

# Digital Epidemiology: designing machine learning approaches to combine Internet-based data sources

Case studies characterizing flu, dengue, Zika, and Ebola epidemics.



**Mauricio Santillana, PhD**

Assistant Professor, Harvard Medical School

Faculty member, Computational Health Informatics Program, Boston Children's Hospital

Associate faculty, Harvard Institute of Applied Computational Sciences



# Digital Epidemiology: designing machine learning approaches to combine Internet-based data sources

Case studies characterizing flu, dengue, Zika, and Ebola epidemics.



**Collaborators:** Sam C. Kou (Harvard Statistics), Shihao Yang (Harvard Statistics), Fred Lu (BCH), Nicholas Brooke (Break Dengue), Matt Biggerstaff (CDC), Julia Gunn (Boston Public Health Commission), Joe Conidi (Boston Public Health Commission), Michael Johansson (CDC), Nick Reich (Umass Amherst), Roni Rosenfeld (CMU), Sarah McGough (Harvard), Alessandro Vespignani (Northeastern Univ), Nathan Kutz (Univ of Washington), Elaine Nsoesie (Univ of Washington), Rumi Chunara (NYU), John Brownstein (Harvard/BCH), Leonardo C. Clemente (Inst. Politécnico Nacional, Mex), and many more



**HARVARD**  
MEDICAL SCHOOL



**HARVARD**  
School of Engineering  
and Applied Sciences

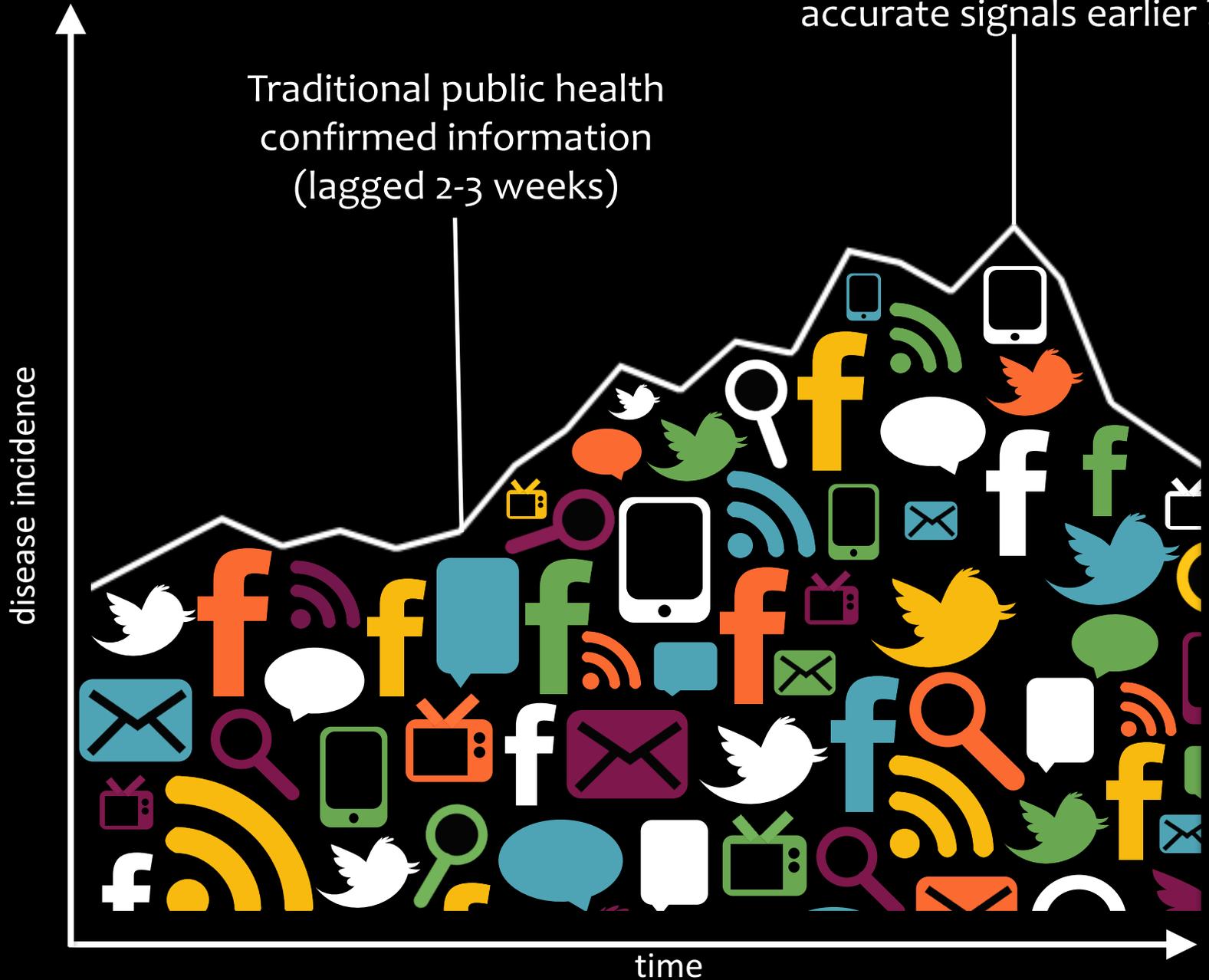


**Boston  
Children's  
Hospital**



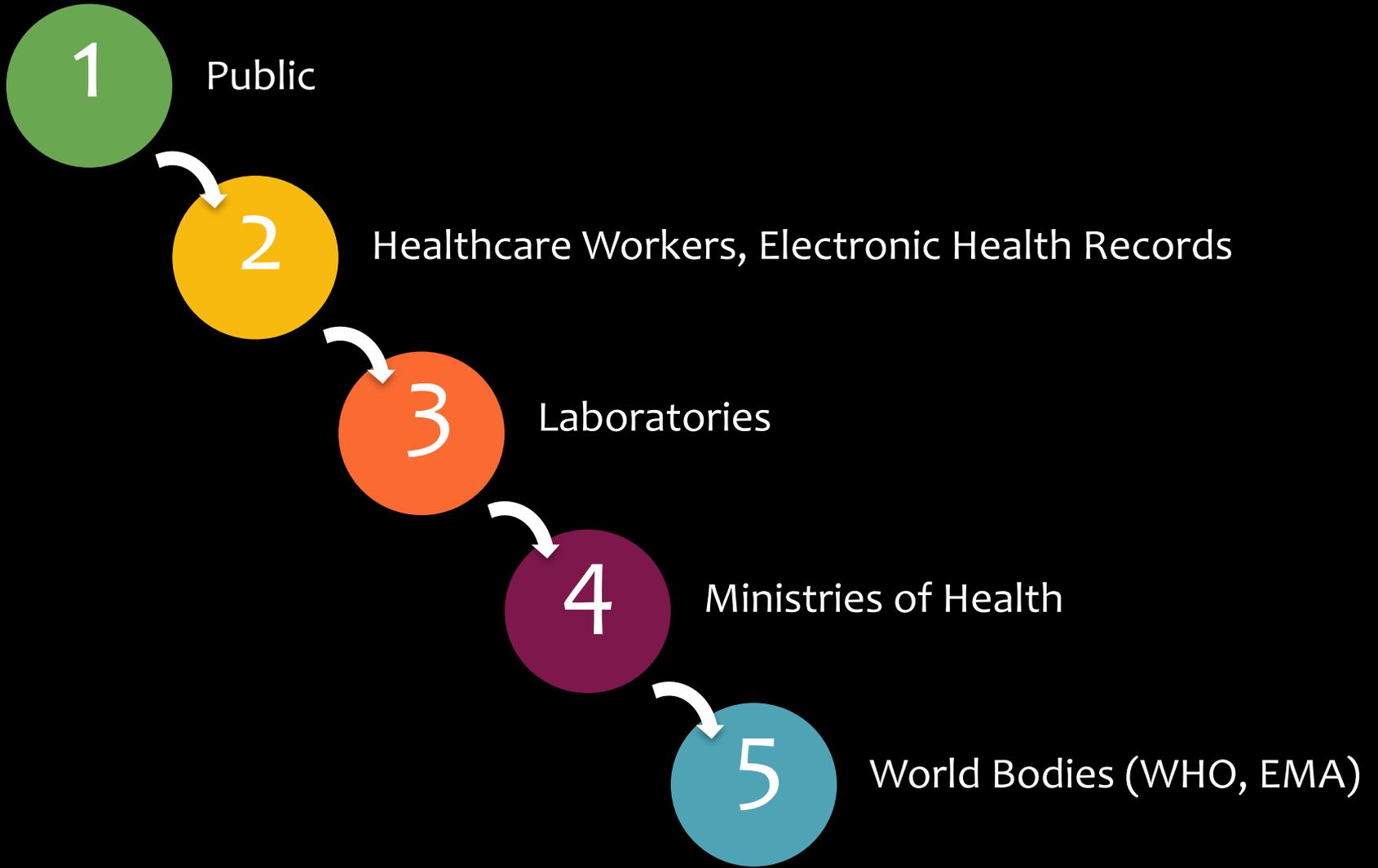
Can Digital disease tracking pick up accurate signals earlier?

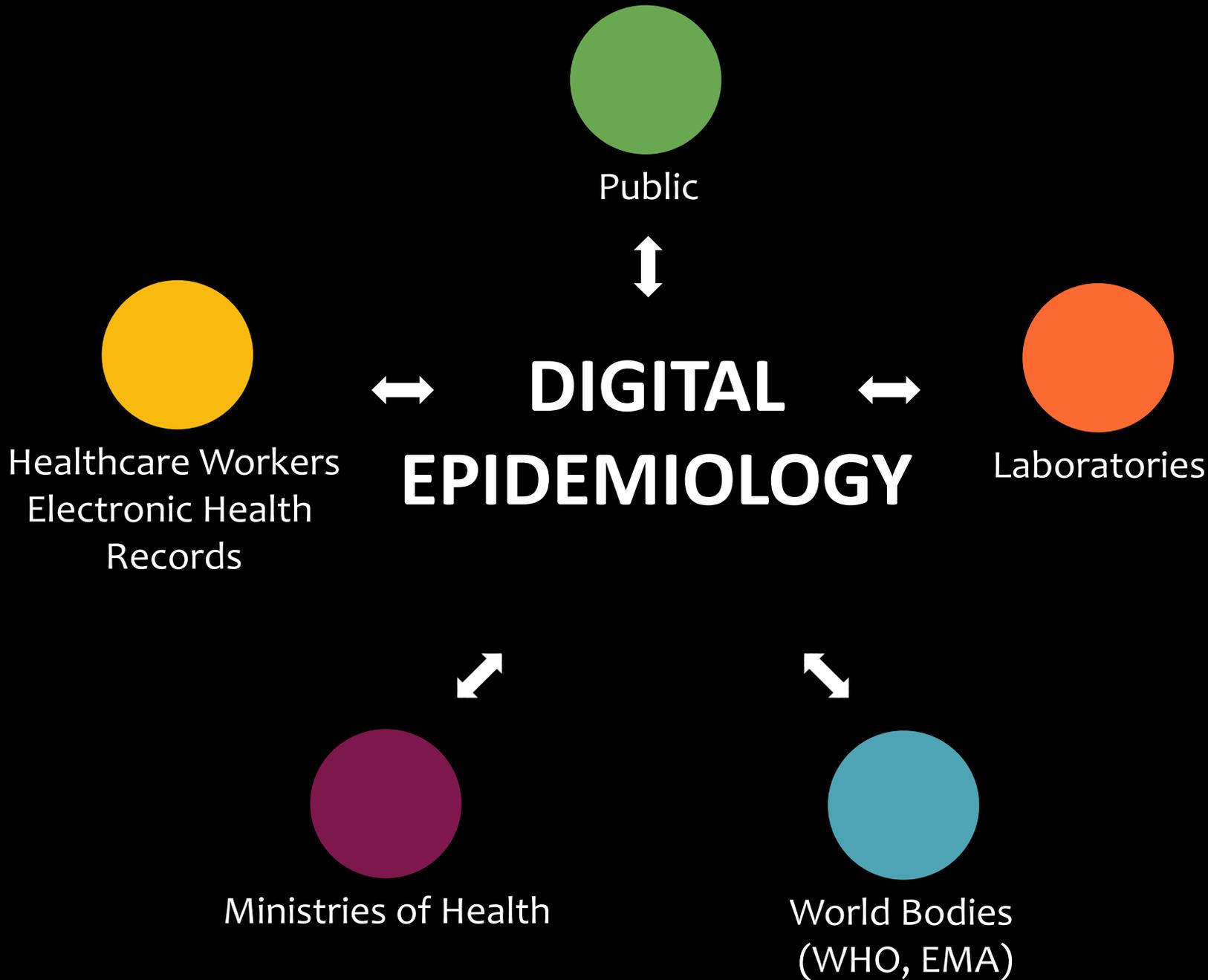
Traditional public health confirmed information (lagged 2-3 weeks)



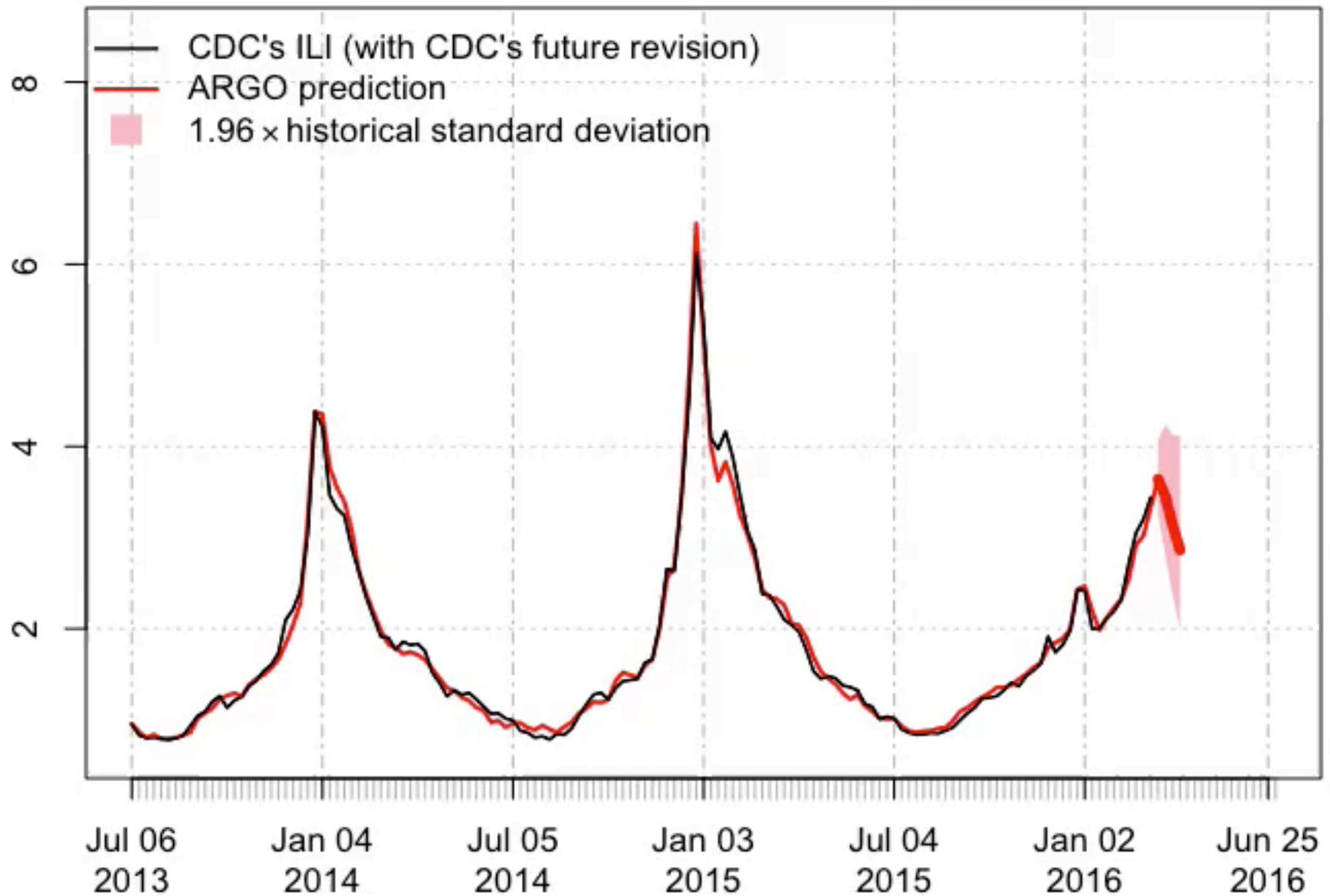


# TRADITIONAL DISEASE REPORTING



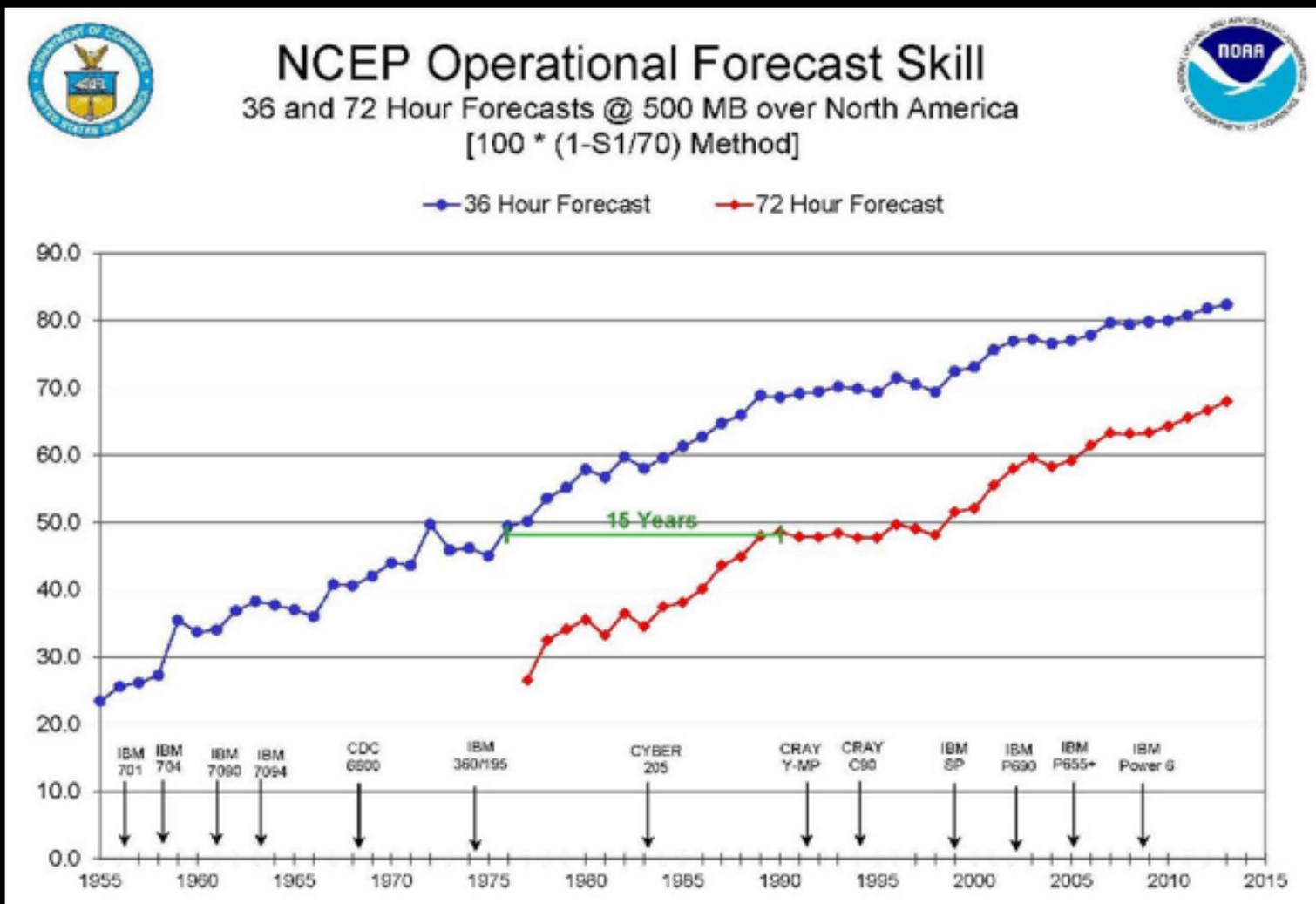


## ARGO Prediction vs. CDC's ILI

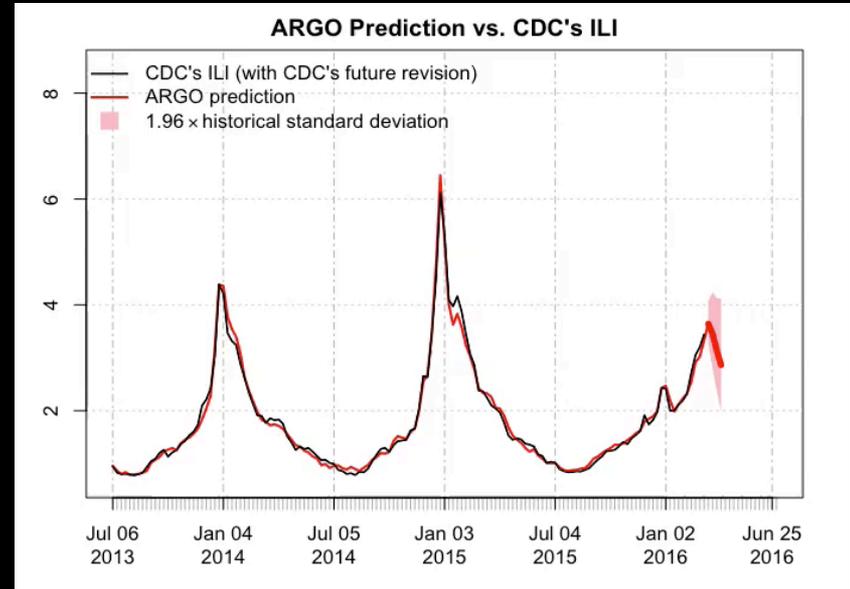
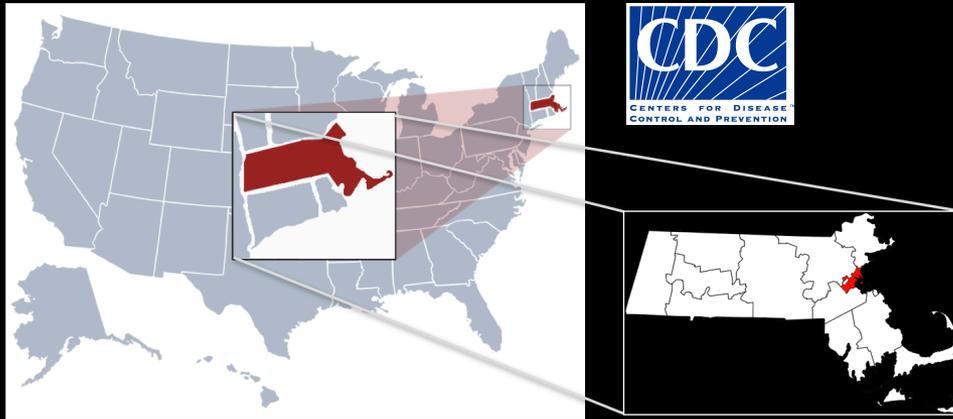


# Real-time tracking vs predictions of disease incidence/risk

## Similarities and differences with weather prediction



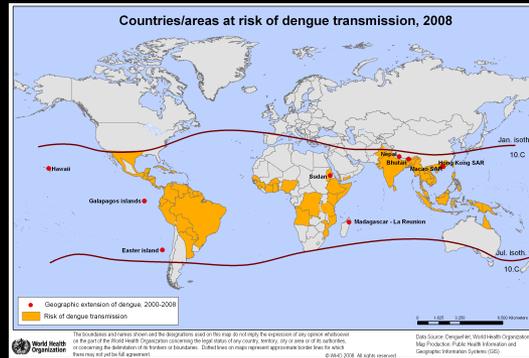
# Part 1. Previous success stories in tracking and forecasting Influenza in data-rich high-income countries: USA



1. Multiple spatial resolutions: National, multi-state, state, city-level
2. Multiple data sources (hybrid systems): traditional healthcare-based, EHR, Google, Twitter, Crowd-sourced disease surveillance.

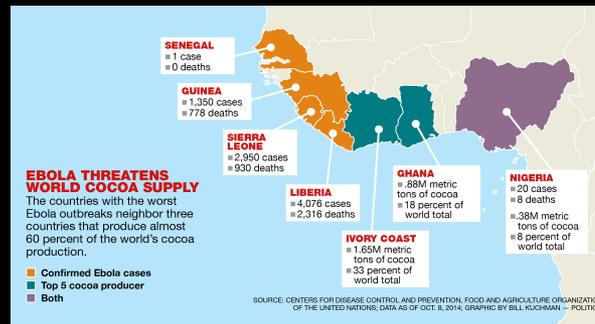
# Part 2. Success stories in tracking and forecasting Flu, Zika, Dengue, Ebola in data-poor medium- to low-income countries.

## Dengue, Zika, and Flu



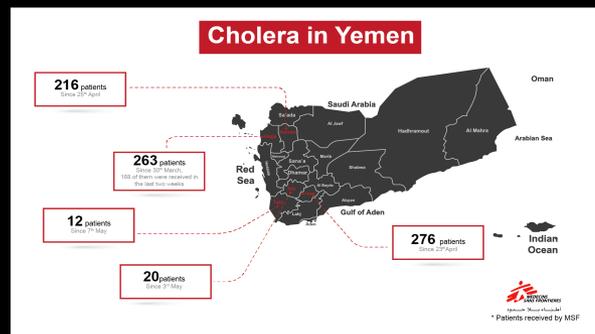
- Latin America (Flu, Zika, Dengue)
- South-east Asia (Dengue)

## Ebola



- West Africa

## Cholera



- Middle East

Seminal work by Google

The promise of big data in public health

# **GOOGLE FLU TRENDS**

Letter

Nature 457, 1012-1014 (19 February 2009) | doi:10.1038/nature07634; Received 14 August 2008; Accepted 13 November 2008; Published online 19 November 2008; [Corrected](#) 19 February 2009

Detecting influenza epidemics using search engine query data

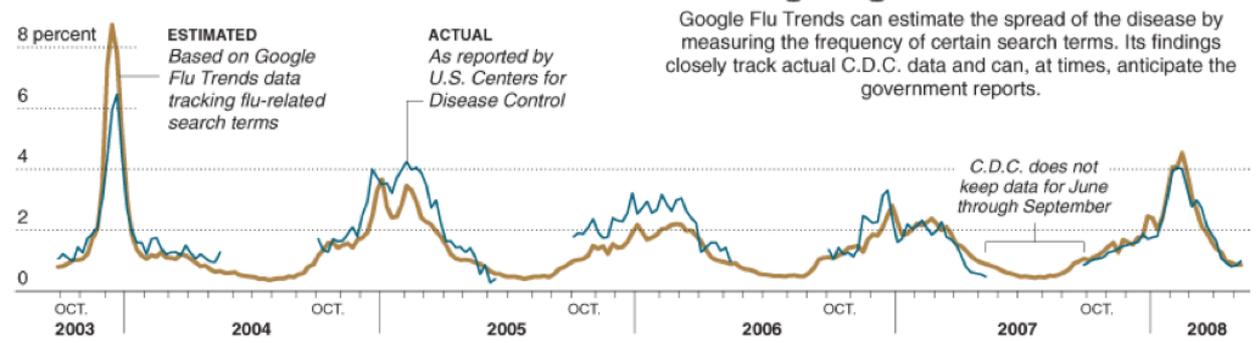
Jeremy Ginsberg<sup>1</sup>, Matthew H. Mohebbi<sup>1</sup>, Rajan S. Patel<sup>1</sup>, Lynnette Brammer<sup>2</sup>, Mark S. Smolinski<sup>1</sup> & Larry Brilliant<sup>1</sup>

- 1. Google Inc., 1600 Amphitheatre Parkway, Mountain View, California 94043, USA
- 2. Centers for Disease Control and Prevention, 1600 Clifton Road, NE, Atlanta, Georgia 30333, USA

Correspondence to: Matthew H. Mohebbi<sup>1</sup> Correspondence and requests for materials should be addressed to J.G. or M.H.M. (Email: [flutrends-support@google.com](mailto:flutrends-support@google.com)).

