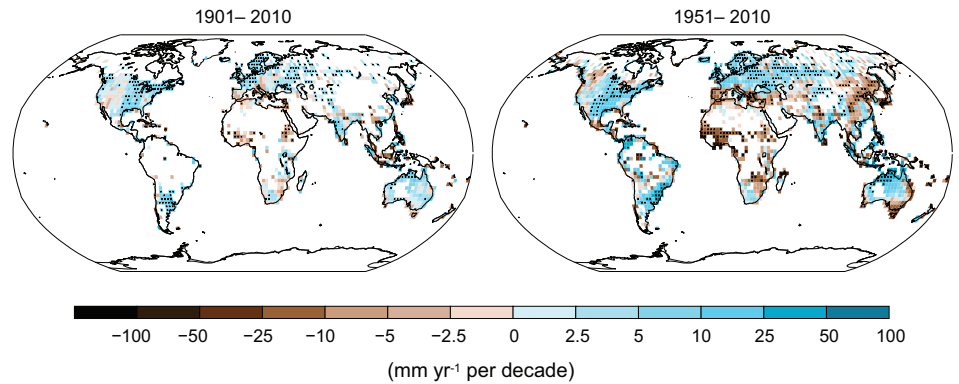
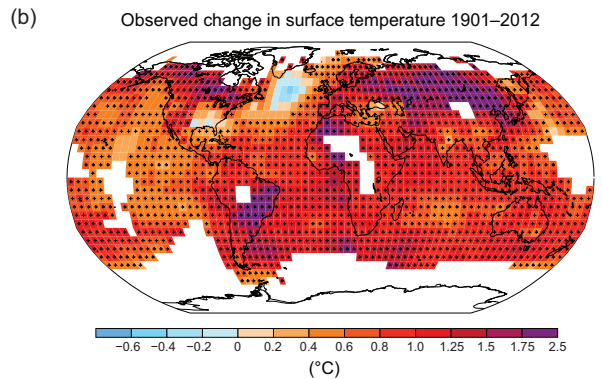
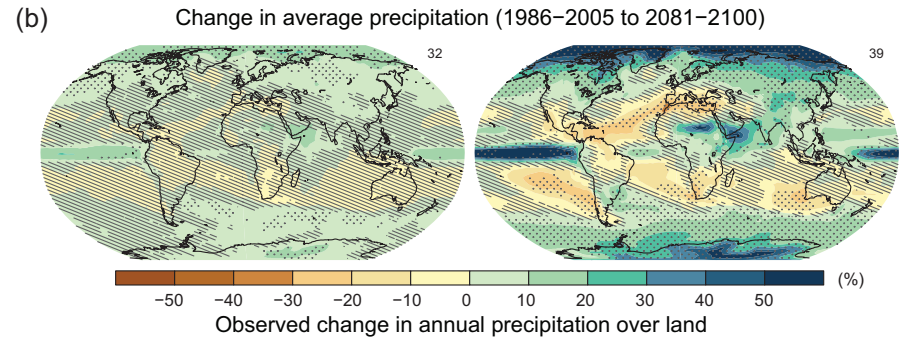
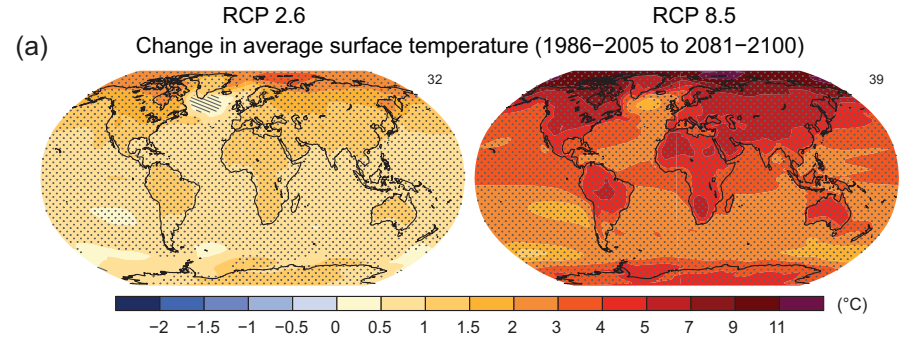
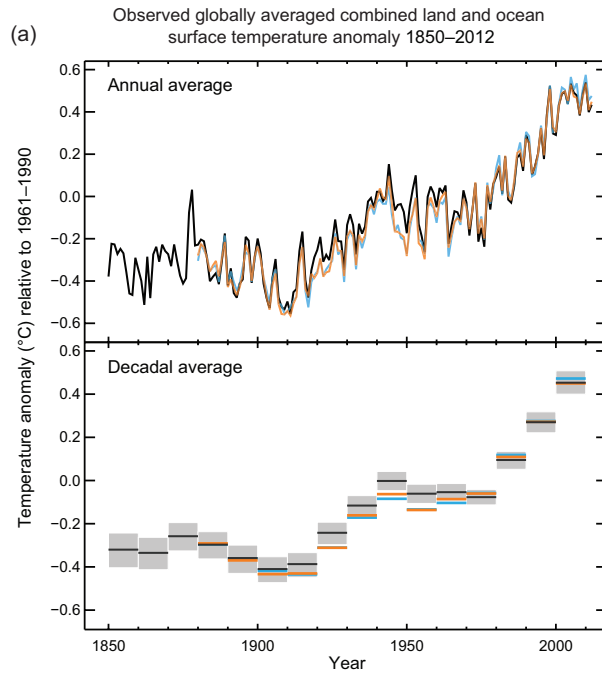


