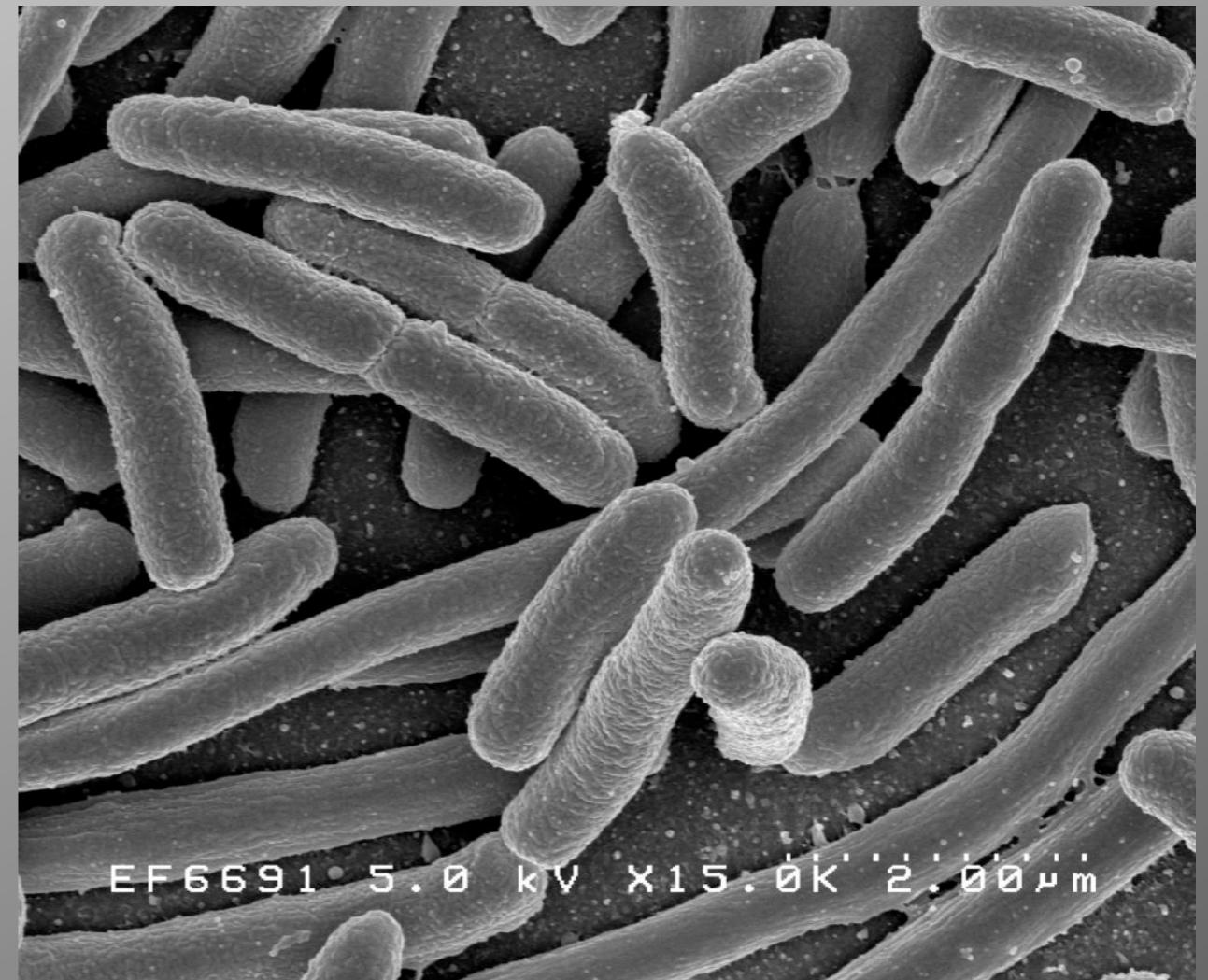


Global Disease Monitoring and Forecasting with Wikipedia + *bonus nerd stuff*

Nicholas Generous
Geoffrey Fairchild
Kyle Hickmann
Alina Deshpande
Sara Y. Del Valle
Reid Priedhorsky

CSCNSI Seminar
June 0x1E, 0x7DF



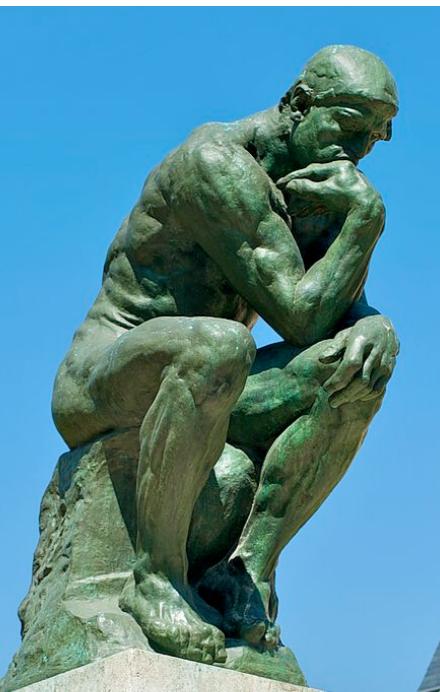
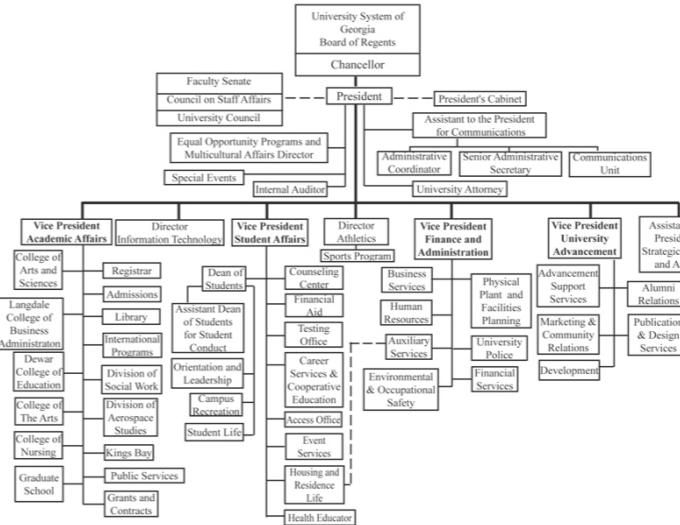
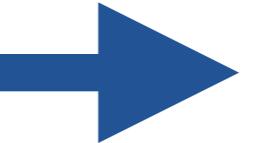
EF6691 5.0 kV x15.0k 2.00 μm

http://en.wikipedia.org/wiki/File:EscherichiaColi_NIAID.jpg

How disease monitoring works

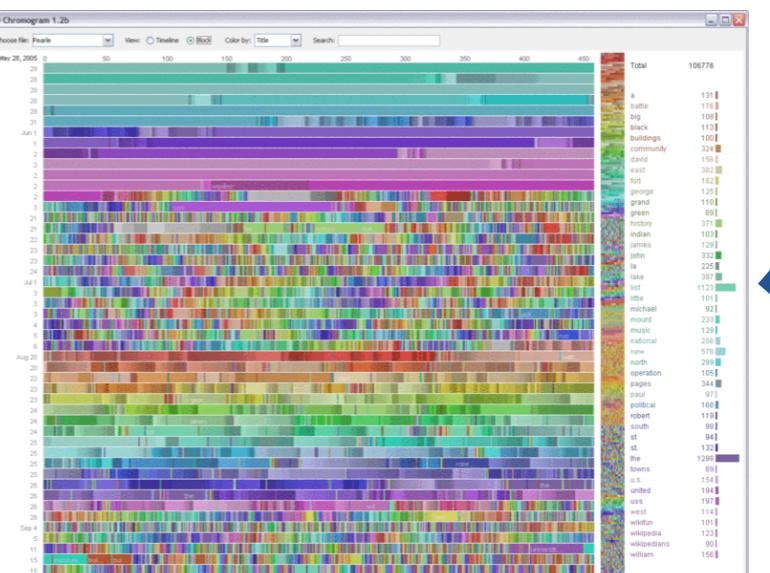
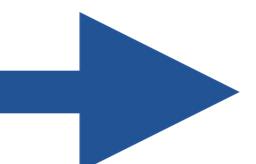
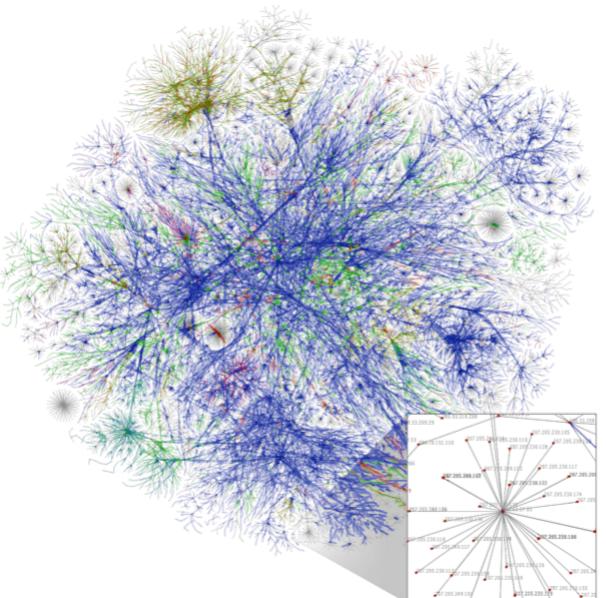
<http://en.wikipedia.org/wiki/File:Blooddraw.jpg>
<http://ww2.valdosta.edu/vsu/org/images/orgchart.jpg>
https://en.wikipedia.org/wiki/File:Paris_2010_-_Le_Penseur.jpg
http://en.wikipedia.org/wiki/File:Internet_map_1024_-_transparent.png
<http://en.wikipedia.org/wiki/File:Viegas-UserActivityonWikipedia.gif>

old:

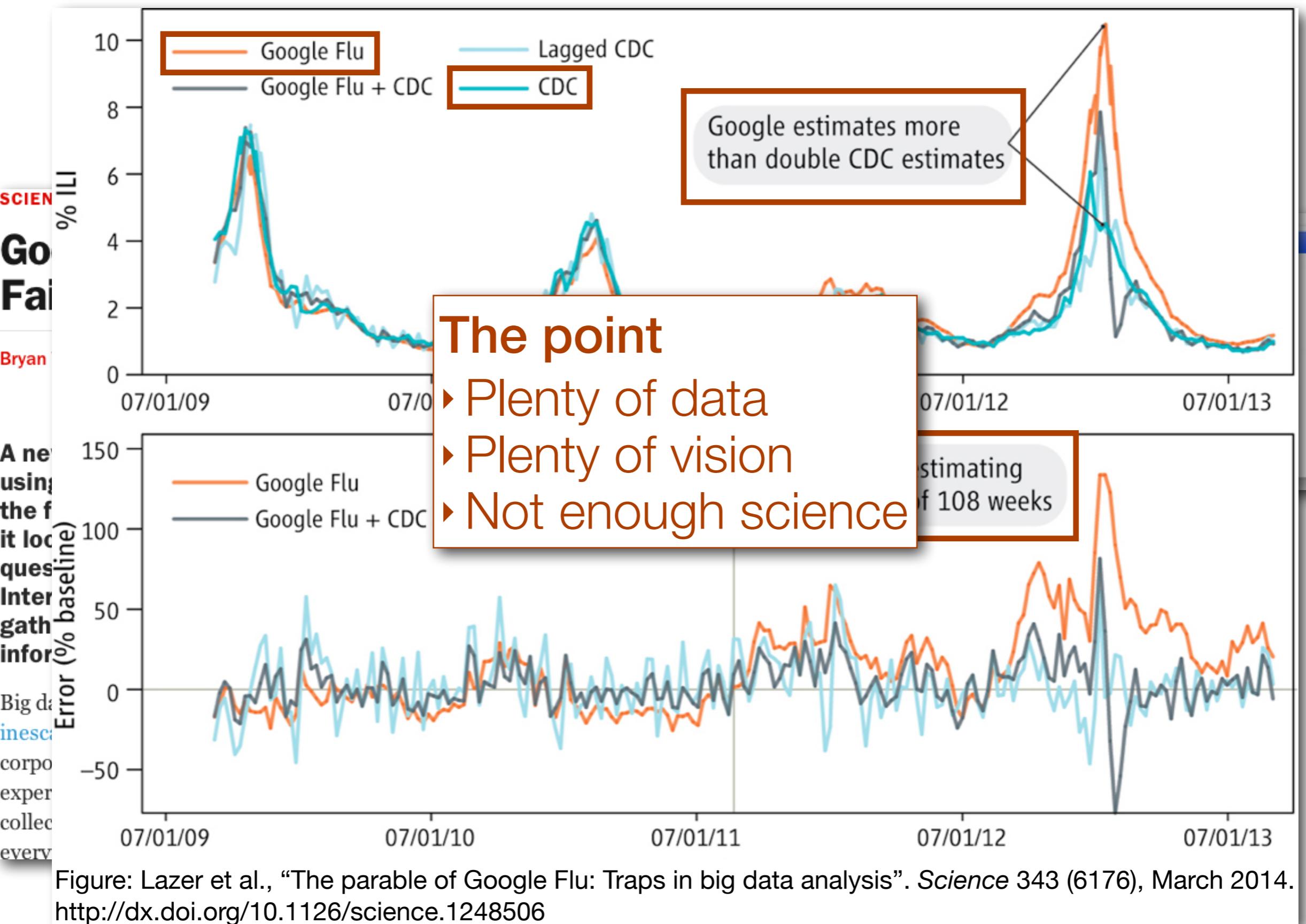


accurate ... costly, slow, & not always available

new:



Fortunately, Google has solved this for us



Related scientific work

Inclusion

- Single-source Internet data
- Biological disease metric

Excluded e.g.'s

- Alerts
- Public perception
- Disease dynamics in model
- Evaluates 3rd party method
- Multiple-source data
- Seasonality metric
- Drug sales metric

McIver & Brownstein

PLOS Comp Bio, April 17, 2014

- Influenza in United States
- Better statistics
- We're better on most other dimensions

The point

- Plenty of science
- Limited vision
- Not enough working systems

Search queries	27	
Baidu	3	avian influenza
Google	20	cancer
Yahoo	2	chicken pox
Yandex	1	cholera
medical sites	3	dengue
Social media	16	dysentery
Twitter	16	gastroenteritis
Server logs	3	gonorrhea
Wikipedia	1	HFMD
medical sites	2	HIV/AIDS
Total	46	influenza 27
		kidney stones
		listeriosis
		malaria
		MRSA
		pertussis
		pneumonia
		RSV
		scarlet fever
		stroke
		suicide
		tuberculosis
		West Nile virus
		Total 55

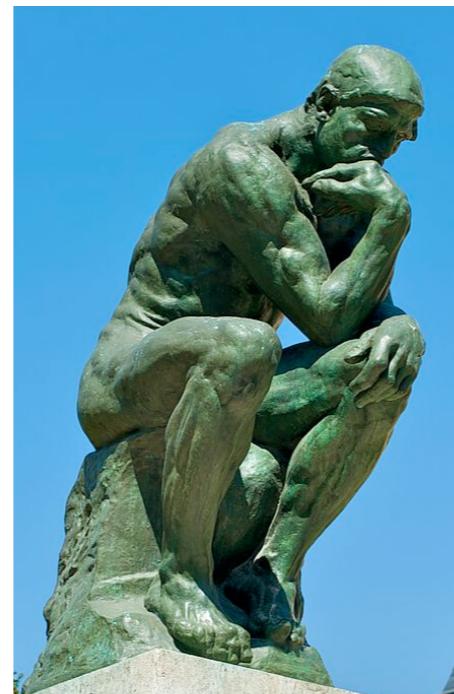
Four challenges for an operational resource

C1. Open data & algorithms

- High quality science
- Continuity & expandability

C2. Breadth

- 100's of disease/country
- Adapt w/o writing code



C3. Transferability

- Missing incidence data
- Transfer w/o re-training

C4. Forecasting

- Simpler = easier to apply

Our goal

**Build an operational disease monitoring
and forecasting system with open data
and open source code.**

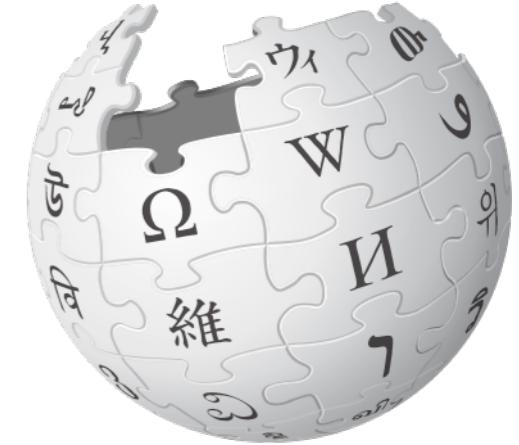
For now: proof-of-concept



WIKIPEDIA
The Free Encyclopedia

What is it?

- ... online encyclopedia
- ... 30 million articles in 287 languages
- ... 850M hits per day (6th)
- ... top result for many search queries



WIKIPEDIA
The Free Encyclopedia

How does it work?

- ... anyone can edit articles
- ... changes live immediately
- ... inverts traditional review/publish

Here be public servants

What data are available?

- ... complete history for every article
- ... hourly requests for each article (“access logs”)
- ... *and lots more*

