

Unsupervised extraction of epidemic syndromes from participatory influenza surveillance self-reported symptoms

Daniela Paolotti
ISI Foundation
Turin, Italy

@danielapaolotti

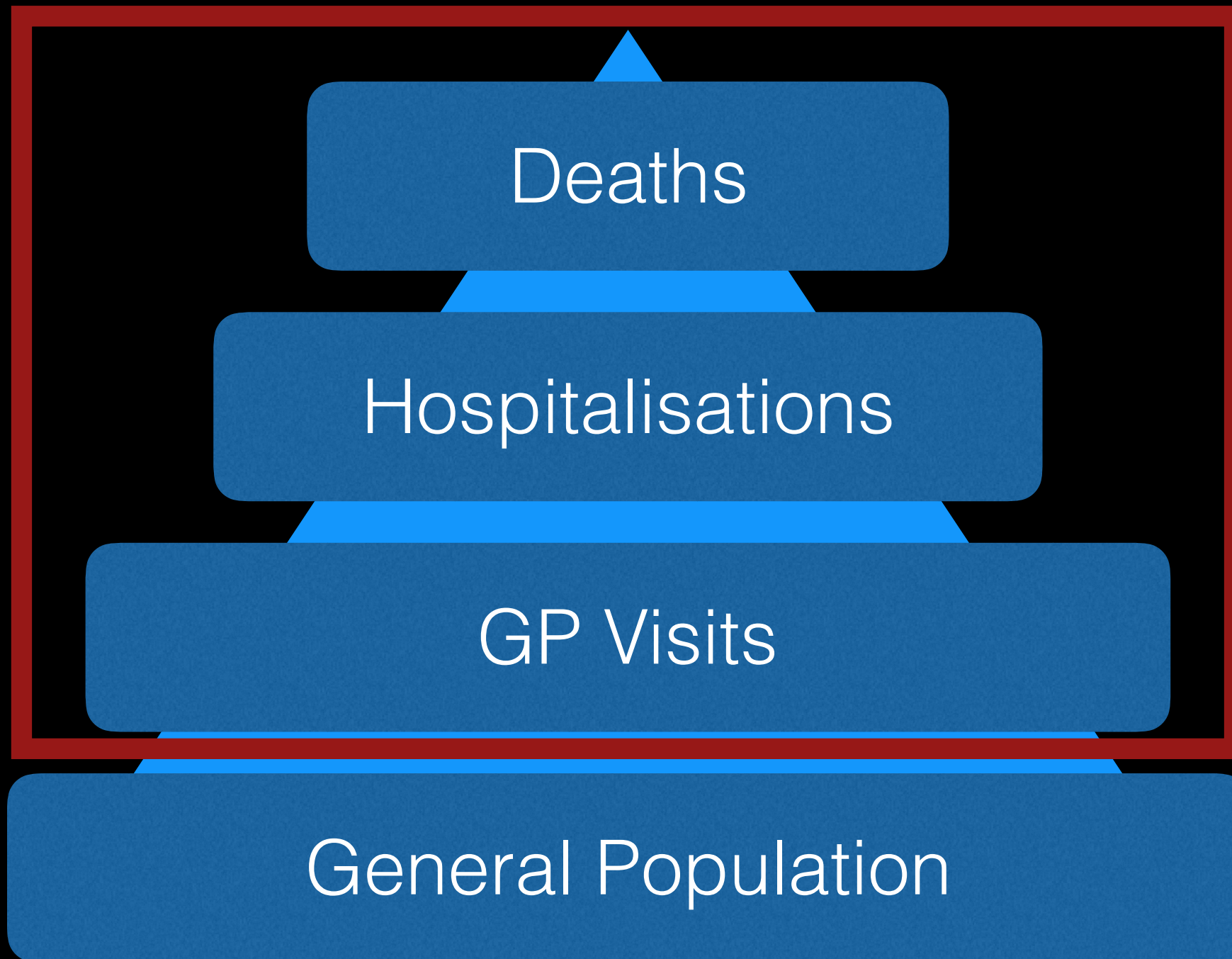
Applied Machine Learning Days
Health Track
January 27th 2020

An iceberg floating in a deep blue ocean under a clear blue sky. The small tip of the iceberg is above the water line, while the much larger, jagged mass is submerged below. The water is a deep, dark blue, and the sky is a lighter, clear blue with a few wispy clouds. The iceberg's surface is textured with various ridges and crevasses.

Traditional Surveillance

Disease Burden

Surveillance



Deaths

The diagram is a funnel shape composed of four horizontal rectangular blocks stacked vertically. The top block is labeled 'Deaths' and is dark blue. The second block is labeled 'Hospitalisations' and is a medium blue. The third block is labeled 'GP Visits' and is a darker blue. The bottom block is labeled 'General Population: Digital Data' and is dark red. The blocks decrease in width from bottom to top, creating a funnel effect. Light blue trapezoidal shapes connect the blocks, pointing upwards. A small dark blue triangle is at the very top.

Hospitalisations

GP Visits

General Population: Digital Data

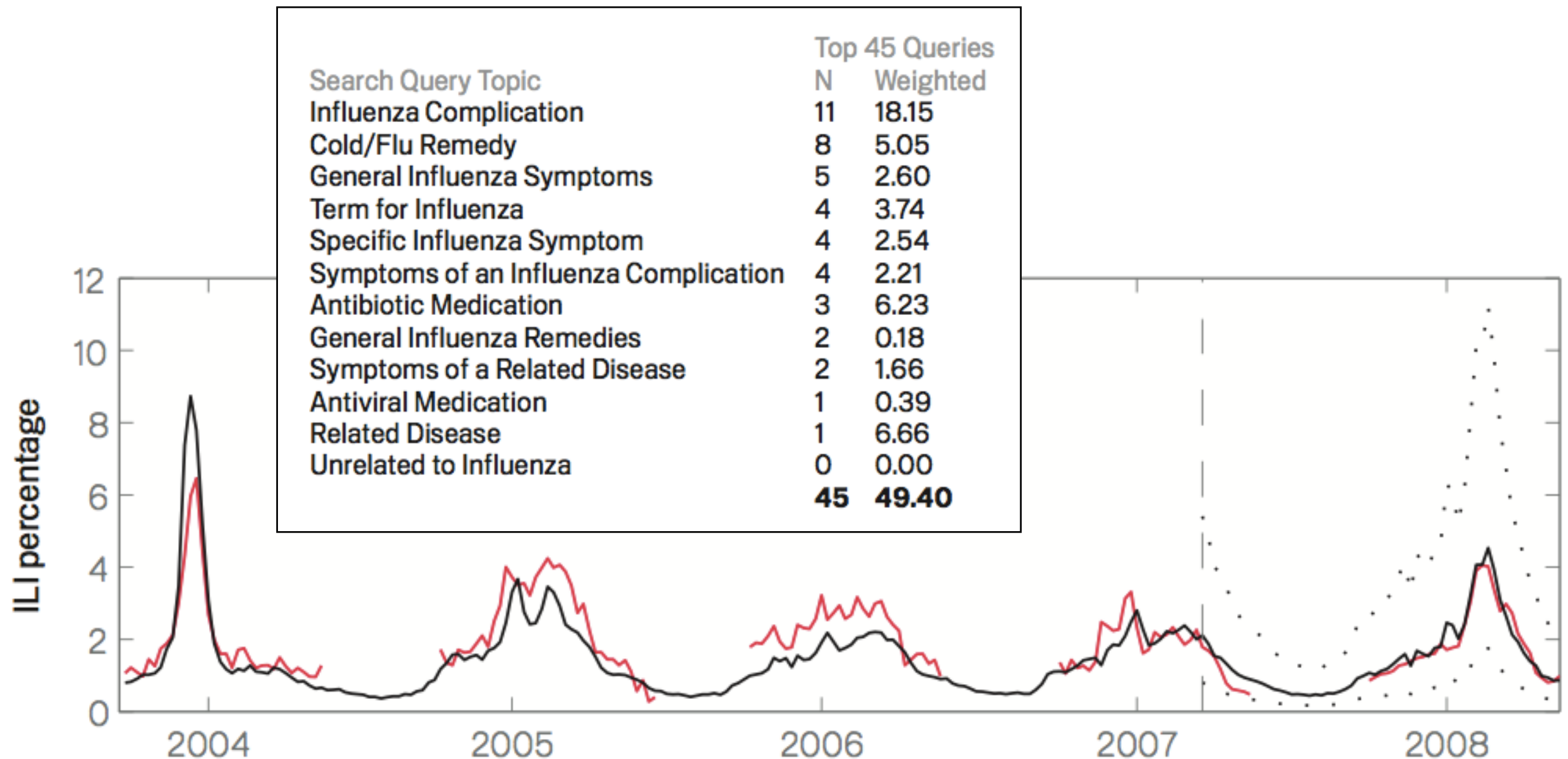
OPPORTUNITIES
DISTRIBUTED
RESEARCH
GENOMICS APPLIED
EXAMPLES
PETABYTES
CURRENTLY USE DESKTOP
TARGET SHARED
DATABASES
GARTNER
CAPACITY
DISK
NETWORKS
SOFTWARE
GROW
USED
MPP
DESCRIBING
INDEXING
CITATION
TYPES
ARCHIVES
TENS
ORGANIZATIONS
SINCE
CAPTURE
SOCIAL
LARGER
NEEDED
HUNDREDS
FORMS
NOW
TERABYTES
CONTINUES
CREATED