# Big data and the rise and spread of infectious disease



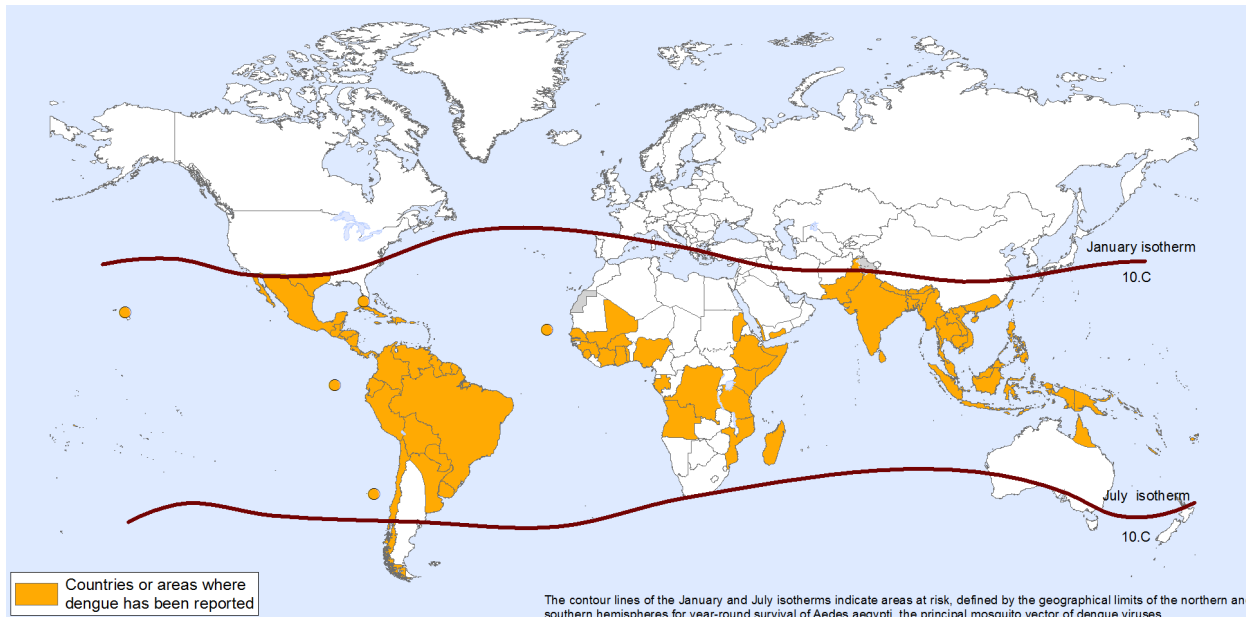
David Taylor, Professor of Tropical Environmental Change  
Department of Geography,  
National University of Singapore

# Climates have changed, are changing & will continue to change .....



Can expect distribution and epidemic potential of climate sensitive infectious diseases to shift ....

Dengue ~ climatically-sensitive vector-borne virus



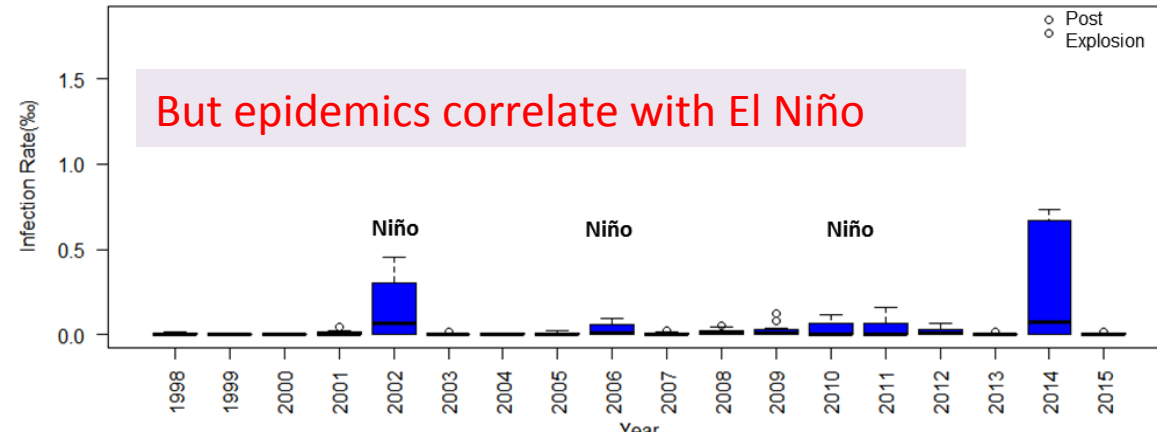
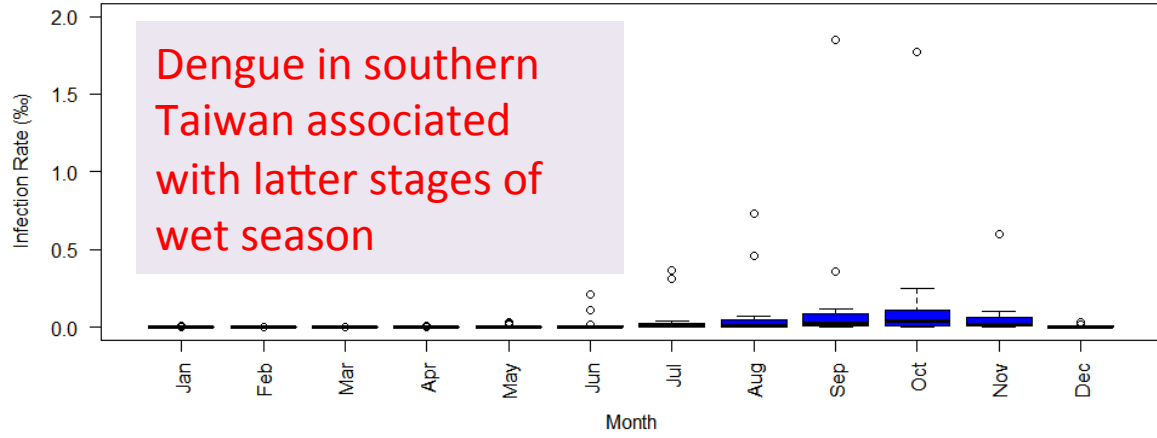
Dengue distribution (2013) From WHO

E.g. Dengue ...

