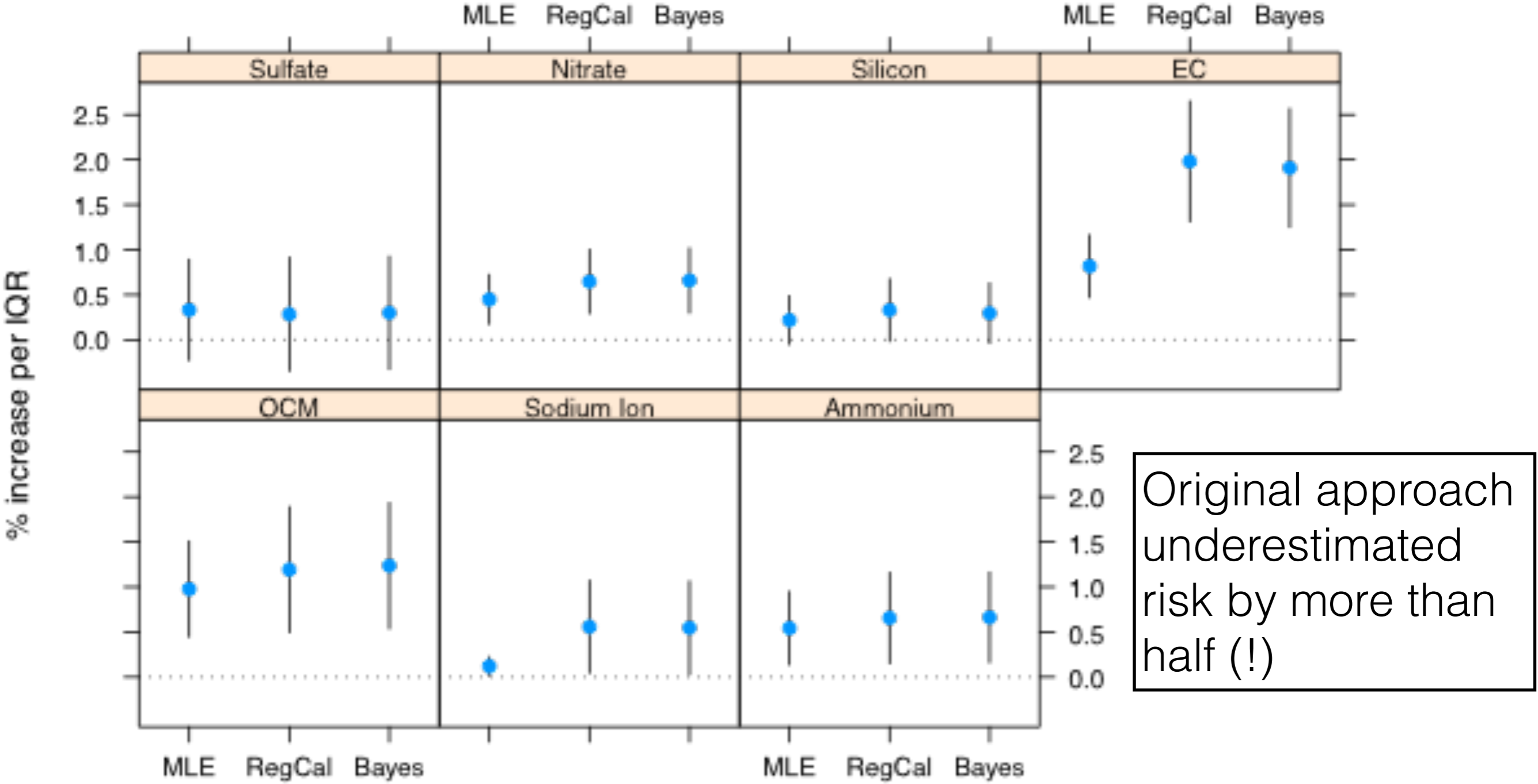
Poisson likelihood
(time series model)



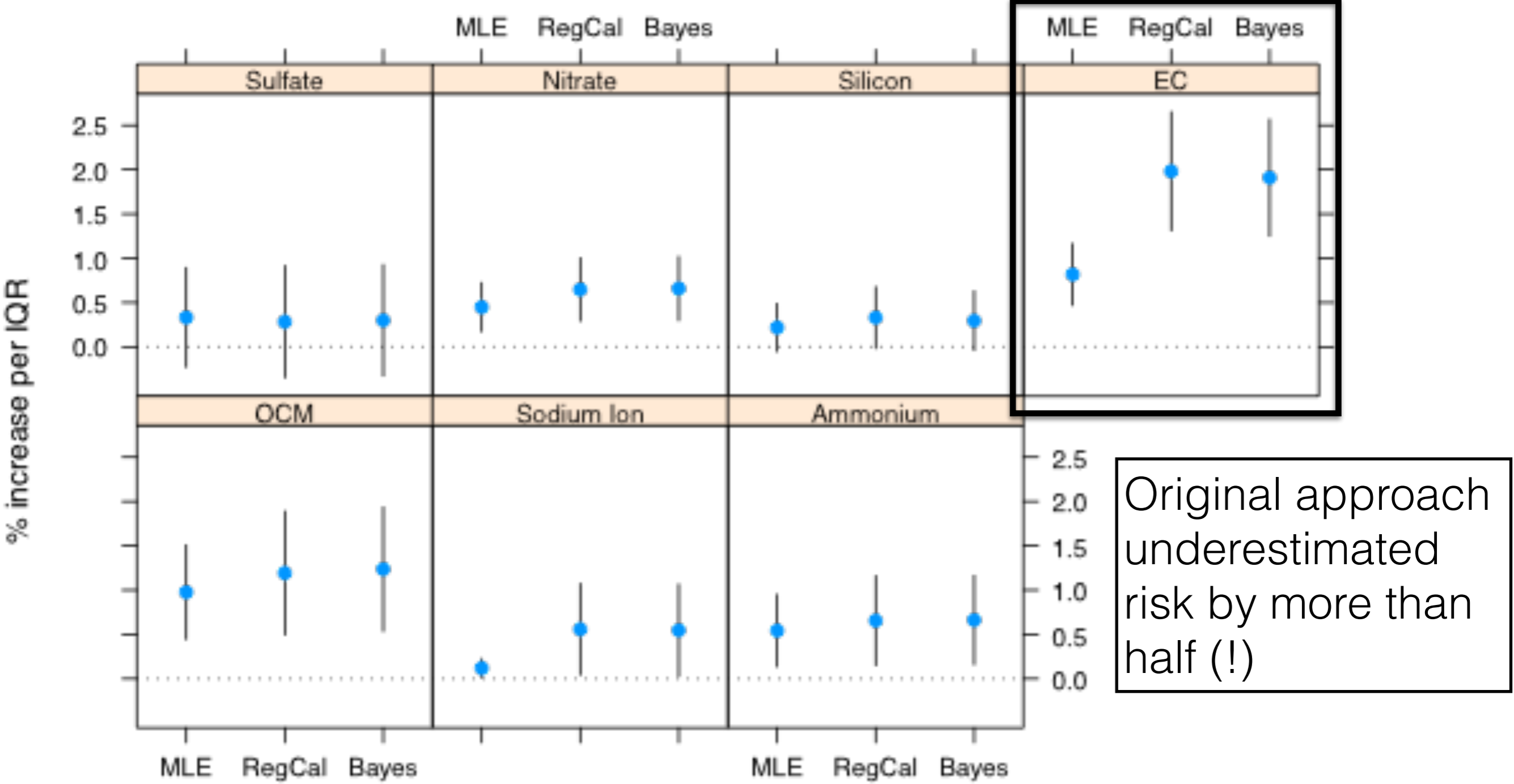