BIG DATA

Detecting influenza epidemics using search engine query data

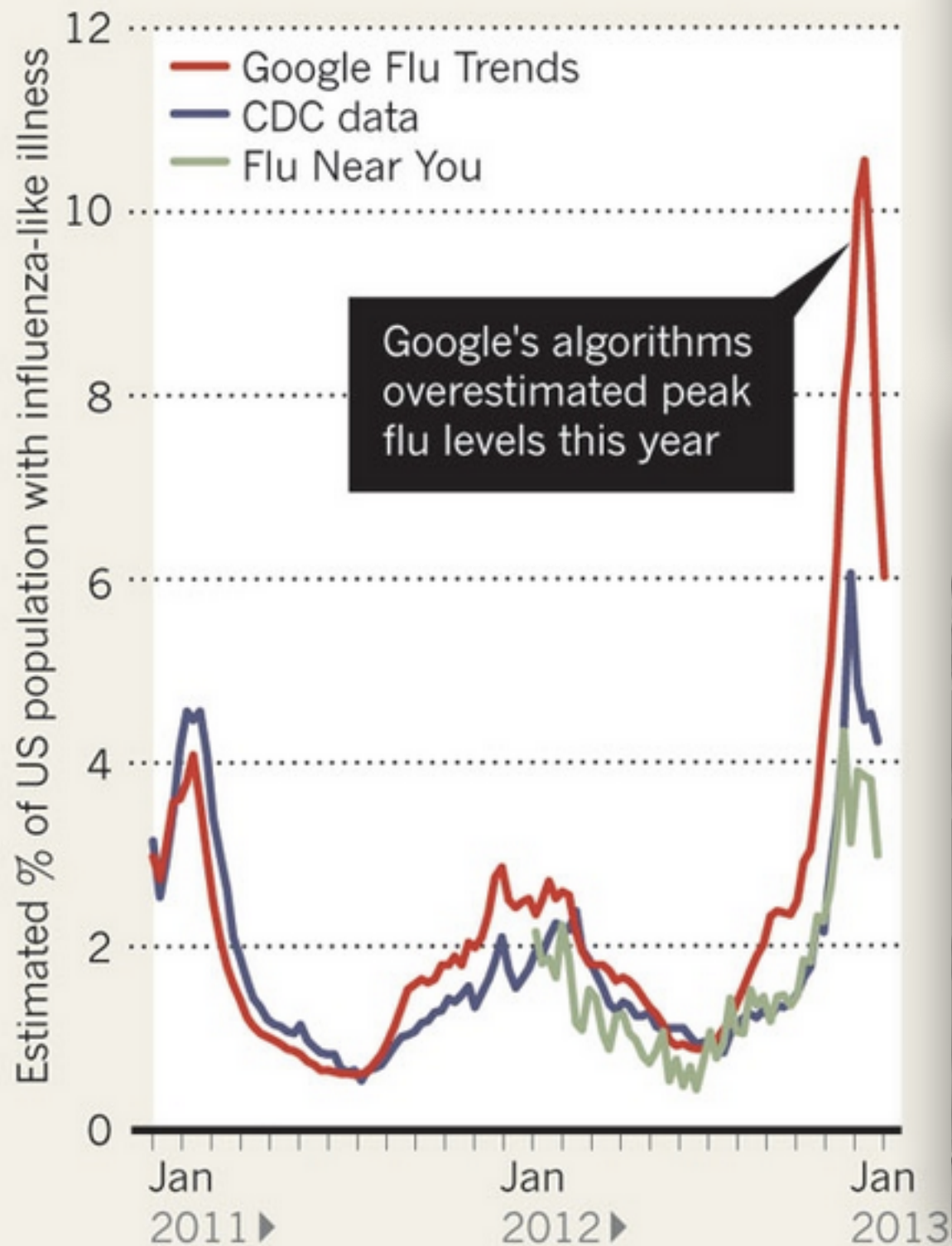
Jeremy Ginsberg¹, Matthew H. Mohebbi¹, Rajan S. Patel¹, Lynnette Brammer², Mark S. Smolinski¹ & Larry Brilliant¹

¹Google Inc. ²Centers for Disease Control and Prevention



FEVER PEAKS

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



When Google got flu wrong

US outbreak foxes a leading web-based method for tracking seasonal flu.

Declan Butler

Science 14 March 2014:
Vol. 343 no. 6176 pp. 1203–1205
DOI: 10.1126/science.1248506

< Prev | Table of Contents

Read Full Text

POLICY FORUM

BIG DATA

The Parable of Google Flu: Traps in Big Data Analysis

David Lazer^{1,2,*}, Ryan Kennedy^{1,3,4}, Gary King³, Alessandro Vespignani^{5,6,3}

⁺ Author Affiliations

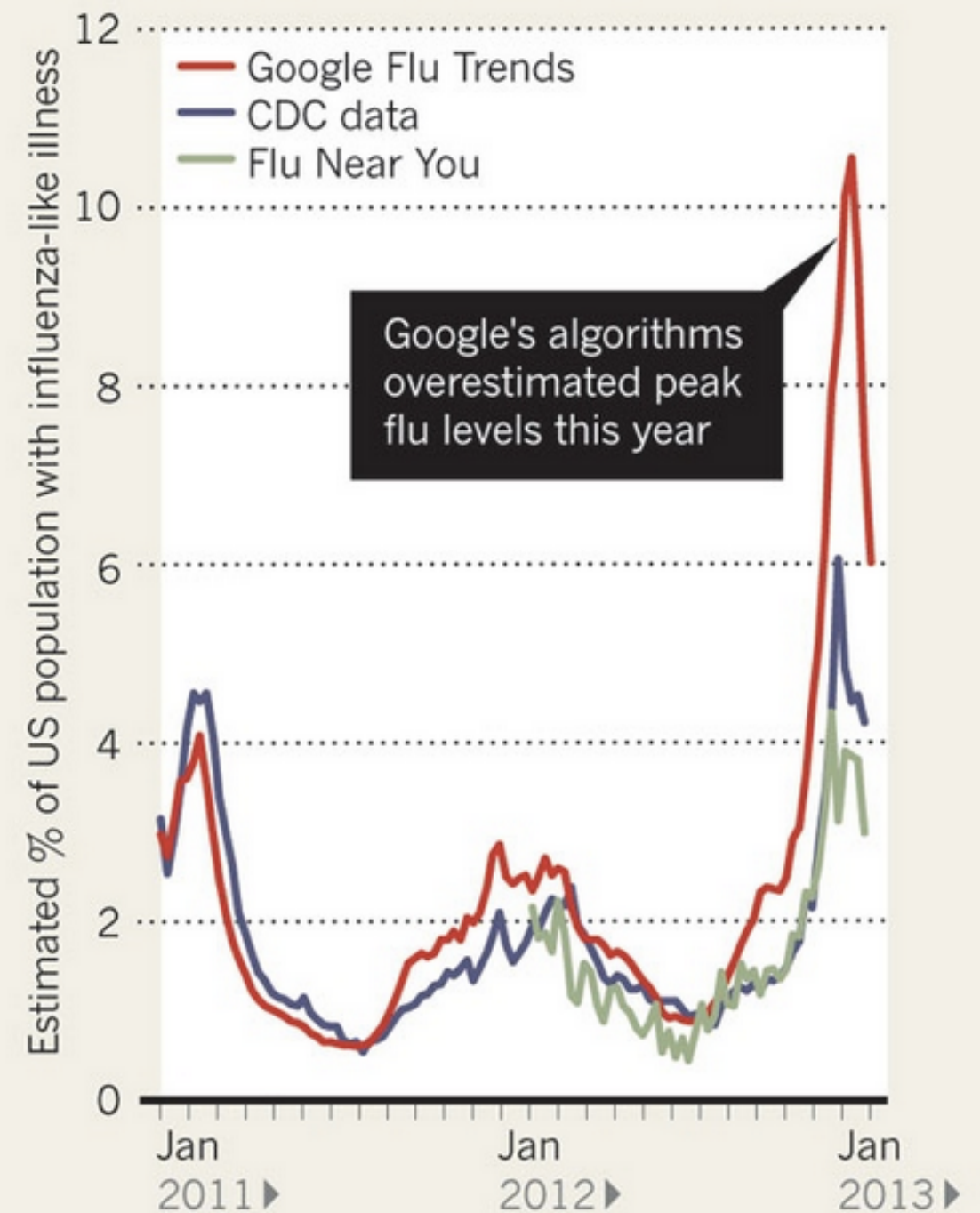
^{*} Corresponding author. E-mail: d.lazer@neu.edu.

In February 2013, Google Flu Trends (GFT) made headlines but not for a reason that Google executives or the creators of the flu tracking system would have hoped. *Nature* reported GFT predicting more than double the proportion of doctor visits for influenza-like illness (ILI) than the Centers for Disease Control and Prevention (CDC), which bases its estimates on surveillance data from laboratories across the United States (1, 2). This happened despite the fact that GFT is often held up as an exemplary use of big data (3). What lessons can we draw from this error?