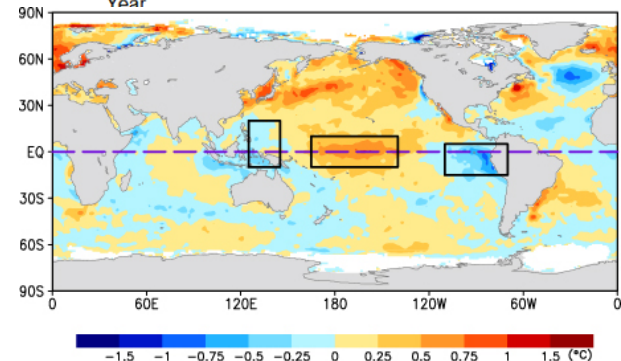
- transmitted by **Aedes** mosquitoes (*Ae. aegypti* and *Ae. Albopictus*)
- Highly **anthrophylic** vectors
- Same vectors transmit **chikungunya, yellow fever** and **zika** (and **Rift Valley fever virus** in Africa!)
- Vectors and diseases are expanding range & severity of infections is increasing
- Not everyone infected shows symptoms!

Recent spread of dengue to higher altitude parts of Asia (e.g. Nepal) suggests that climate change may be a factor

Dengue epidemics in southern Taiwan: a possible link to climate?

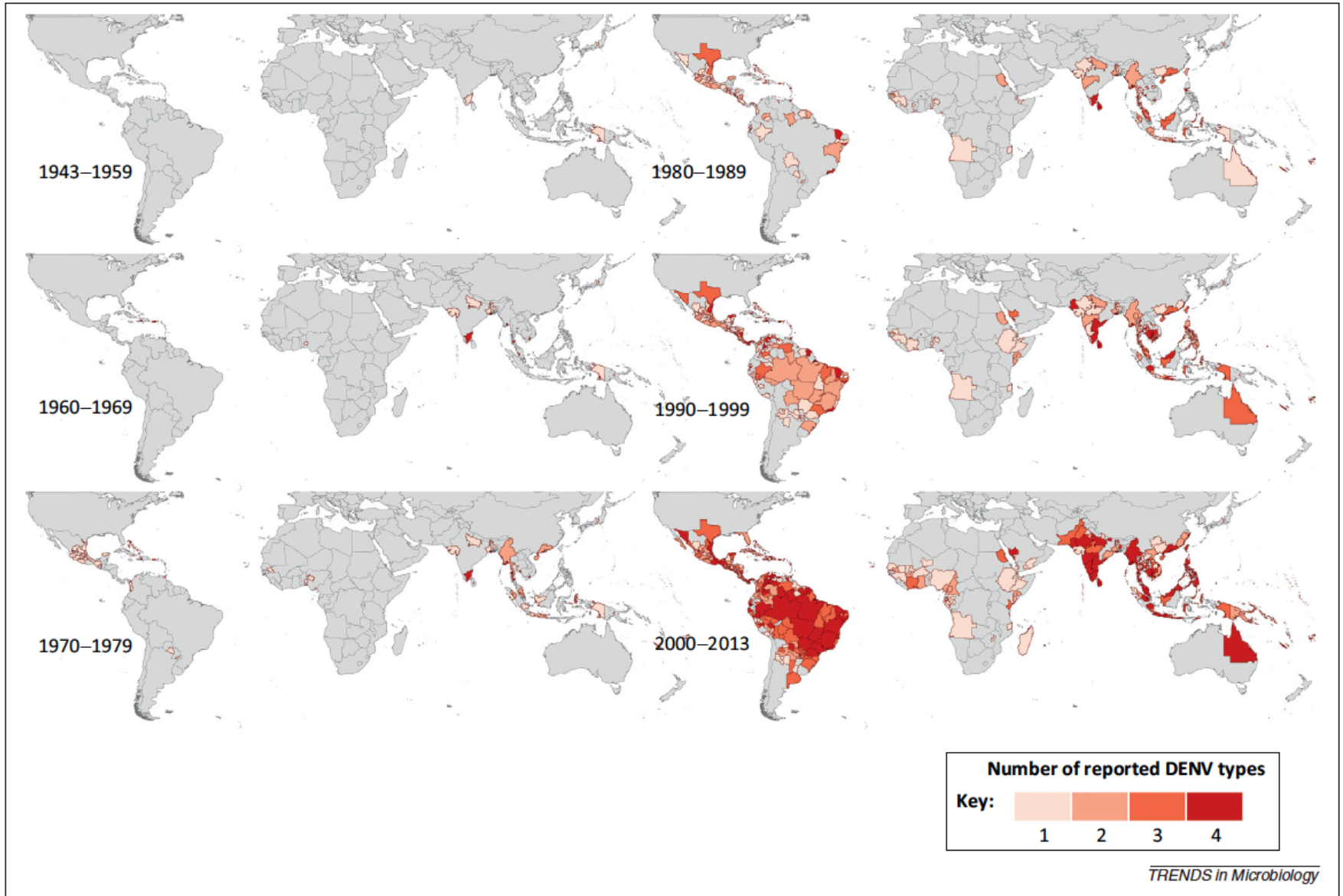


Especially with modoki form of El Niño – modoki El Niño forces south the main typhoon track bringing more rain to Taiwan