Predictive model



Combined Estimates Across 20 Counties

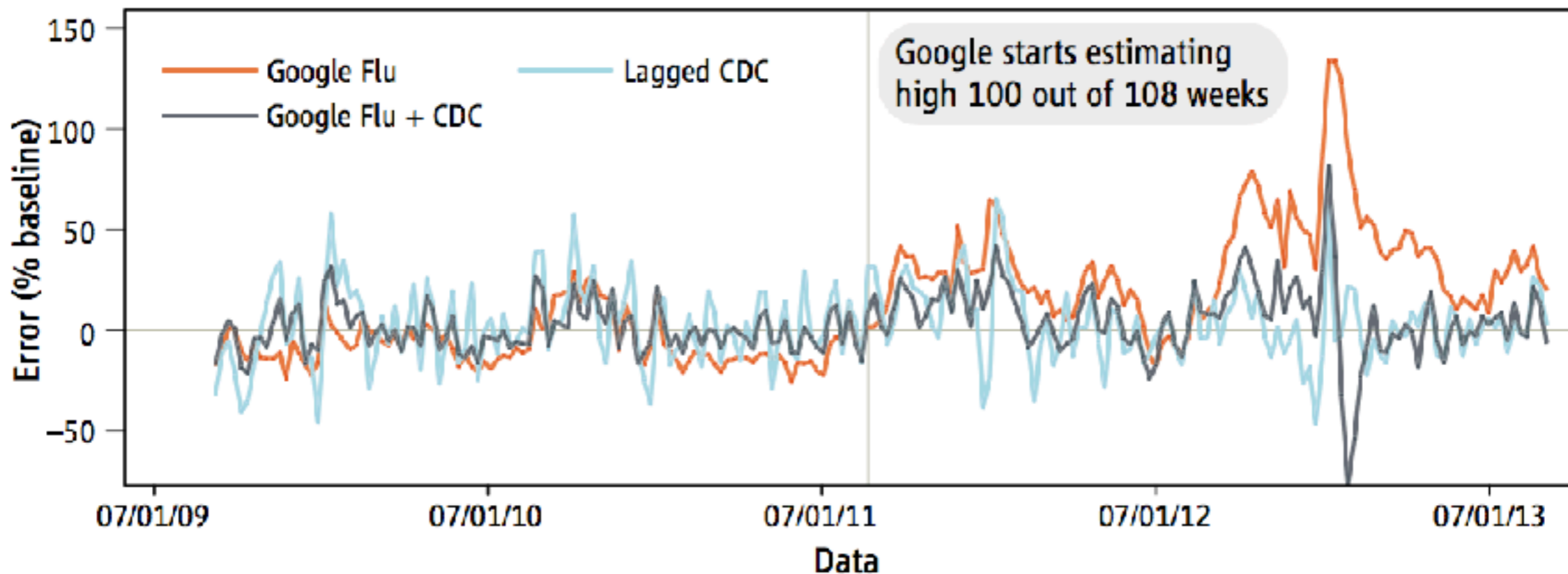