The New York Times

PERCENT OF HEALTH VISITS FOR FLU-LIKE SYMPTOMS Mid-Atlantic region



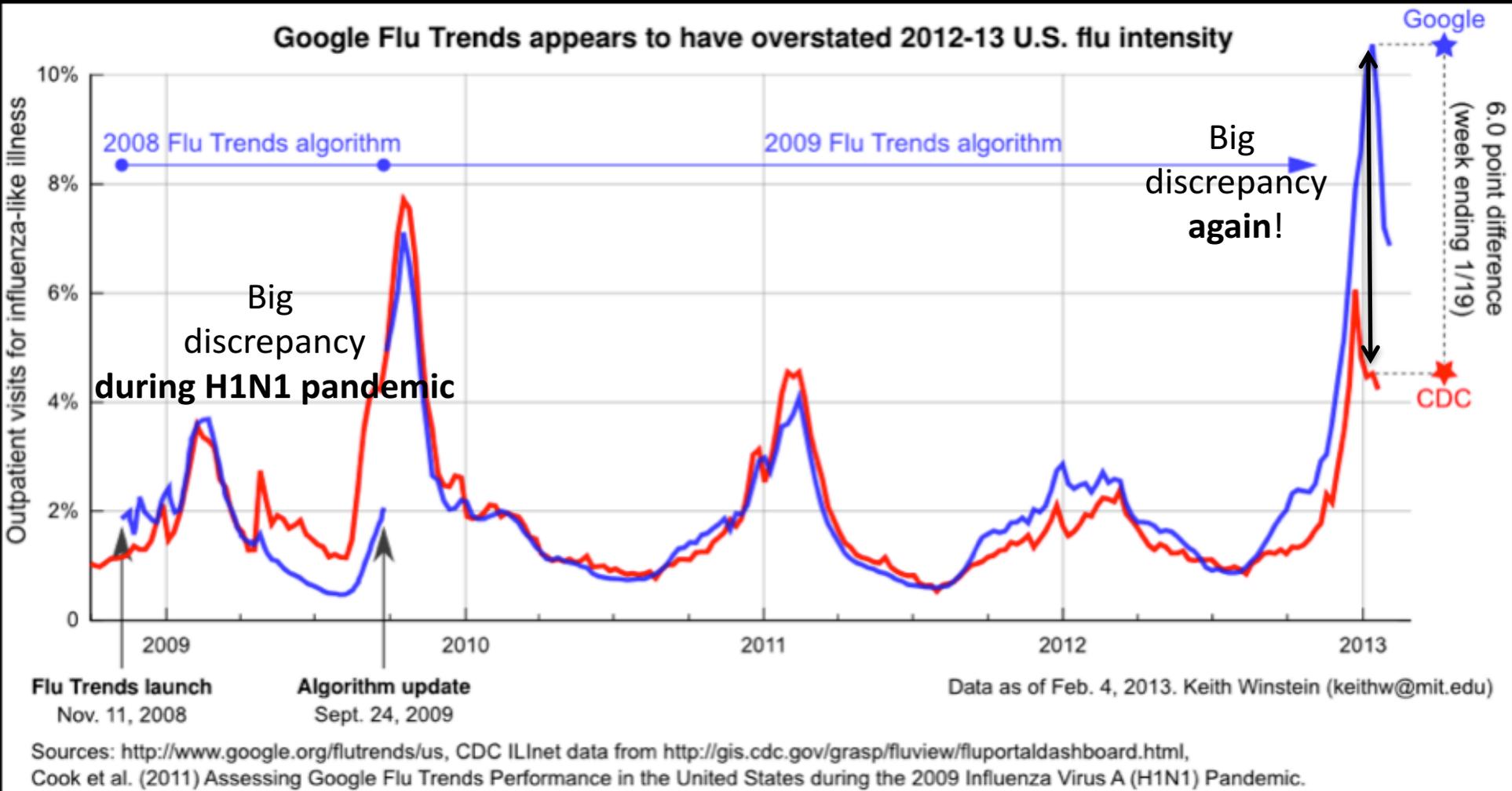
Sources: Google; Centers for Disease Control

Using Google to Monitor the Flu

Google Flu Trends can estimate the spread of the disease by measuring the frequency of certain search terms. Its findings closely track actual C.D.C. data and can, at times, anticipate the government reports.

What next? need to remove (not useful) terms.

Big discrepancies again!



Fixes were reported in: Cook et al. (2011) Assessing Google flu trends performance in the U.S. during the 2009 influenza virus A (H1N1) pandemic. PLoS One

Plot obtained from: <http://blog.keithw.org/2013/02/q-how-accurate-is-google-flu-trends.html>

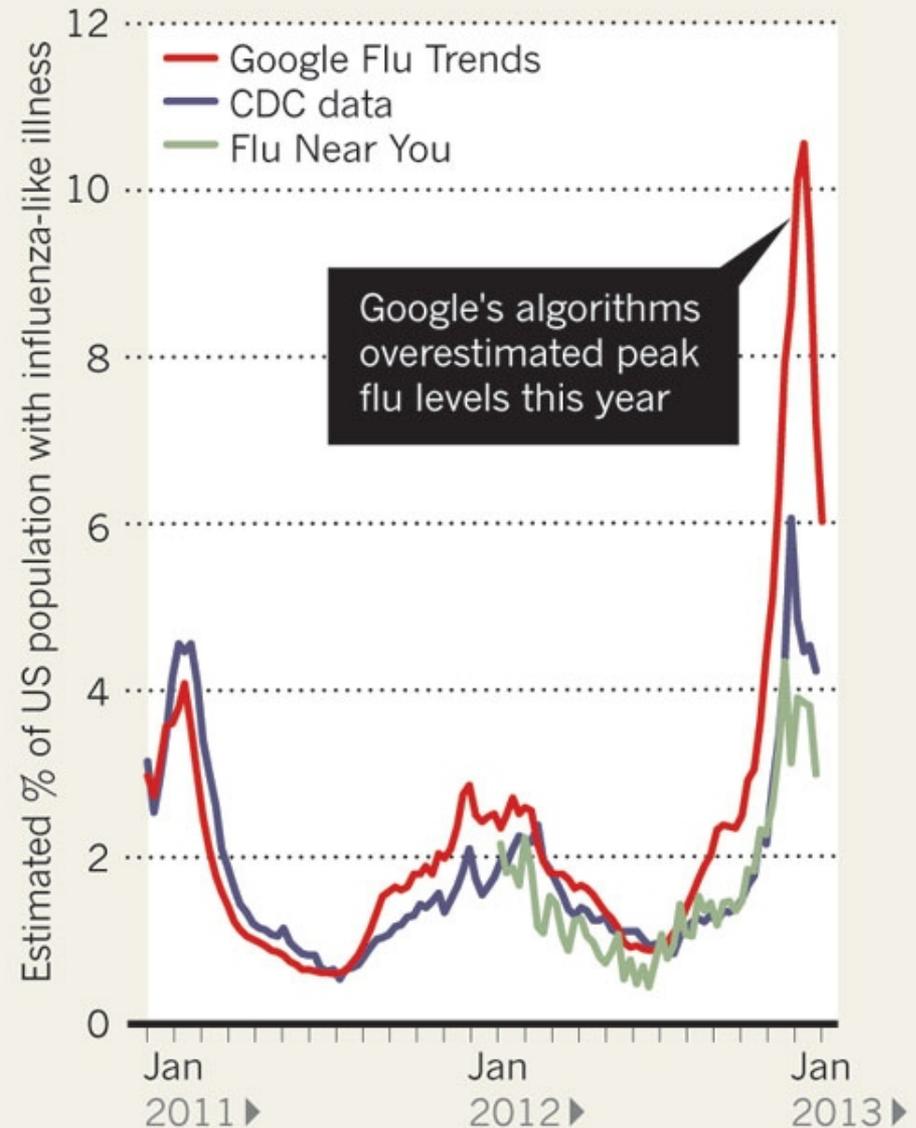


## When Google got flu wrong.

[nature.com/news/when-google-got-flu-wrong](http://nature.com/news/when-google-got-flu-wrong).

### FEVER PEAKS

A comparison of three different methods of measuring the proportion of the US population with an influenza-like illness.



# What Can Digital Disease Detection Learn from (an External Revision to) Google Flu Trends?

Mauricio Santillana, PhD, MS, D. Wendong Zhang, MA, Benjamin M. Althouse, PhD, ScM, John W. Ayers, PhD, MA

© 2014 Published by Elsevier Inc. on behalf of American Journal of Preventive Medicine Am J Prev Med 2014;47(3):341-347 341

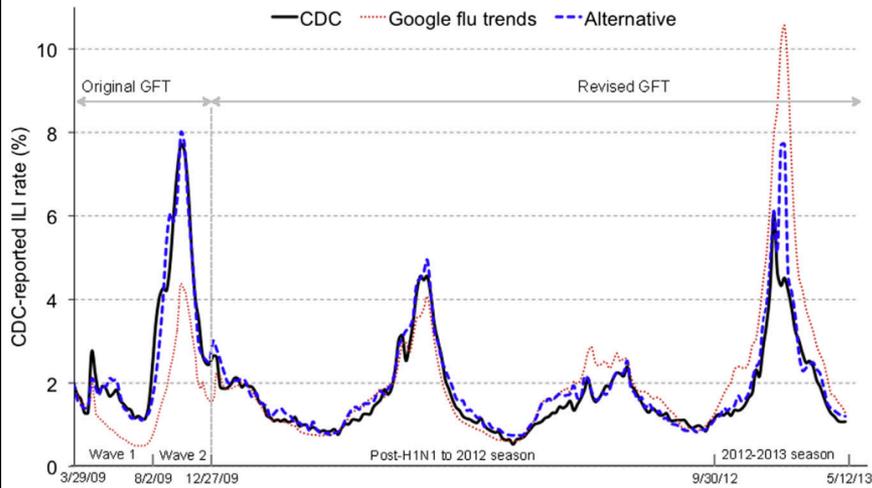


Figure 1. The alternative model outperforms Google Flu Trends

Google incorporated our proposed changes to GFT's engine in Oct 2014

We published a paper proposing improvements to GFT's engine (2014)



## Google Research Blog

The latest news from Research at Google

---

### Google Flu Trends gets a brand new engine

Posted: Friday, October 31, 2014

+1 222
Tweet 161
Like 104

Posted by Christian Stefansen, Senior Software Engineer

Each year the flu kills thousands of people and affects millions around the world. So it's important that public health officials and health professionals learn about outbreaks as quickly as possible. In 2008 we launched [Google Flu Trends](#) in the U.S., using aggregate web searches to indicate when and where influenza was striking in real time. These models [nicely complement](#) other survey systems—they're more fine-grained geographically, and they're typically more immediate, up to 1-2 weeks ahead of traditional methods such as the CDC's official reports. They can also be incredibly helpful for countries that don't have official flu tracking. Since launching, we've expanded Flu Trends to cover 29 countries, and launched [Dengue Trends](#) in 10 countries.

The original model performed surprisingly well despite its simplicity. It was retrained just once per year, and typically used only the 50 to 300 queries that produced the best estimates for prior seasons. We then left it to perform through the new season and evaluated it at the end. It didn't use the official CDC data for estimation during the season—only in the initial training.

SCIENTIFIC REPORTS PDF

Article | OPEN

# Advances in nowcasting influenza-like illness rates using search query logs

Vasileios Lamos , Andrew C. Miller, Steve Crossan & Christian Stefansen

*Scientific Reports* **5**,  
 Article number: 12760 (2015)  
 doi:10.1038/srep12760

Received: 07 May 2015  
 Accepted: 06 July 2015  
 Published online: 03 August 2015

[Download Citation](#)

Applied mathematics  
 Computer science Epidemiology  
 Influenza virus

Google and collaborators published a paper improving our AJPM 2014 methodology in August 2015

We improved last effort by Google team and published our results in PNAS in September 2015

PNAS

 CrossMark  
 ← click for updates

## Accurate estimation of influenza epidemics using Google search data via ARGO

Shihao Yang<sup>a</sup>, Mauricio Santillana<sup>b,c,1</sup>, and S. C. Kou<sup>a,1</sup>

<sup>a</sup>Department of Statistics, Harvard University, Cambridge, MA 02138; <sup>b</sup>School of Engineering and Applied Sciences, Harvard University, Cambridge, MA 02138; and <sup>c</sup>Computational Health Informatics Program, Boston Children's Hospital, Boston, MA 02115

Edited by Wing Hung Wong, Stanford University, Stanford, CA, and approved September 30, 2015 (received for review August 6, 2015)

**Accurate real-time tracking of influenza outbreaks helps public health officials make timely and meaningful decisions that could save lives. We propose an influenza tracking model, ARGO (AutoRegression with Google search data), that uses publicly available online search data. In addition to having a rigorous statistical foundation, ARGO outperforms all previously available Google-search-based tracking models, including the latest version of Google Flu Trends, even though it uses only low-quality search data as input from publicly available Google Trends and Google Correlate websites. ARGO not only incorporates the seasonality in influenza epidemics but also captures changes in people's online search behavior over time. ARGO is also flexible, self-correcting, robust, and scalable, making it a potentially powerful tool that can be used for real-time tracking of other social events at multiple temporal and spatial resolutions.**

CDC's ILI reports have a delay of 1–3wk due to the time for processing and aggregating clinical information. This time lag is far from optimal for decision-making purposes. To alleviate this information gap, multiple methods combining climate, demographic, and epidemiological data with mathematical models have been proposed for real-time estimation of flu activity (18, 21–25). In recent years, methods that harness Internet-based information have also been proposed, such as Google (1), Yahoo (2), and Baidu (3) Internet searches, Twitter posts (4), Wikipedia article views (5), clinicians' queries (6), and crowdsourced self-reporting mobile apps such as Influenzanet (Europe) (26), Flutracking (Australia) (27), and Flu Near You (United States) (28). Among them, GFT has received the most attention and has inspired subsequent digital disease detection systems (3, 8,

APPLIED  
 REMATICS

SCIENTIFIC REPORTS PDF

Article | OPEN

# Advances in nowcasting influenza-like illness rates using search query logs

Vasileios Lamos, Andrew C. Miller, Steve Crossan & Christian Stefansen

*Scientific Reports* **5**,  
 Article number: 12760 (2015)  
 doi:10.1038/srep12760

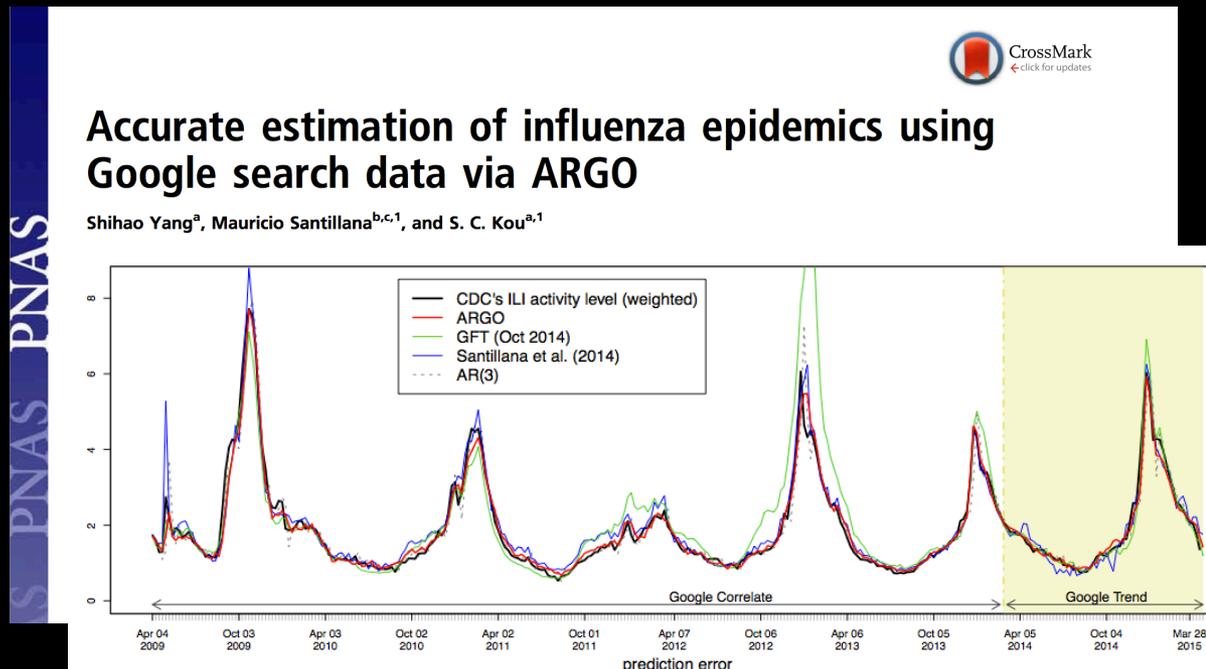
Received: 07 May 2015  
 Accepted: 06 July 2015  
 Published online: 03 August 2015

Download Citation

Applied mathematics  
 Computer science Epidemiology  
 Influenza virus

Google and collaborators published a paper improving our AJPM 2014 methodology in August 2015

We improved last effort by Google team and published our results in PNAS in September 2015



# Google discontinues Flu Trends indefinitely!



## Google Research Blog

The latest news from Research at Google

### The Next Chapter for Flu Trends

Posted: Thursday, August 20, 2015

  7



Instead of maintaining our own website going forward, we're now going to empower institutions who specialize in infectious disease research to use the data to build their own models. Starting this season, we'll provide Flu and Dengue signal data directly to partners including [Columbia University's Mailman School of Public Health](#) (to update their [dashboard](#)), [Boston Children's Hospital/Harvard](#), and [Centers for Disease Control and Prevention \(CDC\) Influenza Division](#). We will also continue to make historical Flu and Dengue estimate data available for anyone to see and analyze.

NEWS

# Google Flu Trends calls out sick, indefinitely

Google will pass along search queries related to the flu to health organizations so they can develop their own prediction models

By [Fred O'Connor](#) | [Follow](#)

IDG News Service | Aug 20, 2015 2:07 PM PT

**MORE LIKE THIS** ::

[Google Begins Tracking Swine Flu in Mexico](#)



[Google's Panicky Flu Estimates Were Dead Wrong](#)

**BIG DATA**

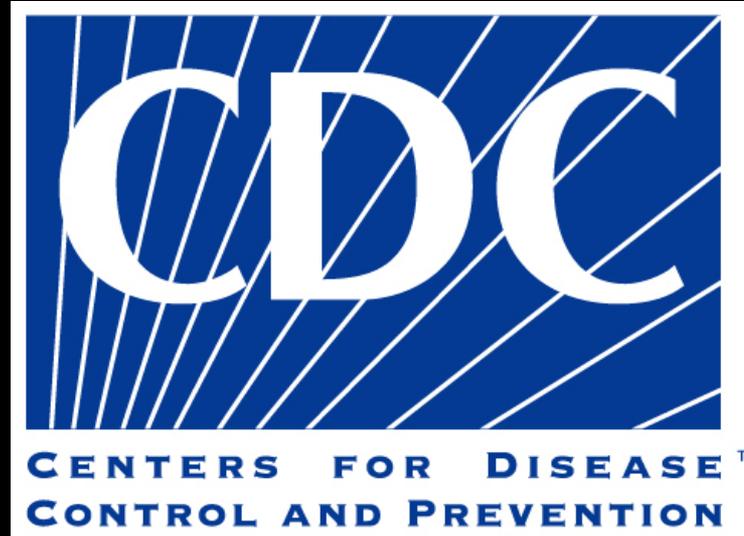
# Google discontinues Flu Trends, starts offering data to researchers

JORDAN NOVET | AUGUST 20, 2015 12:17 PM

TAGS: [GOOGLE](#), [GOOGLE FLU TRENDS](#)

Our team at Boston Children's Hospital now has access to Google's search volumes, as one of the exclusive Google's partners.

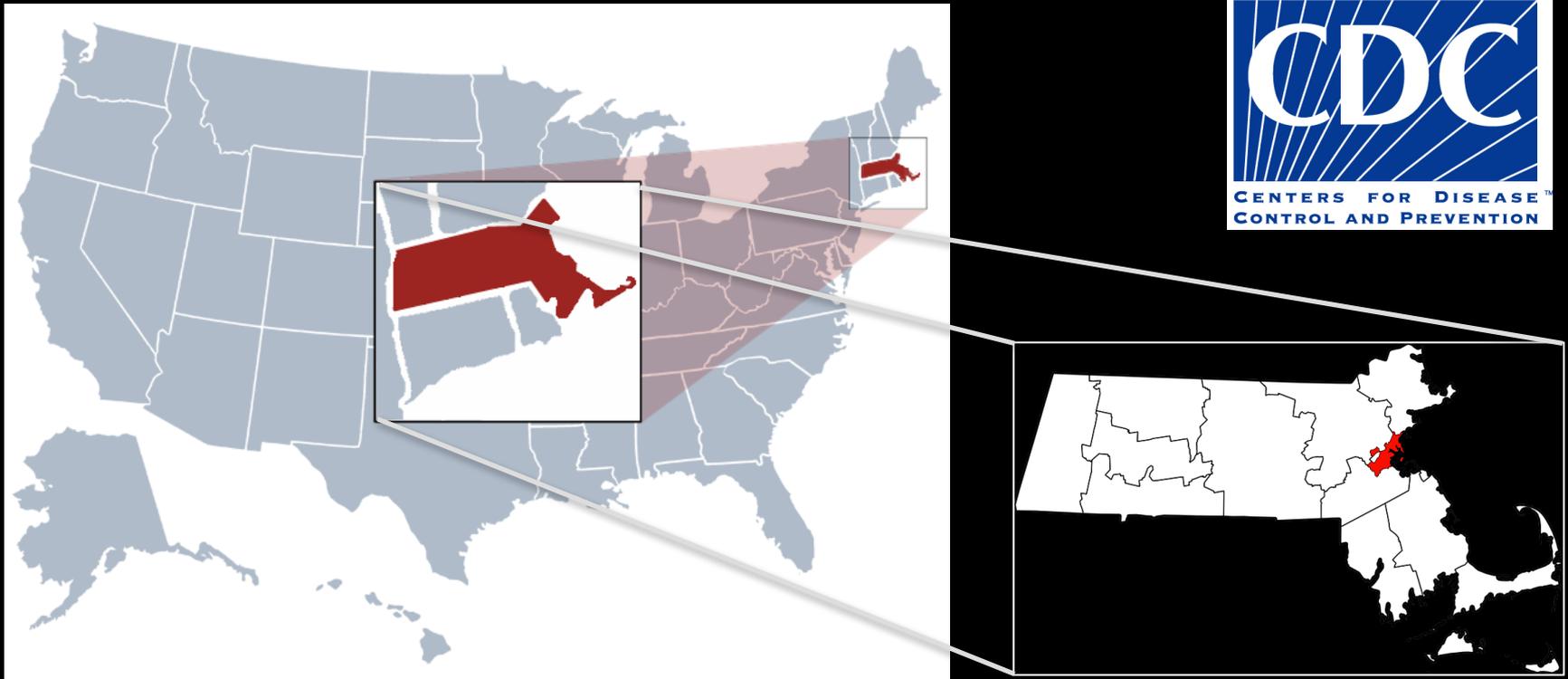
We are helping create a new improved disease forecasting platform supported partially by the Centers for Disease Control and Prevention



In collaboration with the CDC Influenza division, we are extending our work from National and Regional predictions, to state-level and city level (Boston as a pilot)

Grant: *Centers for Disease Control and Prevention's Cooperative Agreement PPHF 11797-998G-15*

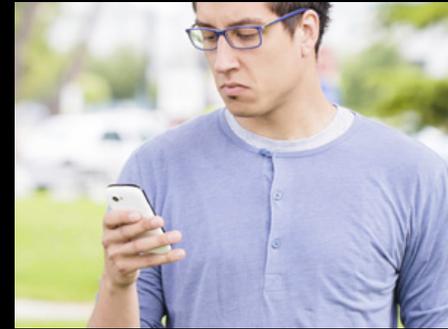
Team members: *Fred Lu, Leonardo C. Clemente*  
CDC liaison and collaborator: *Matt Biggerstaff*



# Beyond Google searches...



What are doctors searching for?

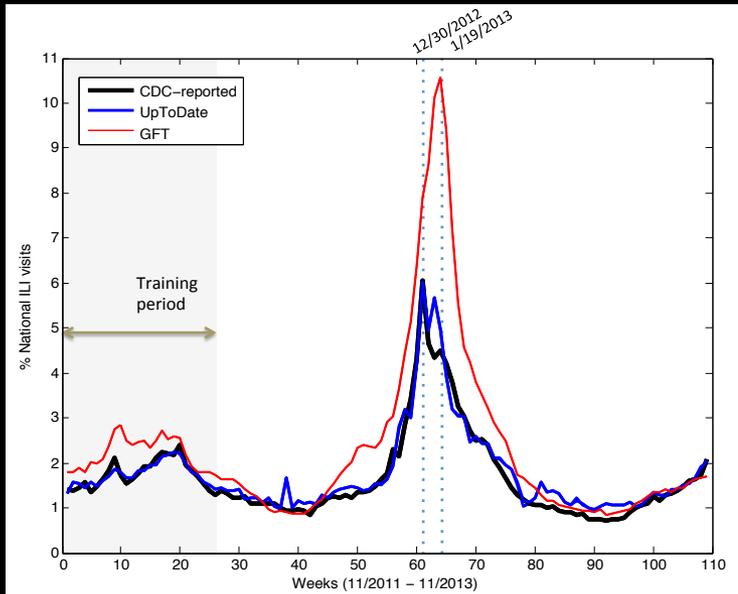


What are people tweeting? What are they reporting on crowd-sourced disease surveillance apps?

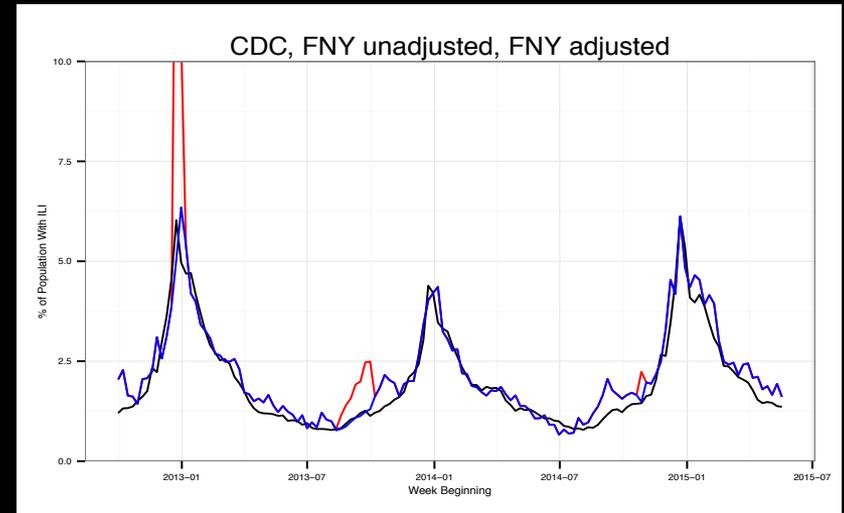


Can we use Electronic Health Records (EHR) to track disease incidence? What lab tests or medications are doctors prescribing?

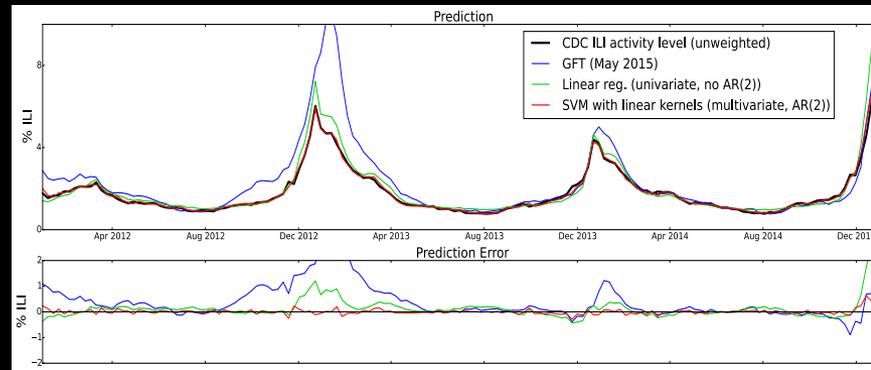
# Beyond Google searches...



What are doctors searching for?

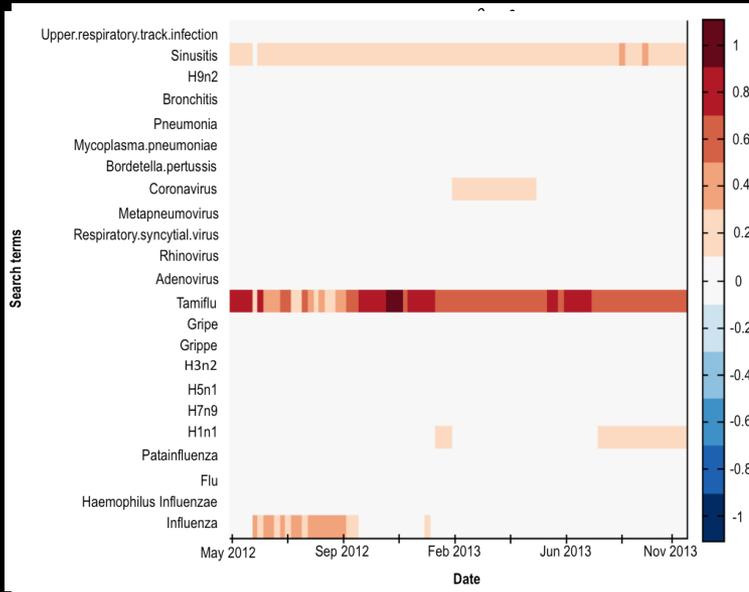


What are people tweeting? What are they reporting on crowd-sourced disease surveillance apps?

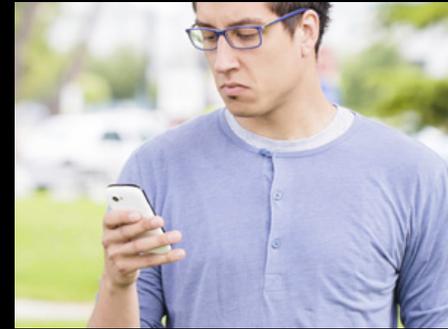


Can we use Electronic Health Records (EHR) to track disease incidence? What lab tests or medications are doctors prescribing?

# Beyond Google searches...



What are doctors searching for?