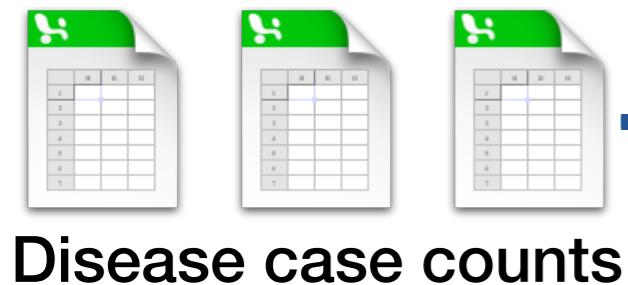
This work

Data pipeline

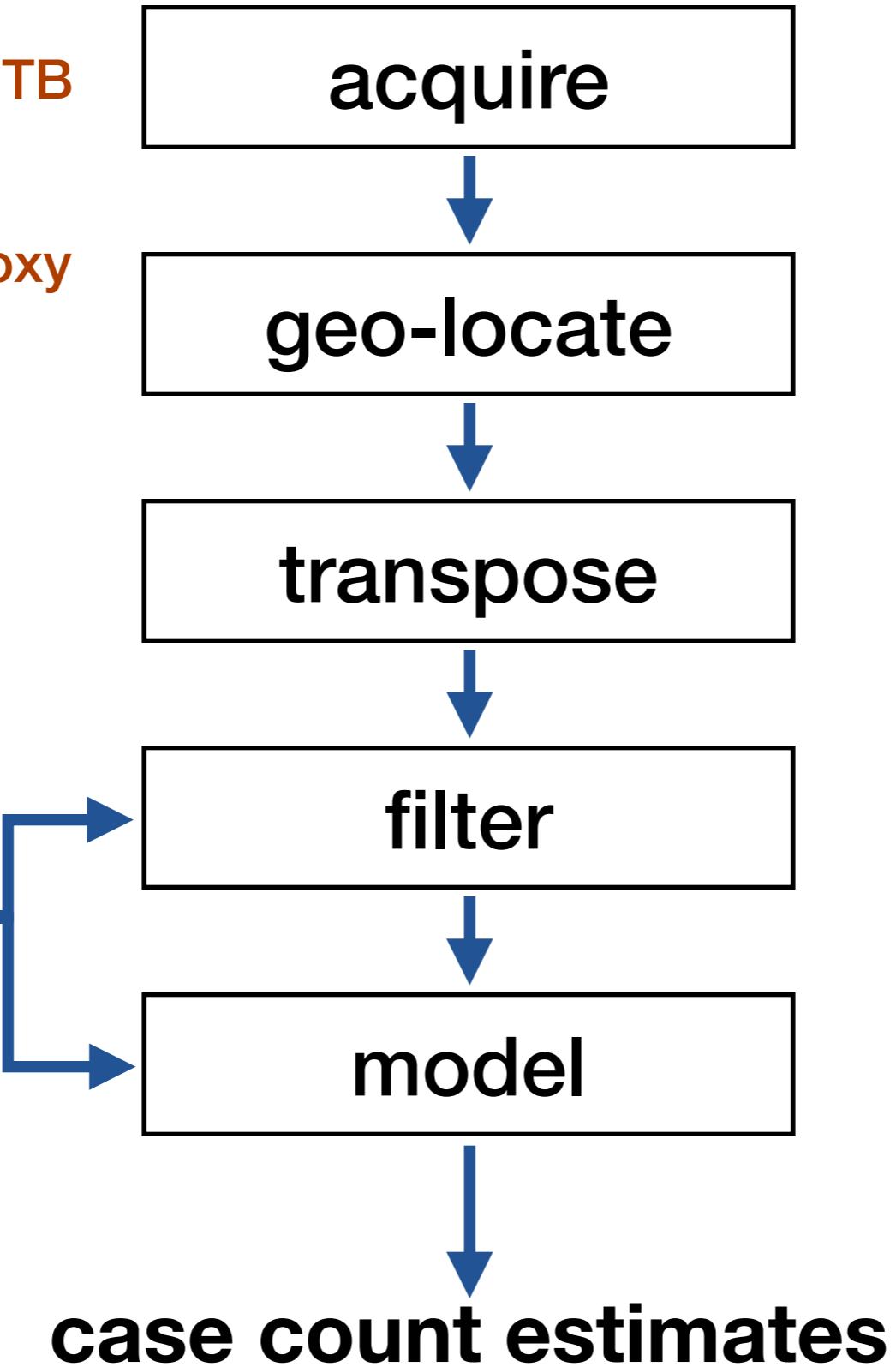


WIKIPEDIA
The Free Encyclopedia

Language as location proxy
‣ Thai = Thailand
‣ English = USA



4.7 TB



Filtering: Candidate articles in English

Article Talk Read View source View history

Influenza

From Wikipedia, the free encyclopedia

"Flu" redirects here. For other uses, see [Flu \(disambiguation\)](#).

"Grippe" redirects here. For other uses, see [Grippe \(disambiguation\)](#).

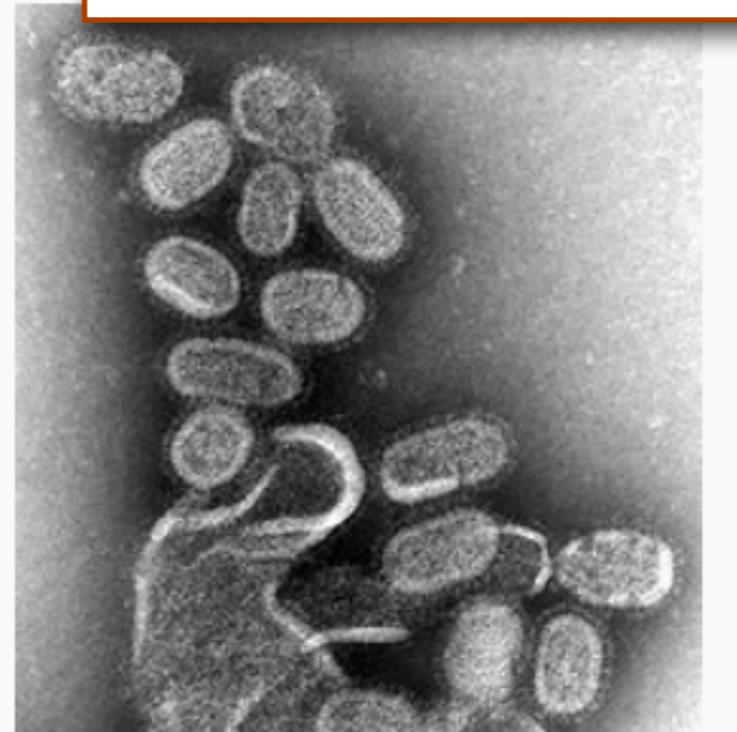
Influenza, commonly known as "the flu", is an infectious disease of birds and mammals caused by RNA viruses of the family [Orthomyxoviridae](#), the influenza viruses. The most common symptoms are chills, fever, runny nose, sore throat, muscle pains, headache (often severe), coughing, weakness/fatigue and general discomfort. Although it is often confused with other influenza-like illnesses, especially the common cold, influenza is a more severe disease caused by a different type of virus.^[1] Influenza may produce nausea and vomiting, particularly in children,^[2] but these symptoms are more common in the unrelated [gastroenteritis](#), which is sometimes inaccurately referred to as "stomach flu" or "24-hour flu".^[3]

Typically, influenza is transmitted through the air by coughs or sneezes, creating [aerosols](#) containing

1. Go to disease article
2. Enumerate linked articles
3. Select articles which are:
 - ▶ symptoms
 - ▶ syndromes
 - ▶ pathogens
 - ▶ conditions
 - ▶ treatments
 - ▶ biological processes
 - ▶ epidemiology

e.g.

- ▶ Influenza
- ▶ Amantadine
- ▶ Swine influenza



TEM of negatively stained influenza virions, magnified approximately 100,000 times

ICD-10 J10, J11

Filtering: Articles in all target languages

The screenshot shows a search results page for 'Influenza' across various Wikipedia editions. A red box highlights the 'Languages' section on the left, which lists numerous language versions. The main content area displays articles for 'Influenza' in English, Polish, and other languages, each with its own header, summary, and edit links. A chart titled 'Percentage of visits for ILI, HHS Region 4, 2013-14 Season through Mar 29, 2014' is overlaid on the Polish Wikipedia page.

4. Translate w/ links

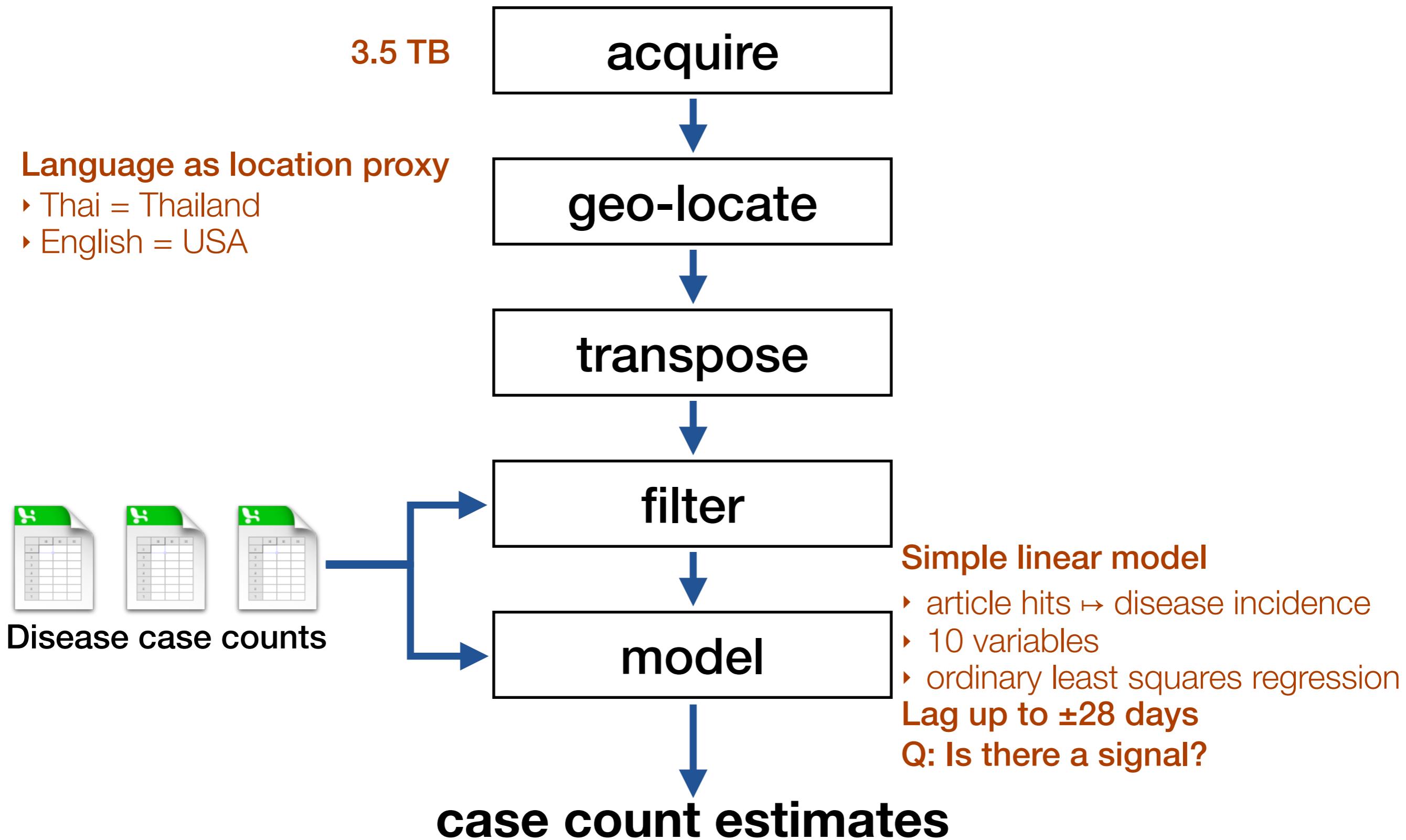
- Grypa
- Amantadyna
- Świńska grypa

5. Compare to incidence

- Correlation r for each article
- Order by decreasing $|r|$

6. Take top 10 articles

Data pipeline



Disease-location contexts analyzed

7 diseases, 9 countries, 14 total contexts

~3 years

Disease	Country	Language	Dates	Resolution
Cholera	Haiti	French	2010-12-05 – 2013-12-05	daily
Dengue	Brazil	Portuguese	2010-03-13 – 2013-03-16	weekly
	Thailand	Thai	2010-10-20 – 2013-11-02	weekly
Ebola	Uganda/DRC	English	2011-01-01 – 2013-12-31	daily
HIV/AIDS	China (PRC)	Chinese	2011-01-01 – 2013-12-01	monthly
	Japan	Japanese	2010-10-09 – 2013-10-12	monthly
Influenza	Japan	Japanese	2010-06-26 – 2013-06-29	weekly
	Poland	Polish	2010-10-22 – 2013-10-22	weekly
	Thailand	Thai	2011-01-09 – 2014-01-12	weekly
	United States	English	2011-01-01 – 2014-01-04	weekly
Plague	United States	English	2011-01-22 – 2014-01-25	weekly
Tuberculosis	China (PRC)	Chinese	2010-12-01 – 2013-12-01	monthly
	Norway	Norwegian	2010-12-01 – 2013-12-01	monthly
	Thailand	Thai	2010-12-01 – 2013-12-01	monthly