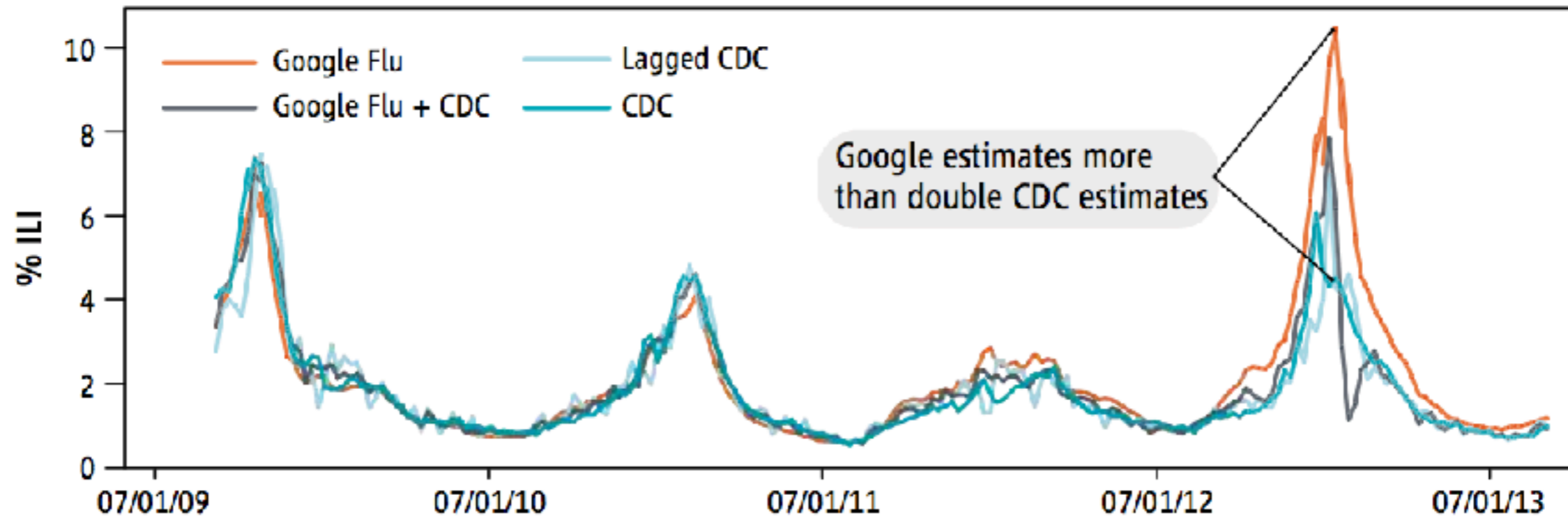


Combined Estimates Across 20 Counties



What could
go wrong?