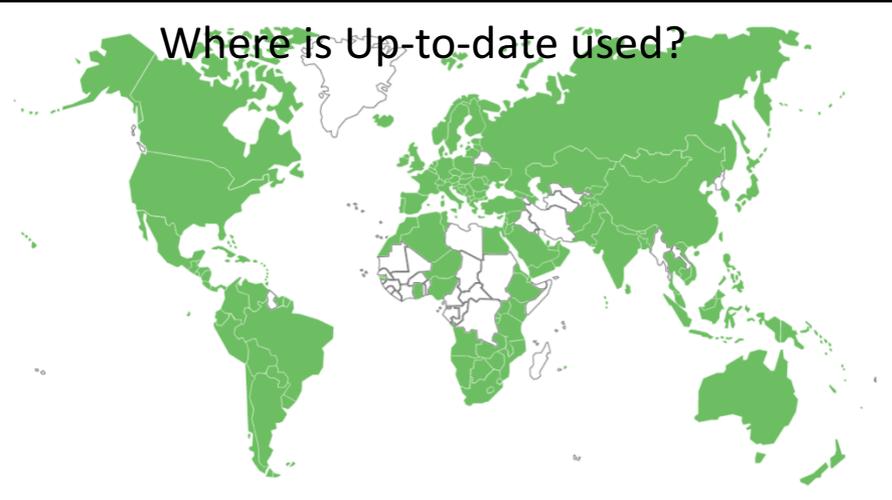
What are people tweeting? What are they reporting on crowd-sourced disease surveillance apps?



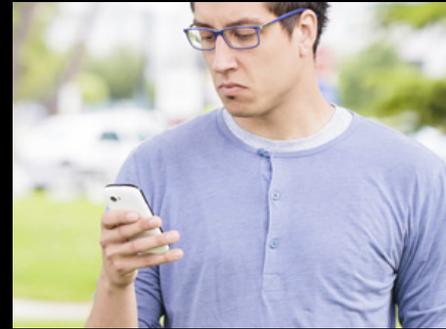
Can we use Electronic Health Records (EHR) to track disease incidence? What lab tests or medications are doctors prescribing?

# Beyond Google searches...

Where is Up-to-date used?



What are doctors searching for?



What are people tweeting? What are they reporting on crowd-sourced disease surveillance apps?



Can we use Electronic Health Records (EHR) to track disease incidence? What lab tests or medications are doctors prescribing?

# Beyond Google searches...

OXFORD JOURNALS

## Clinical Infectious Diseases

### Using Clinicians' Search Query Data to Monitor Influenza Epidemics

Mauricio Santillana,<sup>1,2</sup> Elaine O. Nsoesie,<sup>2,3</sup> Sumiko R. Mekaru,<sup>2</sup> David Scales,<sup>2,4</sup> and John S. Brownstein<sup>1,5</sup>

<sup>1</sup>School of Engineering and Applied Sciences, Harvard University, Cambridge, <sup>2</sup>Children's Hospital Informatics Program, Boston Children's Hospital, <sup>3</sup>Department of Pediatrics, Harvard Medical School, Boston, and <sup>4</sup>Department of Internal Medicine, Cambridge Health Alliance, Massachusetts; and <sup>5</sup>Department of Epidemiology, Biostatistics and Occupational Health, McGill University, Montreal, Quebec, Canada

Search query information from a clinician's database, UpToDate, is shown to predict influenza epidemics in the United States in a timely manner. Our results show that digital disease surveillance tools based on experts' databases may be able to provide an alternative, reliable, and stable signal for accurate predictions of influenza outbreaks.

**Keywords.** digital disease detection; Internet-based disease surveillance; prediction of influenza.

validated traditional surveillance systems and have the potential to provide timely epidemiologic intelligence to inform prevention messaging and healthcare facility staffing decisions.

The potential for the public's search activity to be influenced by anxiety, fears, and rumors raises concerns regarding reliability [10–13]. Although recent revisions to GFT have shown that these concerns can be partially mitigated [13–15], shifting Internet-based surveillance from the entire public to subject-matter experts may maintain timeliness while generating a more reliable and stable signal requiring much less data. A recent small retrospective study using data on queries to a Finnish primary care guidelines database demonstrated, for example, that disease-specific queries for Lyme disease, tularemia, and other infectious diseases correlated well with concurrent confirmed cases [16].

Here, we show that UpToDate ([www.uptodate.com](http://www.uptodate.com)), a physician-authored clinical decision support Internet resource that is used by 700 000 clinicians in 158 countries and almost 90% of academic medical centers in the United States, can be used for syndromic surveillance of influenza. Specifically, we use UpToDate's search query activity related to ILI to design a timely sentinel of influenza incidence in the United States.

What are doctors searching for?

# AJPM American Journal of Preventive Medicine

A Journal of the American College of Preventive Medicine and Association for Prevention Teaching and Research

## Flu Near You: Crowdsourced Symptom Reporting Spanning 2 Influenza Seasons

Mark S. Smolinski, MD, MPH, Adam W. Crawley, MPH, Kristin Baltrusaitis, MA, Rumi Chunara, PhD, MS, Jennifer M. Olsen, DrPH, Oktavia Wijckij, PhD, Mauricio Santillana, PhD, MS, Andre Nguyen, and John S. Brownstein, PhD, MPH

Digital communications technologies have rapidly increased in use for public health disease surveillance. Mobile phones, tablets, digital pens, and satellites are making it possible for surveillance and rapid response teams in even remote areas of the globe to carry out an essential function of public health to protect against outbreaks of infectious disease. To date, public health surveillance has been limited by the capacity of public health authorities to conduct case and contact tracing and a reliance on data provided primarily by the medical system. The increased use of digital communications technology is now making it possible to enable the public to actively be part of the public health surveillance system.

Since 2003, participatory surveillance approaches have leveraged online survey technology with syndromic surveillance of human infectious diseases through volunteer symptom

**Objectives.** We summarized Flu Near You (FNY) data from the 2012–2013 and 2013–2014 influenza seasons in the United States.

**Methods.** FNY collects limited demographic characteristic information upon registration, and prompts users each Monday to report symptoms of influenza-like illness (ILI) experienced during the previous week. We calculated the descriptive statistics and rates of ILI for the 2012–2013 and 2013–2014 seasons. We compared raw and noise-filtered ILI rates with ILI rates from the Centers for Disease Control and Prevention ILINet surveillance system.

**Results.** More than 61 000 participants submitted at least 1 report during the 2012–2013 season, totaling 327 773 reports. Nearly 40 000 participants submitted at least 1 report during the 2013–2014 season, totaling 336 933 reports. Rates of ILI as reported by FNY tracked closely with ILINet in both timing and magnitude.

**Conclusions.** With increased participation, FNY has the potential to serve as a viable complement to existing outpatient, hospital-based, and laboratory surveillance systems. Although many established systems have the benefits of specificity and credibility, participatory systems offer advantages in the areas of speed, sensitivity, and scalability. (*Am J Public Health*. Published online ahead of print August 13, 2015; e1–e7. doi:10.2105/AJPH.2015.302696)

What are people tweeting? What are they reporting on crowd-sourced disease surveillance apps?

# SCIENTIFIC REPORTS

OPEN

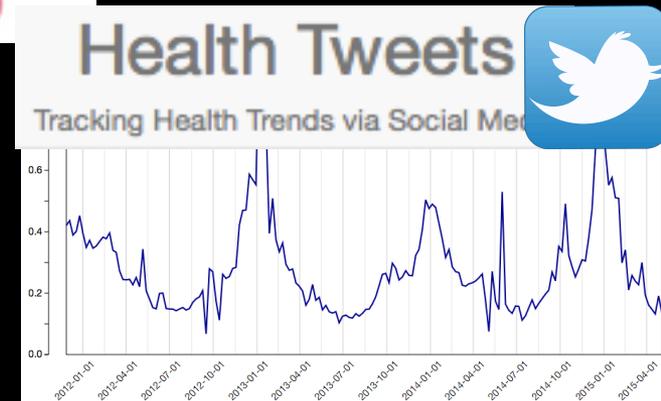
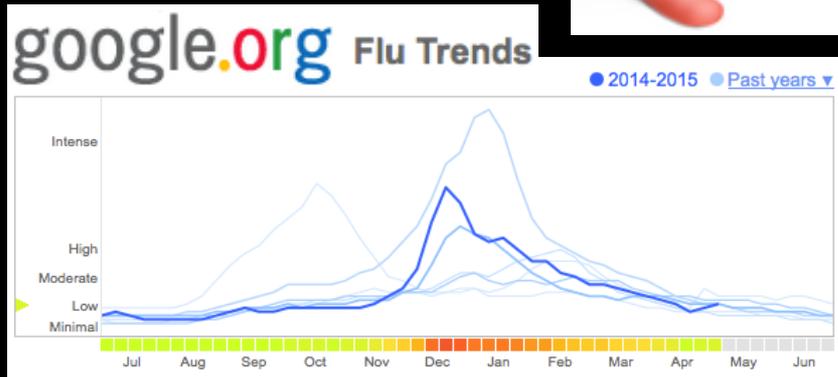
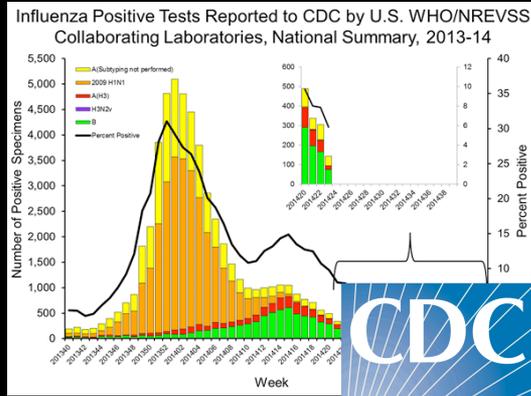
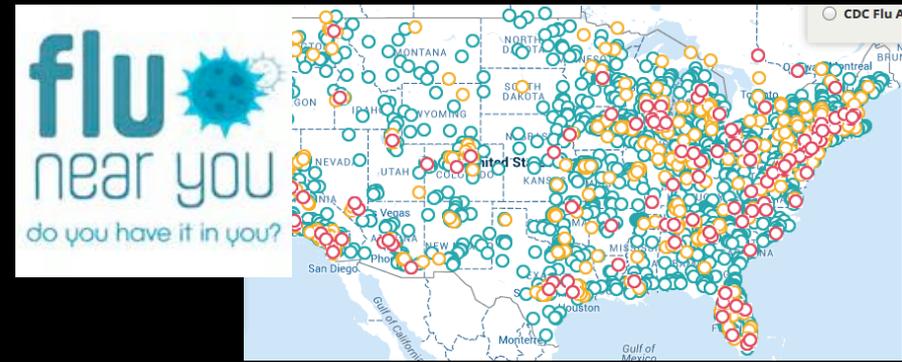
## Cloud-based Electronic Health Records for Real-time, Region-specific Influenza Surveillance

Received: 31 December 2015  
Accepted: 20 April 2016

M. Santillana<sup>1,2,3</sup>, A. T. Nguyen<sup>3</sup>, T. Louie<sup>4</sup>, A. Zink<sup>5</sup>, J. Gray<sup>5</sup>, I. Sung<sup>5</sup> & J. S. Brownstein<sup>1,2</sup>

Can we use Electronic Health Records (EHR) to track disease incidence? What lab tests or medications are doctors prescribing?

# Ensemble approaches yield more accurate and more robust real-time and forecast flu estimates



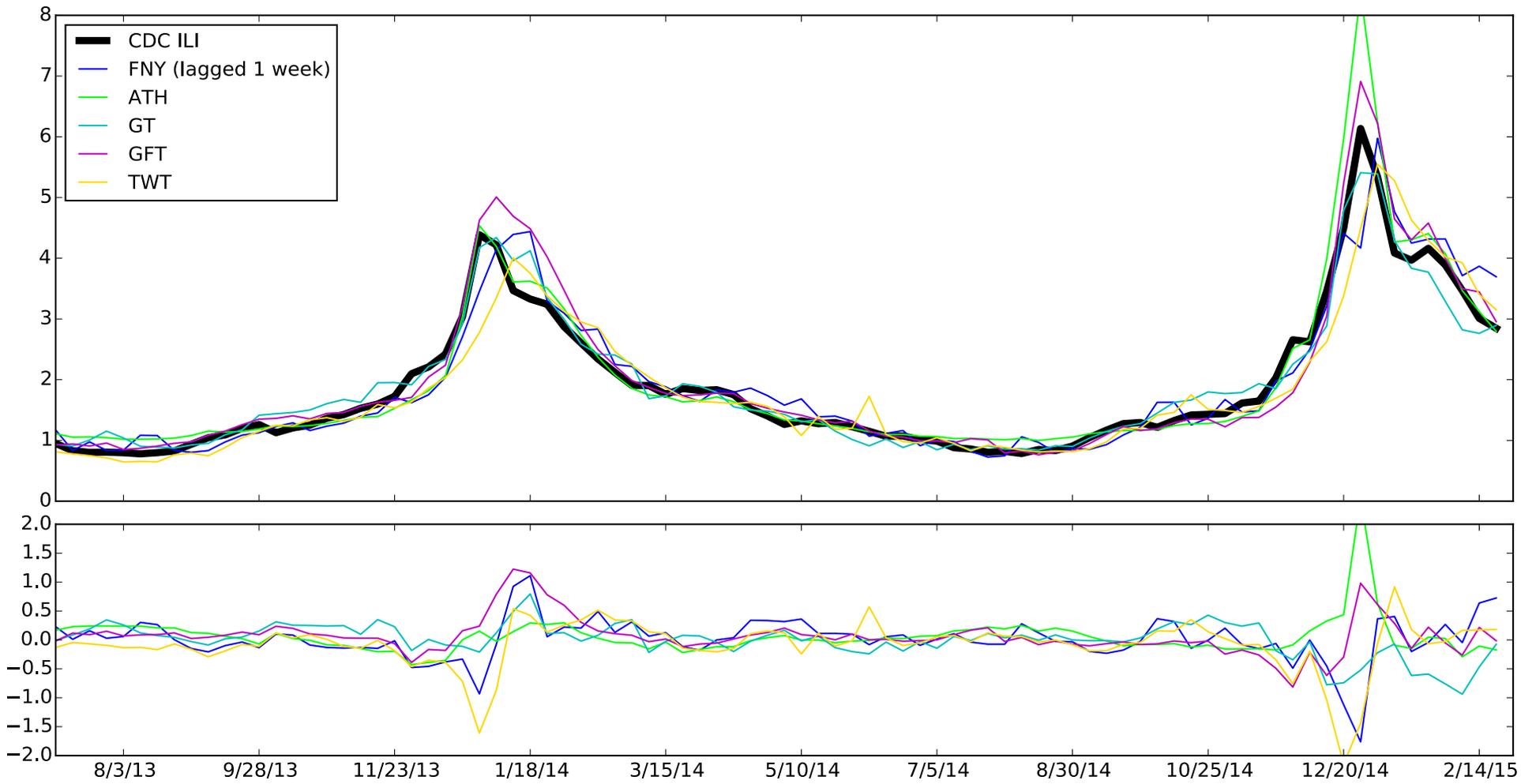
## Performance of individual data sources

	CORR	RMSE (%ILI)	Rel RMSE (%)	RMAE (%)	Hit Rate
FNY	0.948	0.385	15.9	39.3	65.9
ATH	0.977	0.351	14.1	36.7	<b>77.7</b>
GT	0.978	0.245	13.3	42.9	65.9
GFT	0.980	0.333	12.3	35.3	75.3
TWT	0.937	0.414	15.1	50.1	62.4
CDC Baseline	0.930	0.501	18.2	46.7	68.2
CDC Virology	0.923	-	-	-	69.4

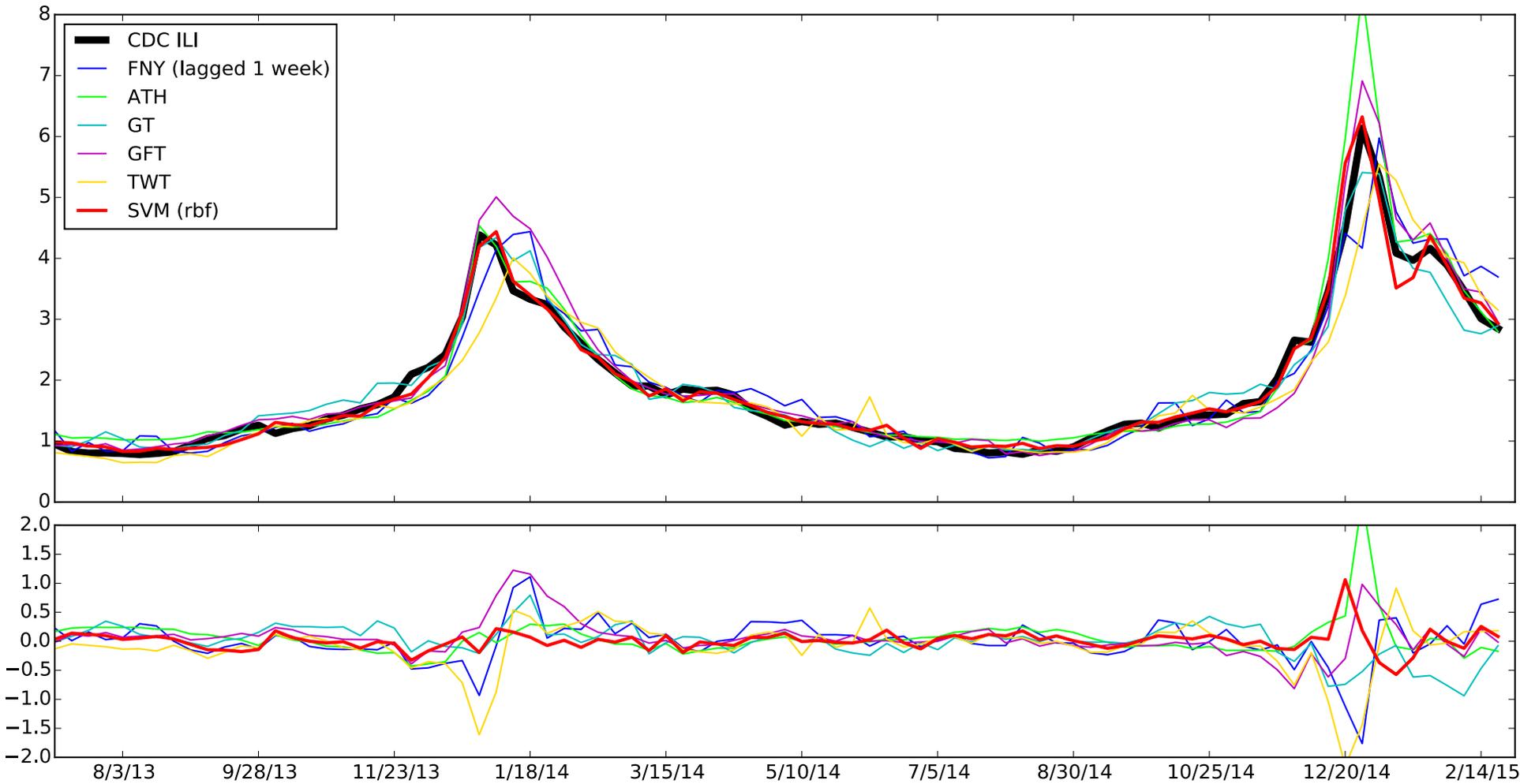
## Performance ensemble

	CORR	RMSE (%ILI)	Rel RMSE (%)	RMAE (%)	Hit Rate
FNY	0.948	0.385	15.9	39.3	65.9
ATH	0.977	0.351	14.1	36.7	<b>77.7</b>
GT	0.978	0.245	13.3	42.9	65.9
GFT	0.980	0.333	12.3	35.3	75.3
TWT	0.937	0.414	15.1	50.1	62.4
CDC Baseline	0.930	0.501	18.2	46.7	68.2
CDC Virology	0.923	-	-	-	69.4
SVM (RBF)	<b>0.989</b>	<b>0.176</b>	<b>8.27</b>	<b>23.6</b>	69.4

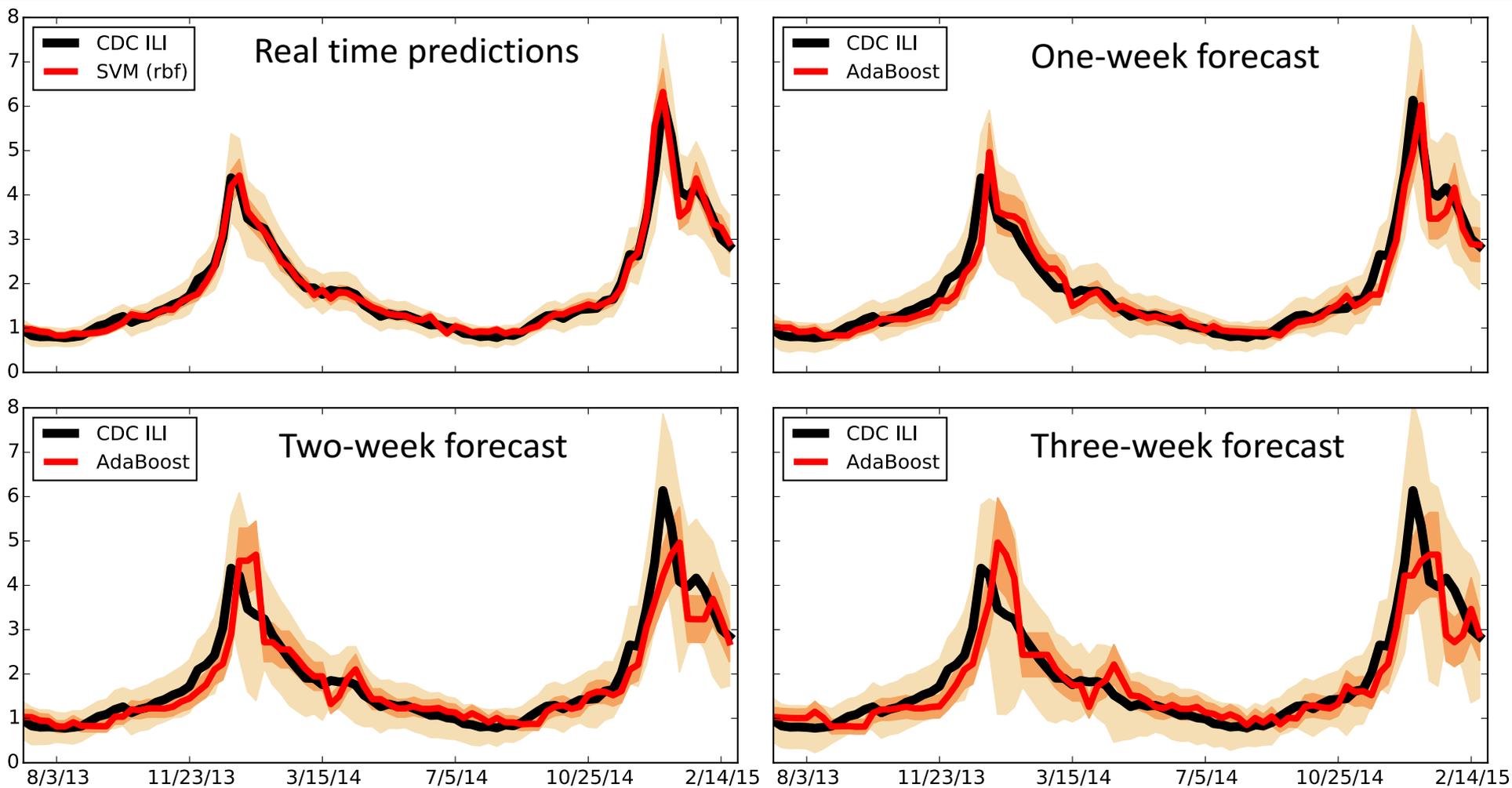
# Performance of individual data sources



# Performance ensemble



# Real time predictions and Forecasts



# Ensemble approaches yield more accurate and more robust real-time and forecast flu estimates

Yang et al. *BMC Infectious Diseases* (2017) 17:332  
DOI 10.1186/s12879-017-2424-7

BMC Infectious Diseases

RESEARCH ARTICLE

Open Access



## Using electronic health records and Internet search information for accurate influenza forecasting

Shihao Yang<sup>1</sup>, Mauricio Santillana<sup>2,3\*</sup>, John S. Brownstein<sup>2,3</sup>, Josh Gray<sup>4</sup>, Stewart Richardson<sup>4</sup> and S. C. Kou<sup>1\*</sup>

### Abstract

**Background:** Accurate influenza activity forecasting helps public health officials prepare and allocate resources for unusual influenza activity. Traditional flu surveillance systems, such as the Centers for Disease Control and Prevention's (CDC) influenza-like illnesses reports, lag behind real-time by one to 2 weeks, whereas information contained in cloud-based electronic health records (EHR) and in Internet users' search activity is typically available in near real-time. We present a method that combines the information from these two data sources with historical flu activity to produce national flu forecasts for the United States up to 4 weeks ahead of the publication of CDC's flu reports.

**Methods:** We extend a method originally designed to track flu using Google searches, named ARGO, to combine information from EHR and Internet searches with historical flu activities. Our regularized multivariate regression model dynamically selects the most appropriate variables for flu prediction every week. The model is assessed for the flu seasons within the time period 2013–2016 using multiple metrics including root mean squared error (RMSE).

**Results:** Our method reduces the RMSE of the publicly available alternative (Healthmap flutrends) method by 33, 20, 17 and 21%, for the four time horizons: real-time, one, two, and 3 weeks ahead, respectively. Such accuracy improvements are statistically significant at the 5% level. Our real-time estimates correctly identified the peak timing and magnitude of the studied flu seasons.

**Conclusions:** Our method significantly reduces the prediction error when compared to historical publicly available Internet-based prediction systems, demonstrating that: (1) the method to combine data sources is as important as data quality; (2) effectively extracting information from a cloud-based EHR and Internet search activity leads to accurate forecast of flu.