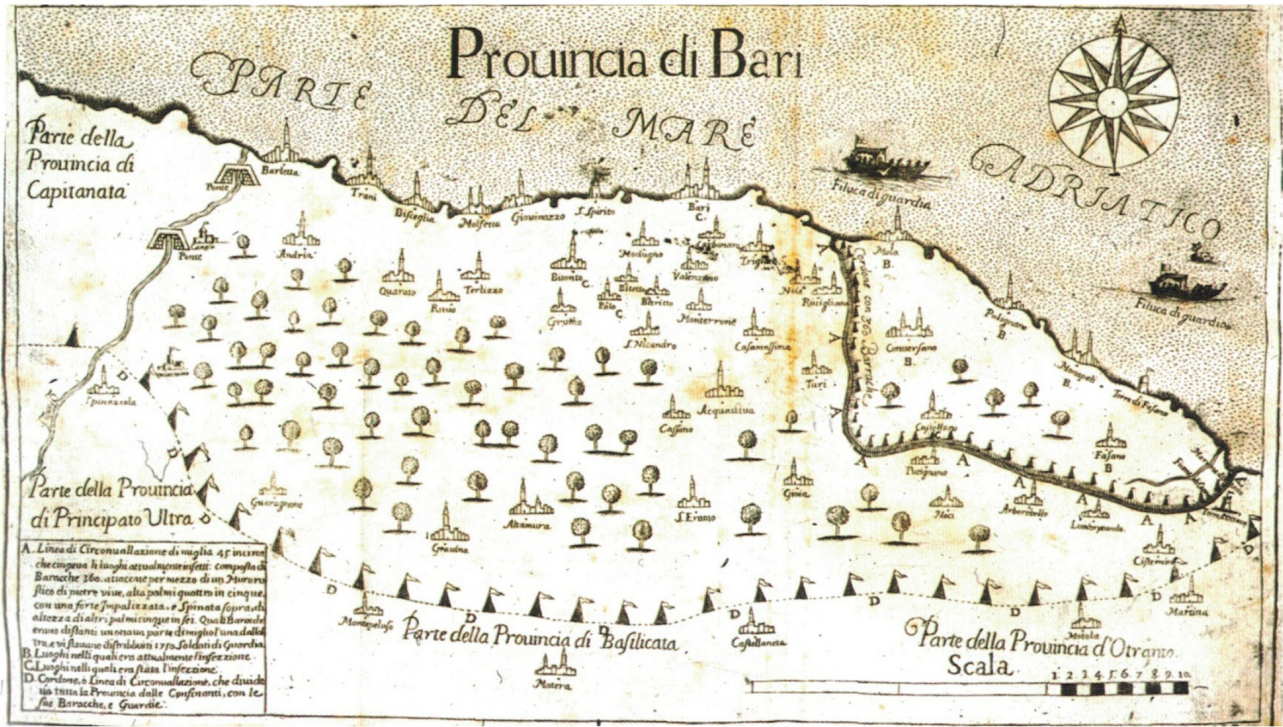
© 2009 Ezilon.com All Right Reserved

# Spread of dengue from late 1970s suggests more than climate change

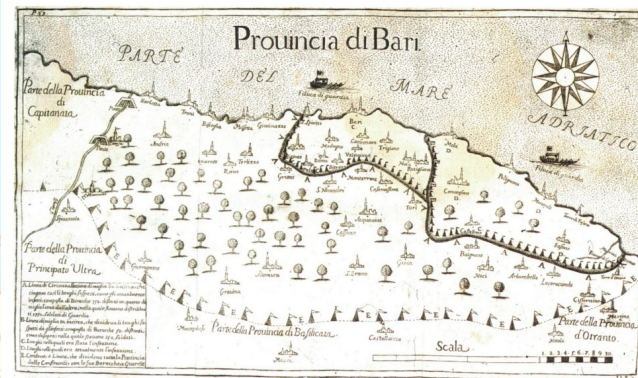


From Messina et al. (2014) Global spread of dengue virus types: mapping the 70 year history, *Trends in Microbiology* 22: 138-146

Rate of spread of new and existing infectious diseases raises questions about conventional disease surveillance and containment techniques ...



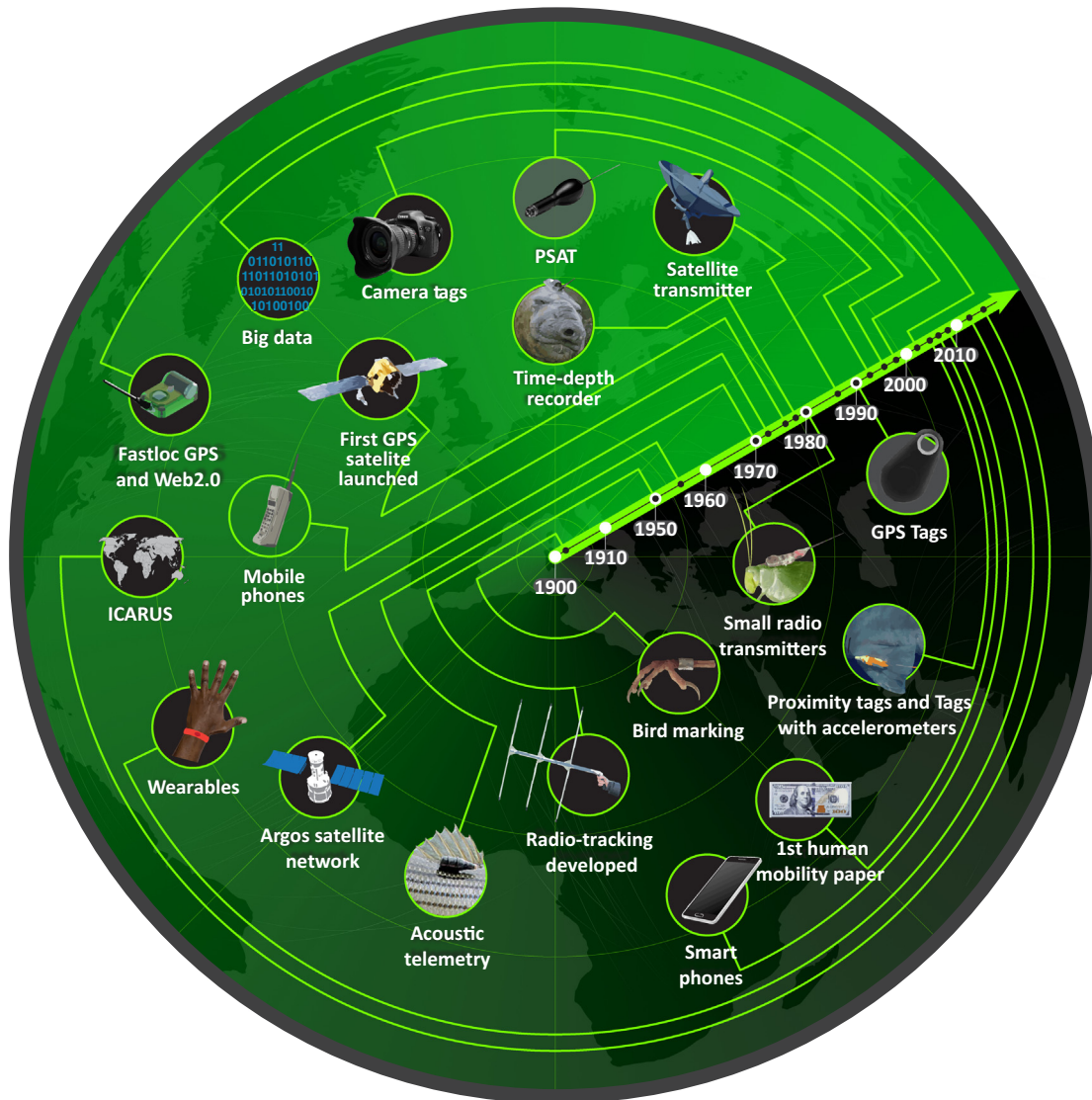
Quarantine – the first public health initiative? First mentioned in 12<sup>th</sup> century, the re-emergence of plague around Bari in the 1690s led to soldiers being used to enforce restrictions on movement ....