~~~~~

# Parable of Google Flu Trends



Manage the Process



# Parable of Personalized Medicine

ARTICLES

nature  
medicine

## Genomic signatures to guide the use of chemotherapeutics

Anil Potti<sup>1,2</sup>, Holly K Dressman<sup>1,3</sup>, Andrea Bild<sup>1,3</sup>, Richard F Riedel<sup>1,2</sup>, Gina Chan<sup>4</sup>, Robyn Sayer<sup>4</sup>,  
Janiel Cragun<sup>4</sup>, Hope Cottrill<sup>4</sup>, Michael J Kelley<sup>2</sup>, Rebecca Petersen<sup>5</sup>, David Harpole<sup>5</sup>, Jeffrey Marks<sup>5</sup>,  
Andrew Berchuck<sup>1,6</sup>, Geoffrey S Ginsburg<sup>1,2</sup>, Phillip Febbo<sup>1-3</sup>, Johnathan Lancaster<sup>4</sup> &  
Joseph R Nevins<sup>1-3</sup>

# Parable of Personalized Medicine

ARTICLES

• Retracted •

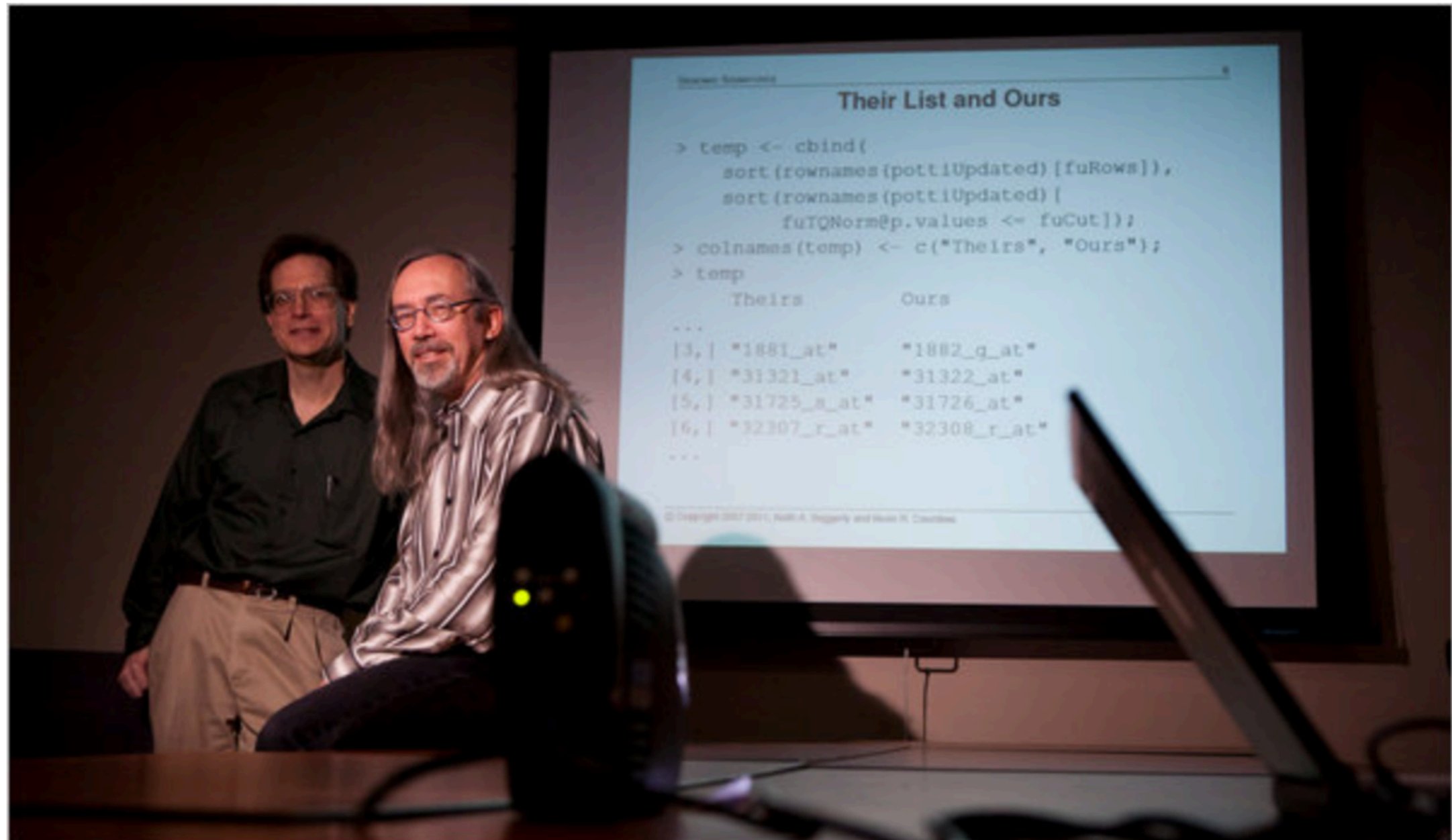
nature  
medicine

## Genomic signatures to guide the use of chemotherapeutics

Anil Potti<sup>1,2</sup>, Holly K Dressman<sup>1,3</sup>, Andrea Bild<sup>1,3</sup>, Richard F Riedel<sup>1,2</sup>, Gina Chan<sup>4</sup>, Robyn Sayer<sup>4</sup>,  