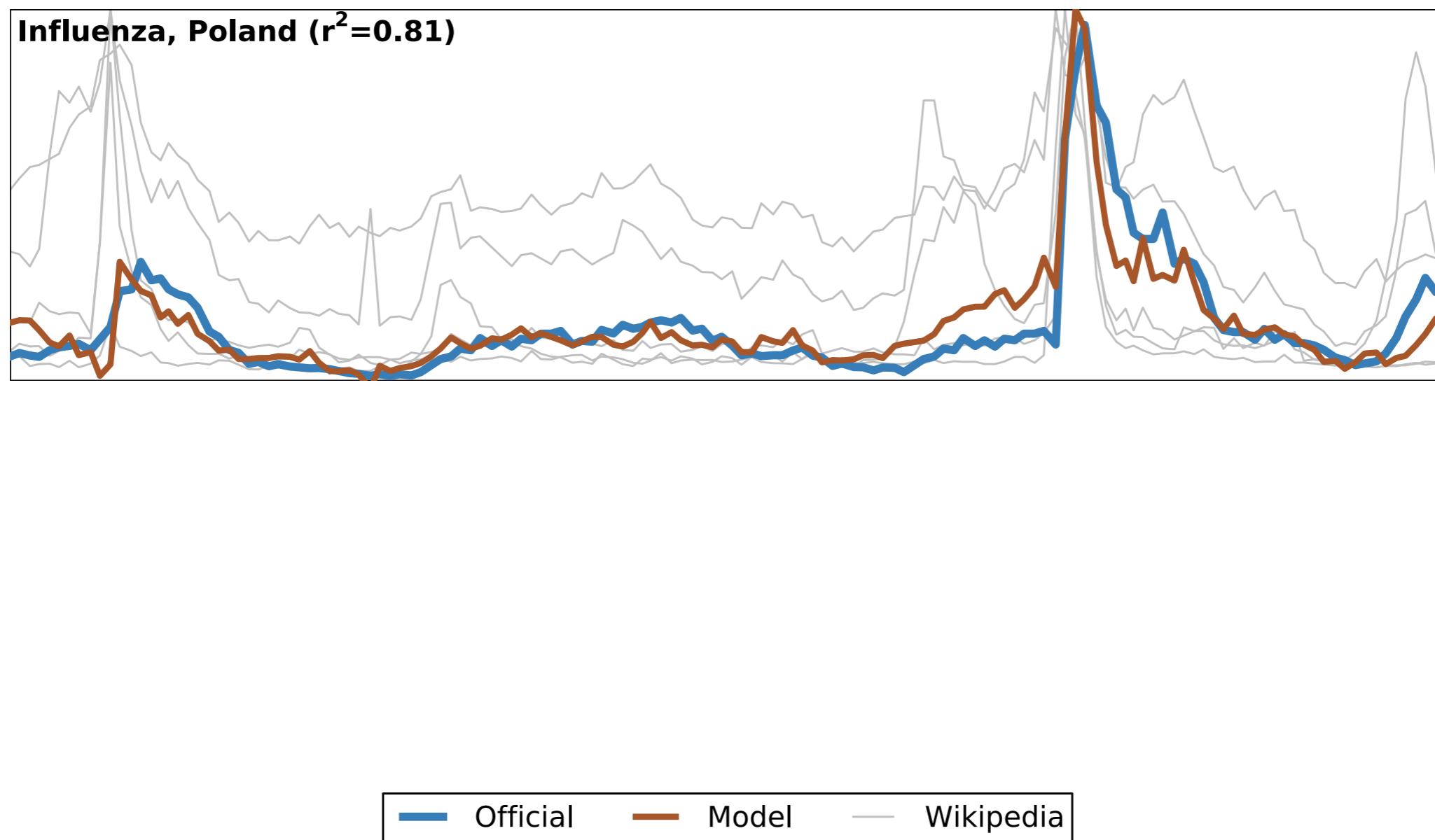
results

Failure modes

- 1. Didn't capture pattern in official data**
- 2. Weak signal-to-noise ratio**

Success (8 of 14 contexts)

↔ ~3 years →



Other successes

1. Forecasting

... best r^2 a few weeks ahead

2. Transferability

... article weights sometimes match across languages

Summary so far

**Wikipedia access logs seem promising
... comparable to other internet data (Twitter, Google)**

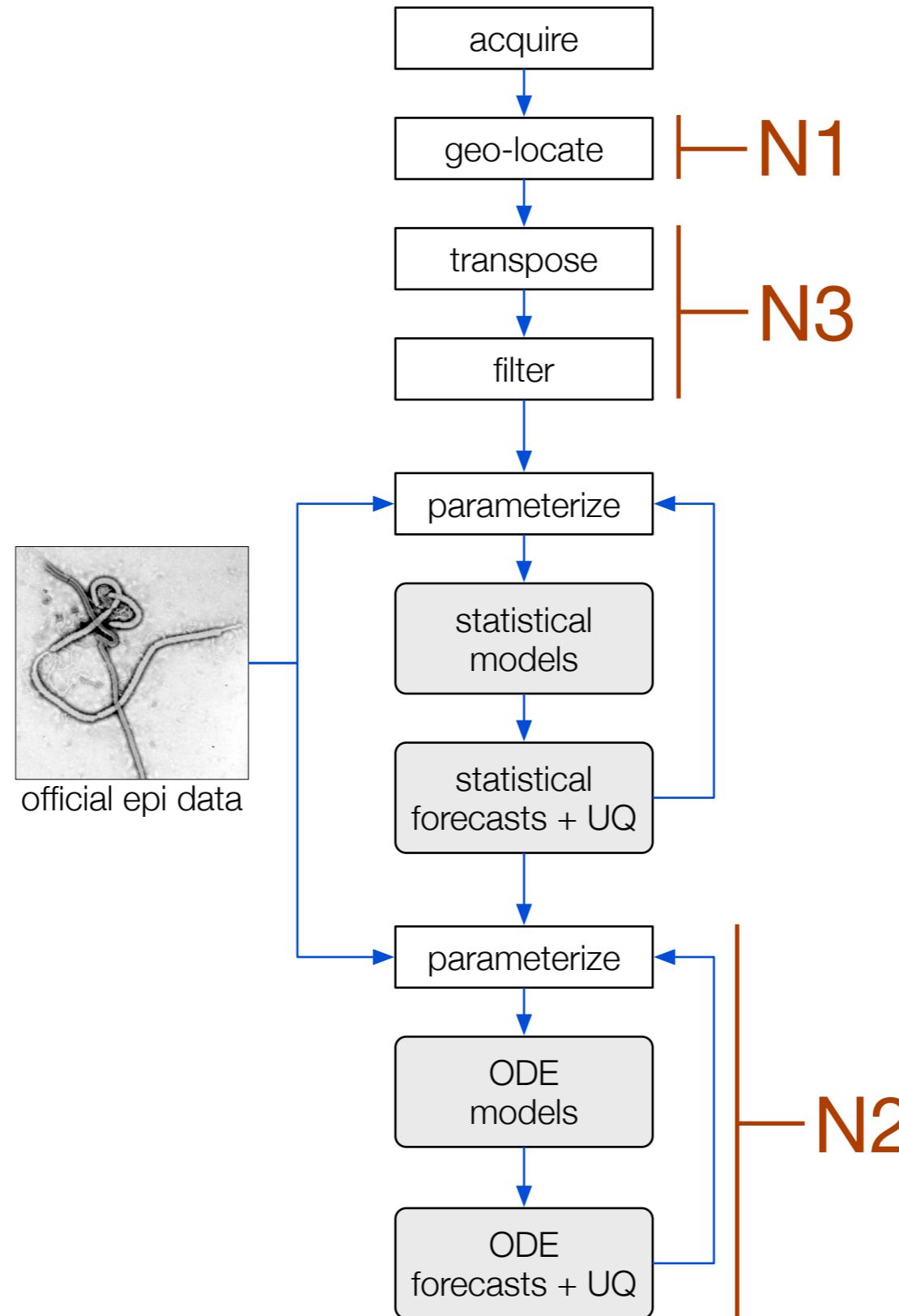
Next steps:

1. Better geo-location
2. Couple with mechanistic models
3. Data-driven article filter

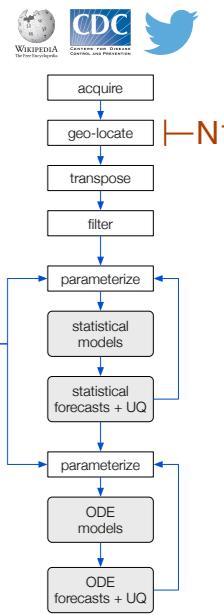
Challenges:

1. Missing ground truth data
2. What is the value-add?

New data flow



C1: Geo-location: Work with WMF



Geo-aggregation of Wikipedia page views: Maximizing geographic granularity while preserving privacy

Reid Priedhorsky, Geoffrey Fairchild, Sara Del Valle
{reidpr,gfairchild,sdelvall}@lanl.gov
Los Alamos National Laboratory

Proposal and request for feedback
LA-UR 15-20145
Draft 2 — January 6, 2015

1 Motivation

Research shows that analyzing internet traces of health-related activity, such as search queries and social media messages, is an effective means of monitoring the spread of disease and has significant promise to tackle problems that traditional disease monitoring tools cannot. The most well-known example of this is Google Flu Trends.¹

Recently, two teams have extended these techniques to use the globally-aggregated Wikipedia page view logs currently available,² with good results.^{3,4} There is even promising evidence that disease forecasting, not simply monitoring, is possible as well using these data. Further, these page view logs are, to our knowledge, the only freely available data source in this class, meaning that the science and utility they can support is significantly greater than proprietary alternatives such as Google queries or Twitter messages.