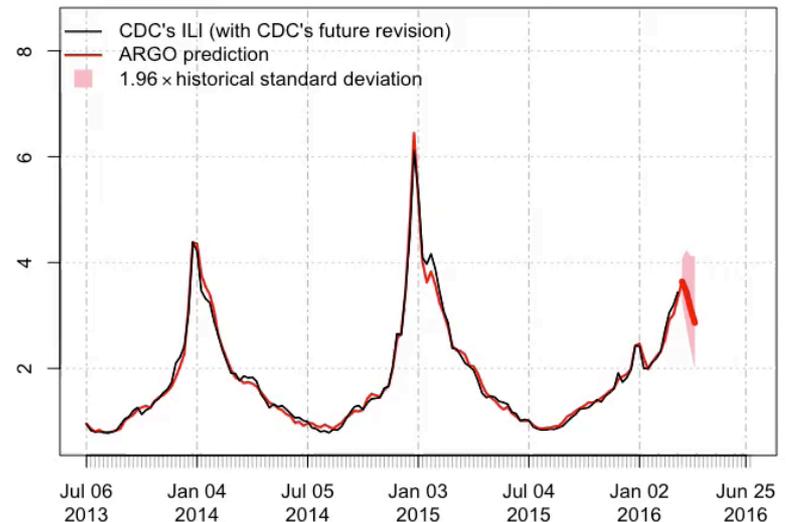
**Keywords:** Influenza-like illnesses reports, Digital disease detection, Dynamic error reduction, Validation test, Autoregression

RESEARCH ARTICLE

## Combining Search, Social Media, and Traditional Data Sources to Improve Influenza Surveillance

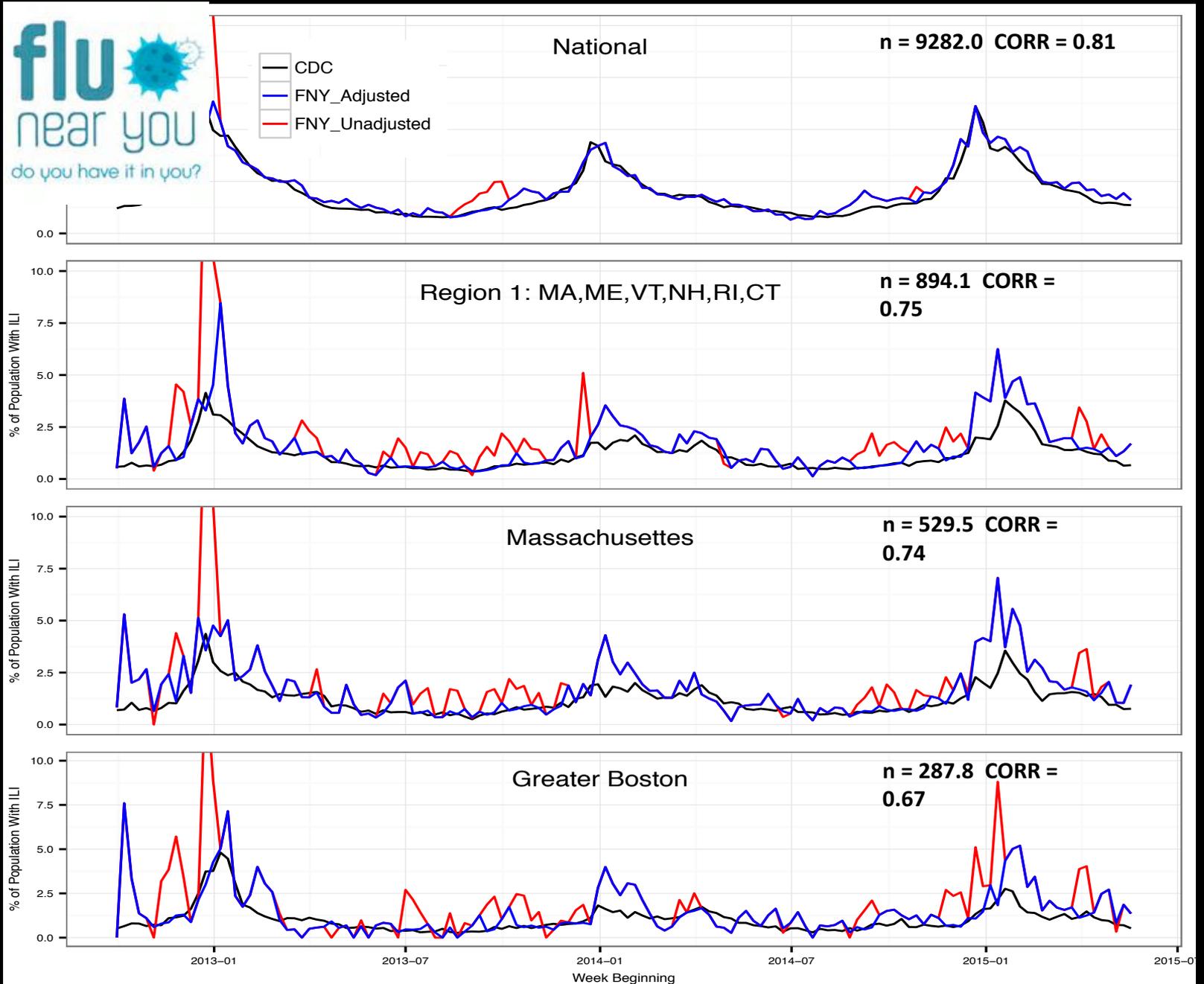
Mauricio Santillana<sup>1,2,3\*</sup>, André T. Nguyen<sup>1</sup>, Mark Dredze<sup>4</sup>, Michael J. Paul<sup>5</sup>, Elaine O. Nsoesie<sup>6,7</sup>, John S. Brownstein<sup>2,3</sup>

ARGO Prediction vs. CDC's ILI

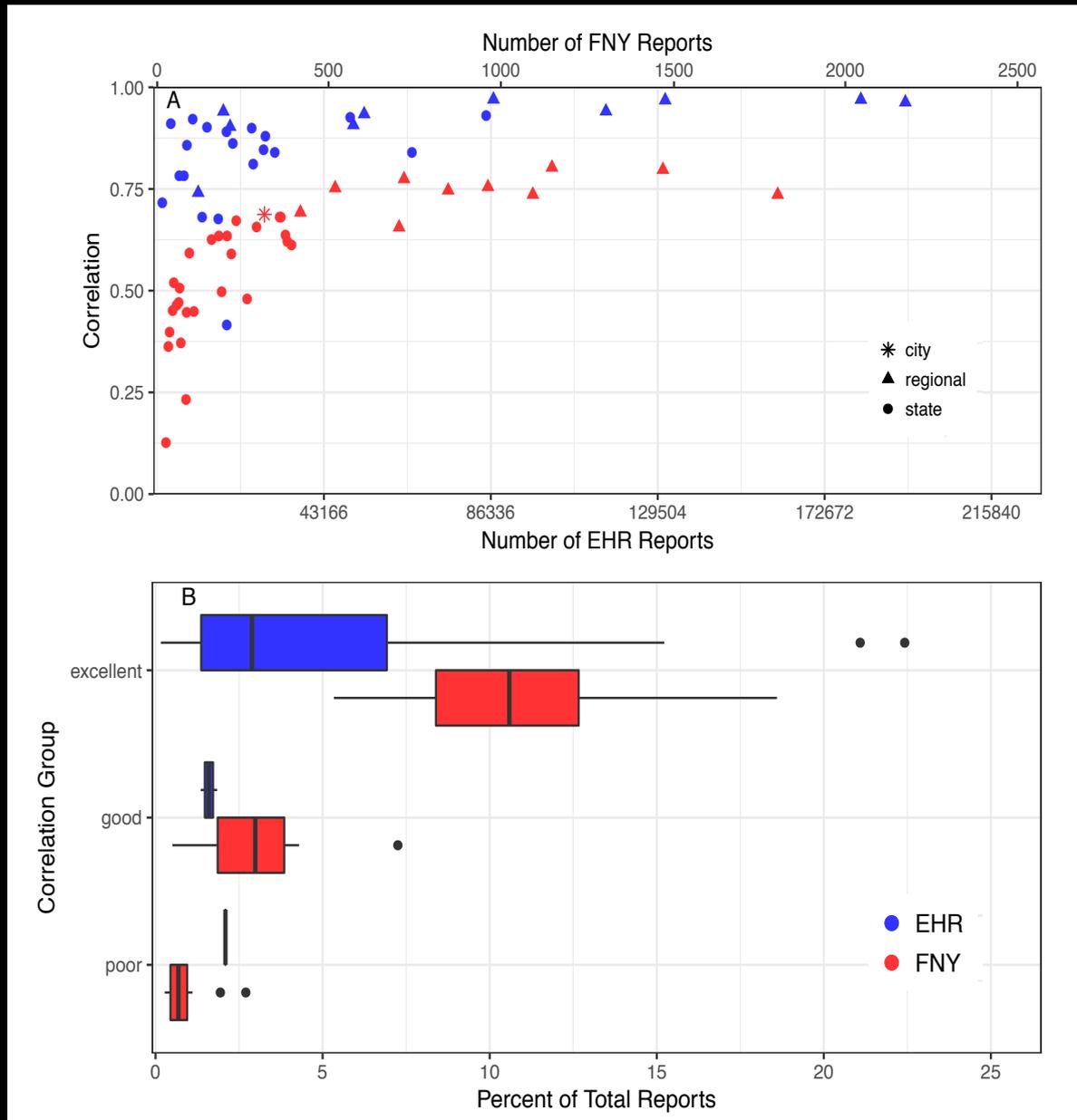


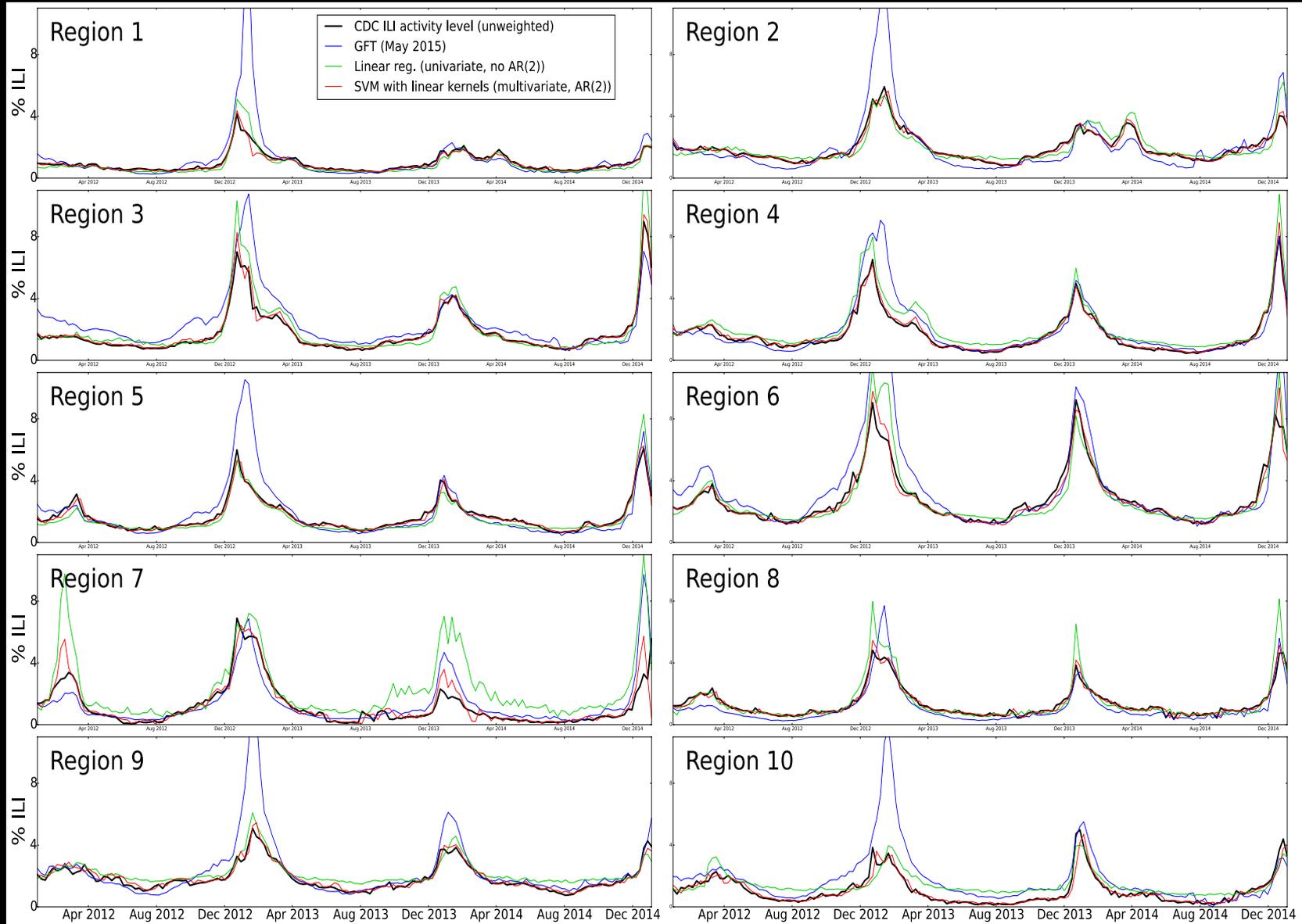
What about estimating flu in a regional level,  
state-level, city-level?

# Correlation of FNY with CDC. Multiple Geographic Scales



# Comparisons between CDC ILINet and FNY and Athenahealth in multiple spatial scales





How about state-level?

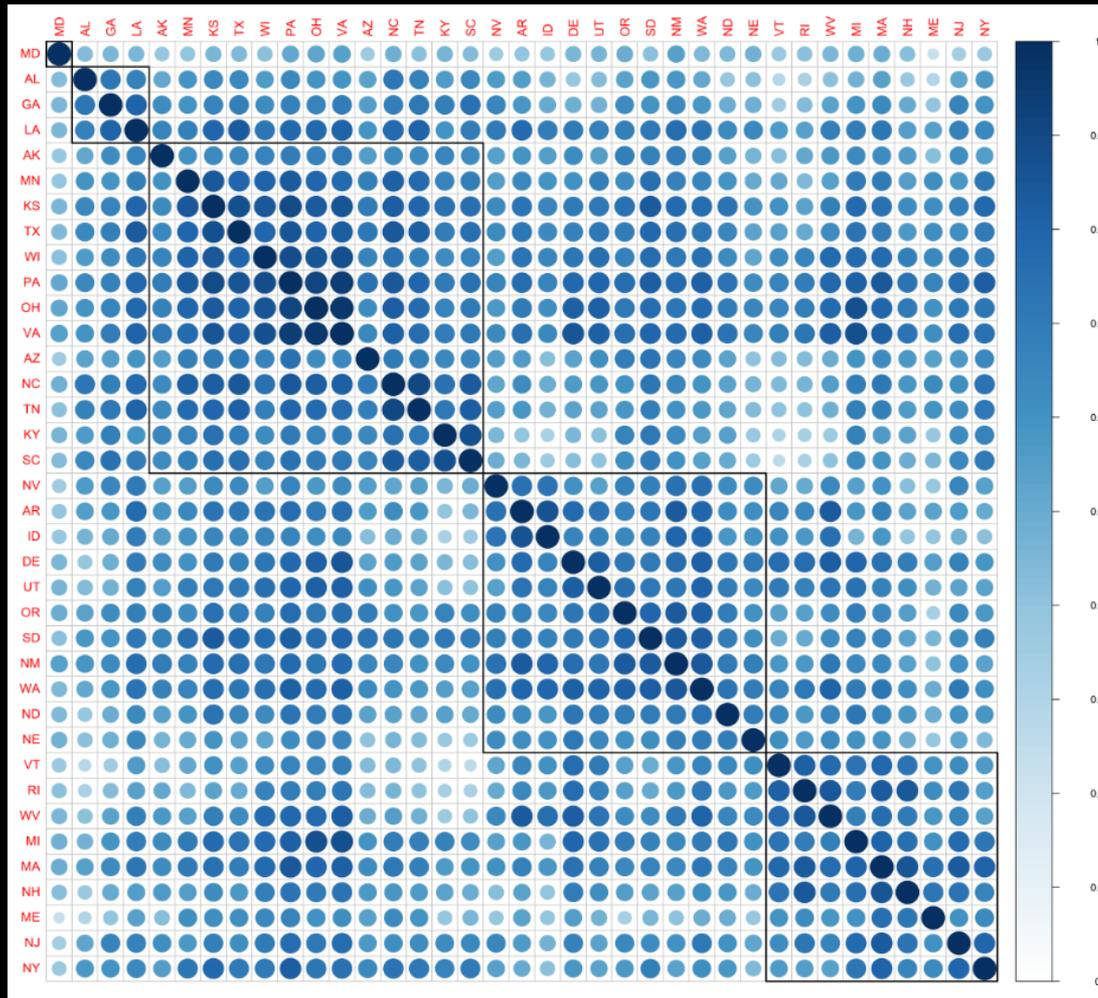
## Spatial-temporal synchronicities



## Flu-related Google search information

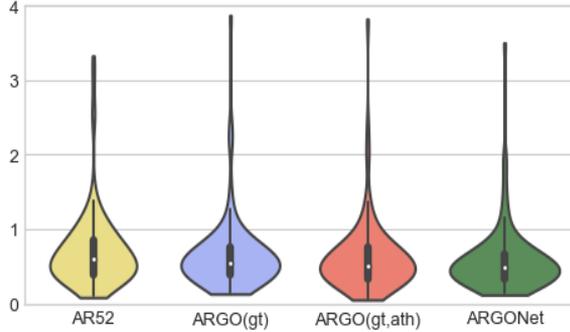


Heat map of pairwise %LI correlations between all states.  
Boxes denote clusters of highly correlated states.

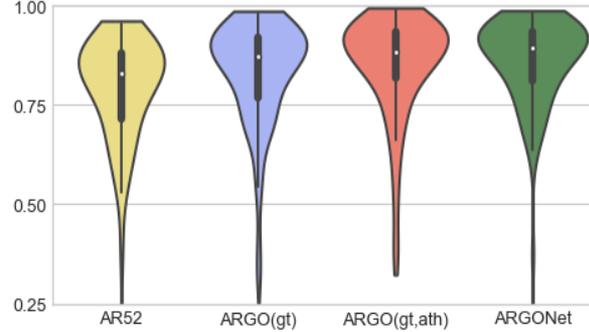


# Clear improvements over previous methodologies

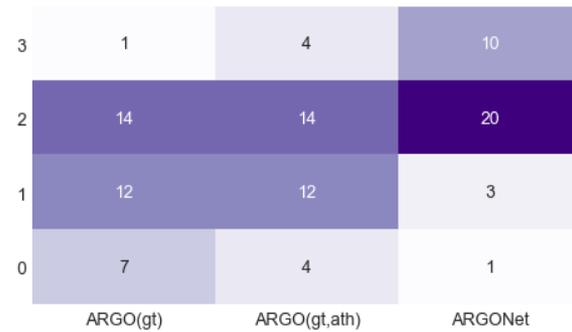
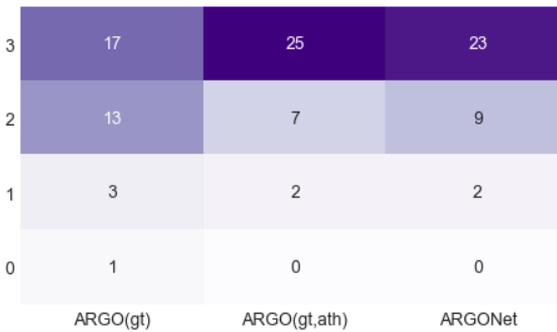
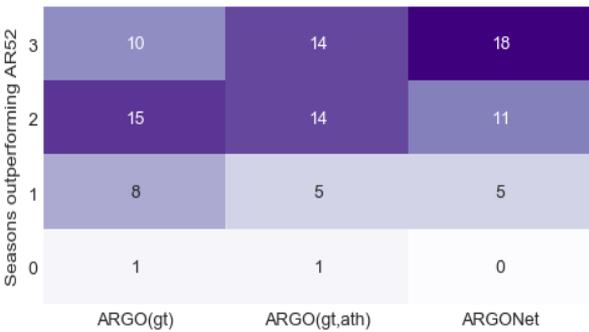
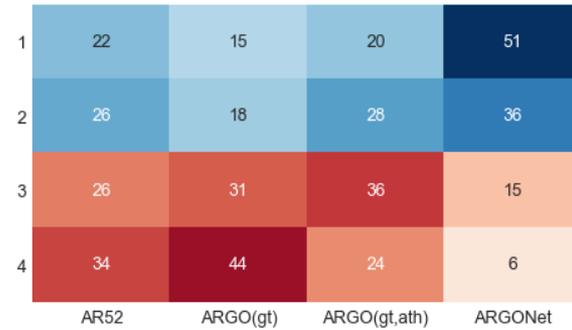
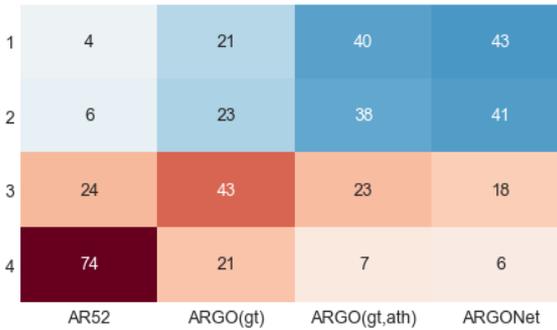
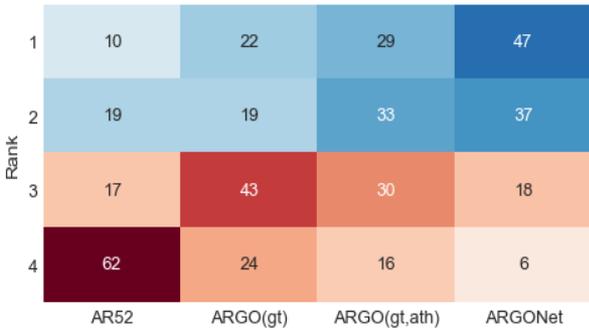
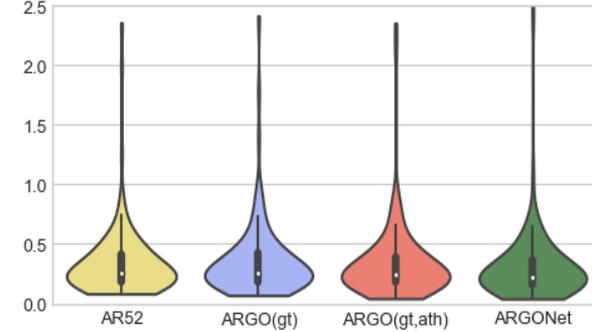
RMSE



Correlation



MAPE

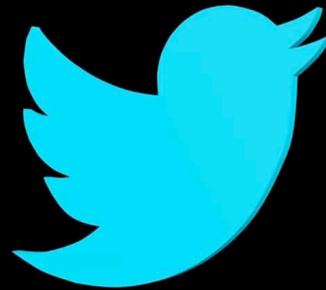


How about city-level?

Refining the spatial resolution...

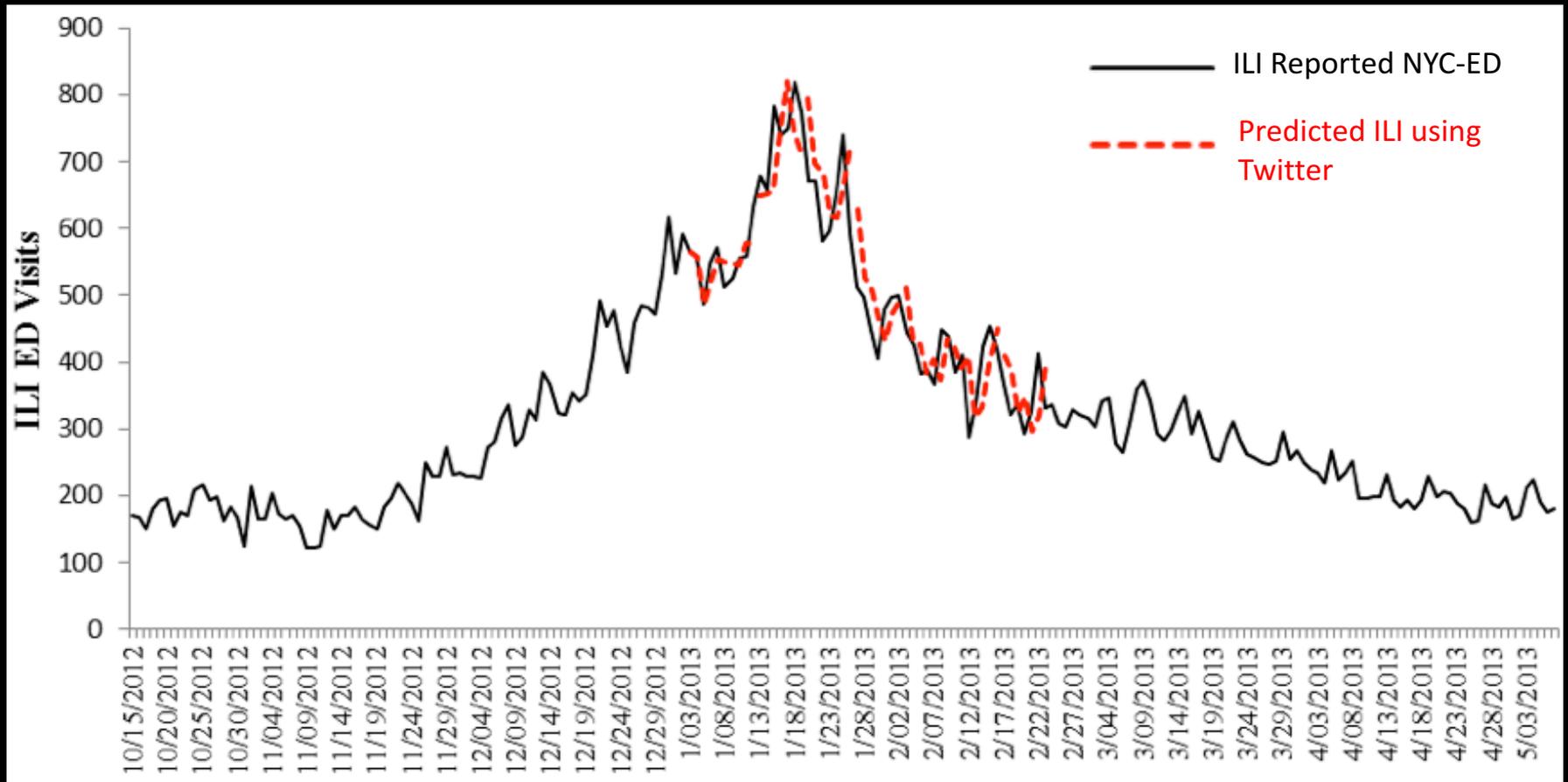


Tracking Flu using twitter  
(Daily analysis in NYC)



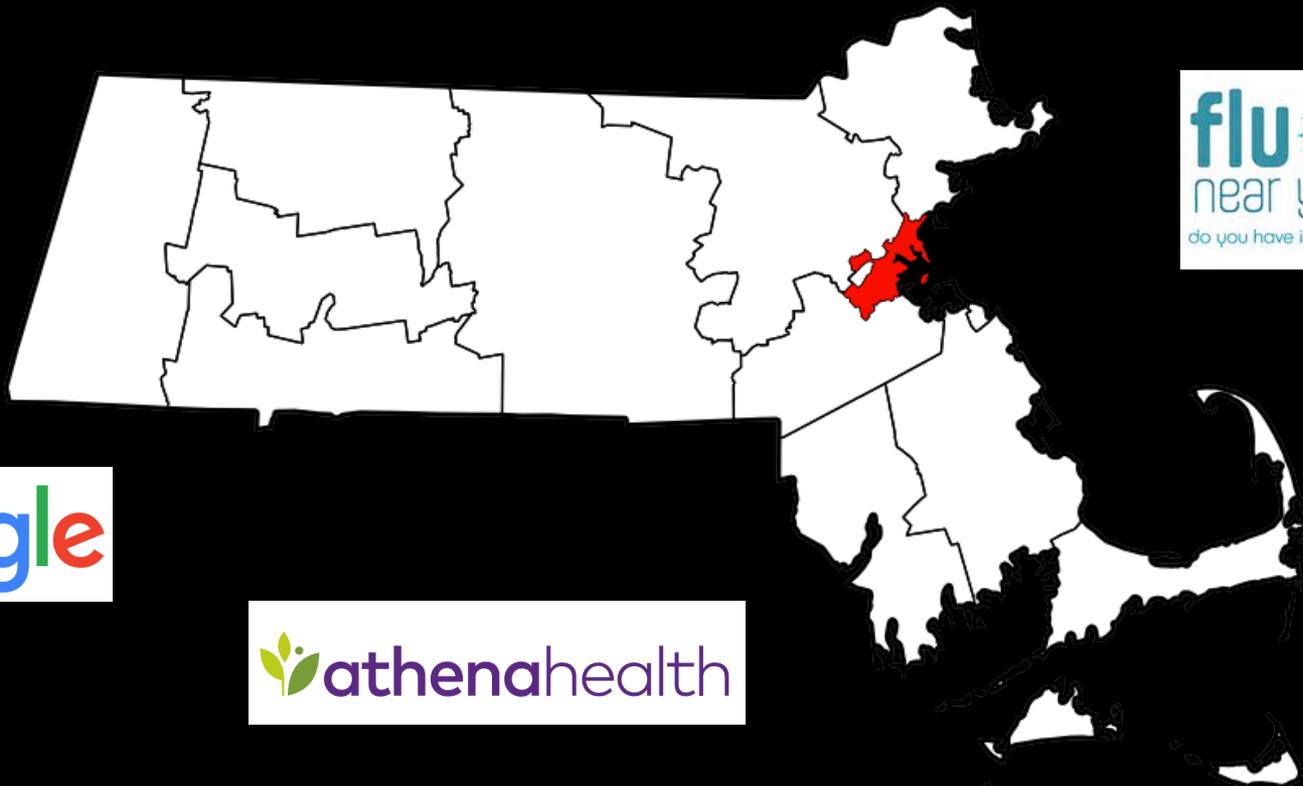
Work with R. Nagar, Q. Yuan, C. Freifeld, A. Nojima, R. Chunara, and J. S. Brownstein

# Daily ILI visits (as reported by the NYC emergency department) compared to predicted ILI using twitter data



# We will extend our methodology to finer spatial resolutions. (Massachusetts and Boston)

Highlights: (a) dynamic-moving training window, (b) automatic feature selection, (c) ensemble approach



Lu F, Hou S, Baltrusaitis K, Shah M, Leskovec J, Sosic R, Hawkins J, Brownstein JS, Conidi G, Gunn J, ..., Santillana M. Accurate influenza monitoring and forecasting in the Boston metropolis using novel Internet data streams. *Journal of Medical Internet Research*. 2018;4 (1) :e4.7

 Sections[Abstract](#)[Introduction](#)[Methods](#)[Results](#)[Discussion](#)[Abbreviations](#)[References](#)[Copyright](#)[↑ Back to top](#)

Published on 09.01.18 in Vol 4, No 1 (2018): Jan-Mar

This paper is in the following e-collection/theme issue:

[Infoveillance, Infodemiology and Digital Disease Surveillance](#) [Infodemiology and Infoveillance](#)

[Article](#)[Cited By \(2\)](#)[Tweetations \(64\)](#)[Metrics](#) Original Paper

# Accurate Influenza Monitoring and Forecasting Using Novel Internet Data Streams: A Case Study in the Boston Metropolis

Fred Sun Lu<sup>1</sup>, AB  ; Suqin Hou<sup>2</sup>, MS  ; Kristin Baltrusaitis<sup>3</sup>, MS  ; Manan Shah<sup>4</sup>  ; Jure Leskovec<sup>4,5</sup>, PhD  ; Rok Susic<sup>4</sup>, PhD  ; Jared Hawkins<sup>1,6</sup>, MMSc, PhD  ; John Brownstein<sup>1,6</sup>, PhD  ; Giuseppe Conidi<sup>7</sup>, MPH  ; Julia Gunn<sup>7</sup>, RN, MPH  ; Josh Gray<sup>8</sup>, MBA  ; Anna Zink<sup>8</sup>, BA  ; Mauricio Santillana<sup>1,6</sup>, MS, PhD 

<sup>1</sup>Computational Health Informatics Program, Boston Children's Hospital, Boston, MA, United States

<sup>2</sup>Harvard Chan School of Public Health, Harvard University, Boston, MA, United States

<sup>3</sup>Department of Biostatistics, Boston University School of Public Health, Boston, MA, United States

<sup>4</sup>Computer Science Department, Stanford University, Stanford, CA, United States

<sup>5</sup>Chan Zuckerberg Biohub, San Francisco, CA, United States

<sup>6</sup>Department of Pediatrics, Harvard Medical School, Boston, MA, United States