Janiel Cragun<sup>4</sup>, Hope Cottrill<sup>4</sup>, Michael J Kelley<sup>2</sup>, Rebecca Petersen<sup>5</sup>, David Harpole<sup>5</sup>, Jeffrey Marks<sup>5</sup>,  
Andrew Berchuck<sup>1,6</sup>, Geoffrey S Ginsburg<sup>1,2</sup>, Phillip Febbo<sup>1-3</sup>, Johnathan Lancaster<sup>4</sup> &  
Joseph R Nevins<sup>1-3</sup>

# “Rock Star” Statisticians

How Bright Promise in Cancer Testing Fell Apart



Michael Stravato for The New York Times

# Deception at Duke

The image shows a screenshot of a video player interface. At the top, a red banner features a stopwatch and the text "60 MINUTES". Below this is a navigation menu with links: HOME, UP NEXT, 60 OVERTIME, NEWSMAKERS, POLITICS, SCIENCE, BUSINESS, and ENTERTA. The main video frame shows a man in a suit standing in front of a backdrop that reads "Deception At Duke" and "Produced By Kyra Darnton". The backdrop also contains some text, including "Six years ago, Duke University...". In the bottom left corner of the video frame, the "60 MINUTES" logo is visible. Below the video frame is a control bar with a play/pause button, a