Other factors behind the rapid spread of some infectious diseases include land cover change, travel & trade, pollution etc



Just as modern technologies have facilitated the rapid spread of infectious diseases, we can also use modern technologies to monitor and anticipate their effects



Timeline of technological advances in animal movement and human mobility.

From Meekan et al. (2017)  
*The ecology of human mobility. Trends in Ecology & Evolution 32: 198-210*



# No agreed definition of Big Data

Generally (e.g. Kitchin 2013) large datasets that are characterised by the *three vs ++*:

- Huge in *volume*
- High in *velocity*
- Diverse in *variety*

And a 4<sup>th</sup> “v”: veracity (?)



And

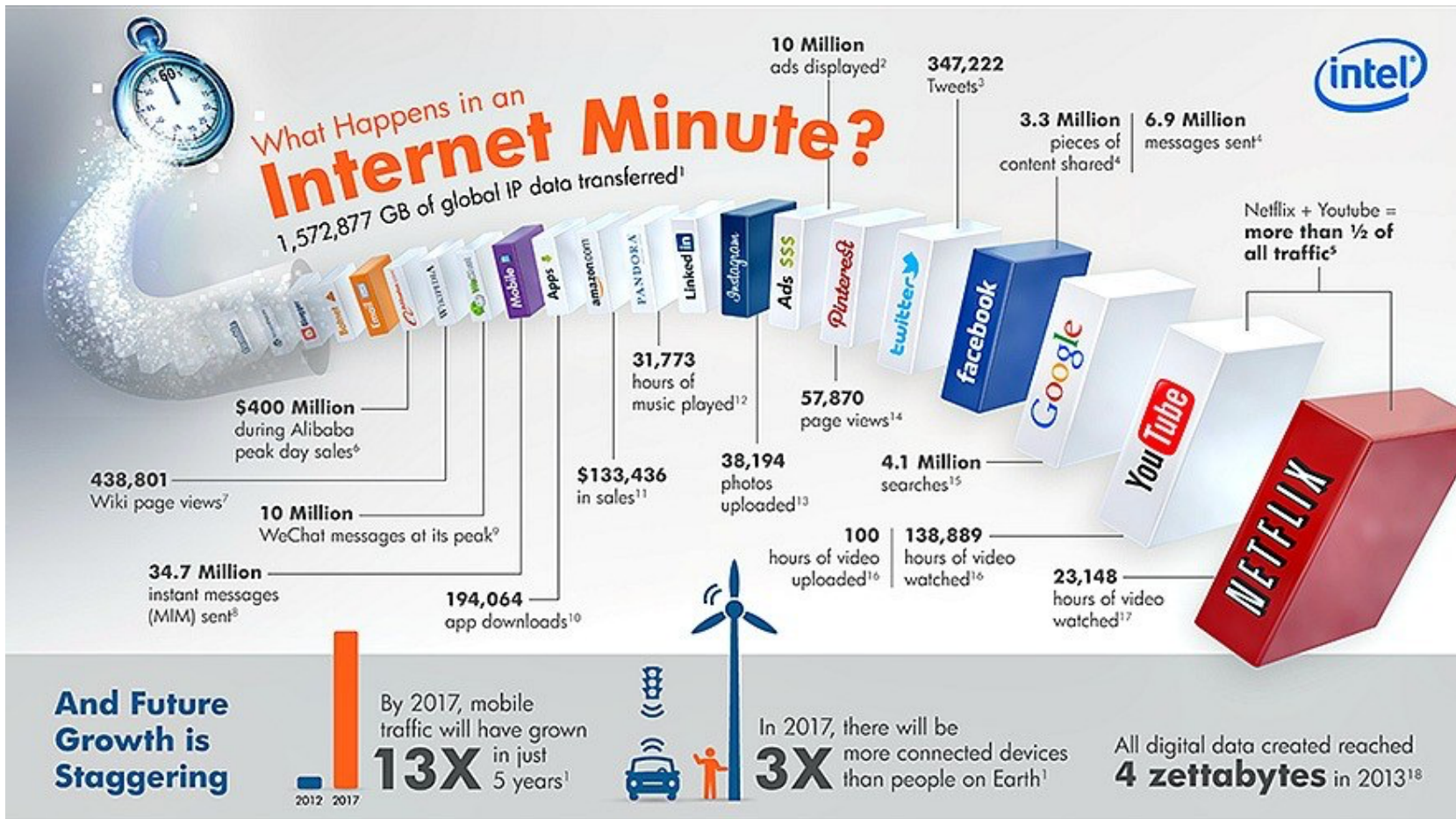
- Exhaustive in **scope** – capture entire populations of data
- Fine-grained in **resolution**
- **Relational** – different datasets can be linked
- **Flexible** – easy to add to and to vary the scale

Also includes the technologies to analyse big data

Different sources of big data:

- 1) **Directed** – generated by digital forms of surveillance
- 2) **Automated** – generated by an inherent, automatic function, e.g. credit card transactions, use of travel cards, tweets, interactions with the internet (clickstream data etc) etc
- 3) **Volunteered** – gifted by users, e.g. crowdsourcing of data, social media posts etc

# The total amount of data produced and consumed is rapidly increasing



We all contribute to the production of data (create “data shadows”), much of it with spatial attributes (“geotags”)