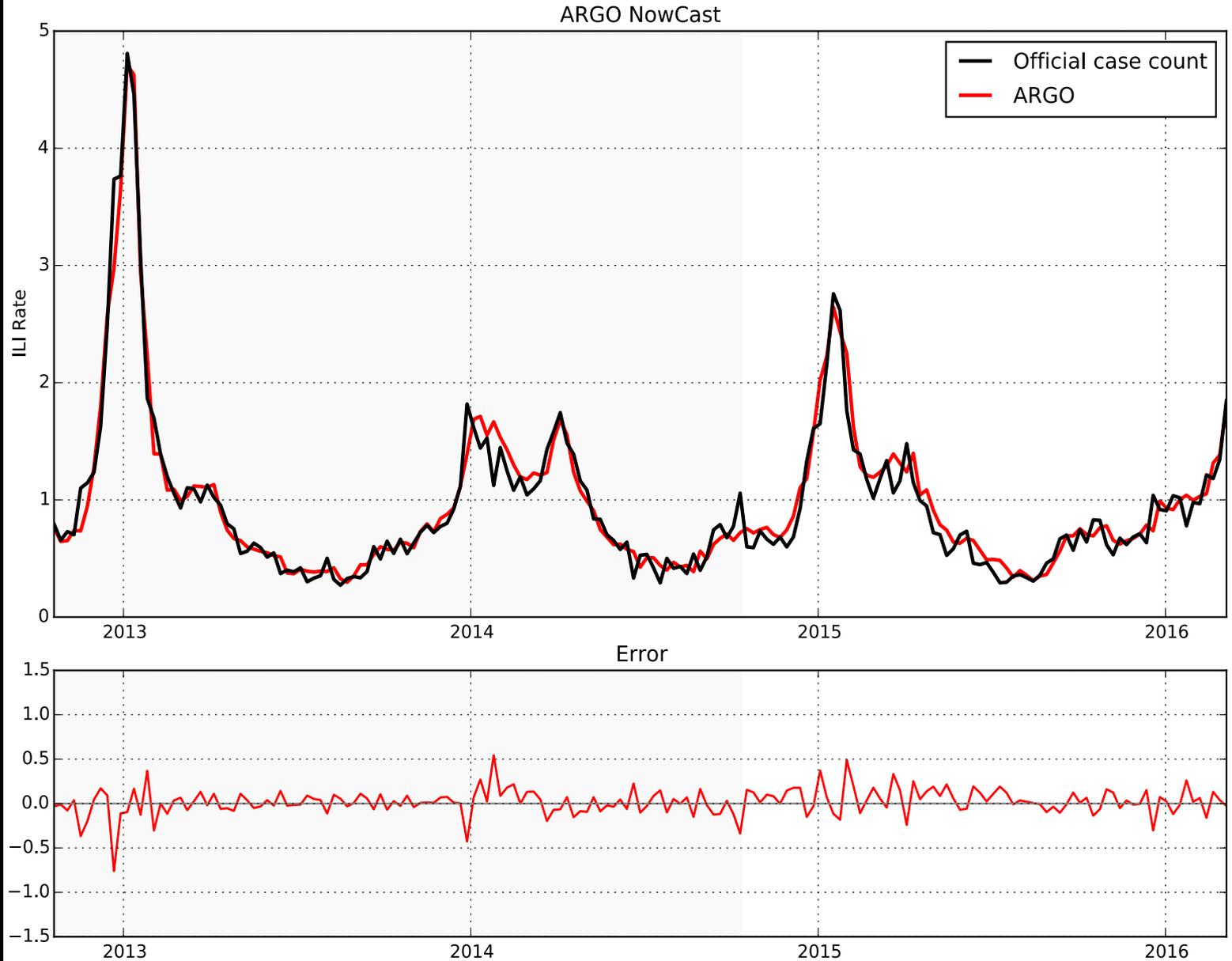
<sup>7</sup>Boston Public Health Commission, Boston, MA, United States

<sup>8</sup>athenaResearch, athenahealth, Watertown, MA, United States

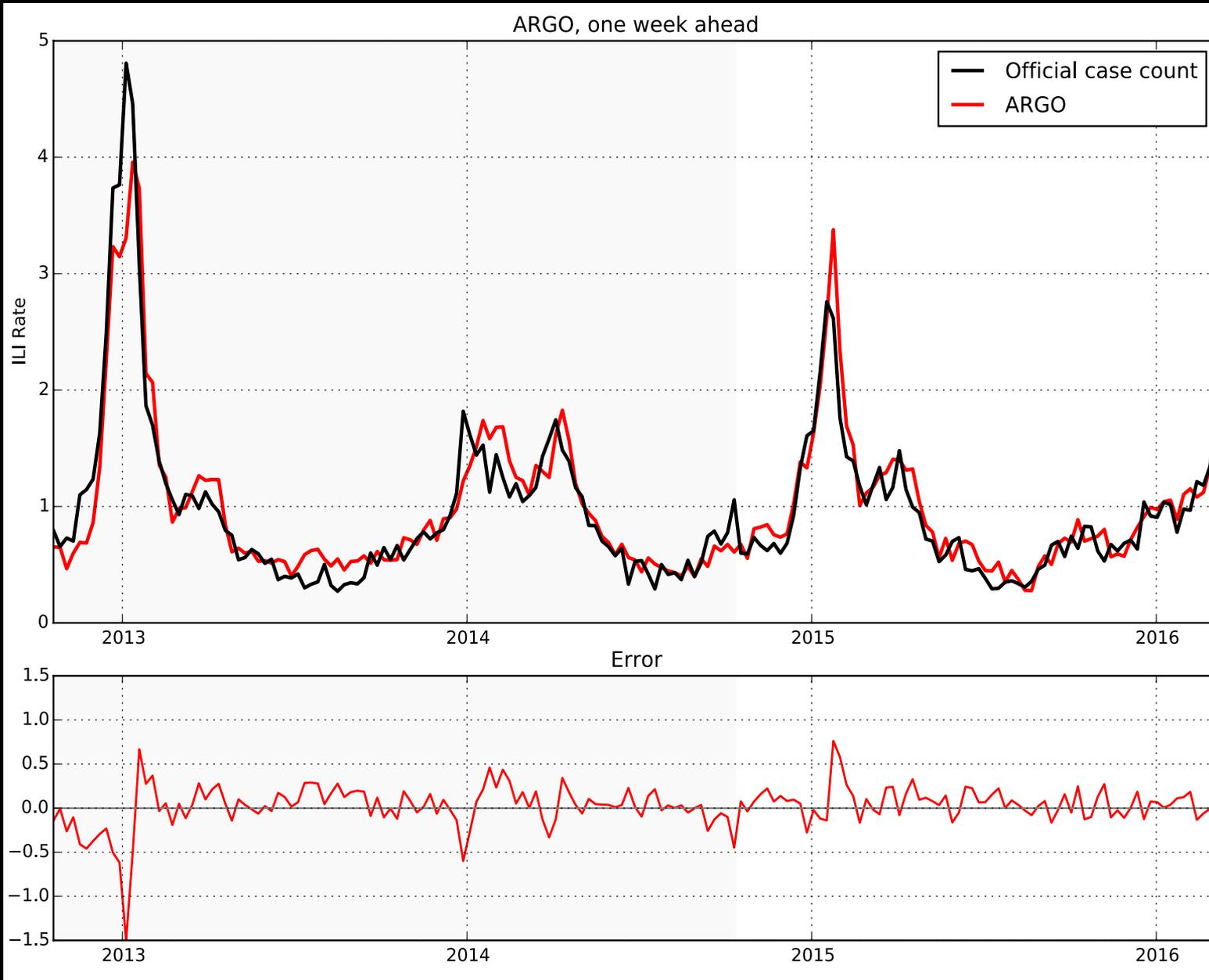
**Corresponding Author:**

Mauricio Santillana, MS, PhD

# Using multiple data sources to track flu in Boston



# Using multiple data sources to forecast flu in Boston

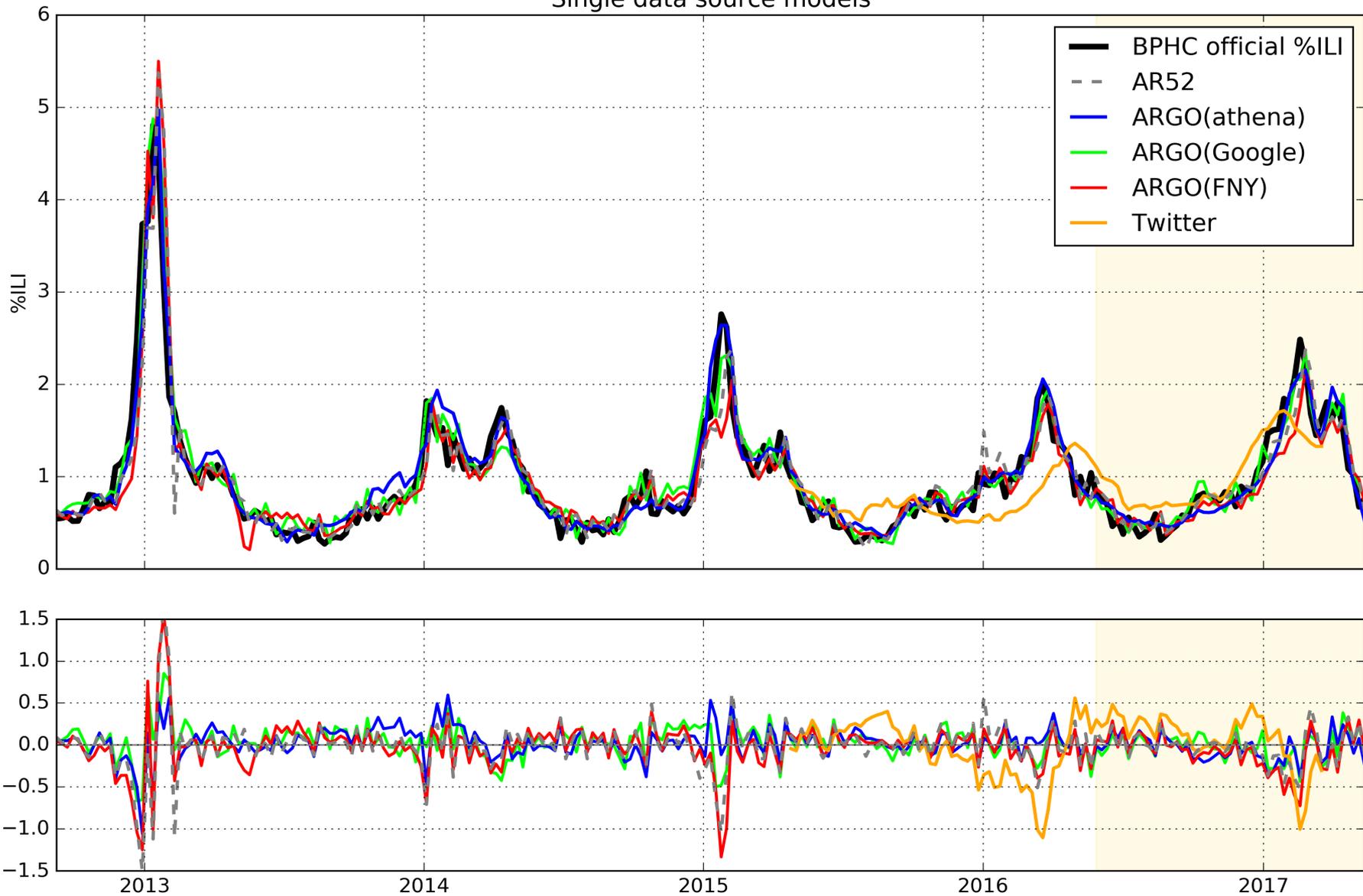


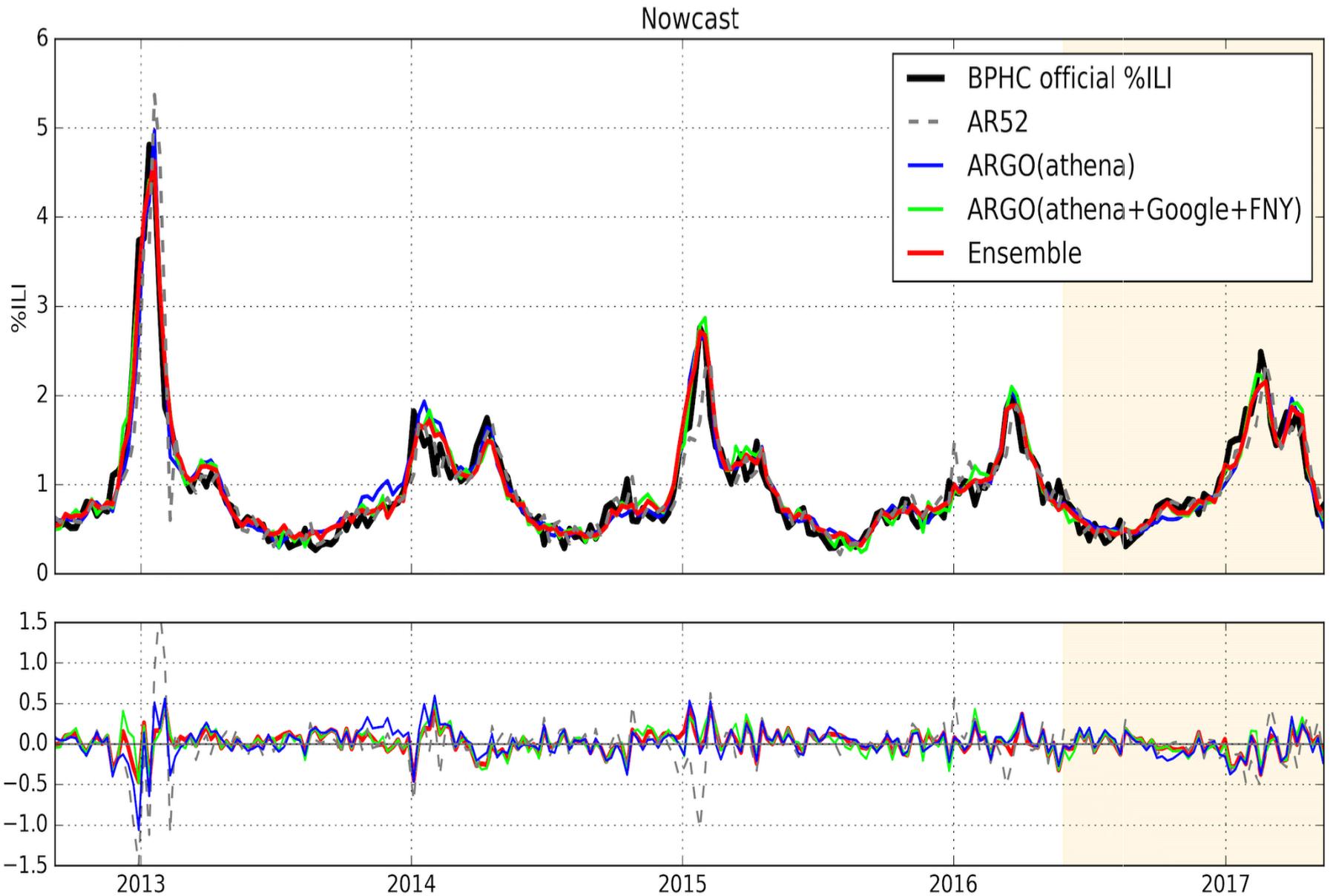
## Boston – out of sample predictions one-week ahead

### Model performance (individual data sources)

Model <sup>a</sup>	Whole period	Flu seasons				Holdout
	2012-16	2012-13	2013-14	2014-15	2015-16	2016-17
<b>Root mean square error</b>						
AR52	0.303	0.577	0.199	0.305	0.217	0.229
ARGO(athena) <sup>b</sup>	0.195	0.306	0.229	0.192	0.133	0.182
ARGO(Google)	0.206	0.312	0.194	0.247	0.161	0.188
ARGO(FNY) <sup>c</sup>	0.299	0.552	0.204	0.343	0.172	0.280
Twitter	—	—	—	0.162	0.427	0.351
GFT <sup>d</sup>	—	0.352	0.271	0.284	—	—
Naive	0.266	0.481	0.208	0.280	0.202	0.219

Single data source models





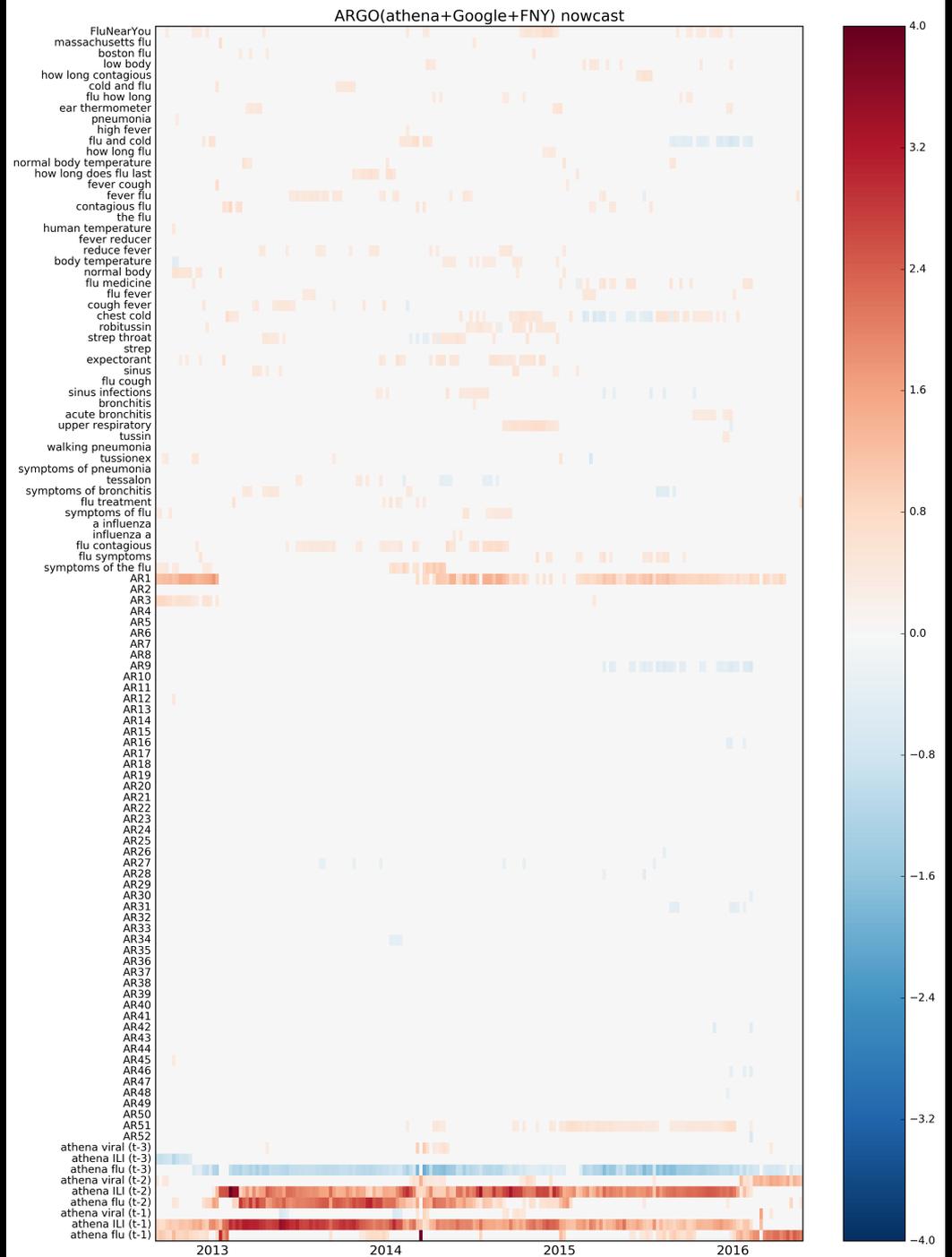
## Model performance (individual data sources)

Model <sup>a</sup>	Whole period	Flu seasons				Holdout
	2012-16	2012-13	2013-14	2014-15	2015-16	2016-17
<b>Root mean square error</b>						
AR52	0.303	0.577	0.199	0.305	0.217	0.229
ARGO(athena) <sup>b</sup>	0.195	0.306	0.229	0.192	0.133	0.182
ARGO(Google)	0.206	0.312	0.194	0.247	0.161	0.188
ARGO(FNY) <sup>c</sup>	0.299	0.552	0.204	0.343	0.172	0.280
Twitter	—	—	—	0.162	0.427	0.351
GFT <sup>d</sup>	—	0.352	0.271	0.284	—	—
Naive	0.266	0.481	0.208	0.280	0.202	0.219

## Model performance (ensemble approaches)

	Whole period (2012-2016)	Flu seasons				Validation 2016-2017
		2012-2013	2013-2014	2014-2015	2015-2016	
<b>RMSE</b>						
AR52	0.303	0.577	0.199	0.305	0.217	0.229
athena	0.208	0.377	<b>0.163</b>	0.219	0.144	0.169
Google	0.222	0.218	0.243	0.303	0.221	0.195
athena+Google	0.186	0.210	0.200	0.224	0.194	<b>0.145</b>
ARGO(athena)	0.195	0.306	0.229	0.192	<b>0.133</b>	0.182
ARGO(Google)	0.206	0.312	0.194	0.247	0.161	0.188
ARGO(athena+Google+FNY)	<b>0.165</b>	<b>0.199</b>	0.192	<b>0.189</b>	0.168	0.156
Ensemble	<b>0.151</b>	<b>0.193</b>	<b>0.170</b>	<b>0.176</b>	<b>0.139</b>	<b>0.150</b>

When combined, what are the strongest predictors?

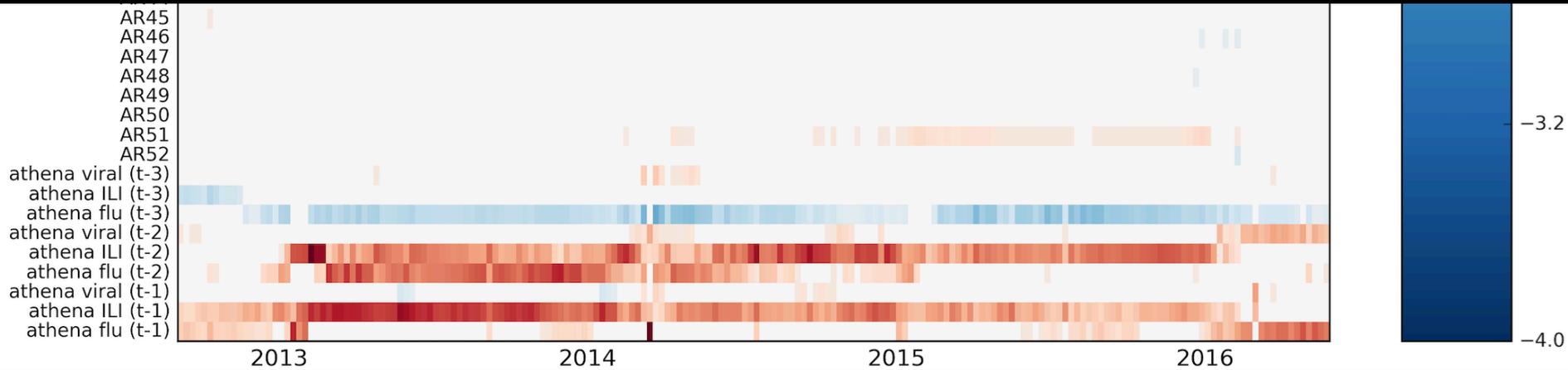


# When combined, what are the strongest predictors?

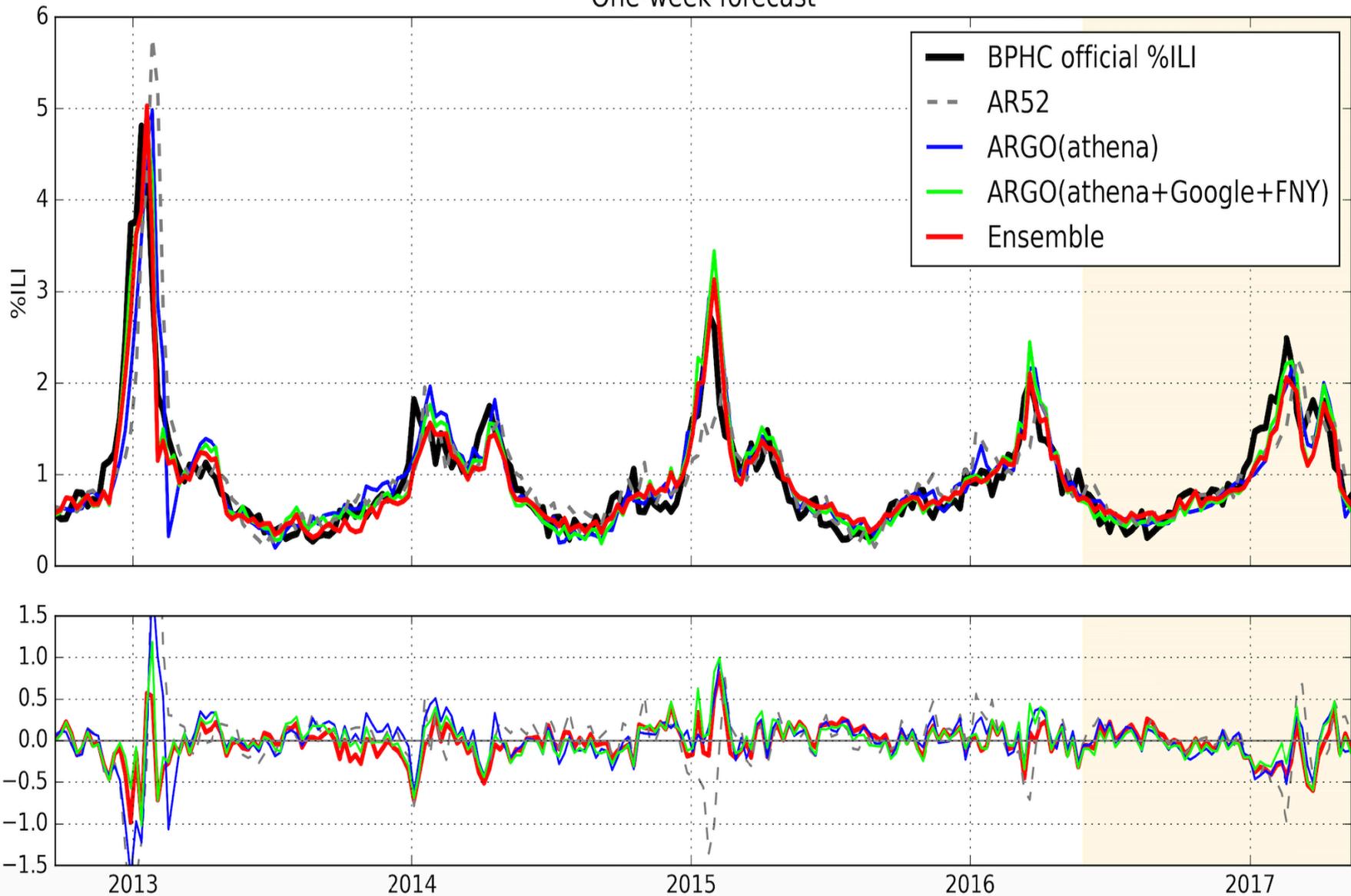
ARGO(athena+Google+FNY) nowcast



When combined, what are the strongest predictors?