progress bar showing "0:52 / 13:46", and a "SHARE" button. Below the control bar are social media sharing options: "23 Comments", "Share this Video:", "Recommend" (473), "Tweet" (49), and "363". At the bottom of the page, the video title "Deception at Duke" is displayed, along with the date and time "February 12, 2012 4:00 PM" and a short description: "Were some cancer patients at Duke University given experimental treatments based on fabricated data? Scott Pelley reports."

HOME UP NEXT 60 OVERTIME NEWSMAKERS POLITICS SCIENCE BUSINESS ENTERTA

Deception At Duke  
Produced By Kyra Darnton

60 MINUTES

0:52 / 13:46 SHARE

23 Comments Share this Video: Recommend 473 Tweet 49 363

**Deception at Duke**  
February 12, 2012 4:00 PM  
Were some cancer patients at Duke University given experimental treatments based on fabricated data? Scott Pelley reports.

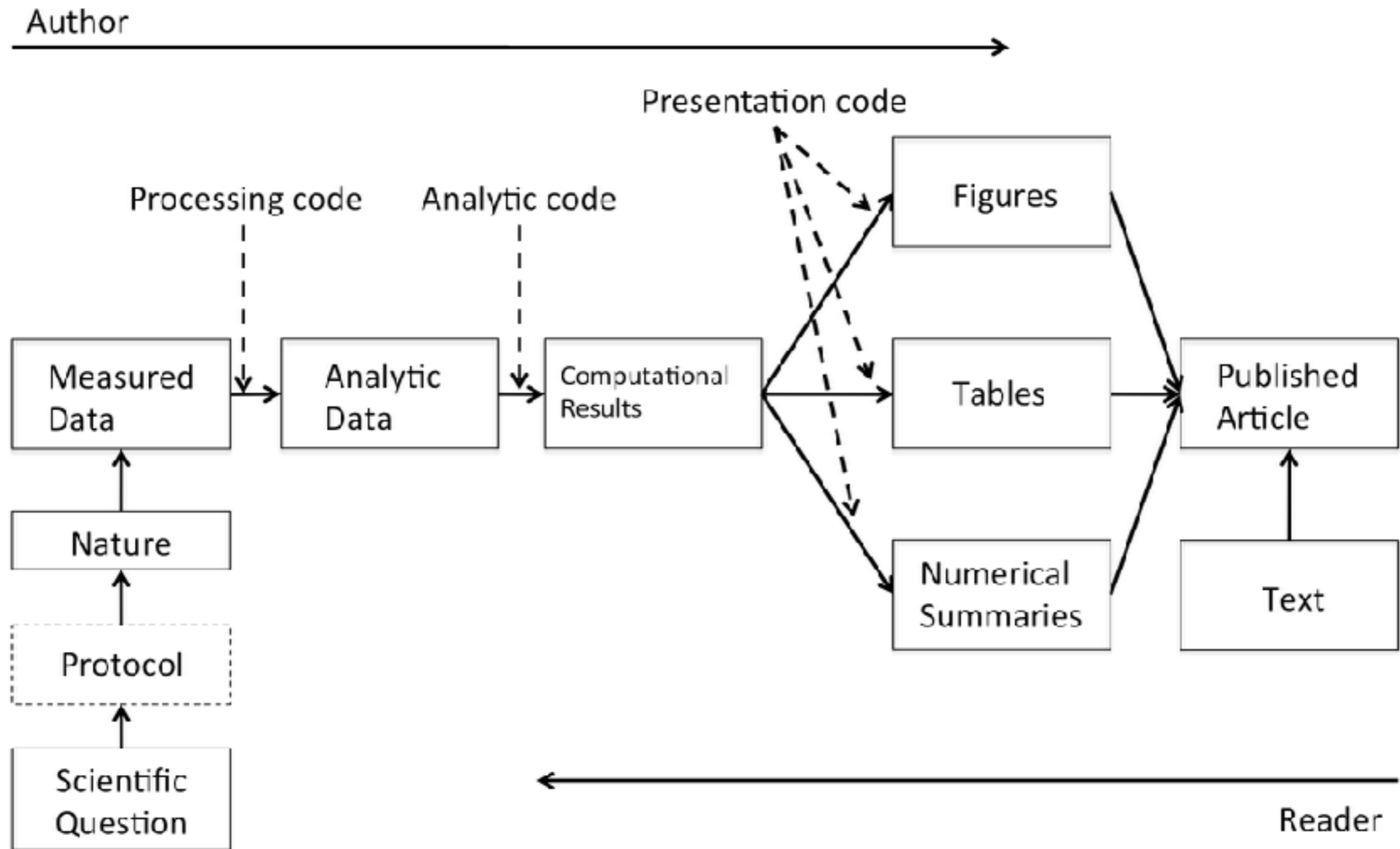


A photograph of a server room showing a massive, tangled mess of network cables. The cables are primarily blue and yellow, with some red and white ones scattered