- Passive data sources don't describe who is well
- Low specificity

FEVER PEAKS

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





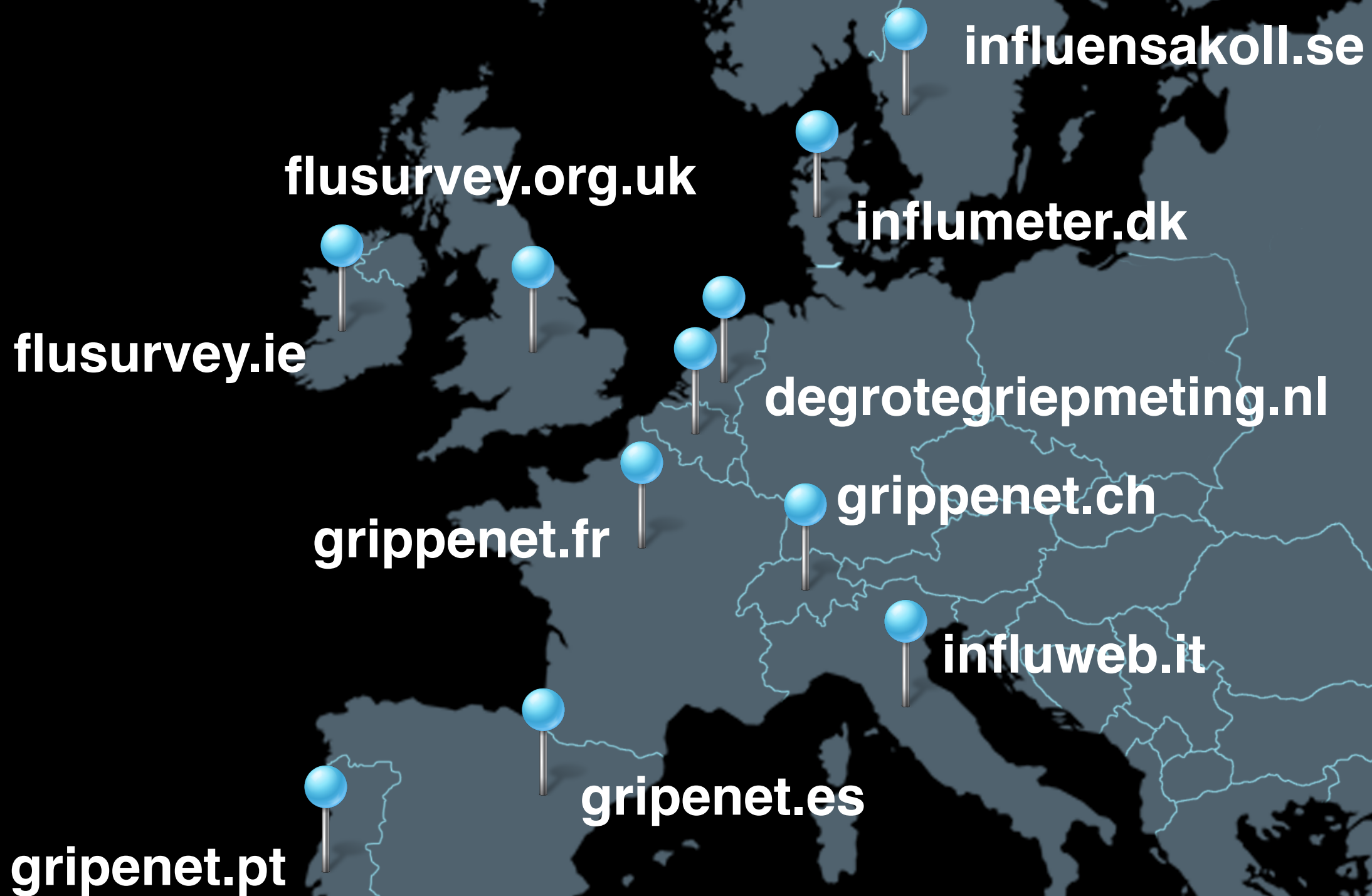
Deaths

Hospitalisations

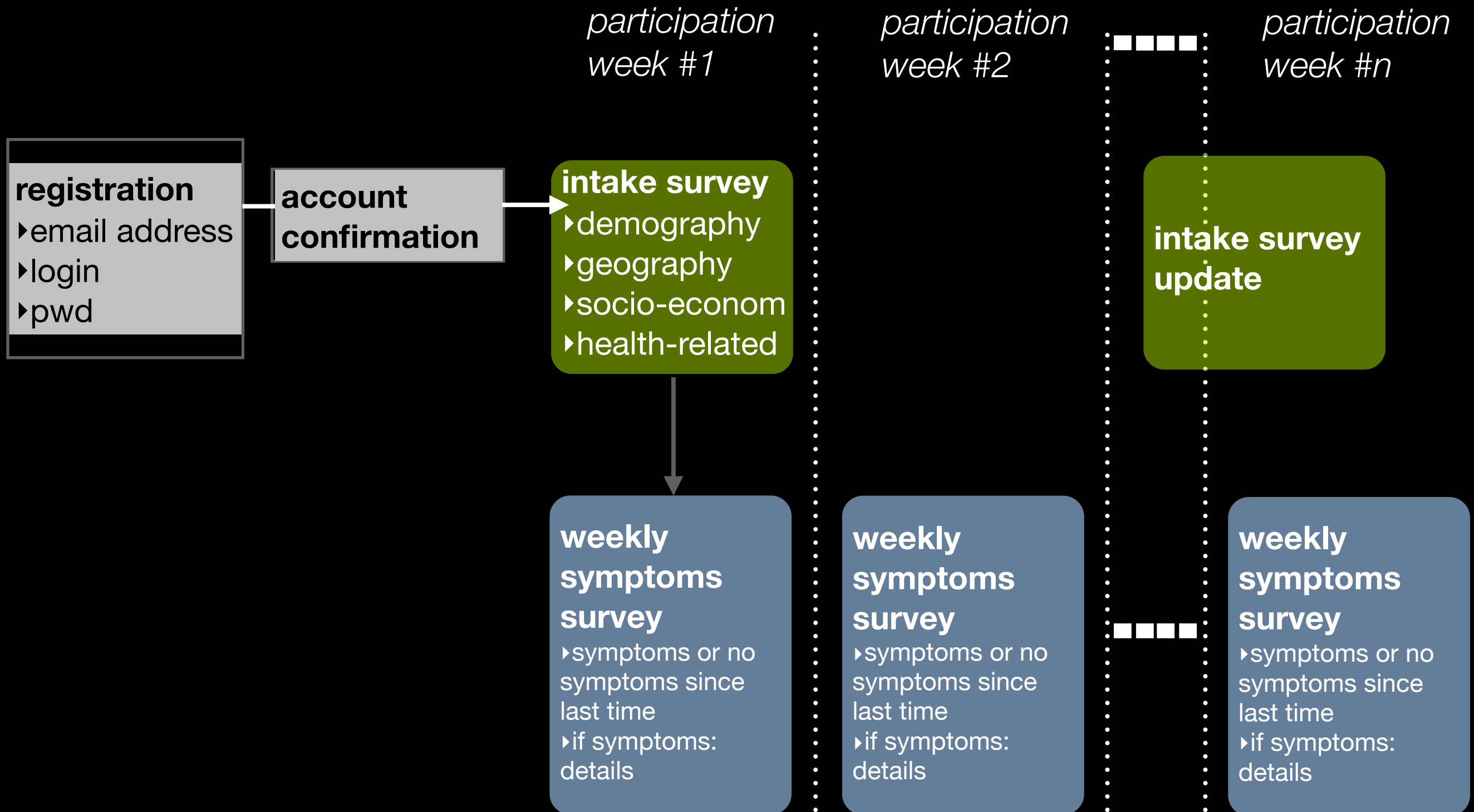
GP Visits

General Population: Influenzanet

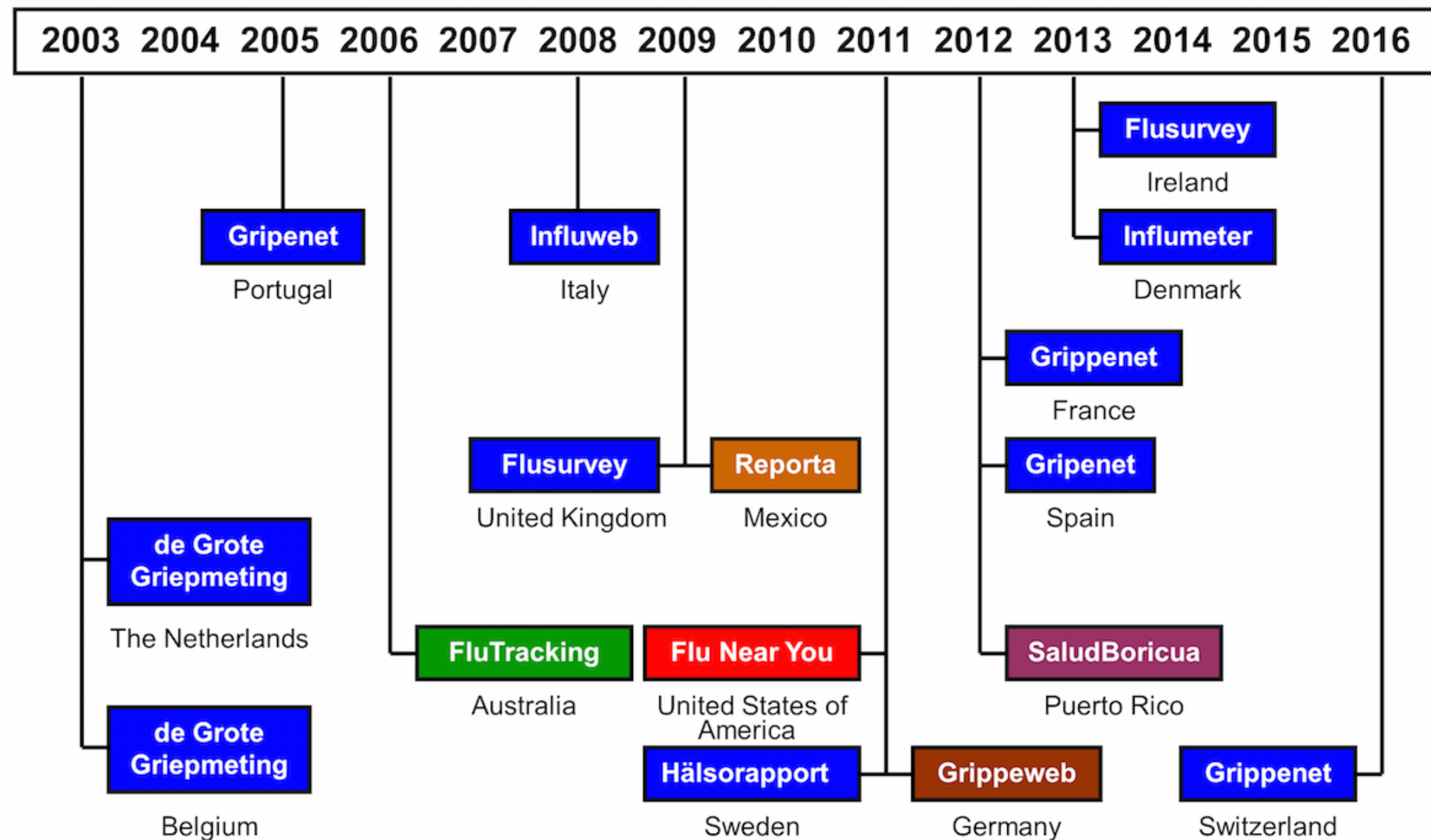
INFLUENZANET.INFO



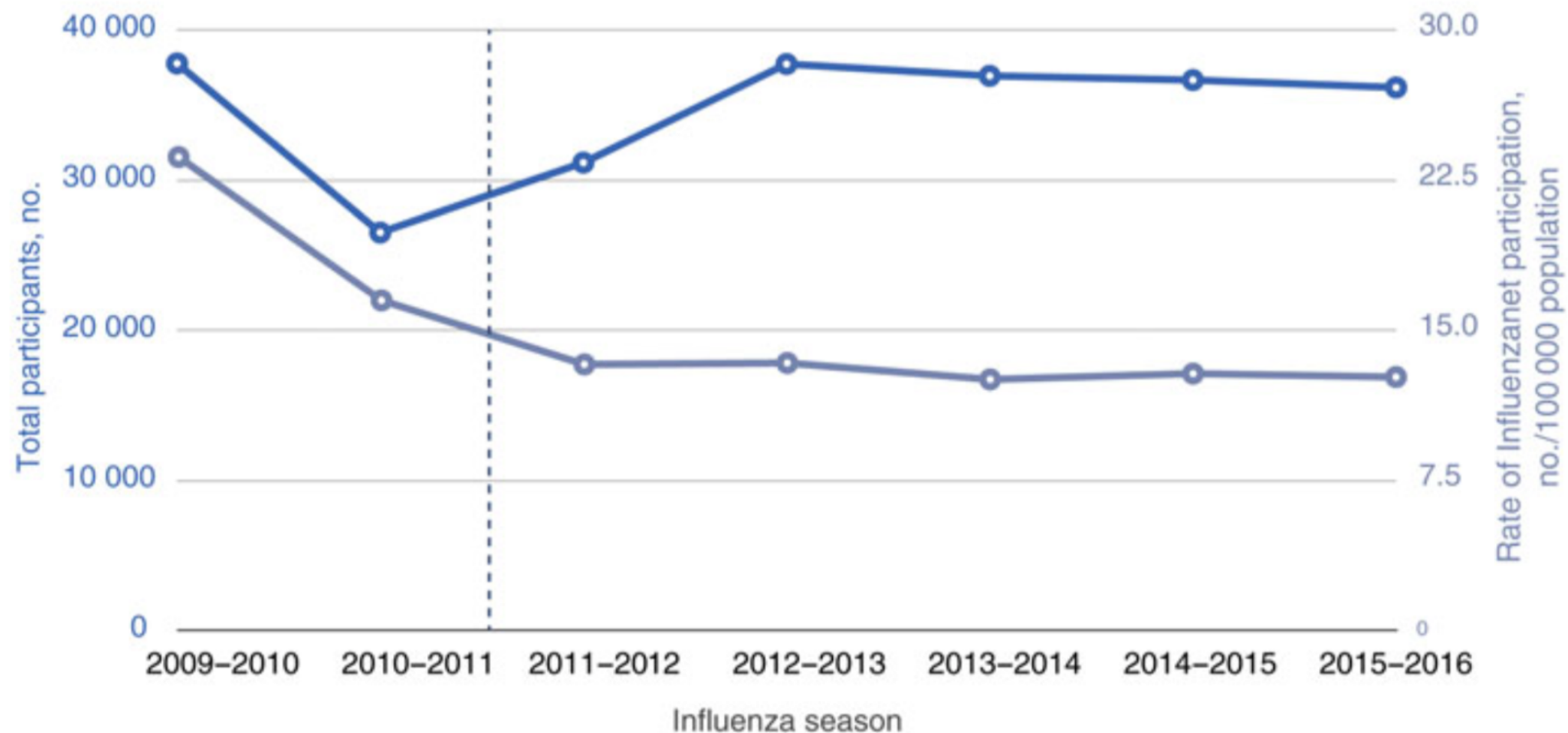
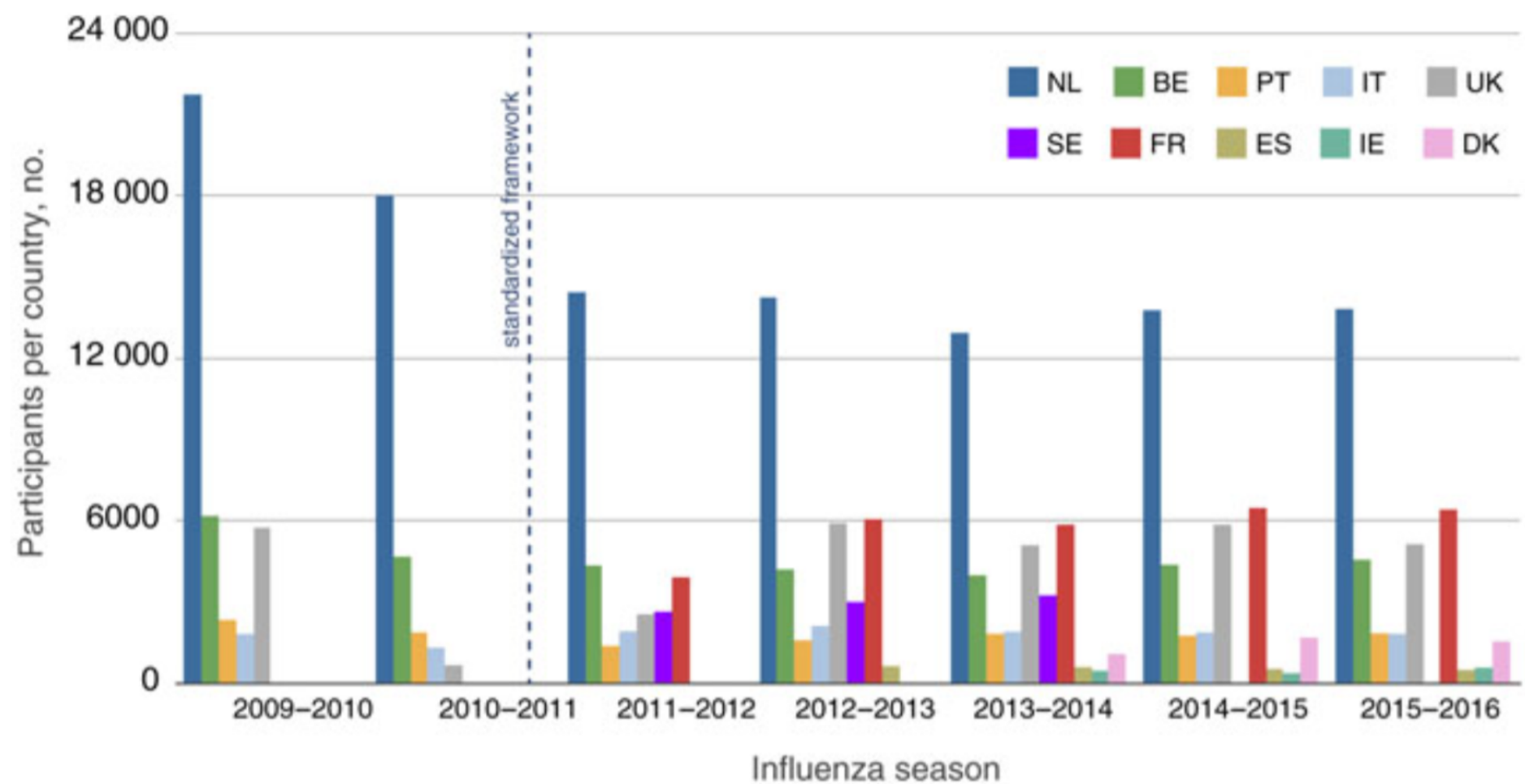
INFLUENZANET study design



INFLUENZANET: a timeline



Koppeschaar CE, Colizza V, Guerrisi C, Turbelin C, Duggan J, Edmunds WJ, Kjelsø C, Mexia R, Moreno Y, Meloni S, Paolotti D, Perrotta D, van Straten E, Franco AO
 Influenzanet: Citizens Among 10 Countries Collaborating to Monitor Influenza in Europe
 JMIR Public Health Surveill 2017;3(3):e66



C. Guerrisi et al,
*Participatory
 Syndromic
 Surveillance of
 Influenza in Europe*
 Journal of Infectious
 Diseases (2016) 214
 (suppl 4): S386-
 S392.

INFLUENZANET.INFO

influweb.it - ISI Foundation & Istituto Superiore di Sanità, Italy

grippenet.fr - INSERM, France

gripenet.pt - Instituto Nacional de Saúde Doutor Ricardo Jorge, Portugal

flusurvey.net - Public Health England

influenzakoll.se - Public Health Agency of Sweden

influmeter.dk - Staten Serum Institute, Denmark

grippenet.ch - Global Health Institute, Geneva

grippeweb.rki.de - Robert Koch Institute, Germany

influweb
|||||

grippenet.fr

gripenet 

flusurvey 

www.influenzakoll.se
Kartlägger förekomst och spridning av influensa i Sverige

influmeter.dk
Kortlægger forekomst og spredning af influenza i Danmark

 **grippe**net

GrippeWeb

[Summary](#)[Primary care data](#)[Severity](#)[Virus characteristics](#)[By country](#)[System](#)[Archives](#)

Week 02/2020 (6-12 January 2020)

- Activity increased compared to week 01/2020, particularly in the southern part of the Region, with two Member States reporting high intensity and six reporting medium intensity. The remainder reported baseline or low intensity levels.