Effective disease monitoring is fundamentally geographic and requires geo-located data sources. Both of the above efforts infer geography at the country level from the wiki language (e.g., views of the Thai Wikipedia are assumed to come from Thailand), whether implicitly or explicitly. Wikimedia traffic statistics generally support this technique.⁵ However, it has several problems that make it unsuitable for operational use, including:

1. It cannot be applied at geographic granularity finer than the country level. However, in many countries, disease monitoring must be carried out at the state or metro-area level in order to be effective.
2. Many or most countries can't be covered at all. For example, this approach is not feasible for Chile, which comprises only 4.4% of Spanish-language requests and whose signal is swamped by other Spanish-speaking countries.

¹http://www.google.org/flu_trends/

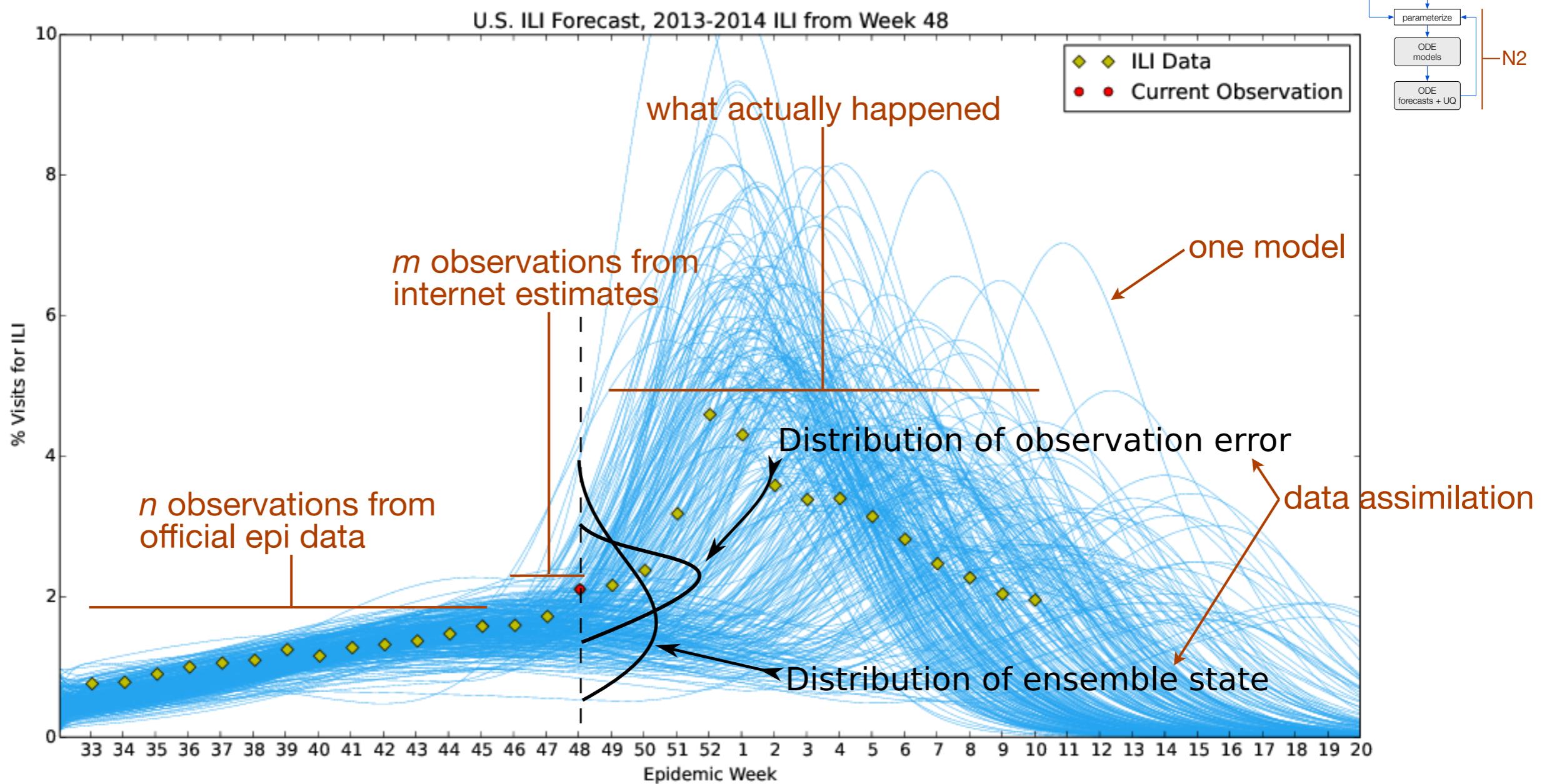
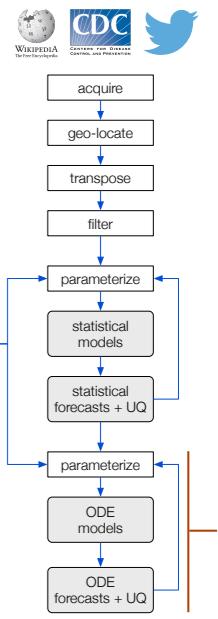
²<http://dumps.wikimedia.org/other/pagecounts-raw/>

³McIver & Brownstein (<http://www.ploscompbiol.org/article/info:doi/10.1371/journal.pcbi.1003581>).

⁴Generous et al. (<http://www.ploscompbiol.org/article/info:doi/10.1371/journal.pcbi.1003892>).

⁵<http://stats.wikimedia.org/wikimedia/squids/SquidReportPageViewsPerLanguageBreakdown.htm>

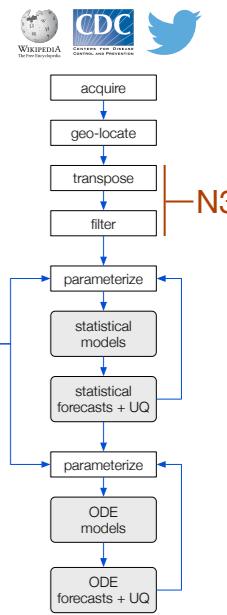
C2: Model coupling via data assimilation



C3: Filter: Raw data file



C3: Filter: Goals for data format



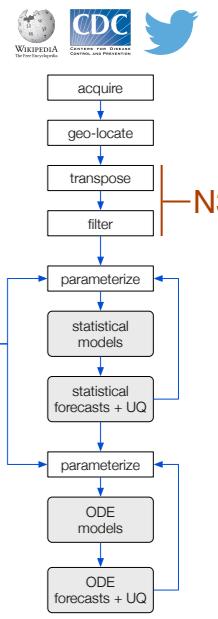
Use cases

- Retrieve ~1 article time series (serial)
- Iterate through all time series (parallel)

Maintenance cases

- Bulk transpose all raw data (parallel)
- Update with 1 day of new data (serial)

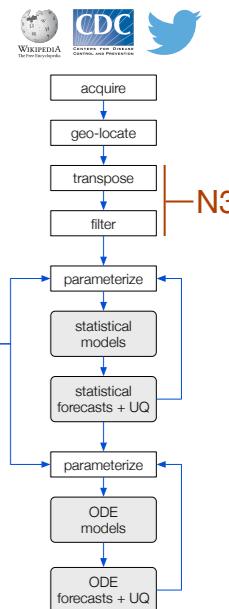
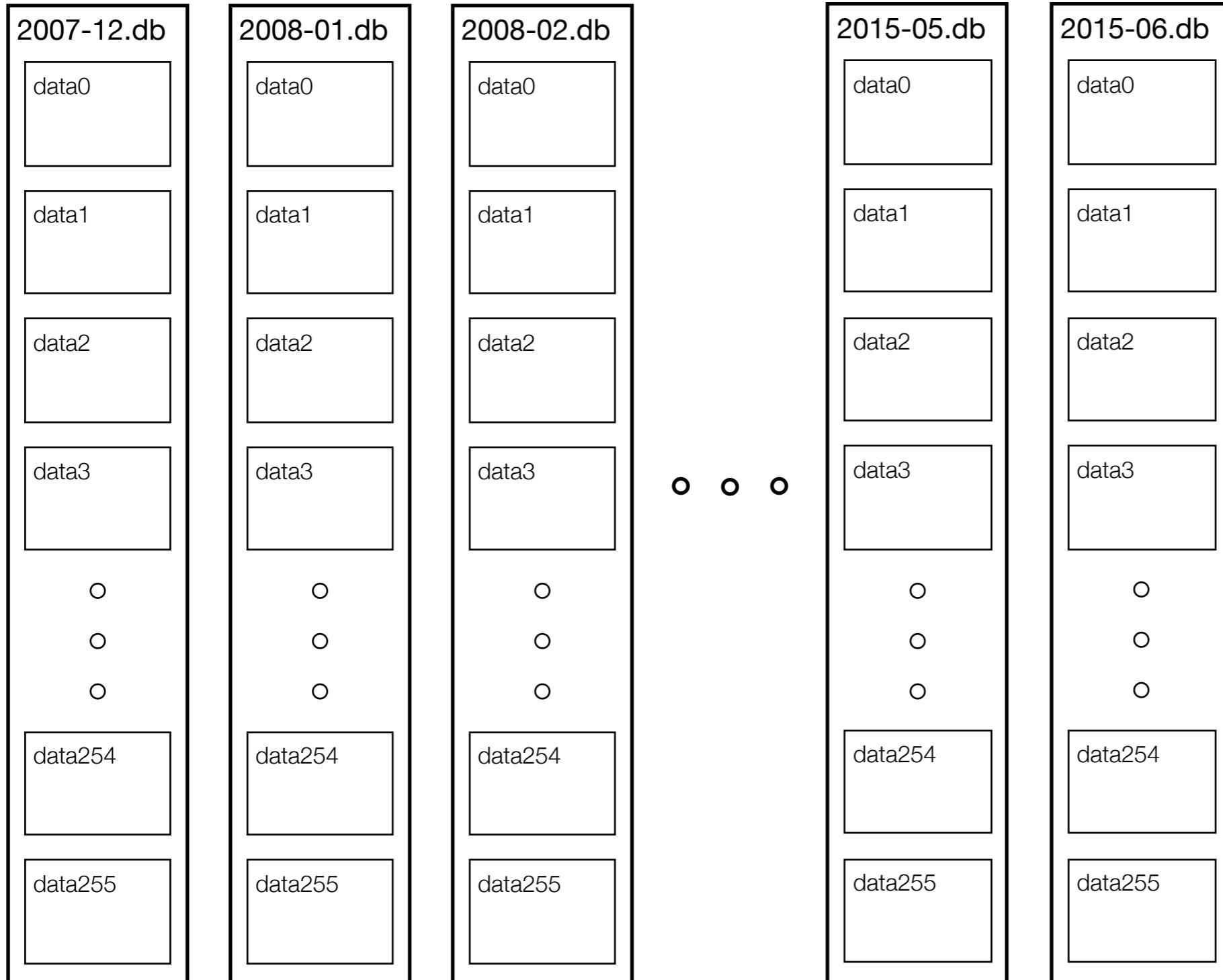
C3: Filter: Data funnel