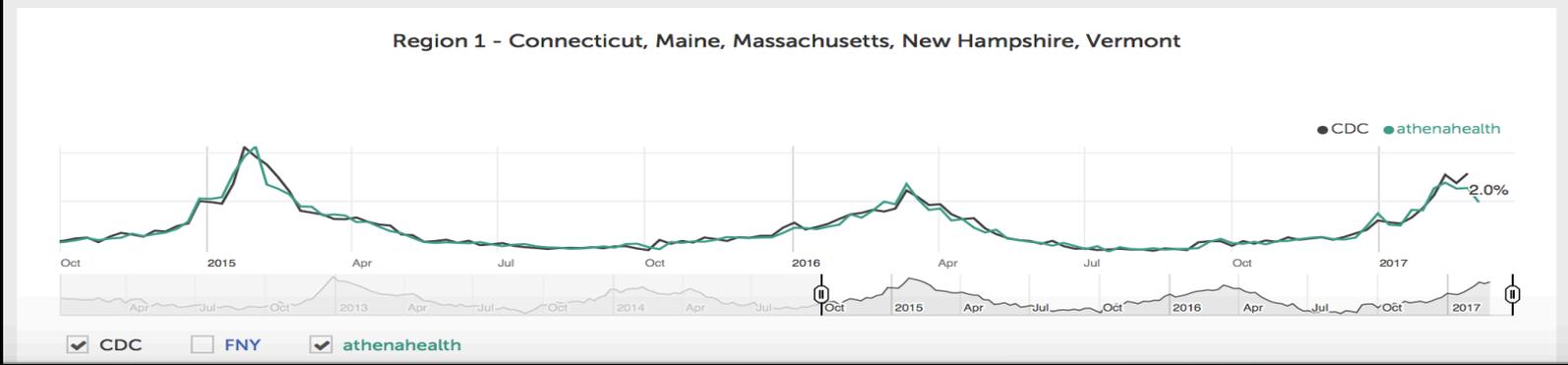
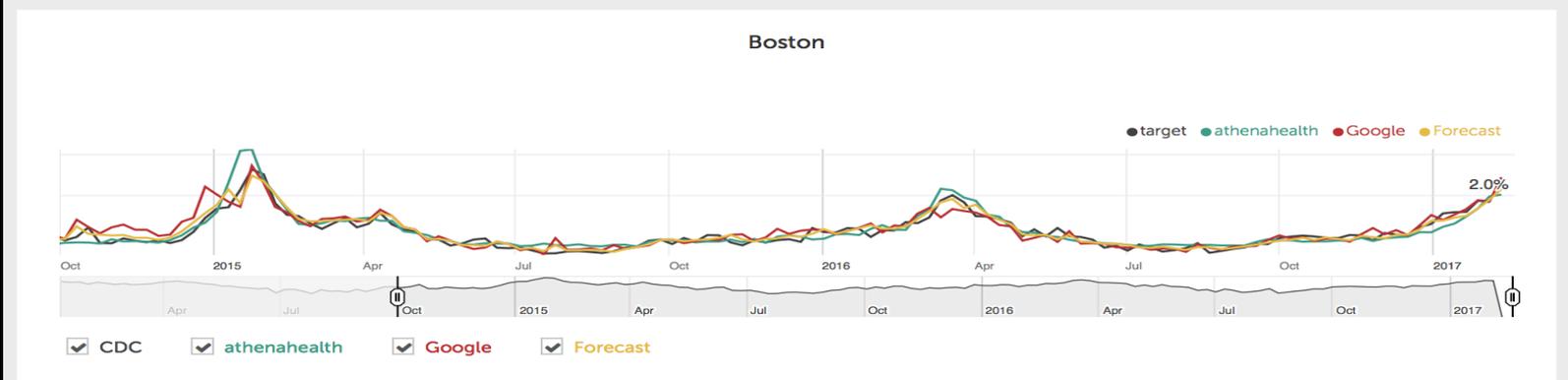
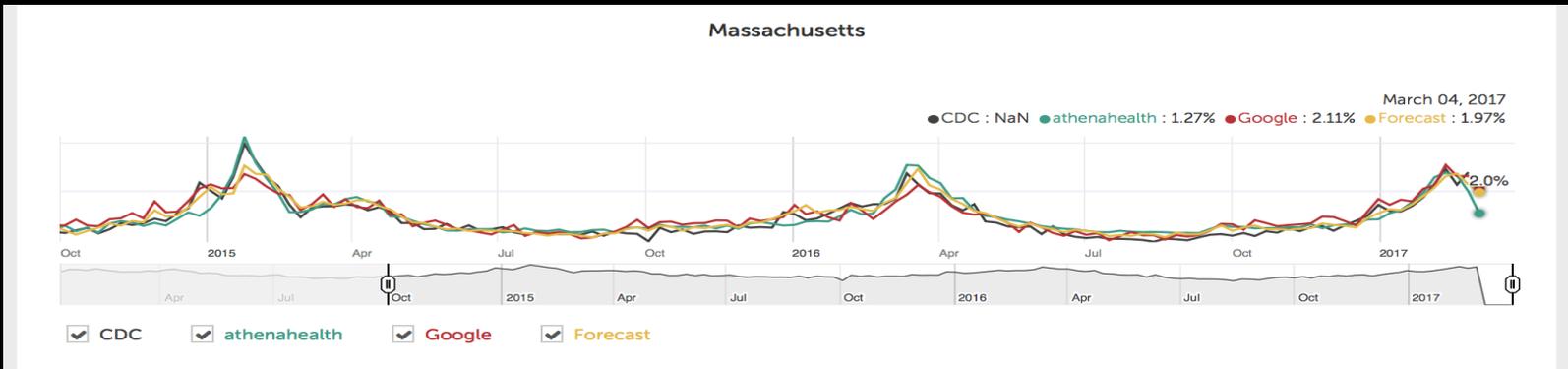


# One-week forecast



# Aim is to display these predictions in a joint CDC-BCH website

## Using multiple data sources to track flu at the state-level in the USA



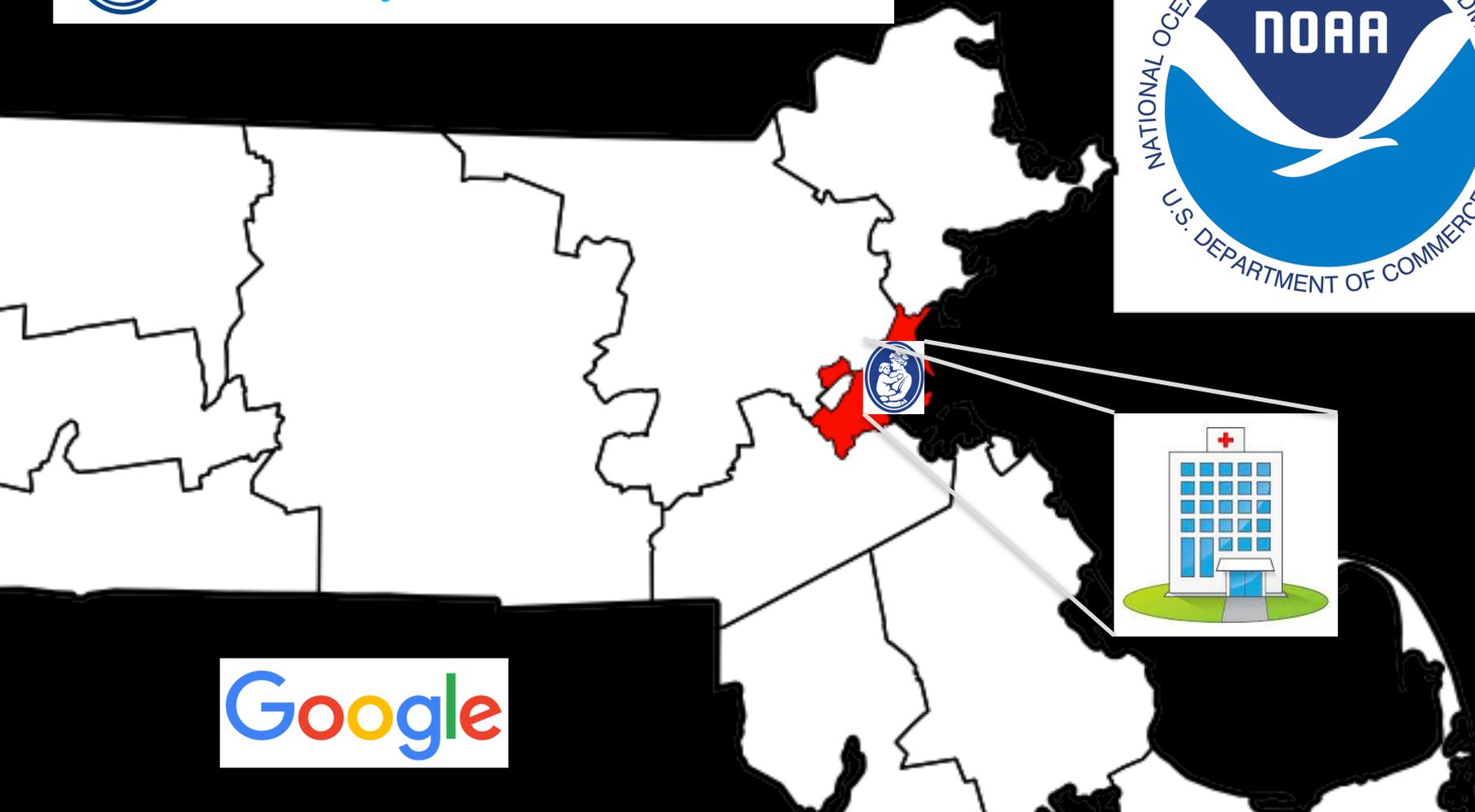
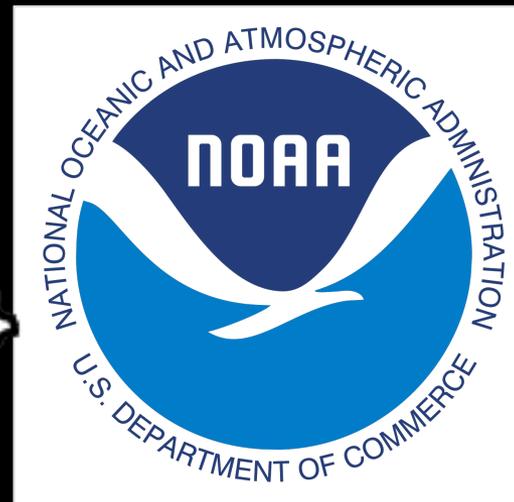
## Hyper-local predictions

Can we predict daily emergency department visits in a hospital?

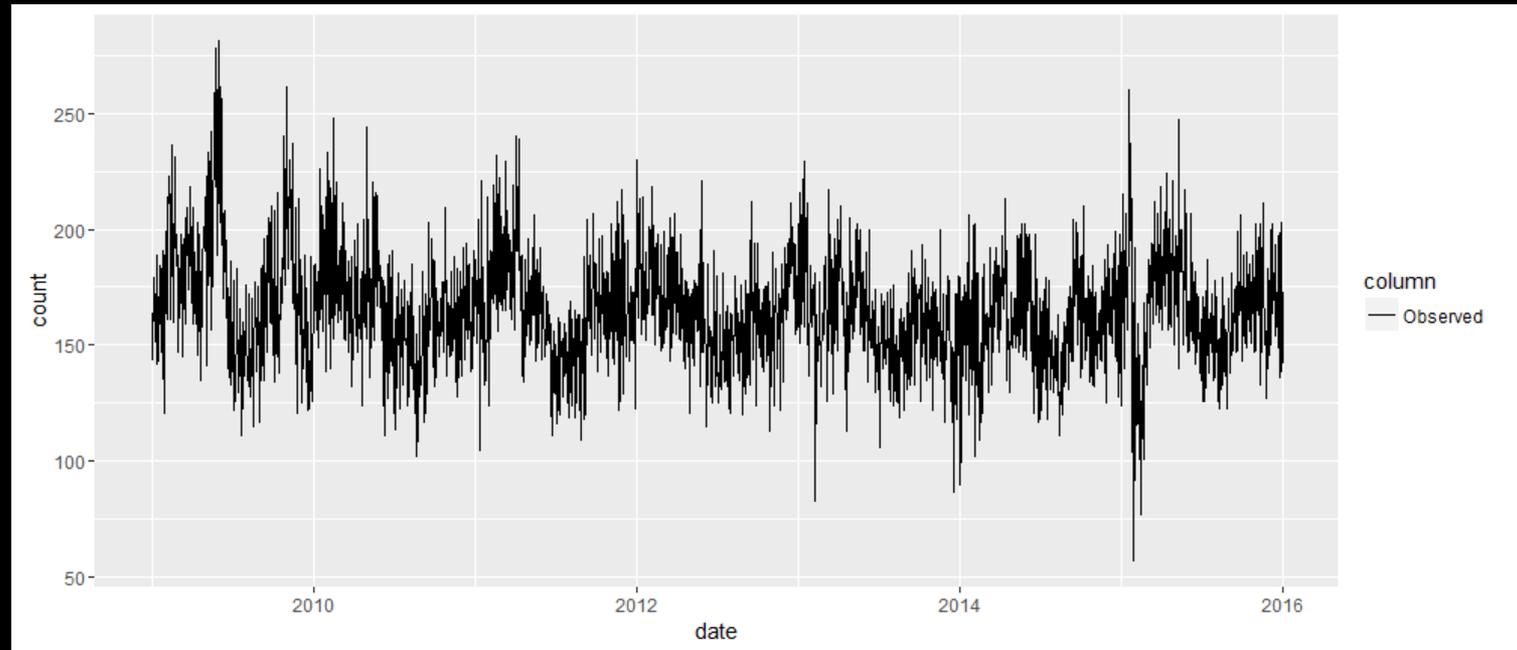


**Boston Children's Hospital**

Until every child is well<sup>SM</sup>

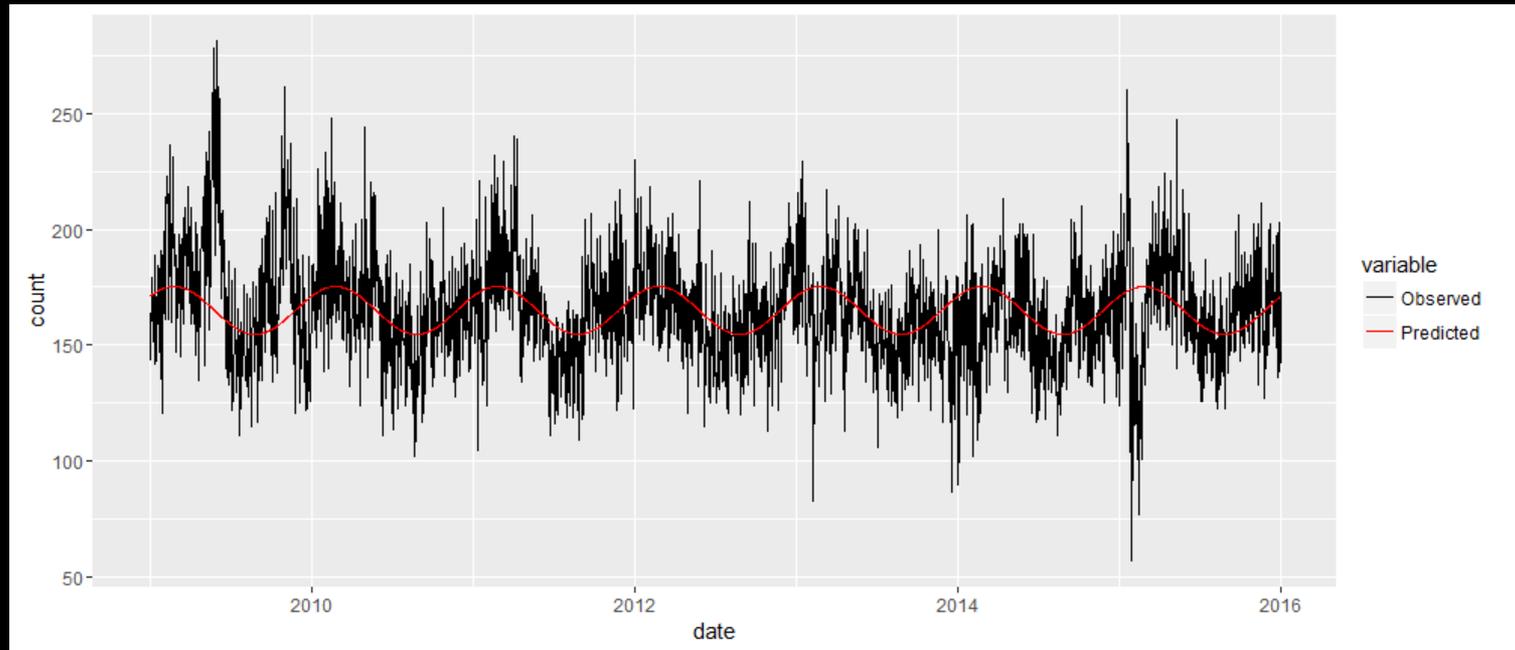


# Daily Visits 2009-2015

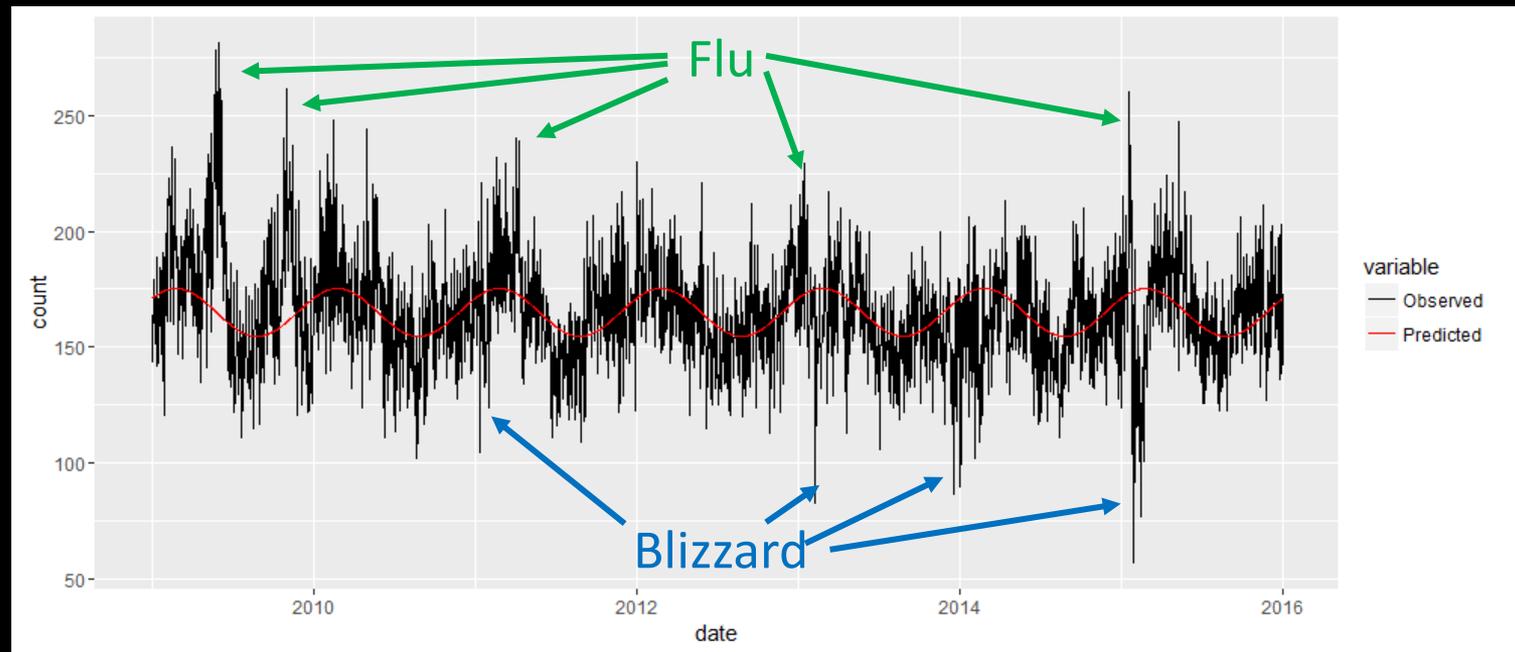


In collaboration with: Sam Tideman, Mauricio Santillana, Jon Bickel, and Ben Reis