- The percentage of samples from sentinel ILI surveillance patients that tested positive for influenza virus increased from 27% in the previous week to 40% this week.
- The majority of reported influenza virus detections from sentinel ILI surveillance across the Region for week 02/2020 were type A (67%): this percentage has decreased from a high of 78% in week 49. The distribution of viruses detected varied between Member States and areas and within sub-regions.
- Data from the 22 countries or regions reporting to the [EuroMOMO](#) project indicated that all-cause mortality was at expected levels for this time of the year.
- ECDC published an [Influenza virus characterization](#) report, summarizing surveillance data in Europe through December 2019

2019/20 season overview

- For the Region as a whole, influenza activity commenced earlier than previous years.
- Influenza activity in the European Region, based on sentinel sampling, first exceeded a positivity rate of 10% in week 47/2019 and has remained over 10% for 8 weeks. There has been an overall increasing trend in the weekly positivity rate for influenza virus detections among sentinel ILI surveillance patients, following a dip in week 52.
- Type A viruses have dominated across the European Region, though several Member States and areas have reported influenza type B virus dominance or co-dominance of types A and B viruses.
- In sentinel sources, both influenza A subtypes, A(H3N2) and A(H1N1)pdm09, are co-circulating and of the influenza B viruses, the vast majority (98%) have been B/Victoria lineage.
- [Influenzanet](#), which uses self-reported symptoms for ILI surveillance in the general population of European countries, is included in the bulletin as a pilot for the 2019/2020 season.
- ECDC and WHO Regional Office published a joint [Regional Situation Assessment](#) for the 2019–2020 influenza season up to week 49/2019, which focused on disease severity and impact on healthcare systems to assist forward planning in Member States.