# Big data comes in different forms, and generally needs “cleaning”

**Structured**

ID	ACTL_SVC_NUM	COMMUTER_CATEGORY	ENTRY_DTTM	EXIT_DTTM	ORIG_LOC_ID_NUM	ORIG_X_ORIG_Y
1	818048348A8A8DF021076A74F9C0E6C03931	adult	11/10/2014 12:00	11/10/2014 12:31	92159	
2	5098A29200A796149347A99183F0412D0F8F0	Senior Citizen	11/10/2014 12:00	11/10/2014 12:20	92089	
3	10C4918D01120E70899A218D93AF3A8A8A70DCE	Senior Citizen	11/10/2014 12:00	11/10/2014 12:36	92091	
4	00A1A8E4C8F481B315FC6D46F8BE773099A1	Student	11/10/2014 12:00	11/10/2014 12:11	94611	
5	AAC06979A1058130813DC454F40C0E4729F0D	adult	11/10/2014 12:00	11/10/2014 12:06	84219	
6	0B0794D034848F1C1500E8434921208418604	222P	11/10/2014 12:00	11/10/2014 12:09	84211	
7	0B0794D034848F1C1500E8434921208418604	Student	11/10/2014 12:00	11/10/2014 12:05	84211	
8	974474D79A94071CE4646A08215D48487CD	adult	11/10/2014 12:00	11/10/2014 12:05	82109	
9	90D8565C2135D04C34F006920C3F88F930CF	21	11/10/2014 12:00	11/10/2014 12:31	83633	
10	90D8565C2135D04C34F006920C3F88F930CF	adult	11/10/2014 12:00	11/10/2014 12:31	83633	
11	CA43317076D0C48F481C009390F08709B	adult				
12	ED825E171E0538778A38E6131C3C38428202C0	adult				
13	288AD830C86C308F41891981C40872AFC09AD8	80				
14	0A0808A8E89F47A8D702C173A70A8080931C4	291P				
15	0A0808A8E89F47A8D702C173A70A8080931C4	291				
16	7CD6796CC678F9317E88D93C79477135D0B8	38				
17	CC7B814075E1081A8A4F70268826982C4A	69				
18	C1F1427E093666044A4A70A81E12764A9F36	62				
19	48C390136A27909D9B281332793CDE91DC33	373				
20	0A0CE80E109270C287A8A38146D7F50704C1	163				
21	62DF866040810F18EC0848E8E575A732065	163				
22	D19F0F01C0208EC0D90A7055F3A9357989	81				
23	975EE20B8879CC08781D0F3461093005C4	83				
24	8C09150679AC746F45C6A8A2C07280A3706CC0	3				
25	7F0DCAAEFC5A40C12FF9764282083290C5C5	3				
26	117F1C207F9F6C8A70C78045D2208E8AC	88				
27	4A8676A686C48A8A8A8A8A8A8A8A8A8A8A8A8	43				

**Semi-Structured**

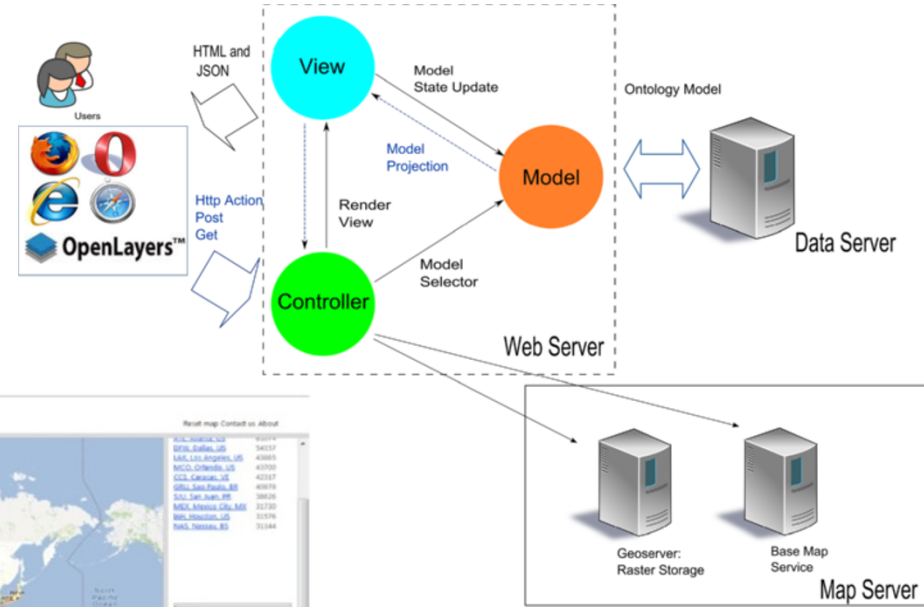
```
{
  "id": "55500355",
  "text": "MASBrowmLEY flying raman monster",
  "user": {
    "id": "8248210166806044",
    "name": "MASBrowmLEY"
  },
  "retweeted_status": {
    "id": "148540708148",
    "text": "Singapore Government"
  },
  "source": "a href='\"http://twitter.com/\"' rel='nofollow'>Twitter Web Client...",
  "coordinates": null,
  "timestamp_ms": "148540708148",
  "entities": {
    "hashtags": null,
    "urls": null,
    "mentions": [
      {
        "id": "55500355",
        "name": "MASBrowmLEY"
      }
    ],
    "media": null
  }
}
```

**Unstructured**  
(texts, photos, videos)

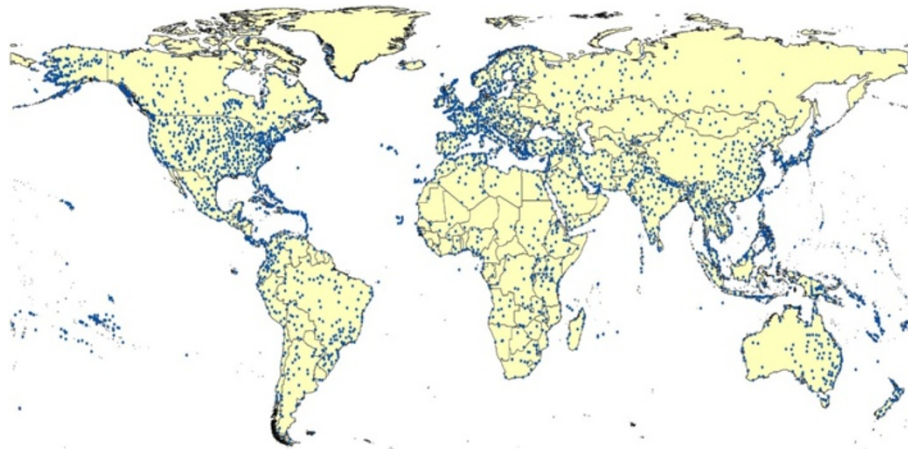
Need to separate signal from noise, identify misinformation (“fake news”) and determine how representative the information is