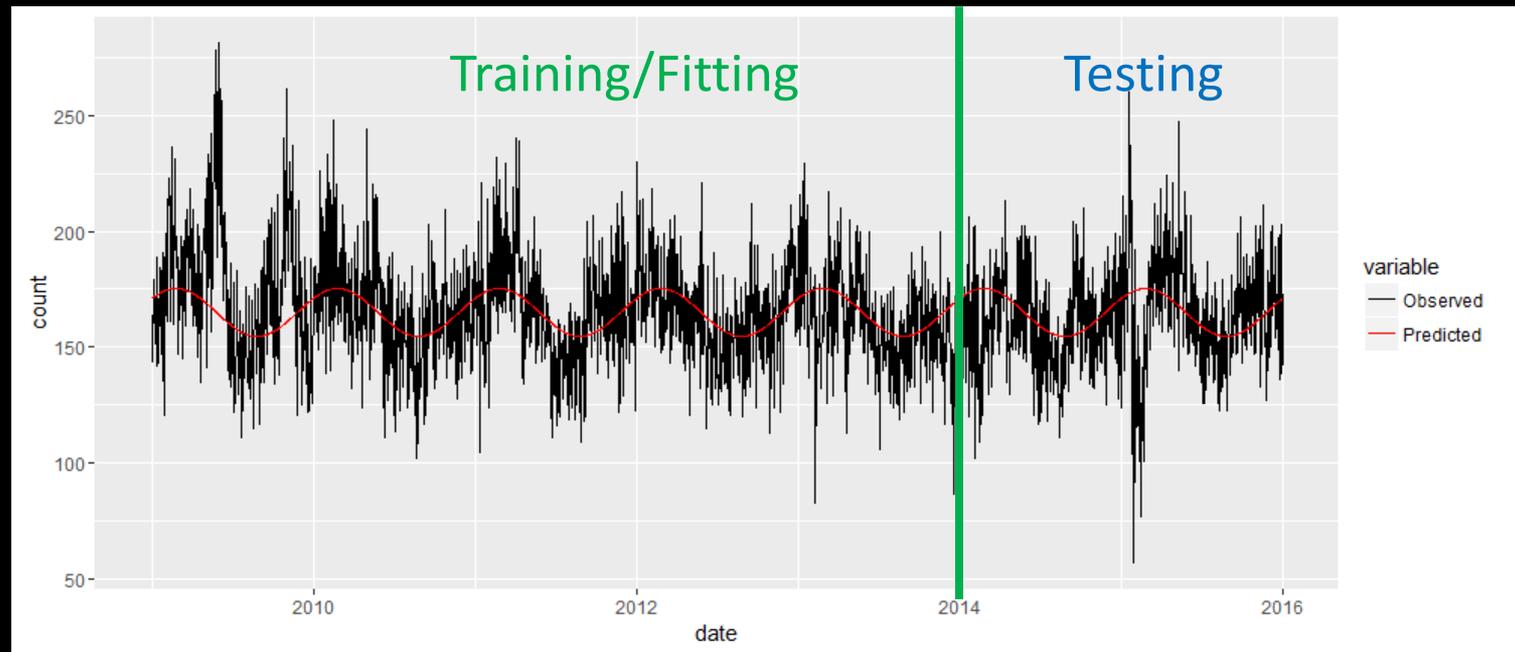
# Seasonal Trend



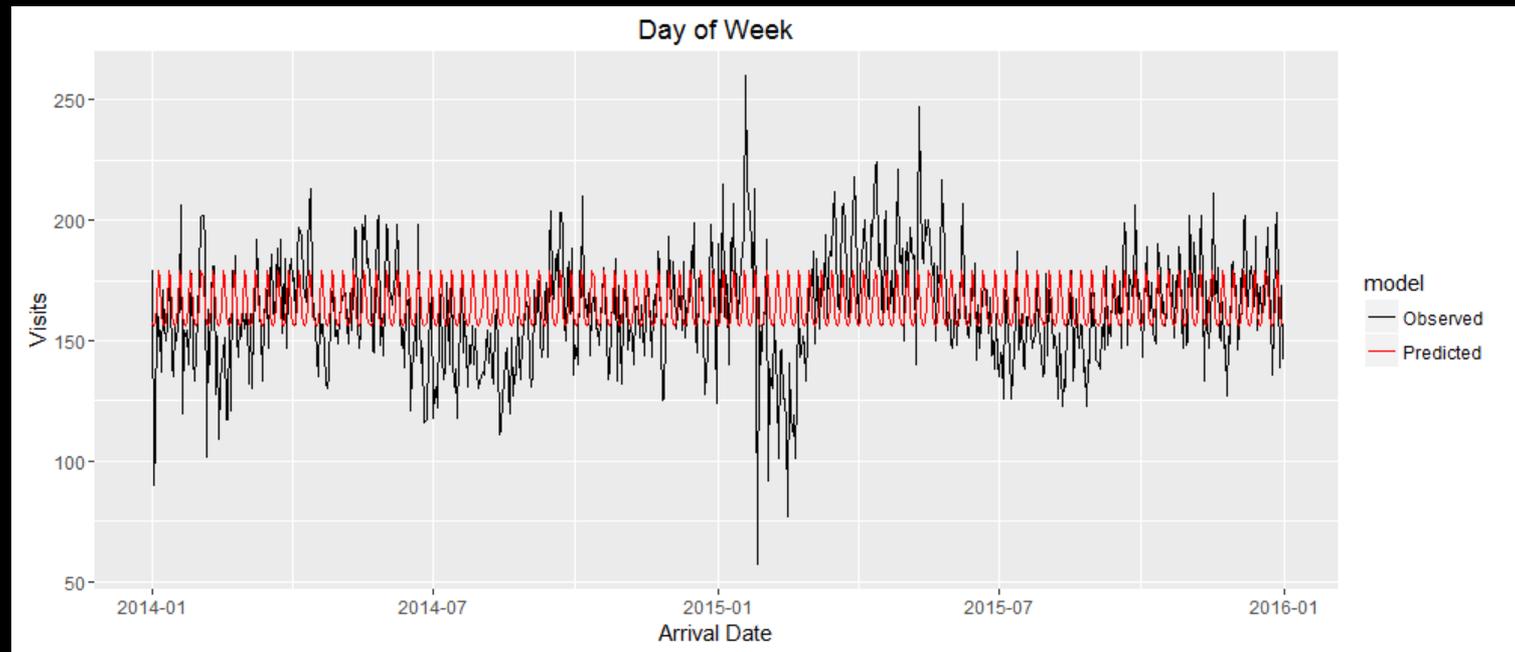
# Noticeable Events



# Split data for modeling



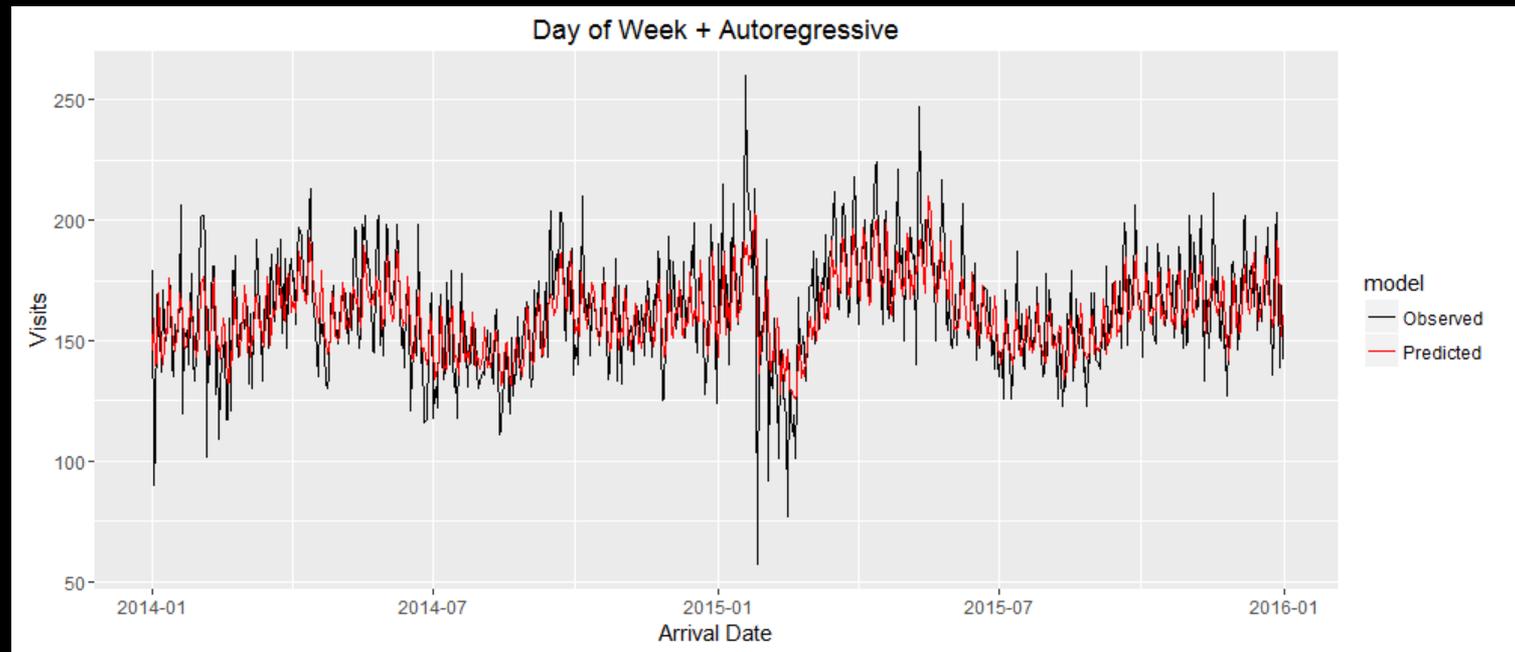
# Current Staffing model = Day of Week



MAPE = 11.0%

Percent of days with bad staffing= 11.2%

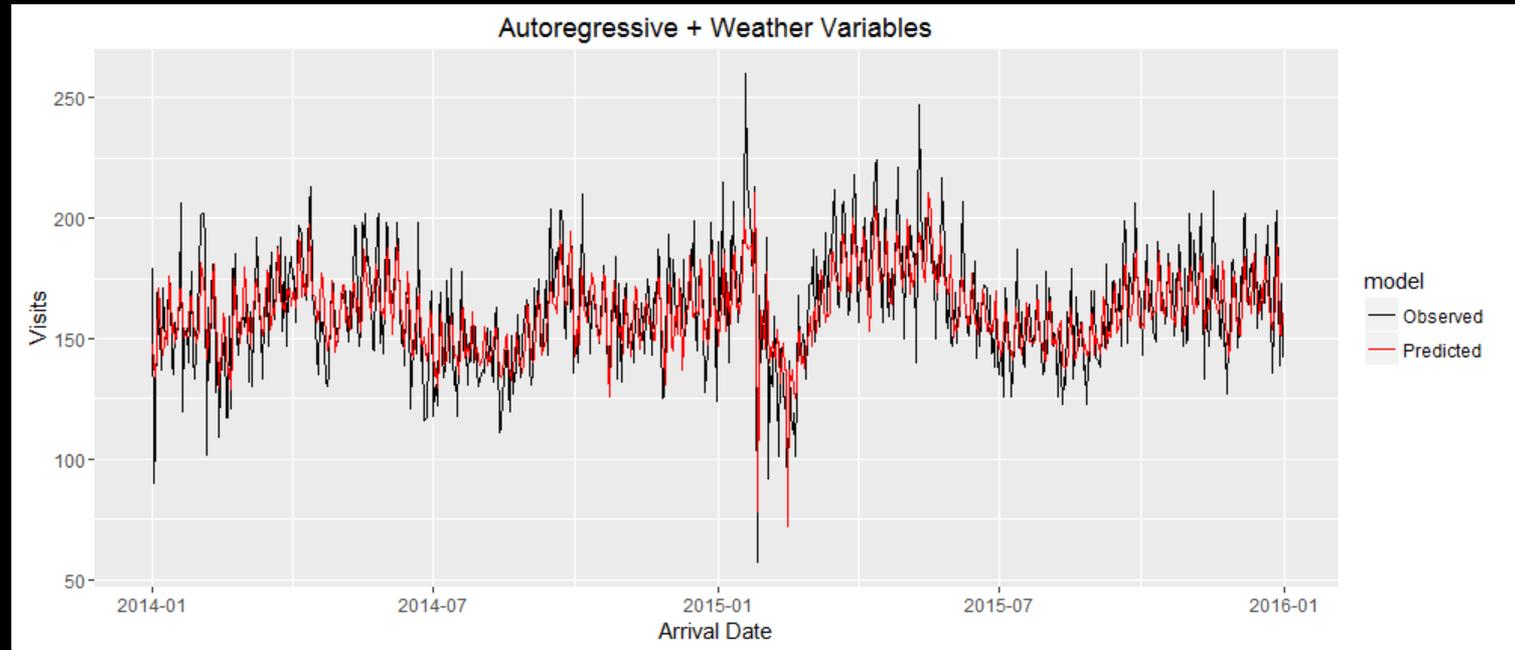
# Add in Autoregression



MAPE = 8.4%

Percent of days with bad staffing= 4.9%

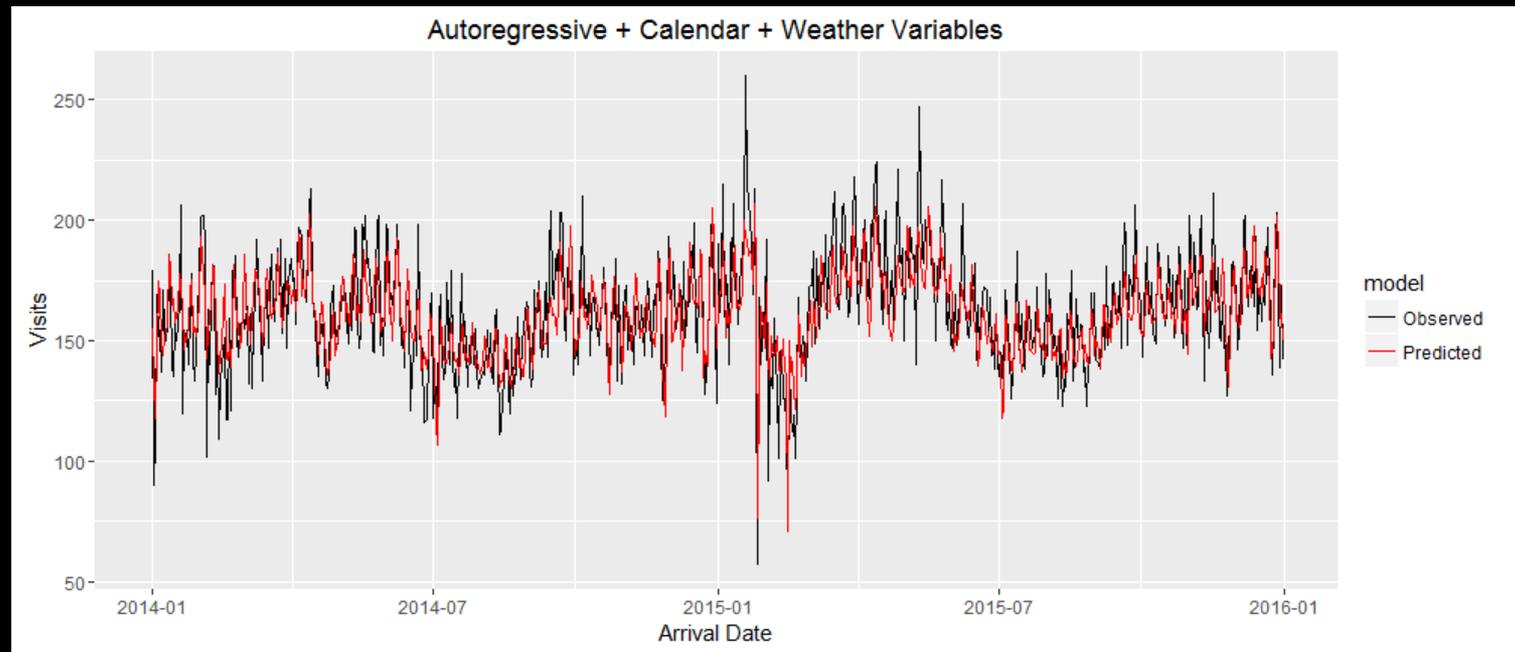
# Add in Weather Data



MAPE = 7.9%

Percent of days with bad staffing= 4.8%

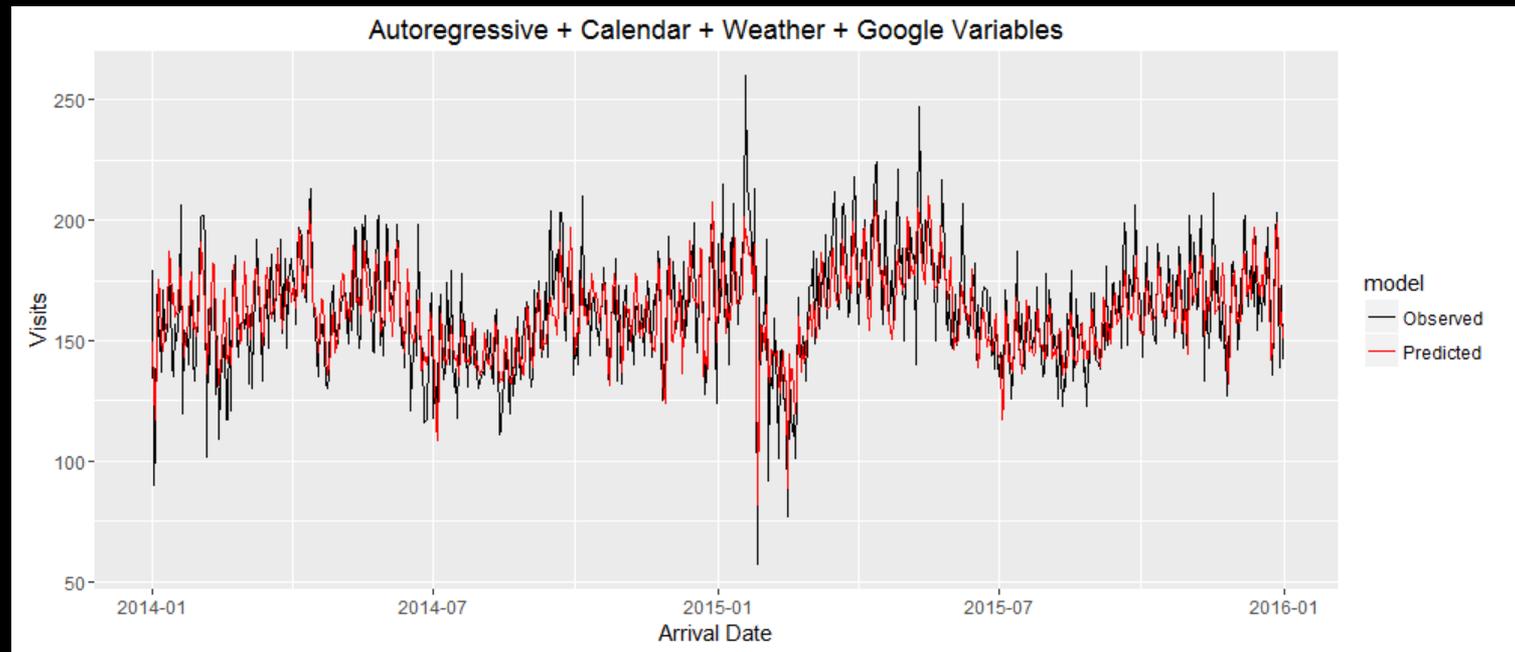
# Add in Calendar Data



MAPE = 7.7%

Percent of days with bad staffing = 3.8%

# Add in Google Data

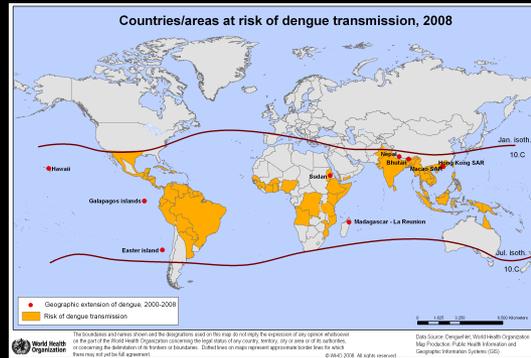


MAPE = 7.6%

Percent of days with bad staffing = 3.3%

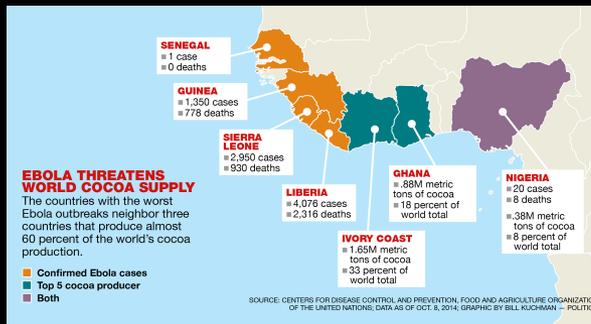
# Part 2. Success stories in tracking and forecasting Flu, Zika, Dengue, Ebola in data-poor medium- to low-income countries.

## Dengue, Zika, and Flu



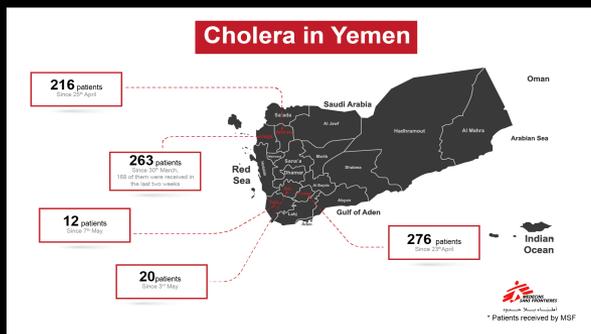
- Latin America (Flu, Zika, Dengue)
- South-east Asia (Dengue)

## Ebola



- West Africa

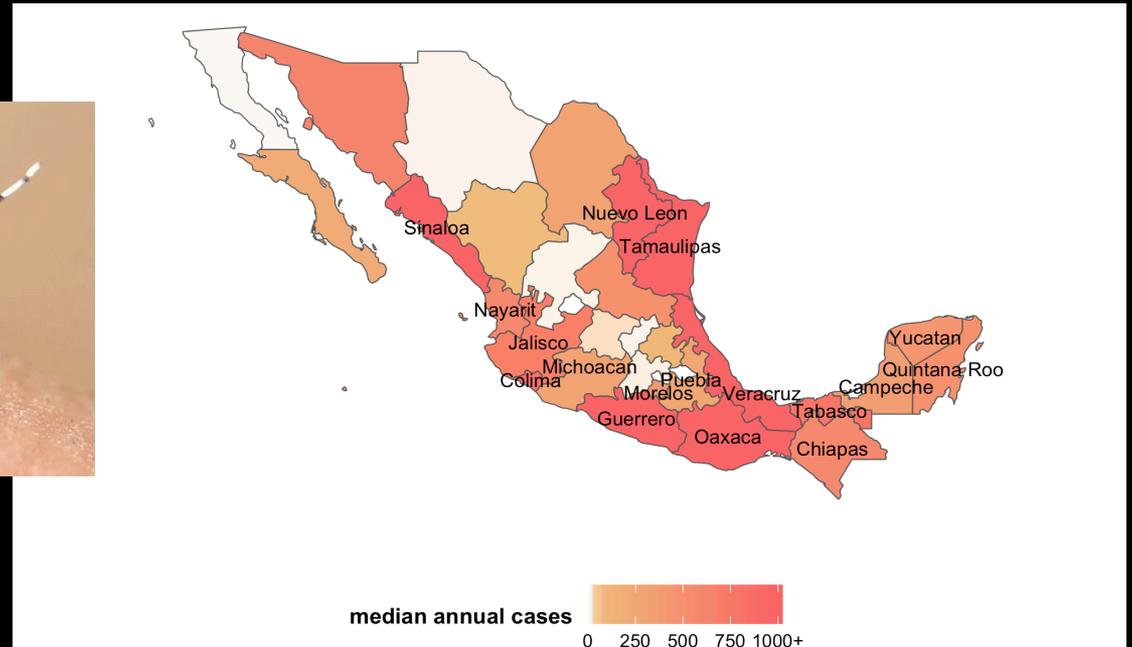
## Cholera



- Middle East

# Forecasting Dengue Incidence in Mexico

Establishing a prediction baseline



Team:

*Mauricio Santillana* (BCH, Harvard),  
*Michael Johansson* (CDC Puerto Rico),  
*Aditi Hota* (Columbia Univ),  
*John Brownstein* (BCH, Harvard),  
*Nick Reich* (Umass Amherst)



Altmetric: 10 Citations: 4

[More detail >>](#)

Article | [OPEN](#)

# Evaluating the performance of infectious disease forecasts: A comparison of climate-driven and seasonal dengue forecasts for Mexico

Michael A. Johansson , Nicholas G. Reich, Aditi Hota, John S. Brownstein & Mauricio Santillana 

*Scientific Reports* **6**, Article number: 33707  
(2016)

doi:10.1038/srep33707

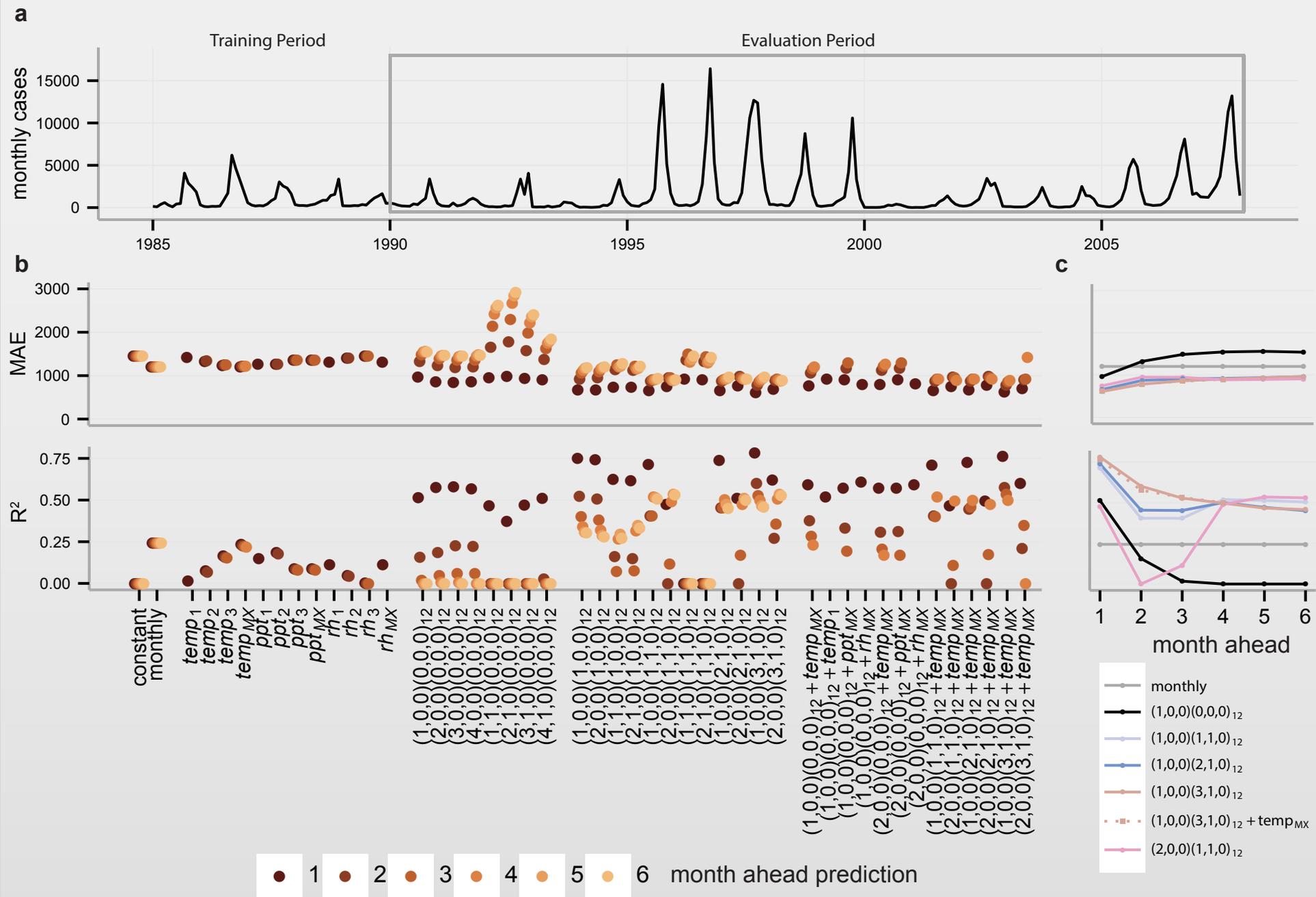
[Download Citation](#)

Received: 17 March 2016

Accepted: 24 August 2016

Published online: 26 September 2016

# Mexico Dengue incidence (Country-level)





Extending use of Google searches to track Dengue in other countries:

Latin America: Mexico, Brazil

Southeast Asia: Thailand, Singapore, Taiwan

 OPEN ACCESS  PEER-REVIEWED

RESEARCH ARTICLE

# Advances in using Internet searches to track dengue

Shihao Yang, Samuel C. Kou , Fred Lu, John S. Brownstein, Nicholas Brooke, Mauricio Santillana 

Published: July 20, 2017 • <https://doi.org/10.1371/journal.pcbi.1005607>

Article

Authors

Metrics

Comments

Related Content



## Abstract

Author summary

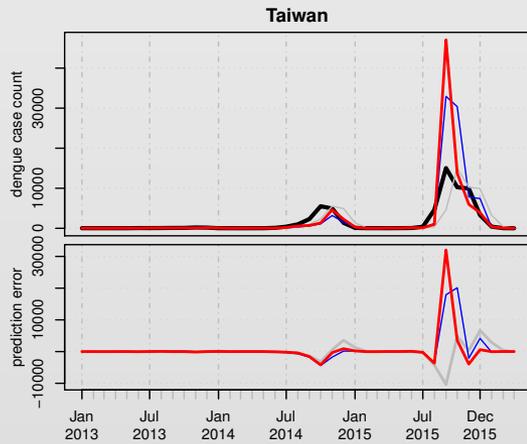
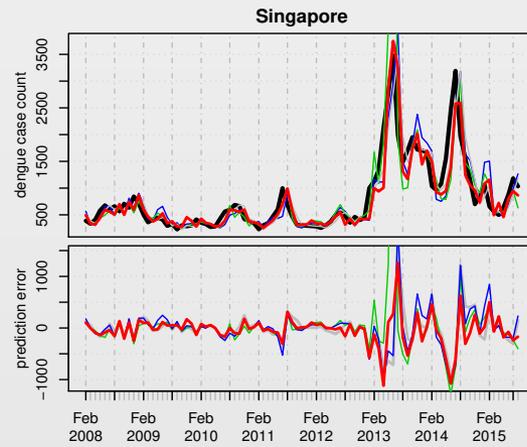
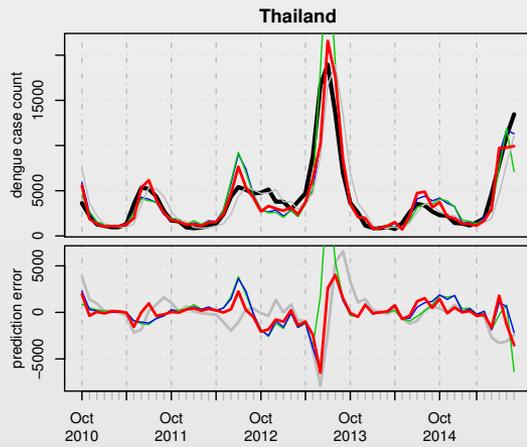
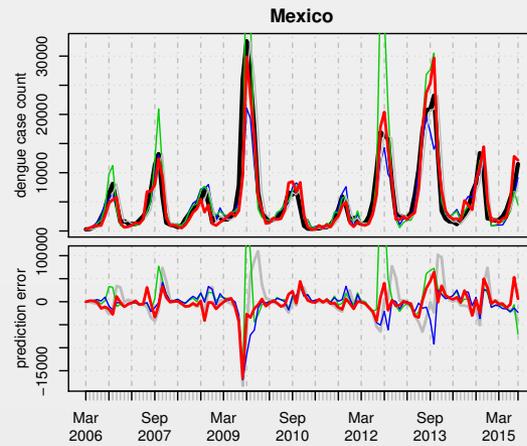
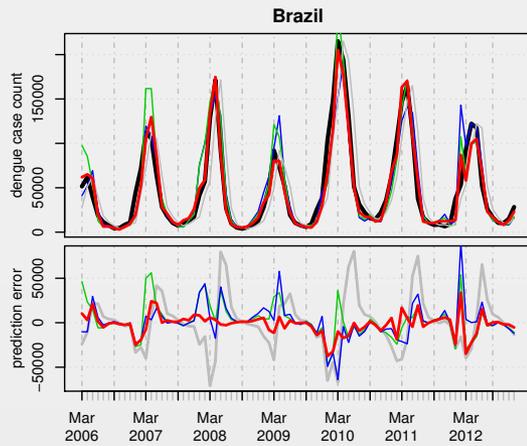
Introduction

Materials and methods

Results

## Abstract

Dengue is a mosquito-borne disease that threatens over half of the world's population. Despite being endemic to more than 100 countries, government-led efforts and tools for timely identification and tracking of new infections are still lacking in many affected areas. Multiple methodologies that leverage the use of Internet-based data sources have been proposed as a



- Target
- ARGO
- SAR
- SAR + GDT
- naive

# 2014 Ebola Outbreak: Media Events Track Changes in Observed Reproductive Number

APRIL 28, 2015 · COMMENTARY

 [Print or Save PDF](#)

 [Citation](#)

 [XML](#)

 [Email](#)

 [Tweet](#)

 [Like](#)

 10

## ■ AUTHORS

[Maimuna S. Majumder](#) [Sheryl Kluberg](#) [Mauricio Santillana](#) [Sumiko Mekaru](#) [John S. Brownstein](#)

## ■ ABSTRACT

In this commentary, we consider the relationship between early outbreak changes in the observed reproductive number of Ebola in West Africa and various media reported interventions and aggravating events. We find that media reports of interventions that provided education, minimized contact, or strengthened healthcare were typically followed by sustained transmission reductions in both Sierra Leone and Liberia. Meanwhile, media reports of aggravating events generally preceded temporary transmission increases in both countries. Given these preliminary findings, we conclude that media reported events could potentially be incorporated into future epidemic modeling efforts to improve mid-outbreak case projections.

## Data-poor environments (Zika)

A more recent contribution on the 2015 Latin American Zika outbreak



JMIR Publications



JMIR Public Health and Surveillance

Published on 01.06.16 in Vol 2, No 1 (2016): Jan-Jun

This paper is in the following e-collection/theme issue:

[Infoveillance, Infodemiology and Digital Disease Surveillance](#) [Infodemiology and Infoveillance](#)

Article

Cited By (0)

Tweetations (29)

Metrics

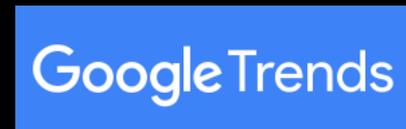
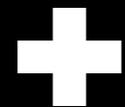
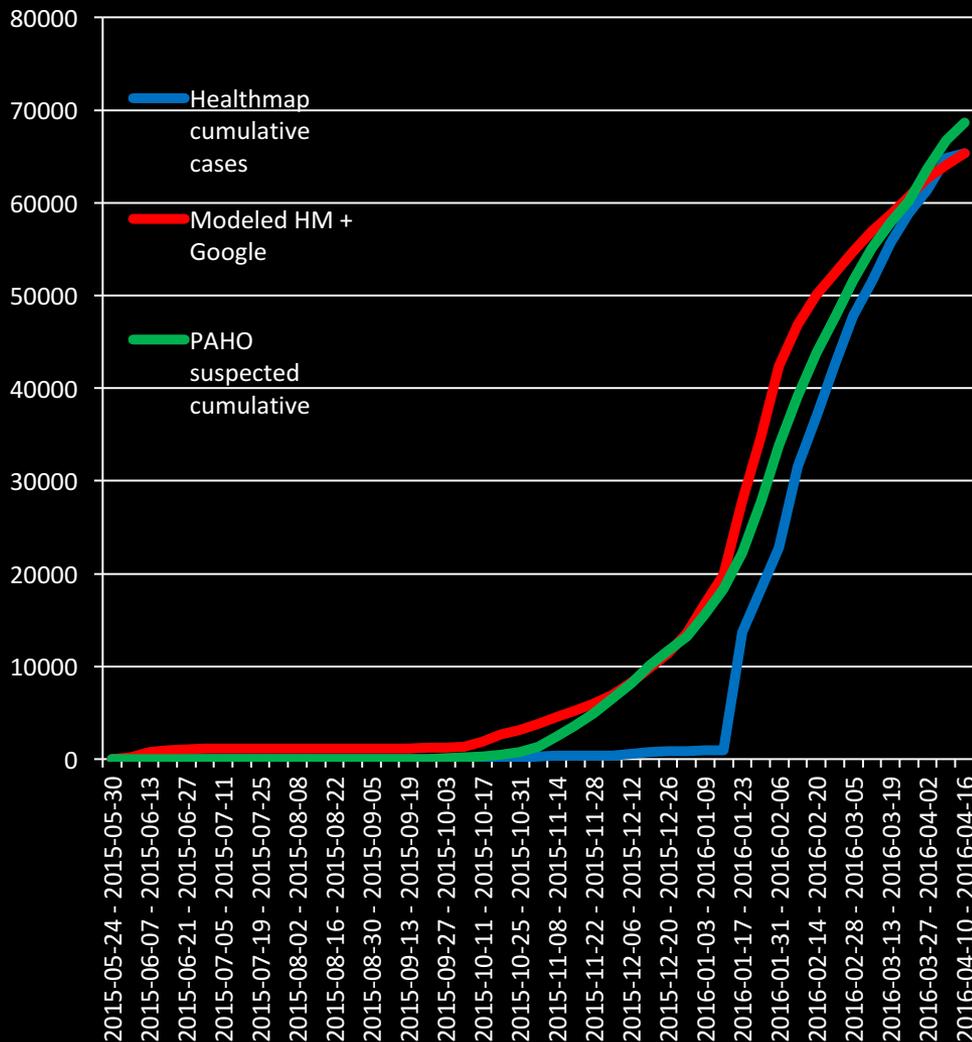
Original Paper

# Utilizing Nontraditional Data Sources for Near Real-Time Estimation of Transmission Dynamics During the 2015-2016 Colombian Zika Virus Disease Outbreak

Maimuna S Majumder<sup>1,2</sup>, MPH ; Mauricio Santillana<sup>1,3,4</sup>, PhD ; Sumiko R Mekaru<sup>1,5</sup>, PhD ; Denise P McGinnis<sup>1</sup>, ScD ;  
Kamran Khan<sup>6,7</sup>, MD ; John S Brownstein<sup>1,4</sup>, PhD

# Data-poor environments (Zika)

## A more recent contribution on the 2015 Latin American Zika outbreak



When we gained access to government-lead disease surveillance information, we found great similarity with the curve we produced ahead of the publication of this information.