Influenza intensity, spread and dominant virus type/subtype

Intensity ▼



2020-W02 ▼

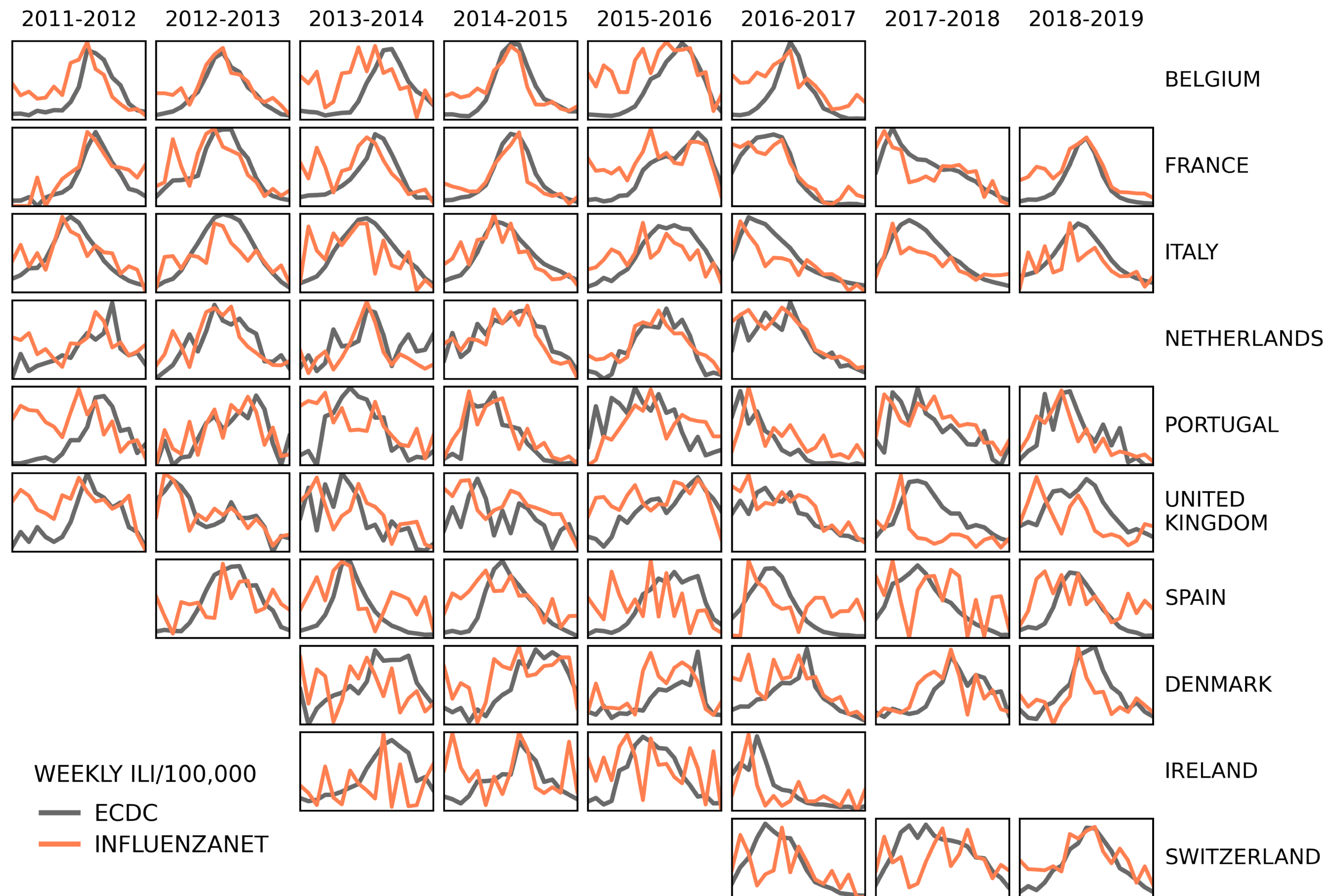


Map type

EU Map



 Export



What is the definition of Influenza-like illness?

Sudden onset of symptoms

AND

at least one of the following four systemic symptoms:

Fever or feverishness, Malaise, Headache, Myalgia

AND

at least one of the following three respiratory symptoms:

Cough, Sore throat, Shortness of breath

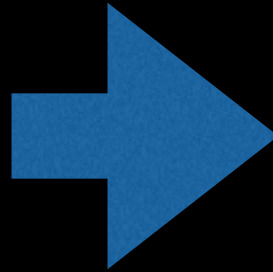
“What if the observed symptoms are the result of a superposition of latent syndromes characterised by an unknown incidence and an unknown composition in terms of symptoms?”

Data

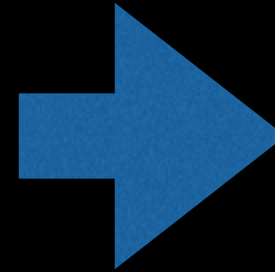
	NL	BE	IT	FR	UK	ES	PT	DK	IE
Number of seasons	6	6	6	6	6	5	6	4	4
Average participants per season	13,450	4,209	1,830	5,757	4,676	526	1,663	1,391	406
Average # surveys per season	206,987	67,420	17,807	68,567	45,543	5,894	17,852	22,782	3,220
Average % of surveys with symptoms	20%	16%	19%	20%	29%	22%	17%	18%	25%
Average of surveys per participant per season	15	16	9	12	9	11	10	16	8

Kalimeri et al, Unsupervised extraction of epidemic syndromes from participatory influenza surveillance self-reported symptoms, Plos Computational Biology 15(4): e1006173

Weekly
Symptoms
Survey



1. Fever
2. Chills
3. Runny/blocked nose
4. Sneezing
5. Sore throat
6. Cough
7. Shortness of breath
8. Headache
9. Muscle/joint pain
10. Chest pain
11. Feeling tired (malaise)
12. Loss of appetite
13. Phlegm
14. Watery, bloodshot eyes
15. Nausea
16. Vomiting
17. Diarrhoea
18. Stomachache
19. Sudden Onset



time
series of daily
symptoms
counts

boolean variables

$$\mathbf{X} = [x_{ij}]$$

matrix whose elements contains the
occurrence of symptom j on day i

Latent Syndromes detection

- it is reasonable to expect that a specific combination of symptoms reported by a user is the symptomatic expression of one or more illnesses, i.e. syndromes, experienced by the user.
- In accordance with this consideration, we postulate that the time series x_{ij} of observed symptoms counts are the result of a linear mixing process driven by K unknown sources, corresponding to the latent syndromes we want to detect.

$$x_{ij} = \sum_{k \in \{1, \dots, K\}} w_{ik} h_{kj} + e_{ij}.$$

Latent Syndromes detection

The mixing equations can be expressed in matrix notation:

$$\mathbf{X} = \mathbf{W}\mathbf{H} + \mathbf{E}$$

$$\mathbf{W} = [w_{ik}], \mathbf{H} = [h_{kj}], \mathbf{E} = [e_{ij}]$$

In this notation, the problem of detecting the unknown K latent sources can then be formulated as a matrix decomposition problem

Non-negative matrix factorization

The specific factorization algorithm we used in this study is a non-negative matrix factorization (NMF) minimizing the Kullback-Leibler loss function:

$$\operatorname{argmin}_{W,H} \sum_{i,j} x_{ij} \log \left(\frac{x_{ij}}{\hat{x}_{ij}} \right) - x_{ij} + \hat{x}_{ij}$$

$$\text{where } \hat{x}_{ij} = \sum_k w_{ik} h_{kj}$$

- this allows a probabilistic interpretation of the decomposition results and, as a consequence, a principled probabilistic way of choosing the intrinsic number K of latent sources or components, based on the model likelihood.

Non-negative matrix factorization

By leveraging on the same probabilistic framework, previously used in the context of semantic analysis of text corpora, we can then interpret the results of the decomposition of X as a mixture of multinomials.

From this probabilistic point of view, by decomposing the matrix X , we are effectively estimating the parameters of a probabilistic model containing a hidden variable which corresponds to the latent component we are looking for and approximating the observed daily proportions of symptoms:

$$\pi(i, j) = x_{ij}/N, \quad N = \sum_{i,j} x_{ij}$$

Non-negative matrix factorization

$$\begin{aligned}\pi(i, j) \approx p(i, j) &= \sum_k p(k) p(i, j|k) \\ &= \sum_k p(k) p(i|k) p(j|k)\end{aligned}$$

mixture of conditionally independent multinomials

$$p(i|k) = w_{ik} / \sum_i w_{ik}, \quad \text{where } \sum_i p(i|k) = 1$$

$$p(j|k) = h_{kj} / \sum_j h_{kj}, \quad \text{where } \sum_j p(j|k) = 1,$$

$$p(k) = \frac{1}{N} \sum_i w_{ik} \sum_j h_{jk}, \quad \text{where } \sum_k p(k) = 1.$$

Non-negative matrix factorization

- The total number of counts N will be proportionally split among K latent components according to $p(k)$.
- These counts will in turn be distributed in each day i according to $p(i|k)$ and finally contribute to the daily symptoms counts according to $p(j|k)$, which describes each component in terms of the expected proportion of symptoms.
- According to this formulation, the total number of counts associated to a latent component k in day i will be given by:

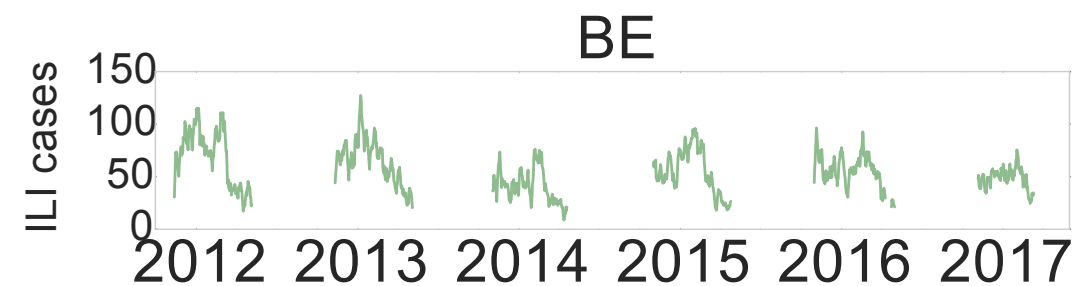
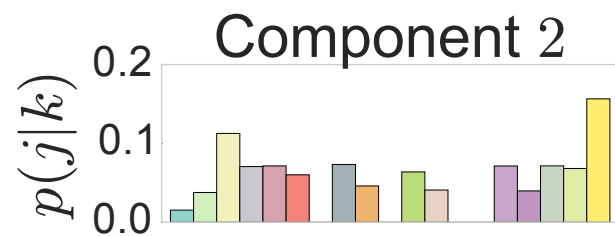
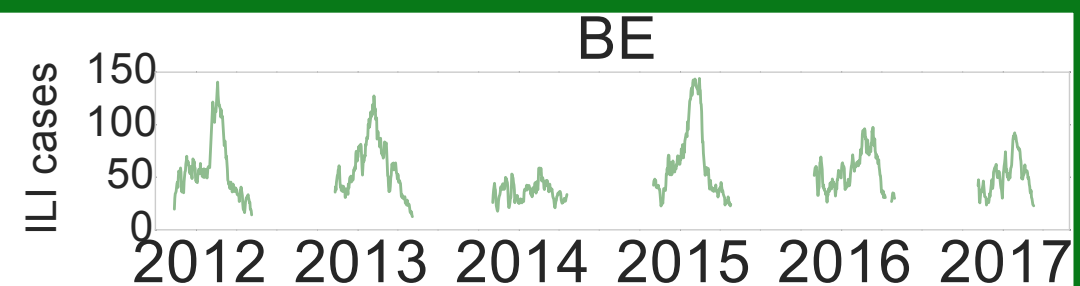
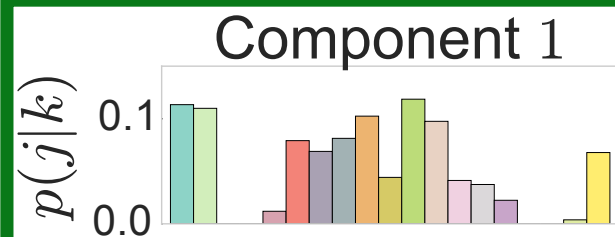
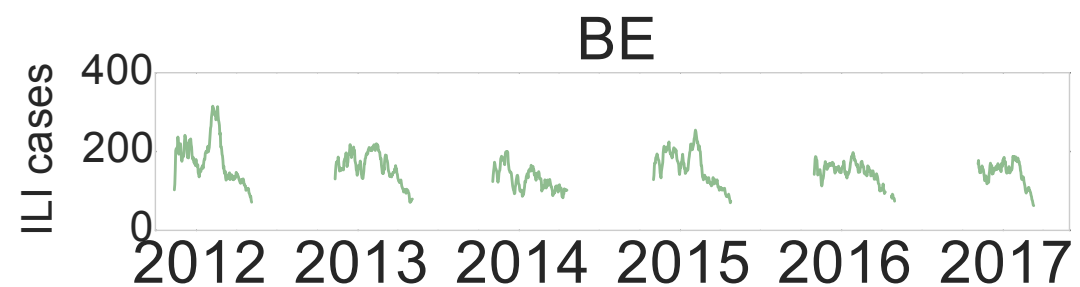
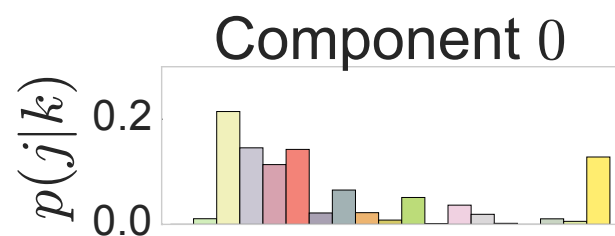
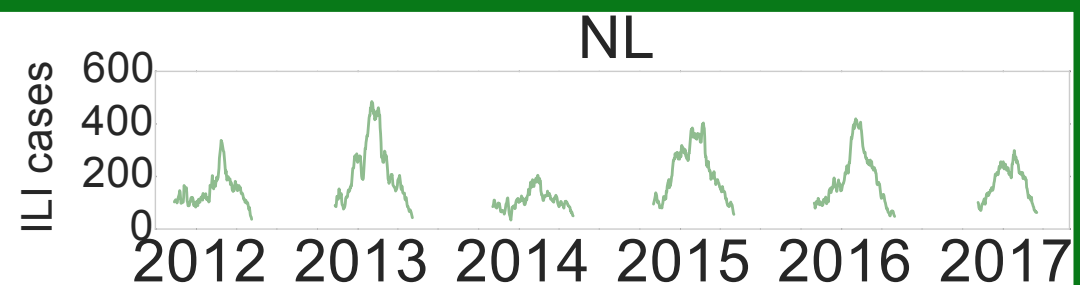
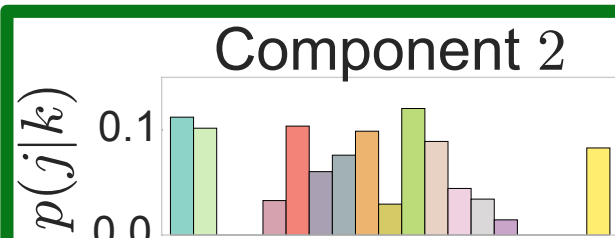
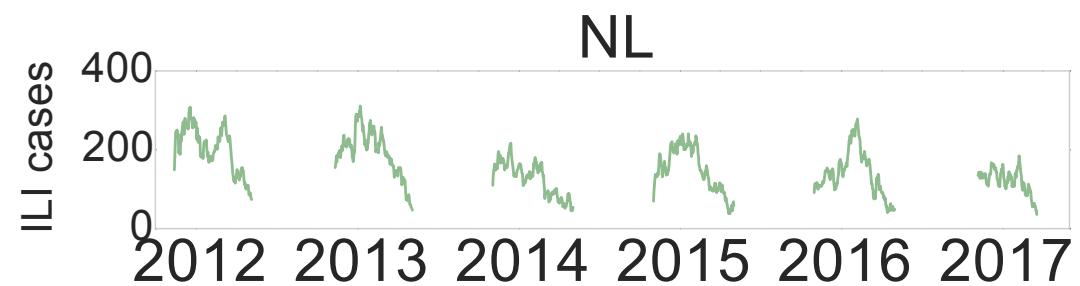
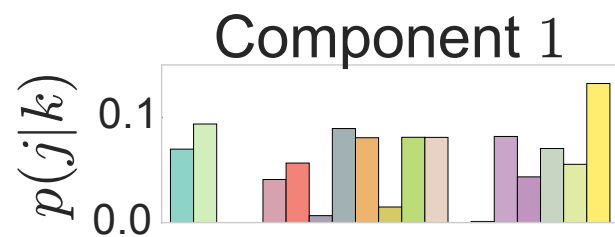
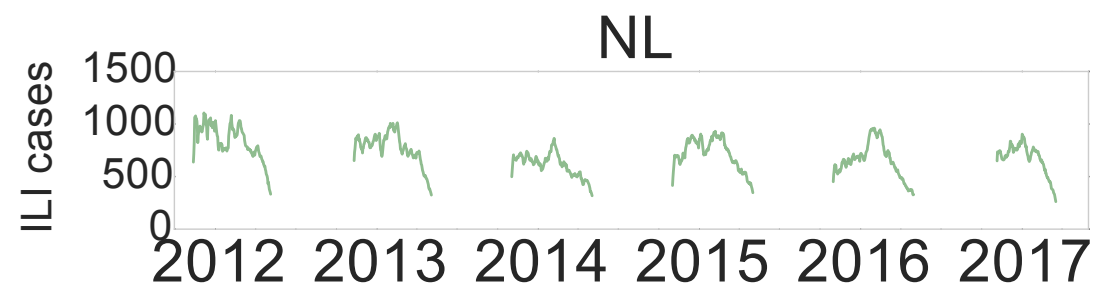
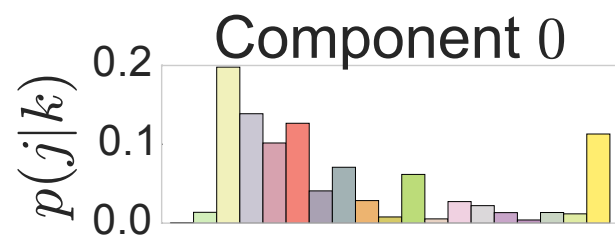
$$y_{ik} = N p(i, k) = N p(k) p(i|k)$$