Big data creates new opportunities (e.g. new sub-disciplines of *infodemiology* and *macroecology*) and brings new insights (e.g. link between travel & spread of infectious disease)

VBD-AIR tool is used to determine the disease importation risk at airports around the world – e.g. dengue for Miami airport, US



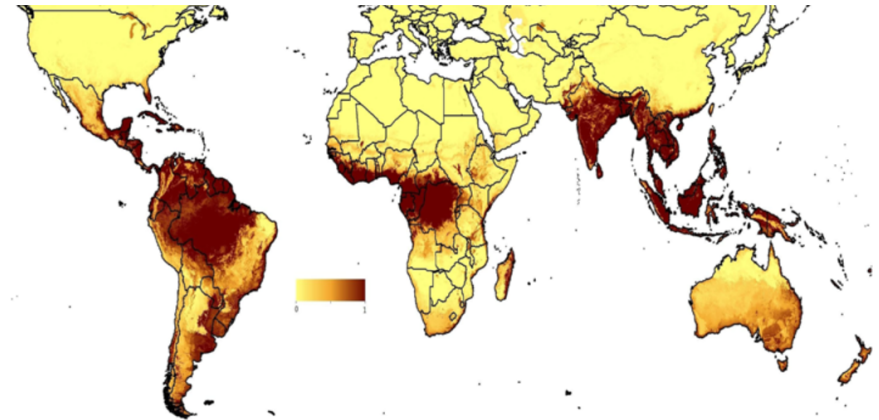
Huang, Z. et al. (2012) Web-based GIS: the vector-borne disease airline importation risk (VBD-AIR) tool. *International Journal of Health Geographies* 11: 33



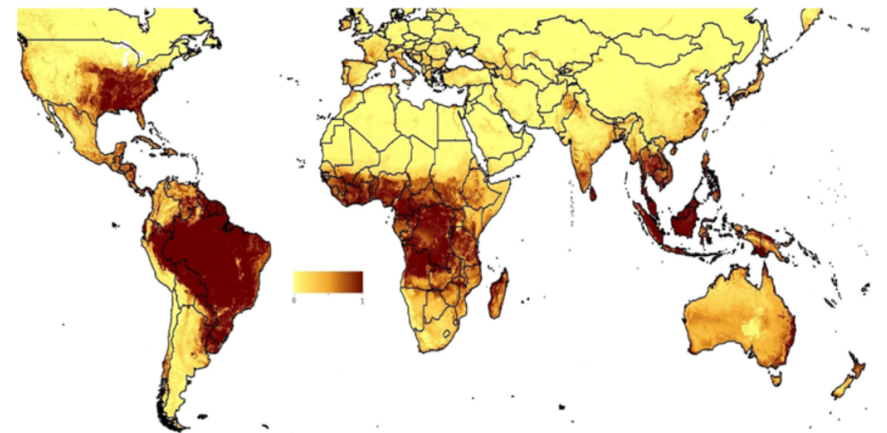
a)



b)



a)



b)

Air network data (a) airports, (b) flight routes for 2011

Information on risk of chikungunya infection (a) and suitability for *Aedes alopictus* vector (b)

Big data have the potential to:

- Augment existing surveillance systems
- Provide an early warning of a disease outbreak
- Provide a basis for research on health and related matters

Big data have the potential to:

- **Augment existing surveillance systems**
- Provide an early warning of a disease outbreak
- Provide a basis for research on health and related matters



# (Mobile) Cellular phone data provide up-to-date information on locations on individuals and their movements

Cellular network comprises a network of Base Transceiver Stations (BTS)

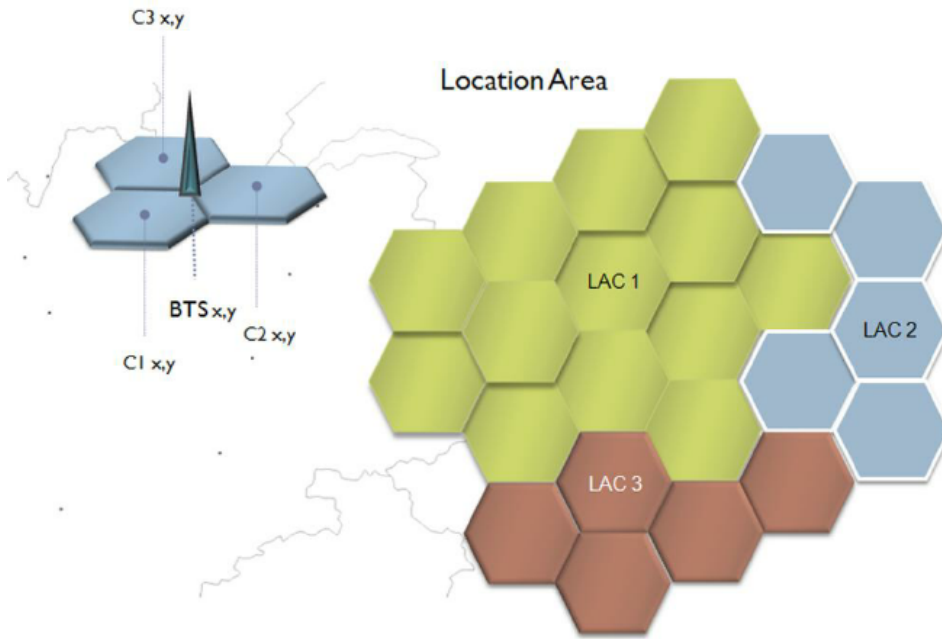
Each cell in the network has a Cell Identifier (ID)

During movement between cells, the network commands the mobile unit to switch to the next cell

The cellular handover records all of the cells through which the mobile unit passes

Cellular network is divided into larger, geo-administrative zones (Location Areas, LA)

Cellular network must know the position of all mobile units at all times to facilitate high connection speeds