throughout. They are haphazardly bundled and draped over server racks, creating a dense, chaotic web of lines. The background shows the dark, perforated metal of server racks. In the center of the image, the text "Bad process!" is written in a bold, white, sans-serif font.

**Bad process!**



# The Data Science Process



# Institute of Medicine Report

REPORT BRIEF  MARCH 2012

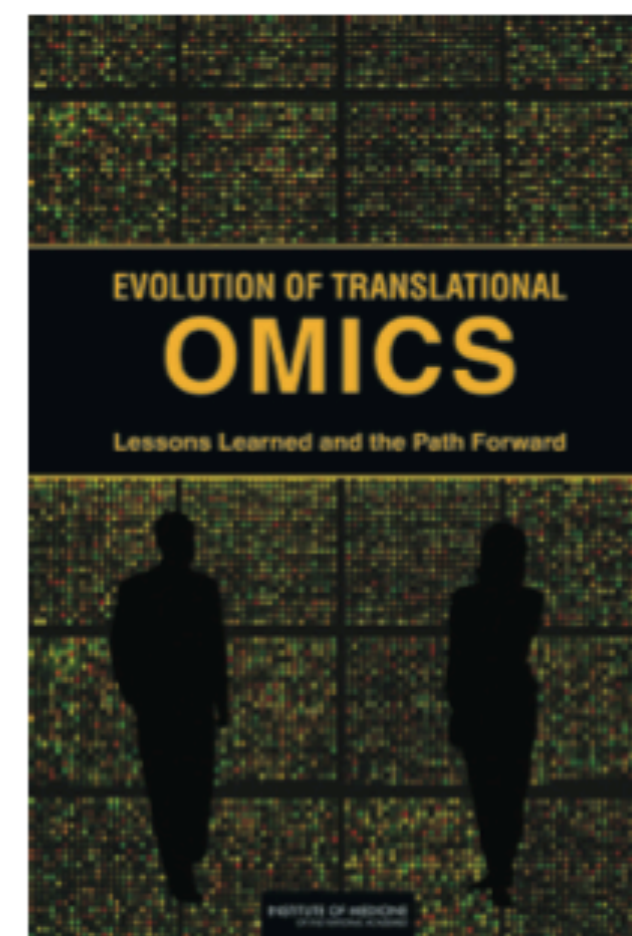
INSTITUTE OF MEDICINE  
OF THE NATIONAL ACADEMIES

Advising the nation • Improving health

For more information visit [www.iom.edu/translationalomics](http://www.iom.edu/translationalomics)

## Evolution of Translational Omics

Lessons Learned and the  
Path Forward



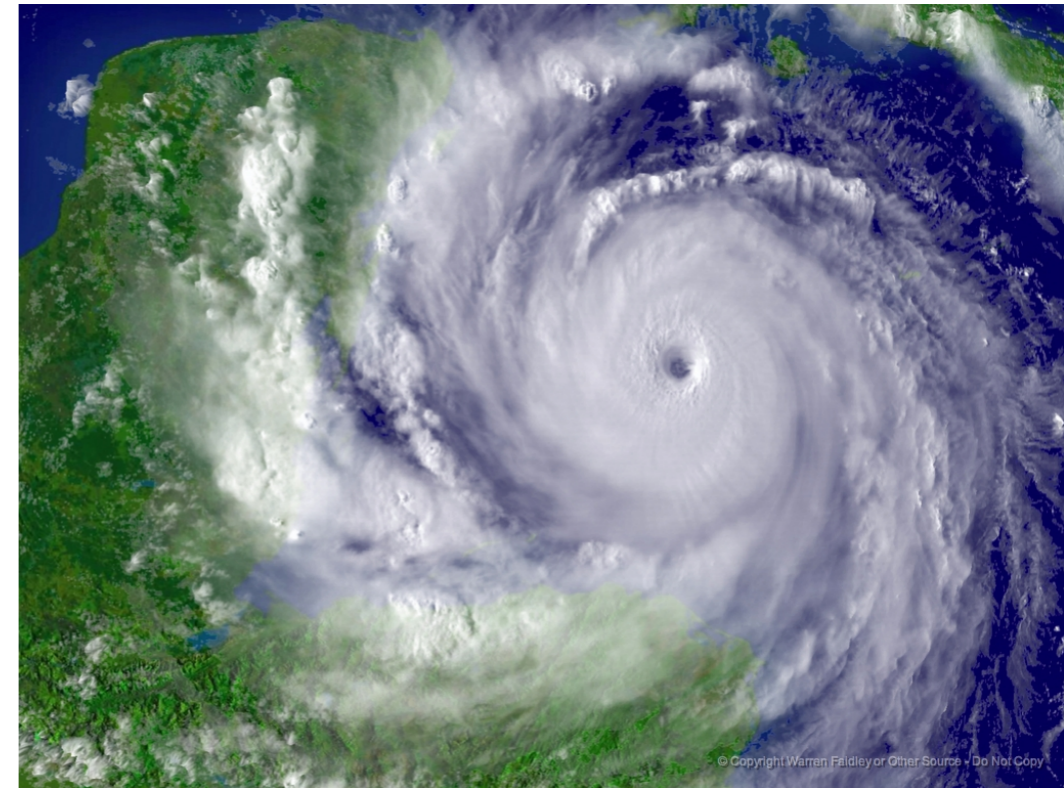
# Institute of Medicine Report

- **Data/metadata** used should be made publicly available
- The **computer code** and fully specified computational procedures used should be made available
- Ideally, the computer code that is released will **encompass all of the steps** of computational analysis, including all data preprocessing steps