# Forecasting Zika using Google searches and Twitter



 OPEN ACCESS  PEER-REVIEWED

RESEARCH ARTICLE

## Forecasting Zika Incidence in the 2016 Latin America Outbreak Combining Traditional Disease Surveillance with Search, Social Media, and News Report Data

Sarah F. McGough , John S. Brownstein, Jared B. Hawkins, Mauricio Santillana 

Version 2  Published: January 13, 2017 • <http://dx.doi.org/10.1371/journal.pntd.0005295>

15 Save	0 Citation
3,819 View	17 Share

<b>Article</b> 	<b>Authors</b>	<b>Metrics</b>	<b>Comments</b>	<b>Related Content</b>
--	----------------	----------------	-----------------	------------------------

<b>Download PDF</b> 	
<b>Print</b>	<b>Share</b>

 Check for updates

- Abstract**
- Author Summary
- Introduction
- Methods
- Results
- Discussion
- Supporting Information

### Abstract

#### Background

Over 400,000 people across the Americas are thought to have been infected with Zika virus as a consequence of the 2015–2016 Latin American outbreak. Official government-led case count data in Latin America are typically delayed by several weeks, making it difficult to track the disease in a timely manner. Thus, timely disease tracking systems are needed to design and

Included in the  
Following Collection

[Zika](#)

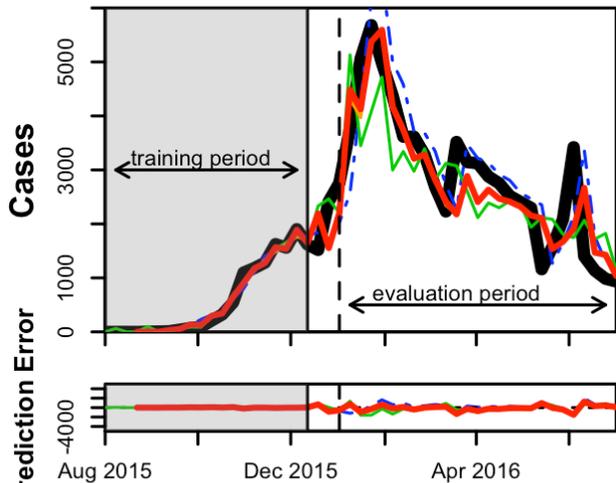
Subject Areas



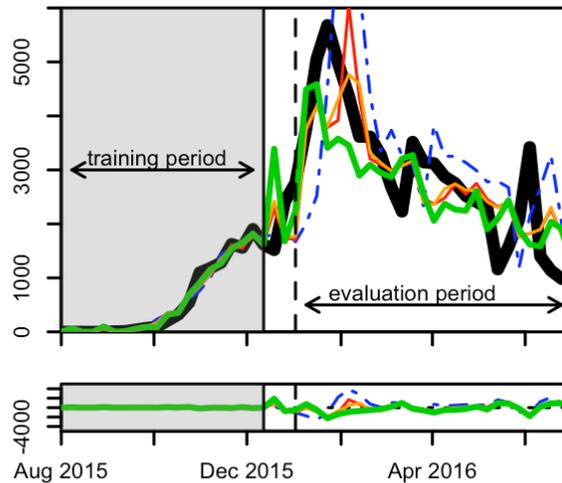
# Forecasting Zika using Google searches and Twitter (with Sarah McGough)

(a) Colombia

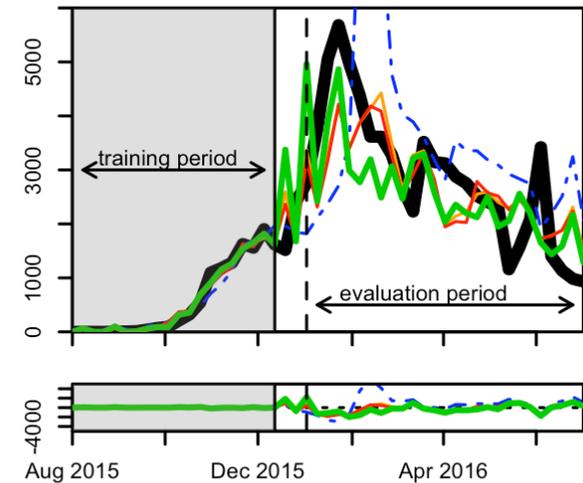
1 Week Ahead



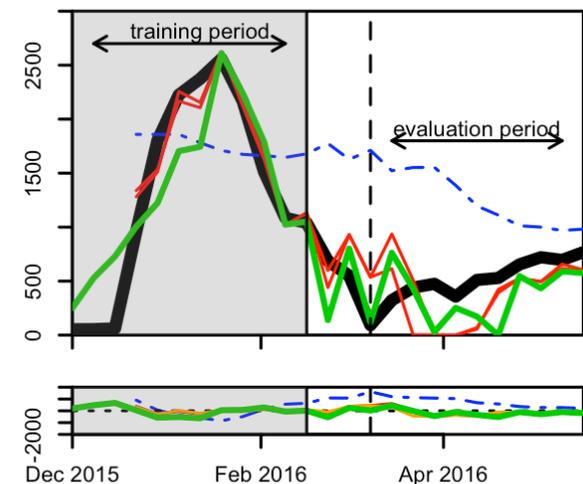
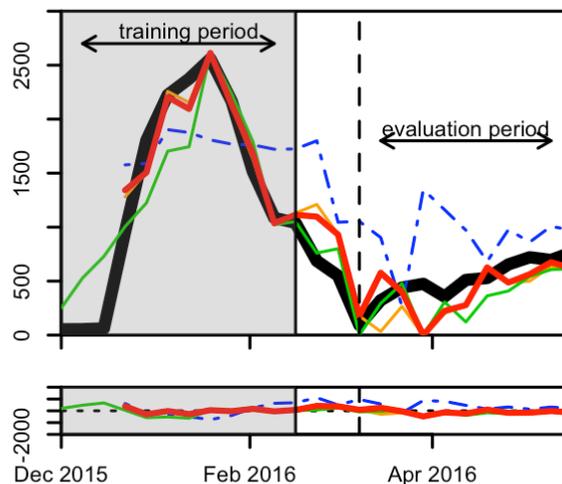
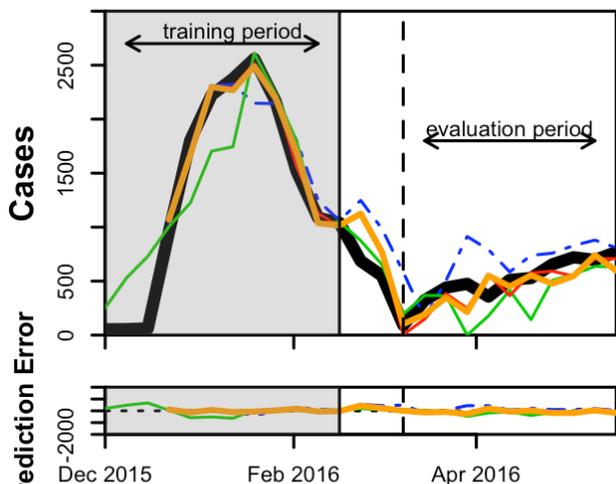
2 Weeks Ahead



3 Weeks Ahead

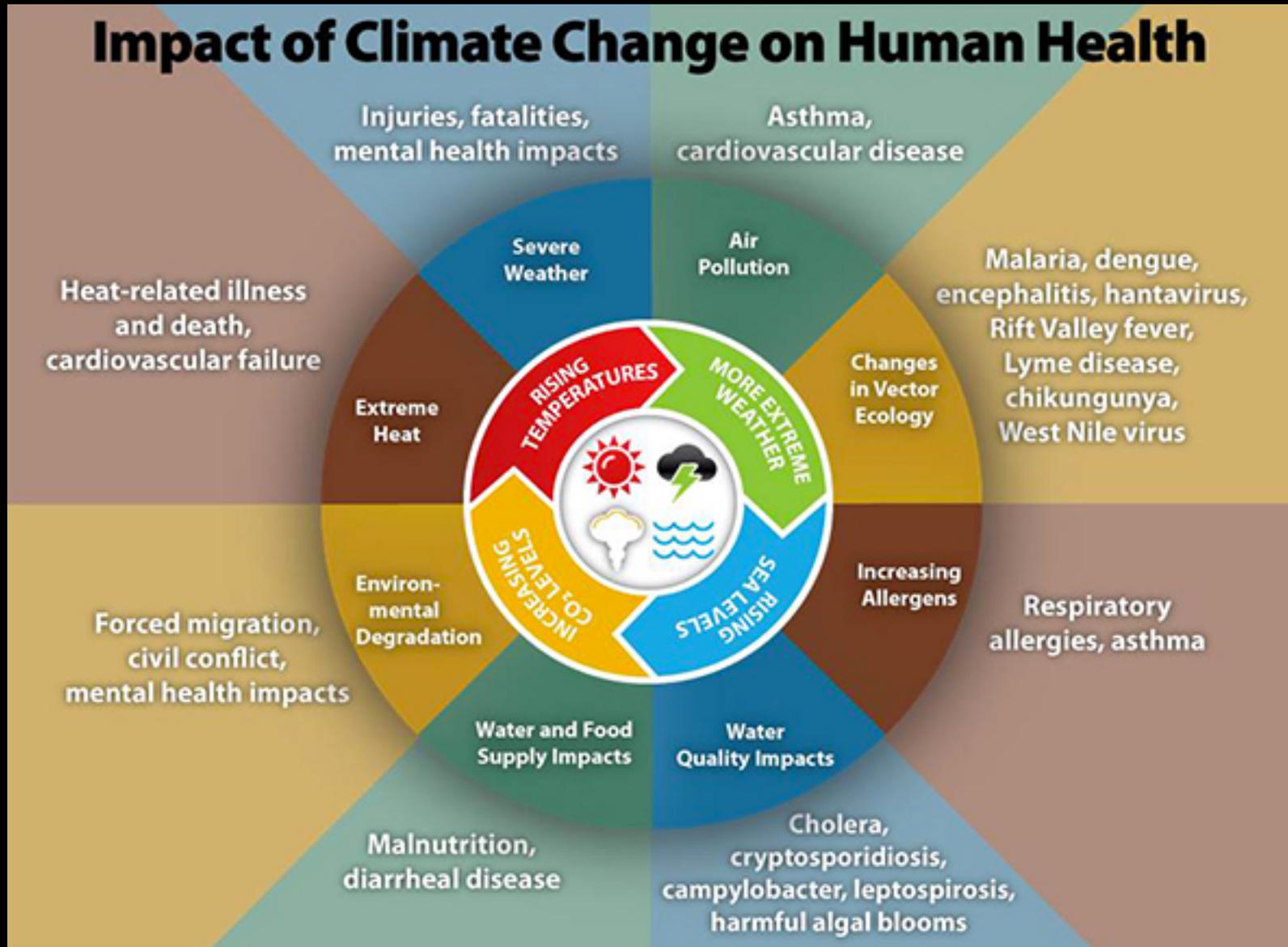


(b) Honduras



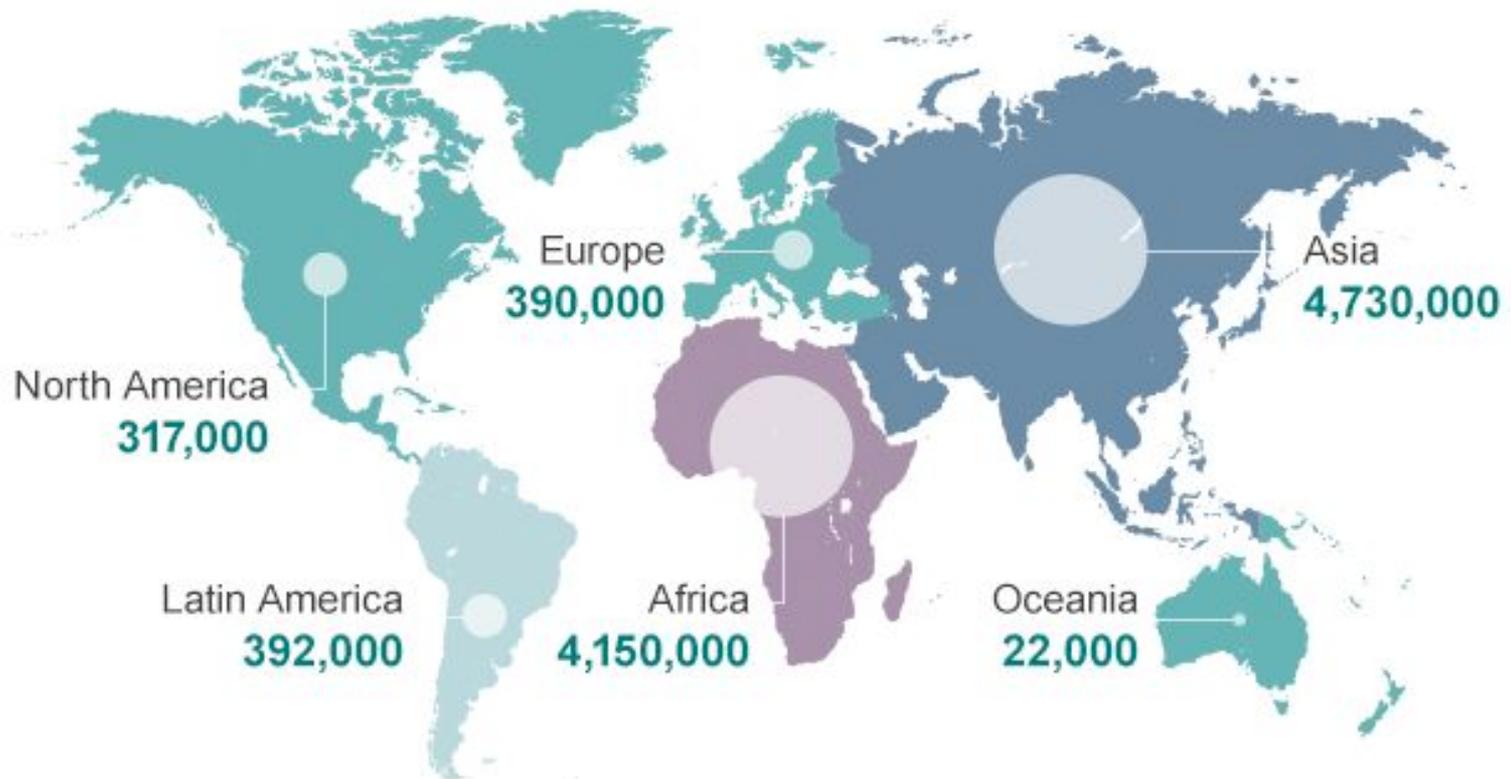
— Observed    - - - AR    — G+T    — ARGO+T    — ARGO+TH

# The influence of weather/climate on public health



**Antibiotic resistance** is now recognized as one of the worlds greatest public health threats, with the potential to render existing antibiotics ineffective in the “not so far” future.

Deaths attributable to antimicrobial resistance every year by 2050



Source: Review on Antimicrobial Resistance 2014

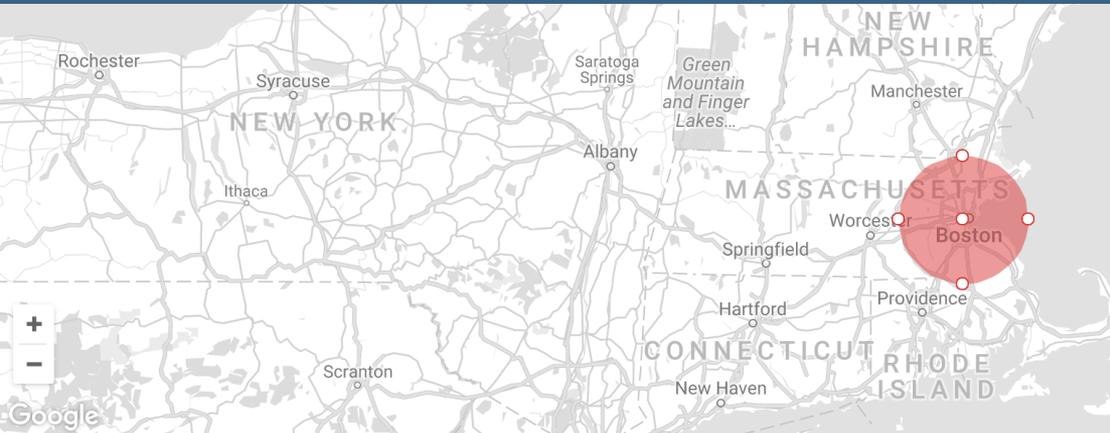


Tool conceived and implemented by *Derek R. MacFadden and John S. Brownstein*

# HealthMap ResistanceOpen

HealthMap ResistanceOpen About

Login



Location

24876 isolates in a 25 mile radius

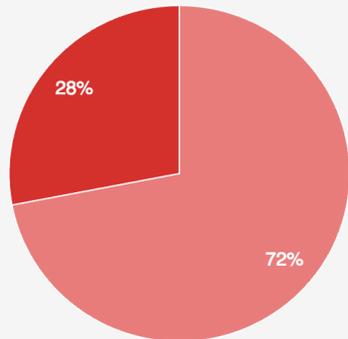
- Years: 2013,2014,2015
- Specimens: Urine,Blood,Respiratory,Sterile,Non-sterile
- Sources: Inpatient,ER,Outpatient

Map data ©2017 Google Terms of Use Report a map error

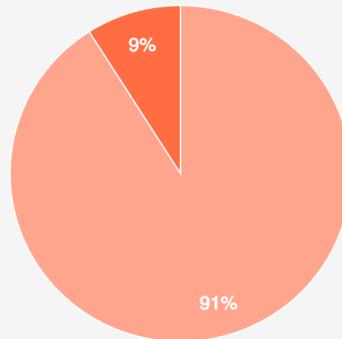
## Antibiotic Resistant Superbugs in Your Area

Search

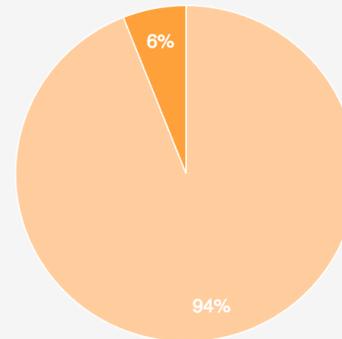
MRSA



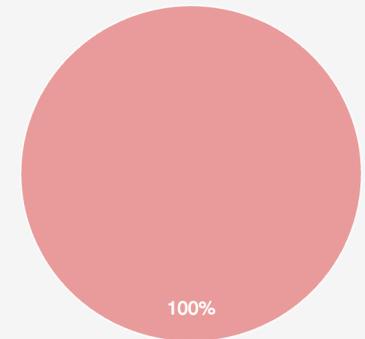
VRE



3rd Gen. Ceph. Resistance

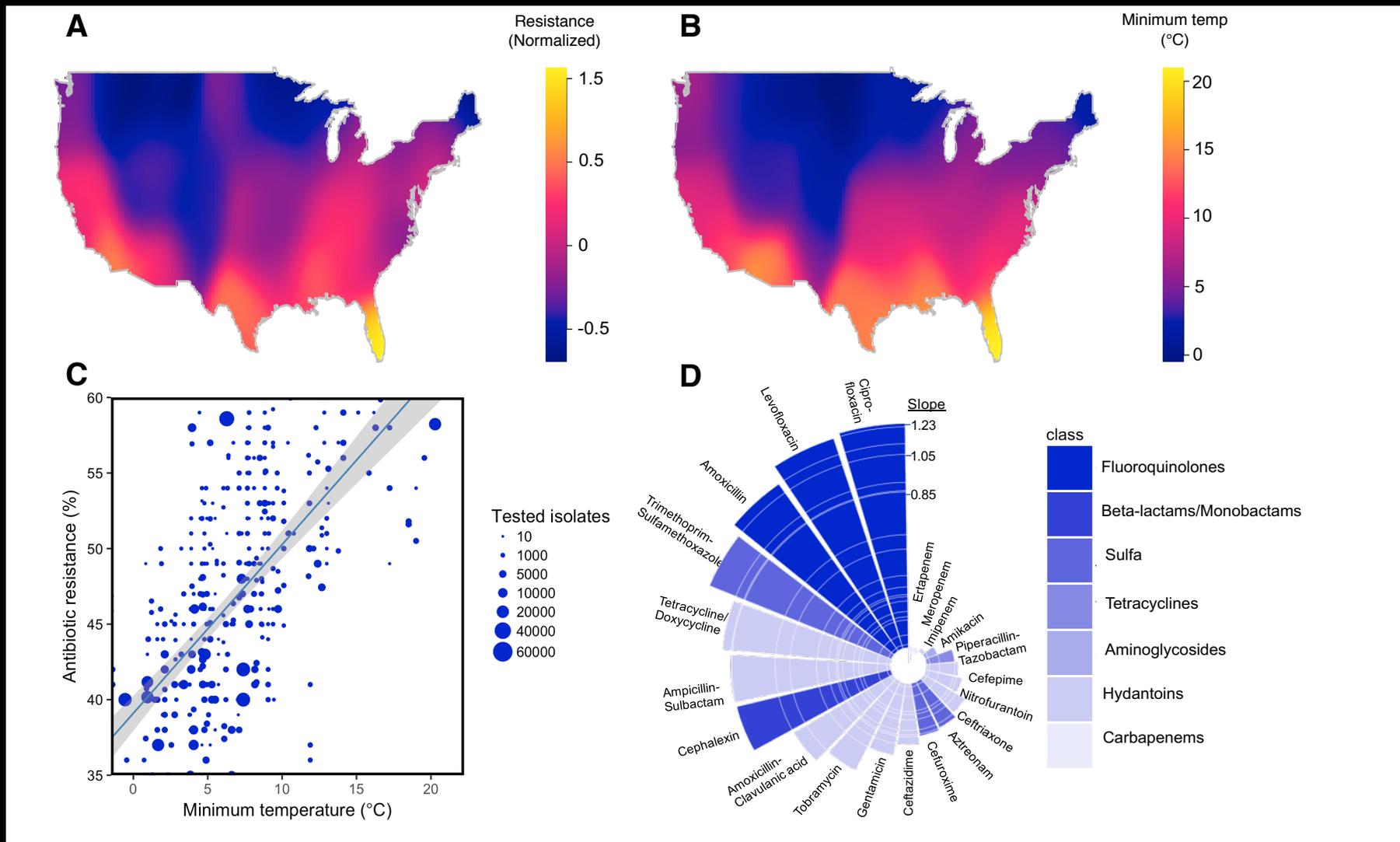


CRE



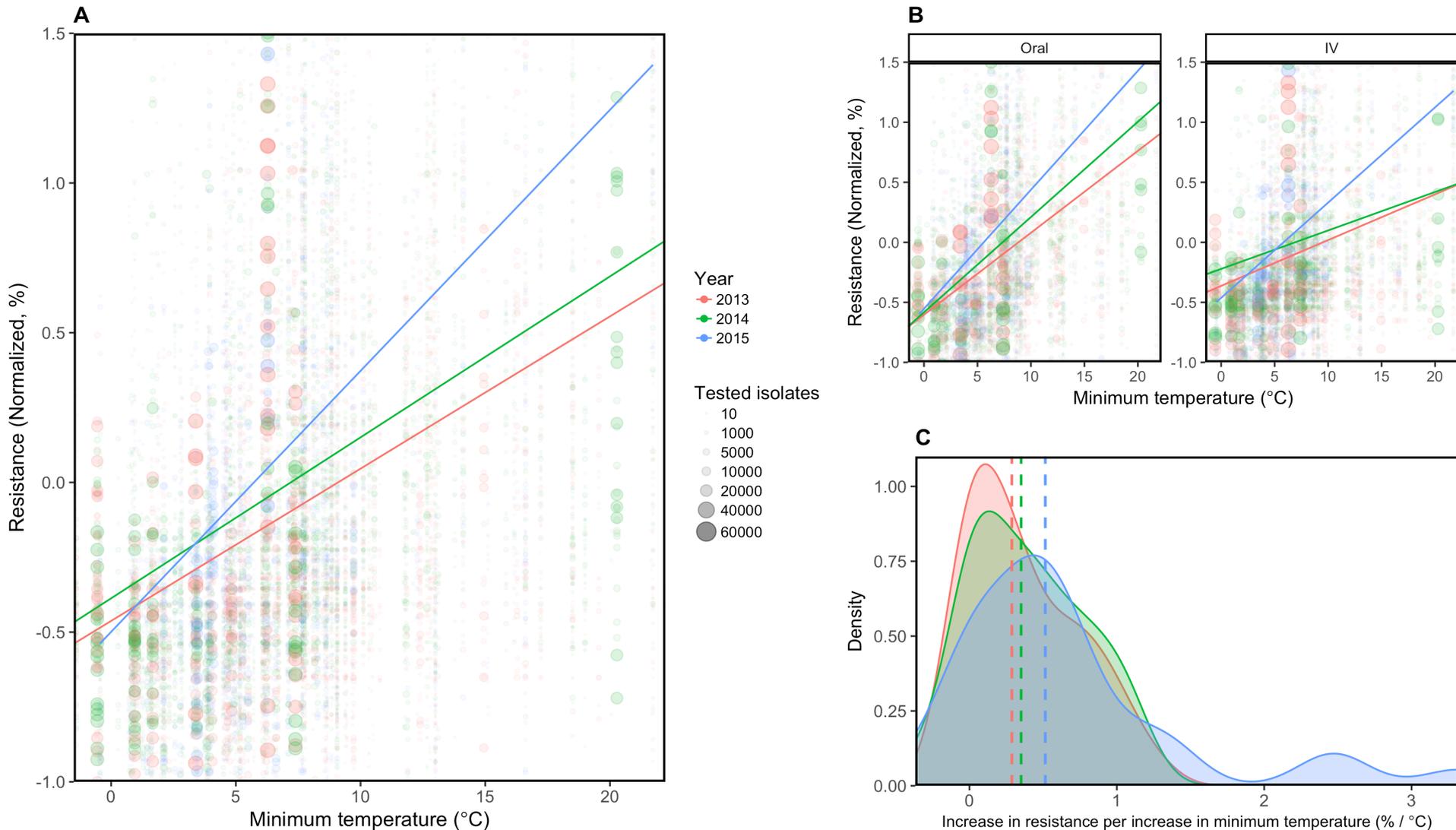
# Antibiotic Resistance Increases with Local Temperature

(Shown in figure *E. coli* resistance for all antibiotics)

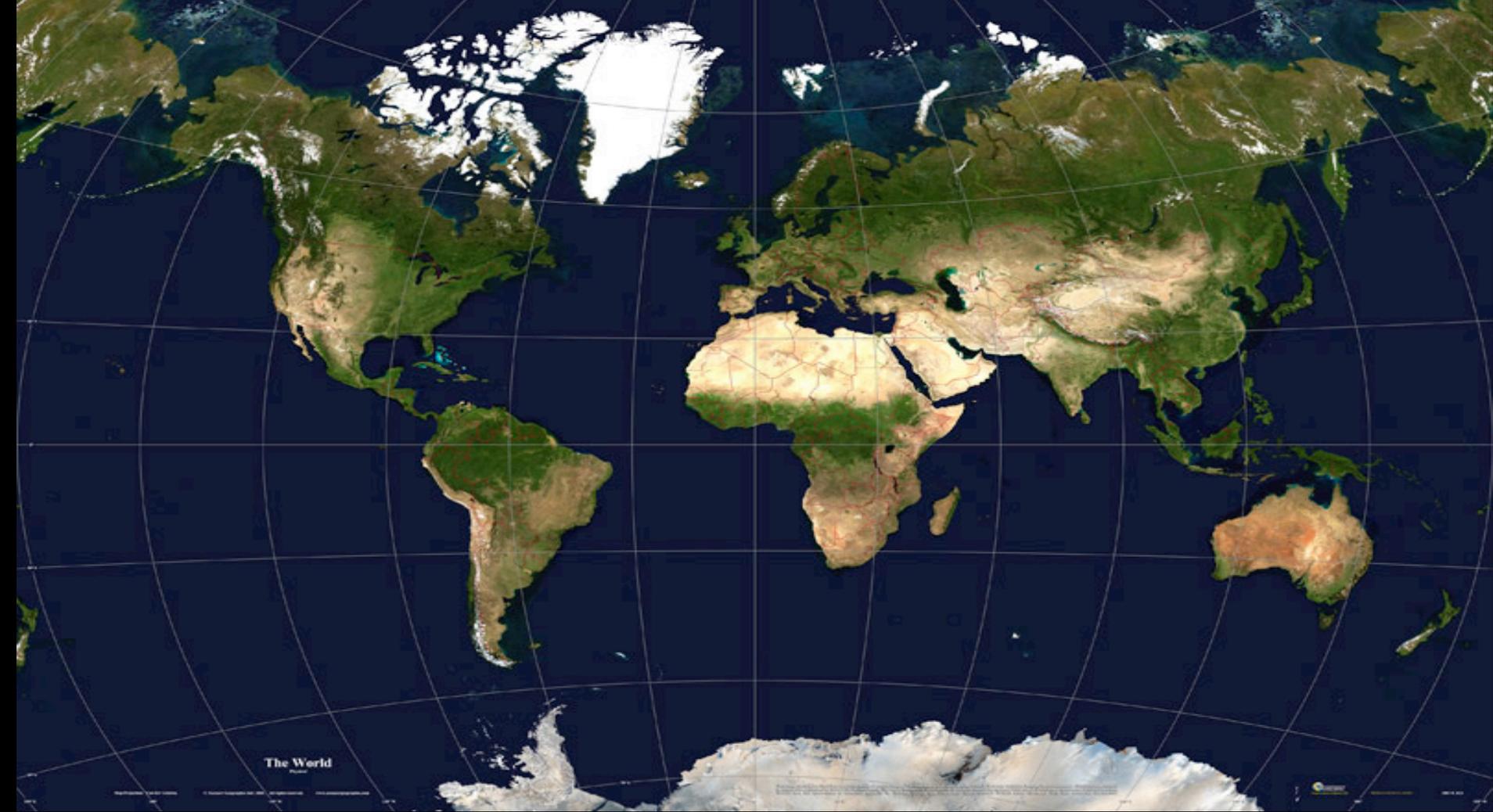


**Team:** Derek R. MacFadden, Sarah F. McGough, David Fisman, John S. Brownstein, and Mauricio Santillana. Nature Climate Change, May 2018

# Antibiotic Resistance Increases with Local Temperature (Shown in figure *E. coli* resistance for all antibiotics)



**Team:** Derek R. MacFadden, Sarah F. McGough, David Fisman, John S. Brownstein, and Mauricio Santillana. *Nature Climate Change*, May 2018



Thank you!

Contact: [msantill@fas.harvard.edu](mailto:msantill@fas.harvard.edu)