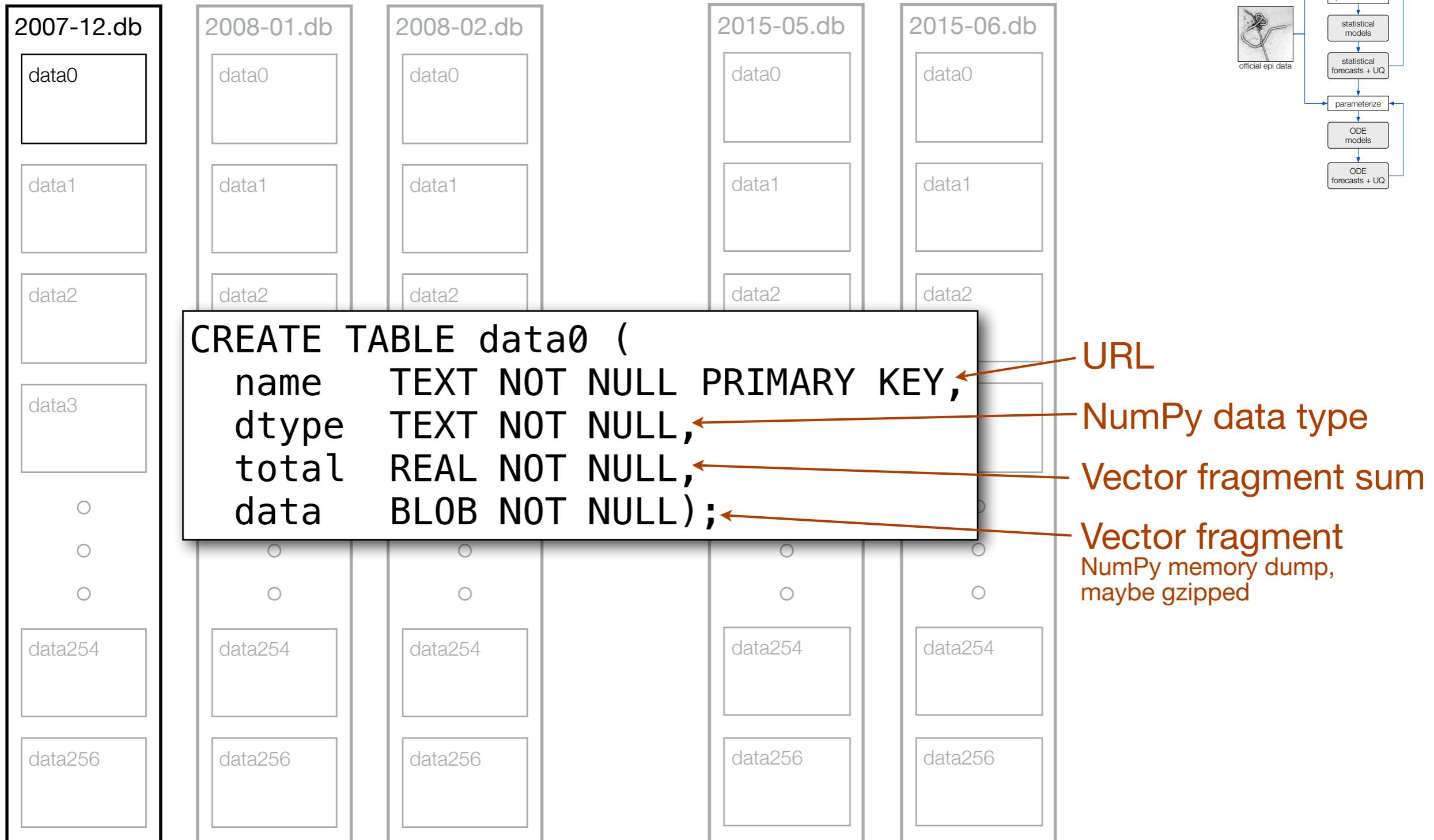
URLs

Total	25G	
Well-formed	20G	<i>egrep</i>
Plausible for any model	37M	60 hits/month
Plausible for specific model	300	correlation
Actually in model	30	linear regression

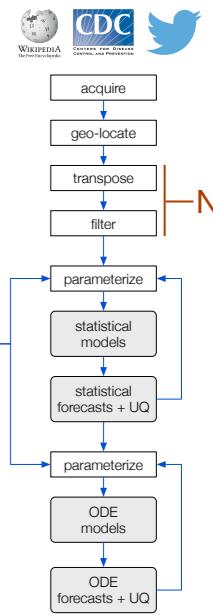
C3: Filter: Data format: Sharded SQLite databases



C3: Filter: Data format: Database schema



C3: Filter: Data format performance



Use cases

- Retrieve ~1 article time series
- Iterate through all time series

1–20 s

7 min (256×350/s) [est.]

Maintenance cases

- Bulk transpose all raw data
- Update with 1 day of new data

24 hours (91× parallel) [est.]

6–12 hours [est.]

Nicholas Generous
 Geoffrey Fairchild
 Kyle Hickmann
 Alina Deshpande
 Sara Y. Del Valle
 Reid Priedhorsky / reidpr@lanl.gov

dx.doi.org/10.1371/journal.pcbi.1003892

OPEN ACCESS Freely available online

PLOS COMPUTATIONAL BIOLOGY

Global Disease Monitoring and Forecasting with Wikipedia

Nicholas Generous*, Geoffrey Fairchild, Alina Deshpande, Sara Y. Del Valle, Reid Priedhorsky

Defense Systems and Analysis Division, Los Alamos National Laboratory, Los Alamos, New Mexico, United States of America

Abstract

Infectious disease is a leading threat to public health, economic stability, and other key social structures. Efforts to mitigate these impacts depend on accurate and timely monitoring to measure the risk and progress of disease. Traditional, biologically-focused monitoring techniques are accurate but costly and slow; in response, new techniques based on social internet data, such as social media and search queries, are emerging. These efforts are promising, but important challenges in the areas of scientific peer review, breadth of diseases and countries, and forecasting hamper their operational usefulness. We examine a freely available, open data source for this use: access logs from the online encyclopedia Wikipedia. Using linear models, language as a proxy for location, and a systematic yet simple article selection procedure, we tested 14 location-disease combinations and demonstrate that these data feasibly support an approach that overcomes these challenges. Specifically, our proof-of-concept yields models with r^2 up to 0.92, forecasting value up to the 28 days tested, and several pairs of models similar enough to suggest that transferring models from one location to another without re-training is feasible. Based on these preliminary results, we close with a research agenda designed to overcome these challenges and produce a disease monitoring and forecasting system that is significantly more effective, robust, and globally comprehensive than the current state of the art.

Citation: Generous N, Fairchild G, Deshpande A, Del Valle SY, Priedhorsky R (2014) Global Disease Monitoring and Forecasting with Wikipedia. PLoS Comput Biol 10(11): e1003892. doi:10.1371/journal.pcbi.1003892

Editor: Marcel Salathé, Pennsylvania State University, United States of America

Received April 18, 2014; Accepted August 21, 2014; Published November 13, 2014

This is an open-access article, free of all copyright, and may be freely reproduced, distributed, transmitted, modified, built upon, or otherwise used by anyone for any lawful purpose. The work is made available under the Creative Commons CC0 public domain dedication.

Data Availability: The authors confirm that all data underlying the findings are fully available without restriction. Data are available in the Supplemental Information. Software is available at <http://github.com/reidpr/quac>.

Funding: This work is supported in part by NIH/NIGMS/MIDAS under grant U01-GM097658-01 and the Defense Threat Reduction Agency (DTRA), Joint Science and Technology Office for Chemical and Biological Defense under project numbers CB3656 and CB10007. Data collected using QUAC; this functionality was supported by the U.S. Department of Energy through the LANL LDRD Program. Computation used HPC resources provided by the LANL Institutional Computing Program. LANL is operated by Los Alamos National Security, LLC for the Department of Energy under contract DE-AC52-06NA25396. Approved for public release: LA-UR-14-22353. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing Interests: The authors have declared that no competing interests exist.

* Email: generous@lanl.gov

Introduction

Motivation and Overview

Infectious disease remains extremely costly in both human and economic terms. For example, the majority of global child mortality is due to conditions such as acute respiratory infection, measles, diarrhea, malaria, and HIV/AIDS [1]. Even in developed countries, infectious disease has great impact; for example, each influenza season costs the United States between 3,000 and 49,000 lives [2] and an average of \$87 billion in reduced economic output [3].