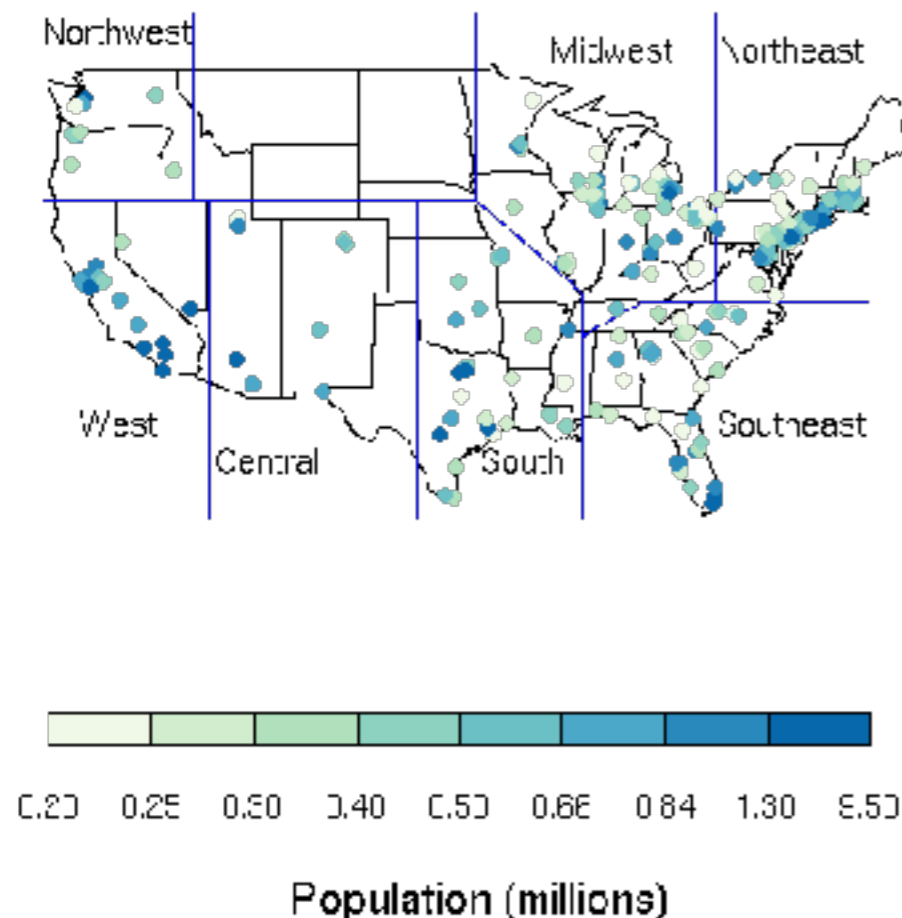
# Air Pollution and Health: A Perfect Storm?

- Estimating small health effects in the presence of much stronger signals
- Results inform substantial policy decisions and affect many stakeholders
- EPA regulations can cost billions of dollars
- Complex statistical methods are needed and subjected to intense scrutiny



# Medicare Cohort Air Pollution Study

## Medicare Air Pollution Study (MCAPS), 1999--2002



### Supplementary Materials

[FREQUENTLY ASKED QUESTIONS](#) about the MCAPS study

[Press release](#) from The Johns Hopkins Bloomberg School of Public Health

Information about the [counties used in the study](#). More information about the counties in MCAPS is also available on our [Google map](#)

Background information particulate matter from the EPA:

- [Particle Pollution in 2003](#)
- [Summary of Counties Violating the PM<sub>2.5</sub> Primary Standards](#)

### Materials for Reproducing Study Results

#### County-specific estimates

Below are tables containing county-specific estimates of the association between PM<sub>2.5</sub> and hospital admissions for various outcomes. The estimates here are the raw regression coefficients (beta) and variances (var) taken from the models described in the paper. **NOTE:** The HTML tables are large and may take a long time to load into your browser.

- Subset of models used for Table 1: [HTML](#) | [Comma separated value \(CSV\)](#) [194K]
- All models, used for Figure 2: [HTML](#) | [Comma separated value \(CSV\)](#) [1.8M]

Comma separated value (CSV) files are better suited for reading into statistical analysis programs.

Please note that the principal findings of the study are estimates of the **national** and **regional** effects of short-term exposure to PM<sub>2.5</sub>. County-specific estimates are provided solely for the purpose of reproducing those findings.



## NMMAPSdata R Package

### Current version: 0.3-4

The NMMAPSdata R package contains daily mortality, air pollution, and weather data originally assembled as part of the National Mortality, Morbidity, and Air Pollution Study (NMMAPS).

There is a [technical report](#) available which contains a brief overview of the package and contains examples of multi-city time series analysis of air pollution and mortality.

- The files [simple.R](#), [seasonal.R](#), and [tdlm.R](#) referenced in the report contain example code and functions for reproducing NMMAPS analyses.

### Database summary information

- Time frame: January 1, 1987 – December 31, 2000
- Causes of death: Total non-accidental, CVD, respiratory, pneumonia, COPD, accidental
  - Age categories: < 65, 65–74, >= 75
- Pollutants: PM<sub>10</sub>, PM<sub>2.5</sub>, CO, O<sub>3</sub>, SO<sub>2</sub>, NO<sub>2</sub>
- Weather: Temperature, dewpoint temperature, relative humidity
- Number of Cities: 108

More detailed information about the database can be found on the iHAPSS website at <http://www.ihapss.jhsph.edu/>.

### Package requirements

- R version 1.9.0 or higher.
- bzip2 compression capability. Most people will *not* have to worry about this since R comes with bzip2 compression capability by default. However, on some Unix-like systems it is possible that the version of R was compiled without it. NMMAPSdata will give an error when the package is loaded if bzip2 capability is not present.
- Approximately 380MB of disk space to store the package.

For **Unix**, **Linux**, and **Mac OS X** users, there is a source package available.

The success of a data analysis depends on the process, not the result.



The Future

# The Future is Bright!

- A tremendous infrastructure on which to build
- Cheap hardware and Moore's Law has made powerful computing available to all with the cloud
- Advanced software has abstracted complex details of data analysis
- The Internet allows analyses to be deployed to the entire world
- Volume of data is increasing dramatically
- A never-ending supply of difficult (but interesting) questions!

# The Future is Bright!

- The future will favor those trained in data science
- Problems and data are coming in too fast
- A perfect training for interdisciplinary work
- Many leadership opportunities

# Johns Hopkins Biostatistics

- Largest / oldest / best school of public health
- PhD program in Biostatistics (~10 per class)
- ScM program in Biostatistics (8-12 per class)
- Rigorous training with focus on science
- **Applications:** Environmental health, genomics, personalized medicine, medical imaging, wearable computing, clinical trials, infectious disease





[About](#) [Education](#) [Research](#) [Prototyping](#) [Partnerships](#)

# The Johns Hopkins Data Science Lab

The Data Science Lab at Johns Hopkins is about all things data science. We produce courses, develop software, prototype apps, conduct research, and generally spread the word about data science. We believe that the intelligent application of data science skills can have a profound impact across all areas of our lives.

[MEET THE TEAM](#)

[SEE EXAMPLES](#)



# Great Students





# Great Faculty

Roger Peng

Brian Caffo

Jeff Leek



Thank You!