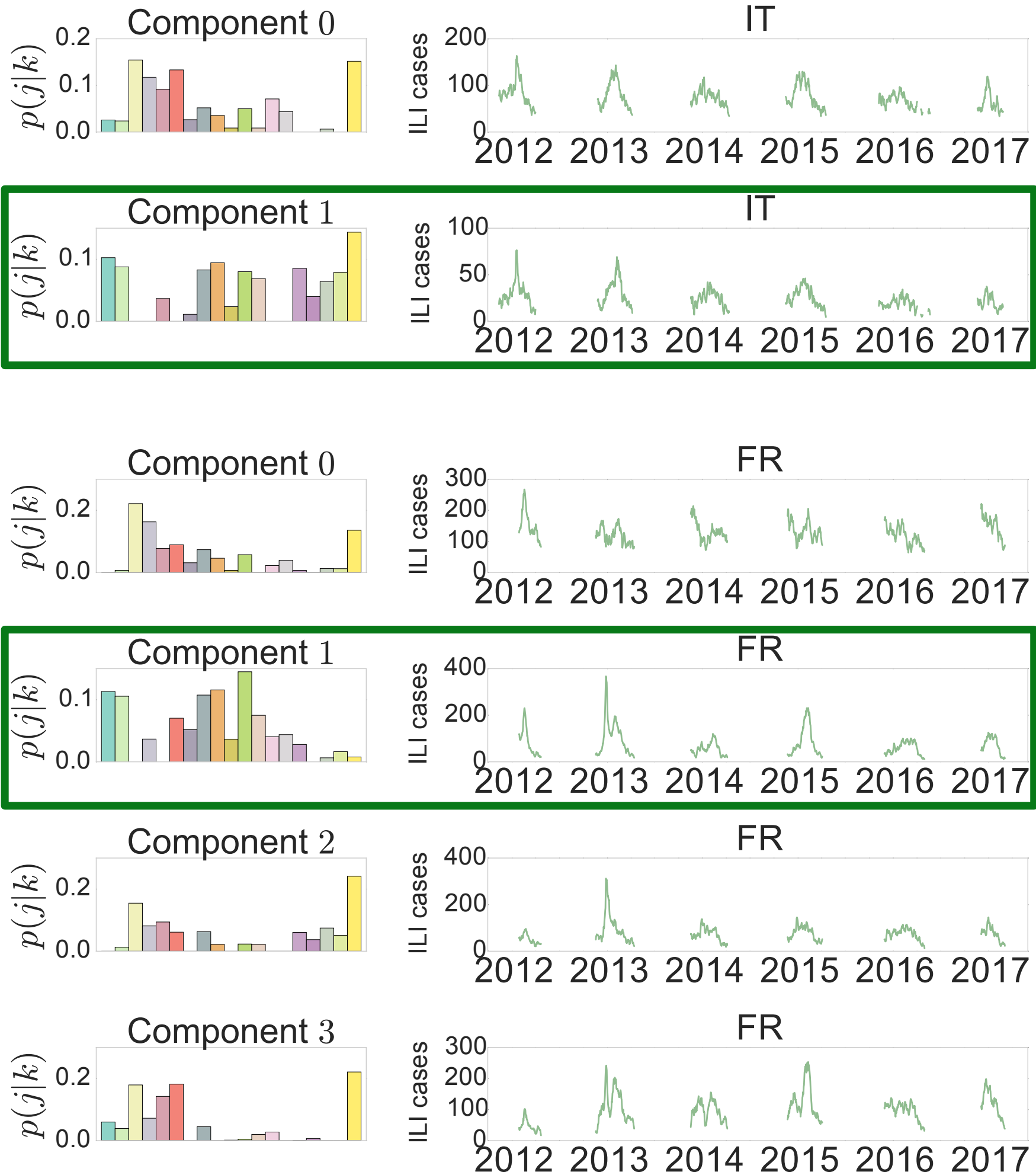
Model selection

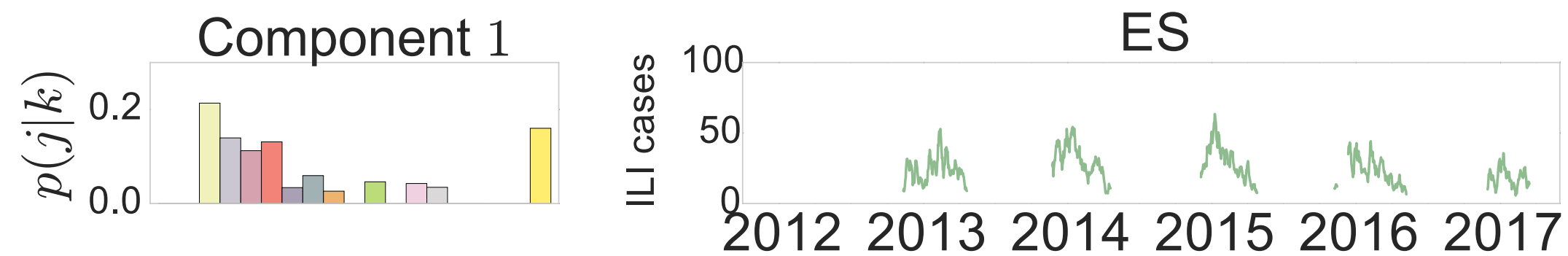
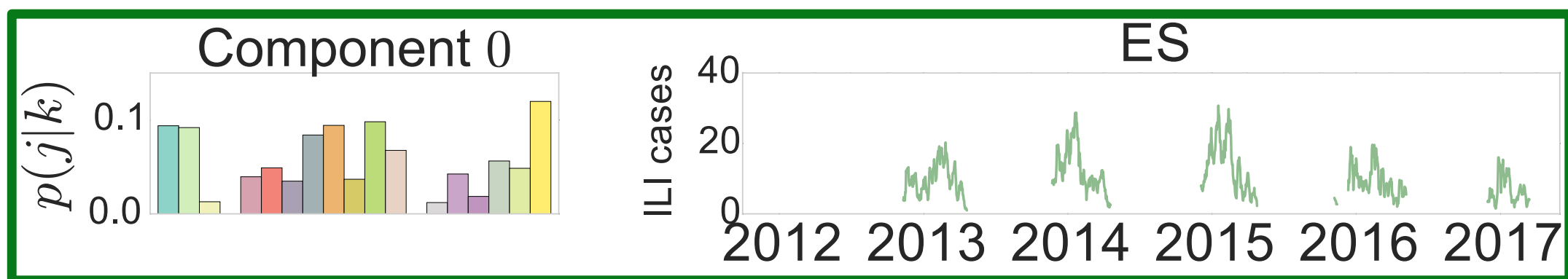
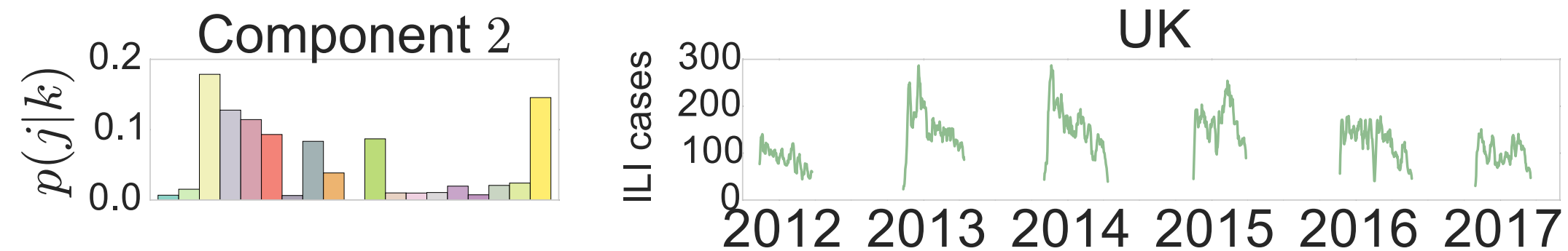
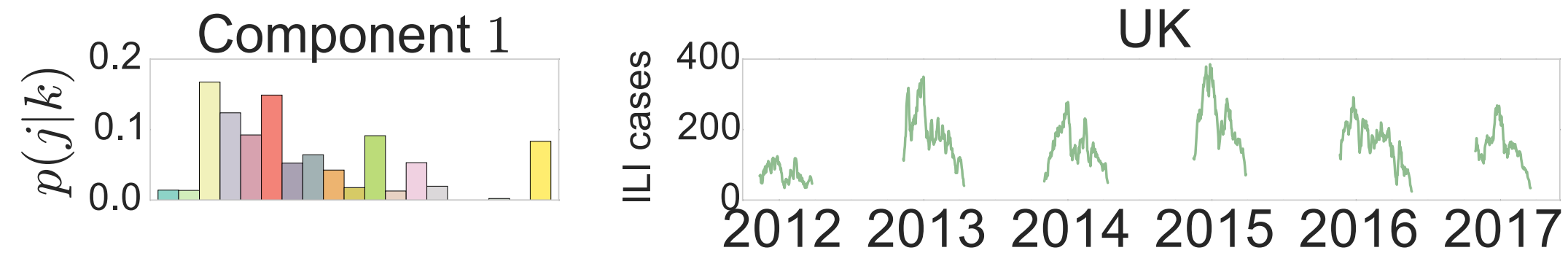
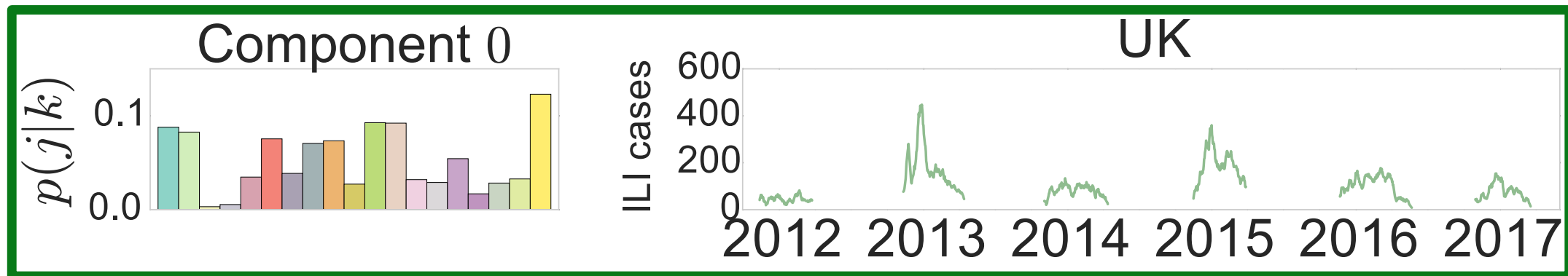
Given a set of candidate models obtained by minimizing the loss function by using an increasing number of hidden components K , we would like to select the best one in terms of its ability to correctly describe the observed phenomenon. The expected value of the Kullback-Liebler loss can be estimated in the asymptotic limit $N \rightarrow \infty$ leading to an approximated model selection criterion by means of the Akaike Information Criterion:

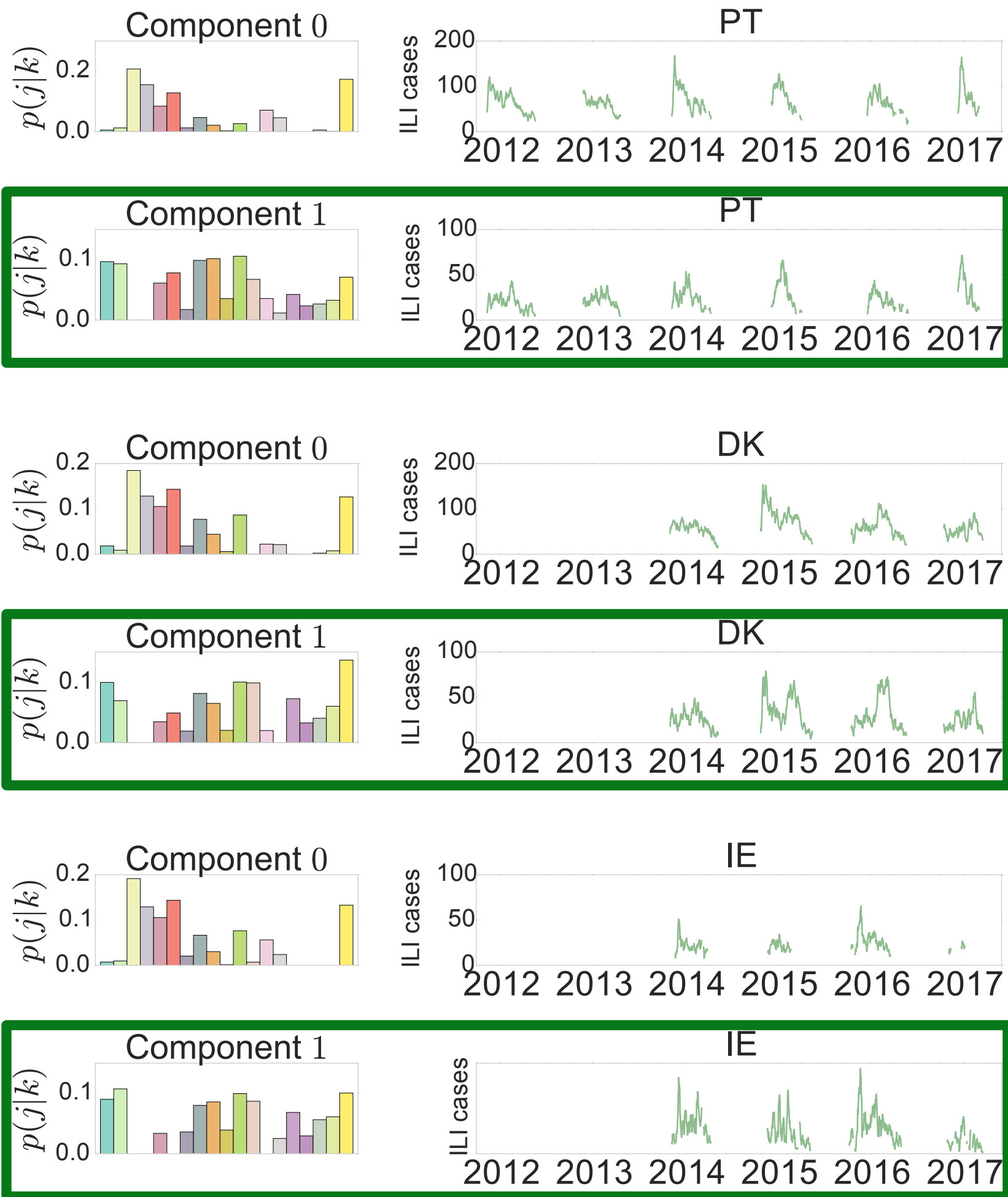
$$AIC_c = -2L(K) + 2P + 2\frac{P(P+1)}{N-P-1},$$

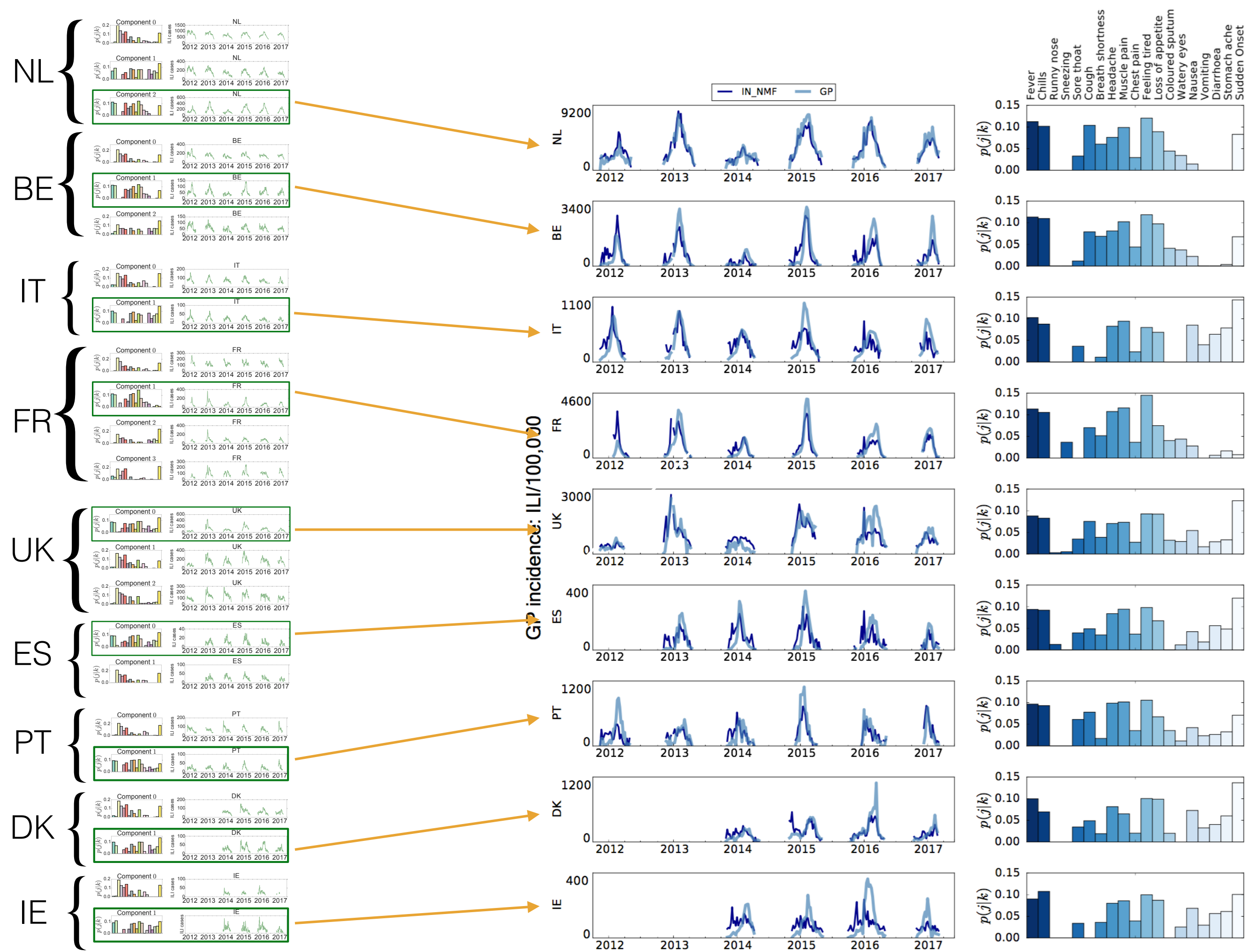
where $L(K)$ is the log-likelihood of the model with K latent components, P is the number of effective parameters of the model: $P = K \setminus, (I + J - 2) - 1$ and N is the the total number of counts. The best model in the set will be the one minimising AIC_c .











	NL	BE	IT	FR	UK	ES	PT	DK	IE
(i) Correlation between IN_ECDC and IN_NMF for the seasons 2011-2017	0.91	0.92	0.86	0.83	0.92	0.86	0.84	0.90	0.82
(ii) Correlation between IN_NMF and GP for the seasons 2011-2017	0.88	0.80	0.69	0.79	0.74	0.65	0.66	0.71	0.38
(iii) Correlation between IN_ECDC and GP for the seasons 2011-2017	0.79	0.72	0.80	0.86	0.75	0.67	0.63	0.68	0.23
(iv) Correlation between IN_NMF prediction for 2016-2017 and GP for the season 2016-2017	0.85	0.82	0.69	0.80	0.60	0.84	0.80	0.76	0.60
(v) Correlation between IN_ECDC and IN_NMF for the season 2016-2017	0.85	0.82	0.86	0.93	0.67	0.59	0.88	0.80	0.71

Kalimeri et al, Unsupervised extraction of epidemic syndromes from participatory influenza surveillance self-reported symptoms, Plos Computational Biology 15(4): e1006173

What's next - I?

- Virological confirmation is needed to estimate more accurately the scaling factor
- extension of the method to other countries and syndromes
- Assess the validity of the method for detection of new emerging diseases

Kalimeri et al, Unsupervised extraction of epidemic syndromes from participatory influenza surveillance self-reported symptoms, Plos Computational Biology 15(4): e1006173

INFLUENZANET.INFO

Daniela Paolotti, Daniela Perrotta (IT)

Vittoria Colizza, Clement Turbelin (FR)

Carl Koppeschaar, Ronald Smalenburg (NL, BE)

Yamir Moreno (ES)

Ricardo Mexia (PT)

Richard Pebody (UK)

Edward Van Straten (SE)

Charlotte Kjelsø (DK)

Jim Duggan (IE)

Antoine Flahaut (CH)

Udo Bucholz (DE)

influweb
|||||

 **grippenet**.fr

deGroteGriepMeting.nl
Het virus in kaart gebracht voor Nederland en België

gripenet 

flusurvey 

www.influensakoll.se
Kartlägger förekomst och spridning av influensa i Sverige

influmeter.dk
Kortlægger forekomst og spredning af influenza i Danmark

 **flusurvey.ie**

 **grippenet**

GrippeWeb

- *D. Perrotta*
- *M. Delfino*
- *K. Kalimeri*
- *C. Cattuto*

The ISI logo consists of the letters 'ISI' in a bold, orange, sans-serif font.

ISI Foundation
& ISI Global Science
Foundation

Thank you!

*–Daniela Paolotti
@danielapaolotti*