Effective and timely disease surveillance — that is, detecting, characterizing, and quantifying the incidence of disease — is a critical component of prevention and mitigation strategies that can save lives, reduce suffering, and minimize impact. Traditionally, such monitoring takes the form of patient interviews and/or laboratory tests followed by a bureaucratic reporting chain; while generally considered accurate, this process is costly and introduces a significant lag between observation and reporting.

These problems have motivated new surveillance techniques based upon internet data sources such as search queries and social media posts. Essentially, these methods use large-scale data mining techniques to identify health-related activity traces within the data streams, extract them, and transform them into some useful metric. The basic approach is to train a statistical estimation model against ground truth data, such as ministry of health disease incidence records, and then apply the model to generate estimates when the true data are not available, e.g., when forecasting or when the true data have not yet been published. This has proven effective and has spawned operational systems such as Google Flu Trends (http://www.google.org/flu_trends/). However, four key challenges remain before internet-based disease surveillance models can be reliably integrated into an decision-making toolkit:

C1. Openness. Models should afford review, replication, improvement, and deployment by third parties. This guarantees a high-quality scientific basis, continuity of operations, and broad applicability. These requirements imply that model algorithms — in the form of source code, not research papers — must be generally available, and they also imply that complete input data must be available. The latter is the key obstacle, as terms are dictated by the data owner rather than the data user; this motivated our exploration of Wikipedia access logs. To our knowledge, no models exist that use both open data and open algorithms.

C2. Breadth. Dozens of diseases in hundreds of countries have sufficient impact to merit surveillance; however, adapting a

reidpr / quac

Code Network Pull Requests 0 Issues 7 Wiki Graphs Settings

Clone in Mac ZIP HTTP SSH Git Read-Only git@github.com:reidpr/quac.git Read+Write access

branch: master Files Commits Branches Tags

quac / 5 commits

null merge of initial icwsm2013 branch to make future merges work

Reid Priedhorsky authored 5 days ago latest commit fbfb05a8ac

File	Time	Description
geo	5 days ago	Initial commit. [reidpr]
sphinx	5 days ago	remove bugs from limitations.rst (closes #7) [reidpr]
tok	5 days ago	initial revision of icwsm2013 branch [reidpr]
.gitignore	5 days ago	Initial commit. [reidpr]
LICENSE	5 days ago	add license file [reidpr]
README	5 days ago	initial revision of icwsm2013 branch [reidpr]
README.mac	5 days ago	initial revision of icwsm2013 branch [reidpr]
collect	5 days ago	Initial commit. [reidpr]

github.com/reidpr/quac

acknowledgements:

- Susan M. Mniszewski
- Funding:
 - NIH/NIGMS: U01-GM097658-01
 - DTRA
 - U.S. Dept. of Energy via LANL/LDRD
 - LANL is operated by LANS, LLC for DOE under contract DE-AC52-06NA25396.



A close-up photograph of a Shiba Inu dog. The dog is a light tan or cream color with darker tan markings on its ears, around its eyes, and on its paws. It is sitting on a light-colored couch. In the background, there is a small wooden shelf with some decorative items, including a small framed picture and some pink flowers. The dog's expression is neutral, and it is looking slightly to the right of the camera.

wow

so results

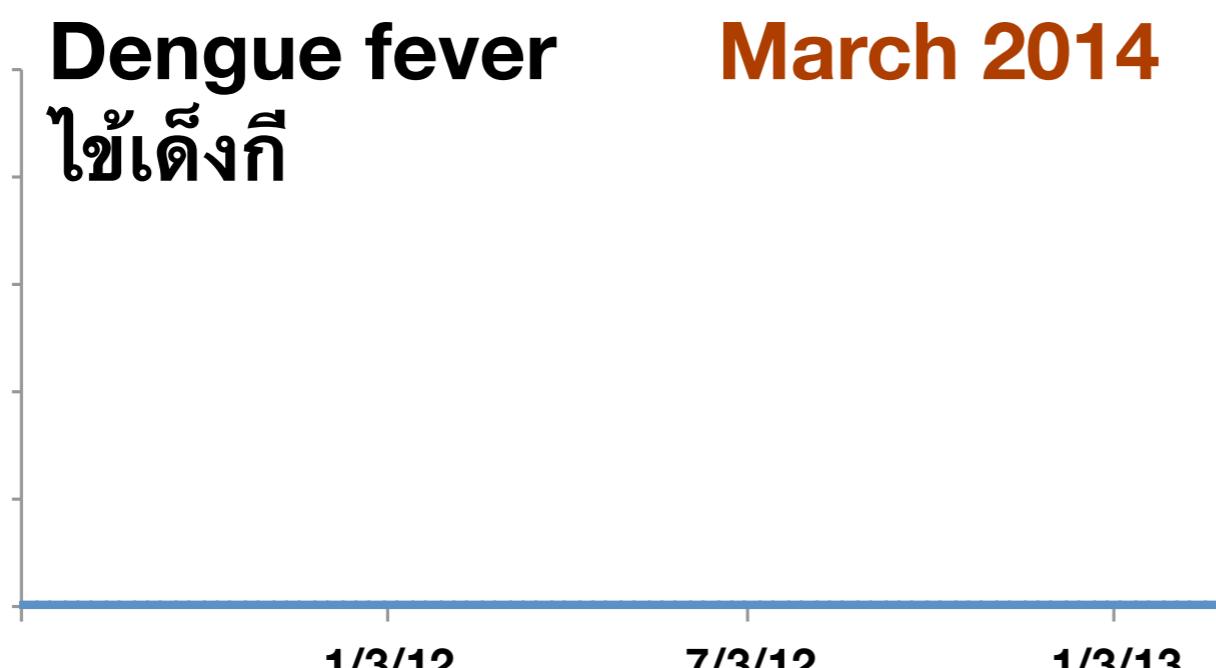
such Big Data Methods™

easy win

much funding

A story about dengue in Thailand

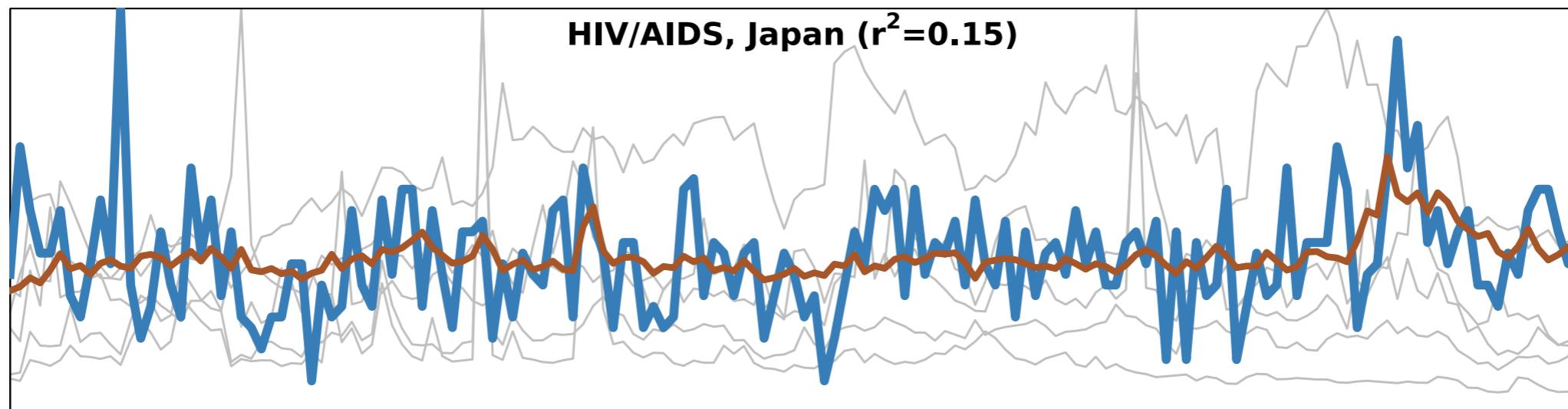
The point: Know your data!



July 23, 2013

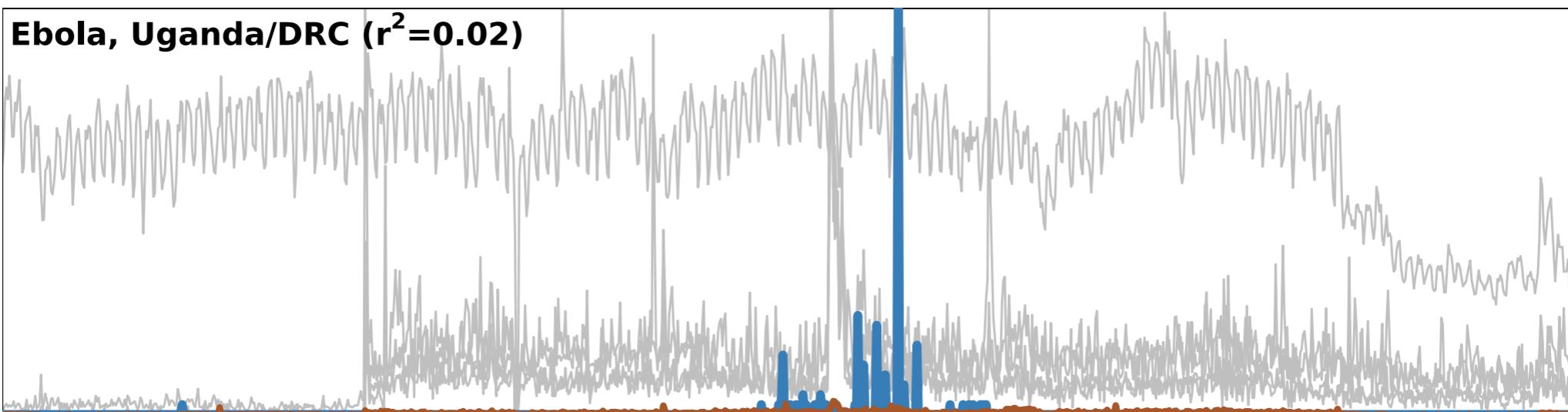
Failed: Didn't capture data pattern

↔ ~3 years →



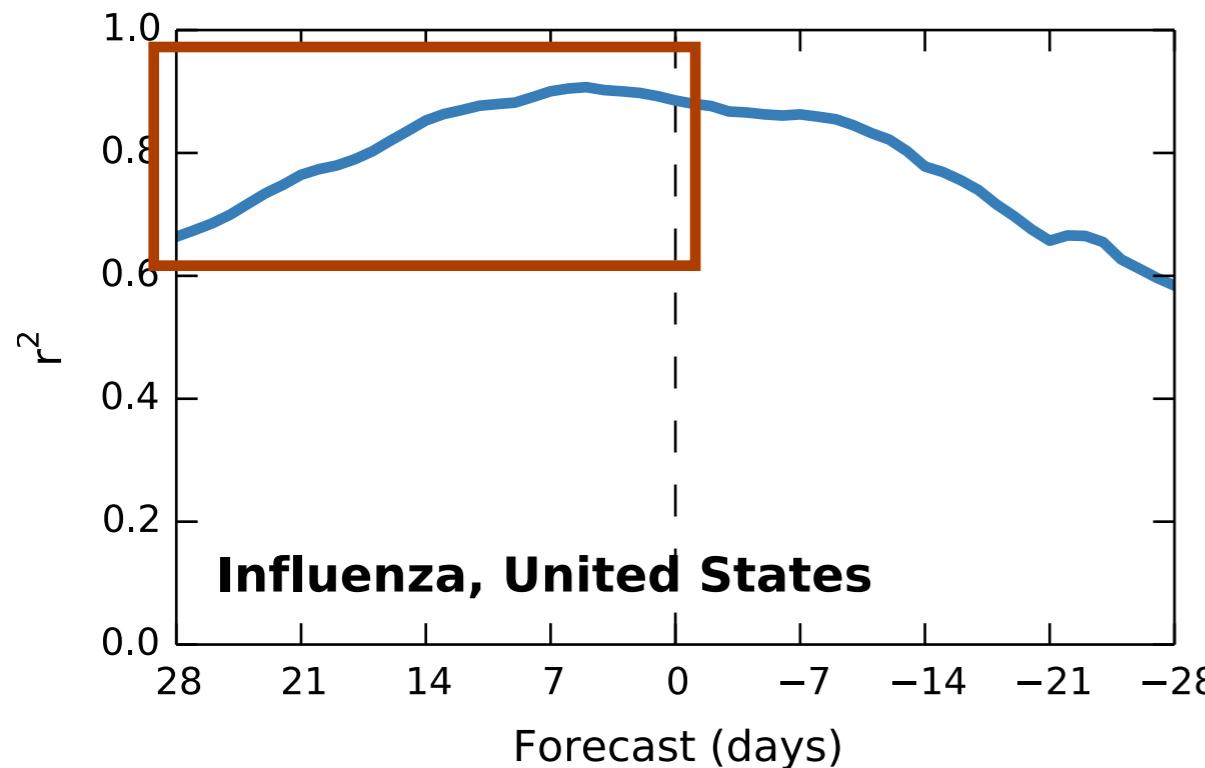
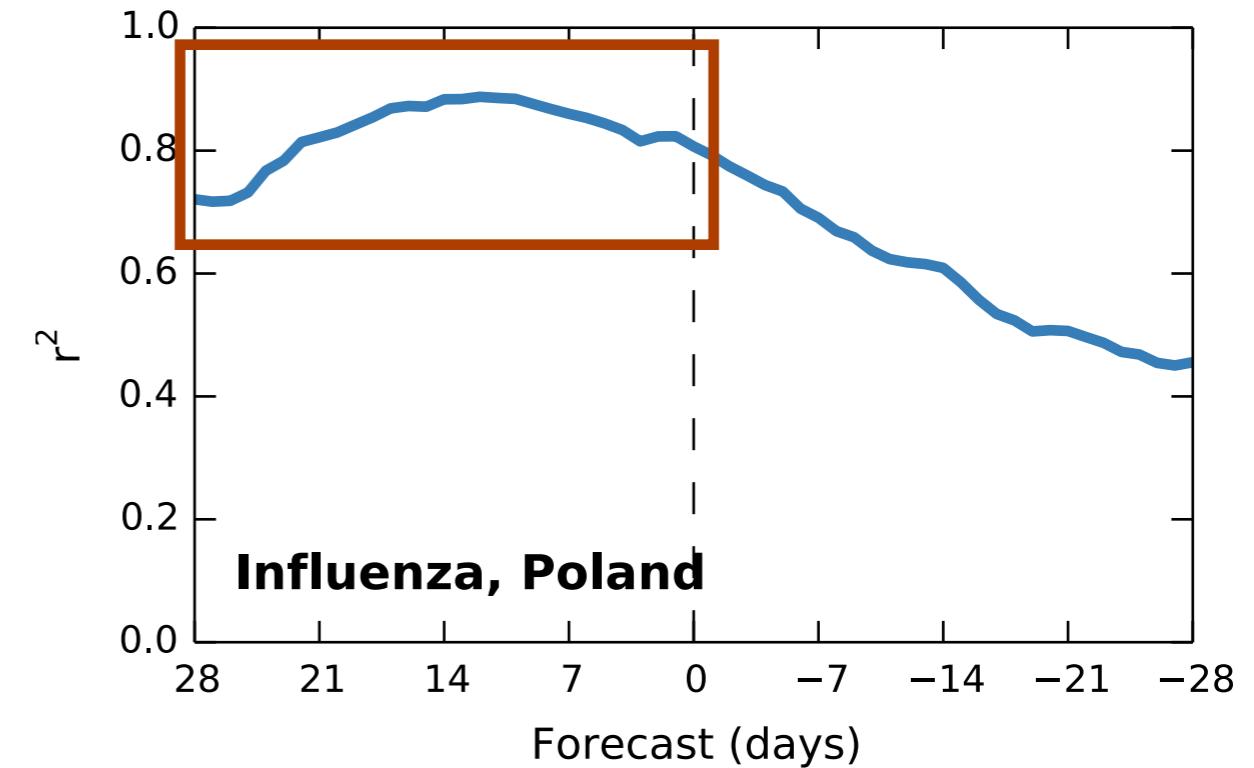
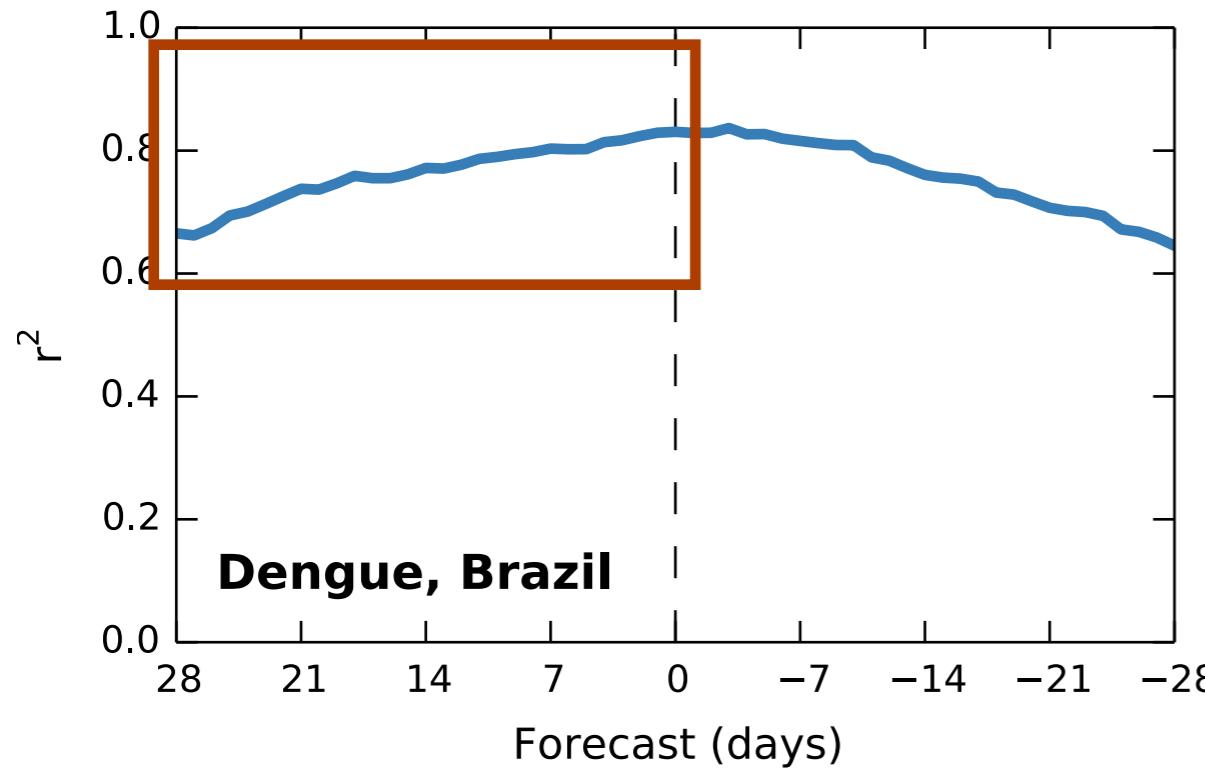
Failed: Weak signal-to-noise ratio (SNR)

~3 years



Success

Seasonality? News? Observing infectors?



Other successes

- ▶ Dengue: Thailand
- ▶ Influenza: Japan, Thailand
- ▶ Tuberculosis: China, Thailand

McIver & Brownstein vs. Generous et al.

	McIver	Generous
First use of Wikipedia	yes	yes
Peer reviewed	yes	yes
Contexts	1	14
Analysis period	5 years	3 years
Language as geo-proxy	mentioned	tested
Statistics	better	worse
Goodness of fit	r, MAE	r^2
Compare to GFT	yes	no
Negative results	no	included
Forecasting	no	28 days
Translatability	no	tested
Focus	result	data source
Research agenda	no	proposed