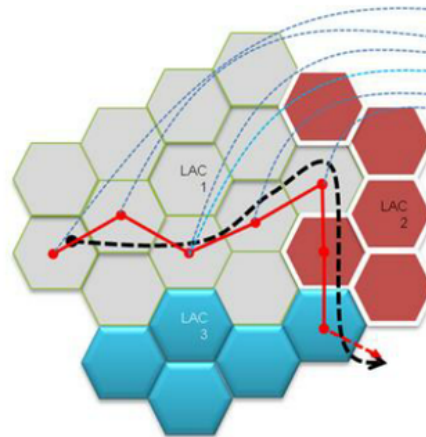


# Mobile phone data collected either actively or passively

**Active data collection** involves provoking production of localisation information for a specific mobile device (e.g. through software or a GPS device included in the phone)

Can be complemented with semantic information

Generally requires agreement of mobile unit owner to participate



DATEEVT	TYPEVT	PARAM	Cellid
30/08/2004 13:06:56	NEWCID		40883
30/08/2004 13:12:37	NEWCID		41721
30/08/2004 13:12:54	NEWCID		40883
30/08/2004 13:34:03	INCCAL	0145294812	40883
30/08/2004 13:48:47	NEWCID		41721
30/08/2004 13:49:01	NEWCID		40883
30/08/2004 17:12:59	RECSMS	888	23123
30/08/2004 17:13:13	LECSMS		14040
30/08/2004 17:13:22	CRESMS		59885
30/08/2004 17:14:13	NEWCID		41702
30/08/2004 17:14:46	ENVSMS	+33689094877	23754
30/08/2004 17:45:22	INCCAL	0608414297	59885
30/08/2004 17:47:01	NEWCID		44960
30/08/2004 17:51:33	NEWCID		23753
30/08/2004 17:51:47	OUTCAL	0608414297	44704
30/08/2004 21:37:25	URLWAP	http://www.orange.fr	3546

Cell ID decoding

Cellid	X	Y	Ville	Adresse
40883	600661	2425302	PARIS	55, rue Vergniaud
41721	600661	2425302	PARIS	55, rue Vergniaud
40883	600661	2425302	PARIS	55, rue Vergniaud
40883	600661	2425302	PARIS	55, rue Vergniaud
41721	600661	2425302	PARIS	55, rue Vergniaud
40883	600661	2425302	PARIS	55, rue Vergniaud
23123	598395	2429344	PARIS	aerg des Invalides
14040	597686	2430483	Paris	16, rue de Berri
59885	598845	2430203	PARIS	11, rue d'Anjou
41702	598243	2430246	PARIS	17, av Malignon
23754	598602	2430447	PARIS	11, rue des Saussaies
59885	598845	2430203	PARIS	11, rue d'Anjou

Stops calculation

deb	fin	duree (sec)	cellid	X	Y	Ville	NomDuSite	Adresse
30/08/2004 09:09:40	30/08/2004 16:27:29	26269	41721	600680	2425340	PARIS	BARRAULT	55, rue Vergniaud
30/08/2004 18:46:07	30/08/2004 18:46:55	48	44448	598850	2430681	PARIS	HAUSSMANN	89, bd Haussman
30/08/2004 18:46:55	30/08/2004 18:48:46	111	44489	598880	2431008	PARIS	ST AUGUSTIN	38/42, rue du Rocher
30/08/2004 18:48:46	30/08/2004 18:51:34	168	49603	598829	2431764	PARIS	ROME	7, rue Marlotte
30/08/2004 18:51:34	30/08/2004 18:52:55	81	3546	599158	2432186	PARIS	LA FOURCHE	70b, av de clichy
30/08/2004 18:52:55	30/08/2004 18:53:08	13	15988	599158	2431808	PARIS	LA FOURCHE	27, rue Lécluse
30/08/2004 18:53:08	30/08/2004 23:44:46	21098	3546	599158	2432186	PARIS	LA FOURCHE	70b, av de Clichy

Data collection & decoding using cell tracing on a mobile phone (dotted line real trajectory, solid line = cell change based trajectory)

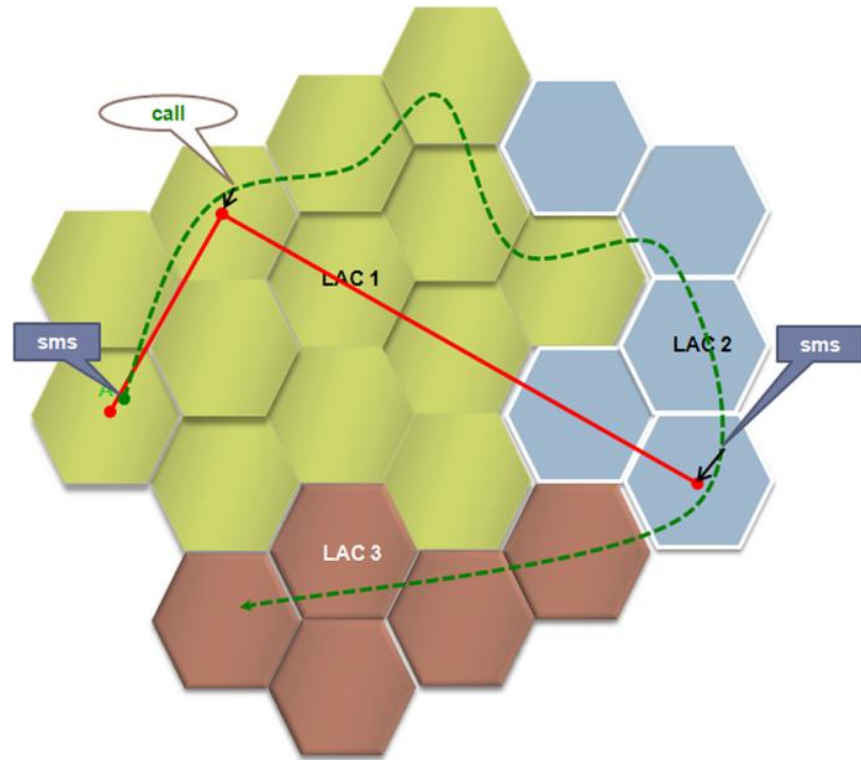
**Passive data collection** utilises billing data, or Call Detail Records

Mobile phone operators collect a large mass of data for billing purposes/system management

Mass of data available

Main draw-back is the lack of semantic information and the need for validation

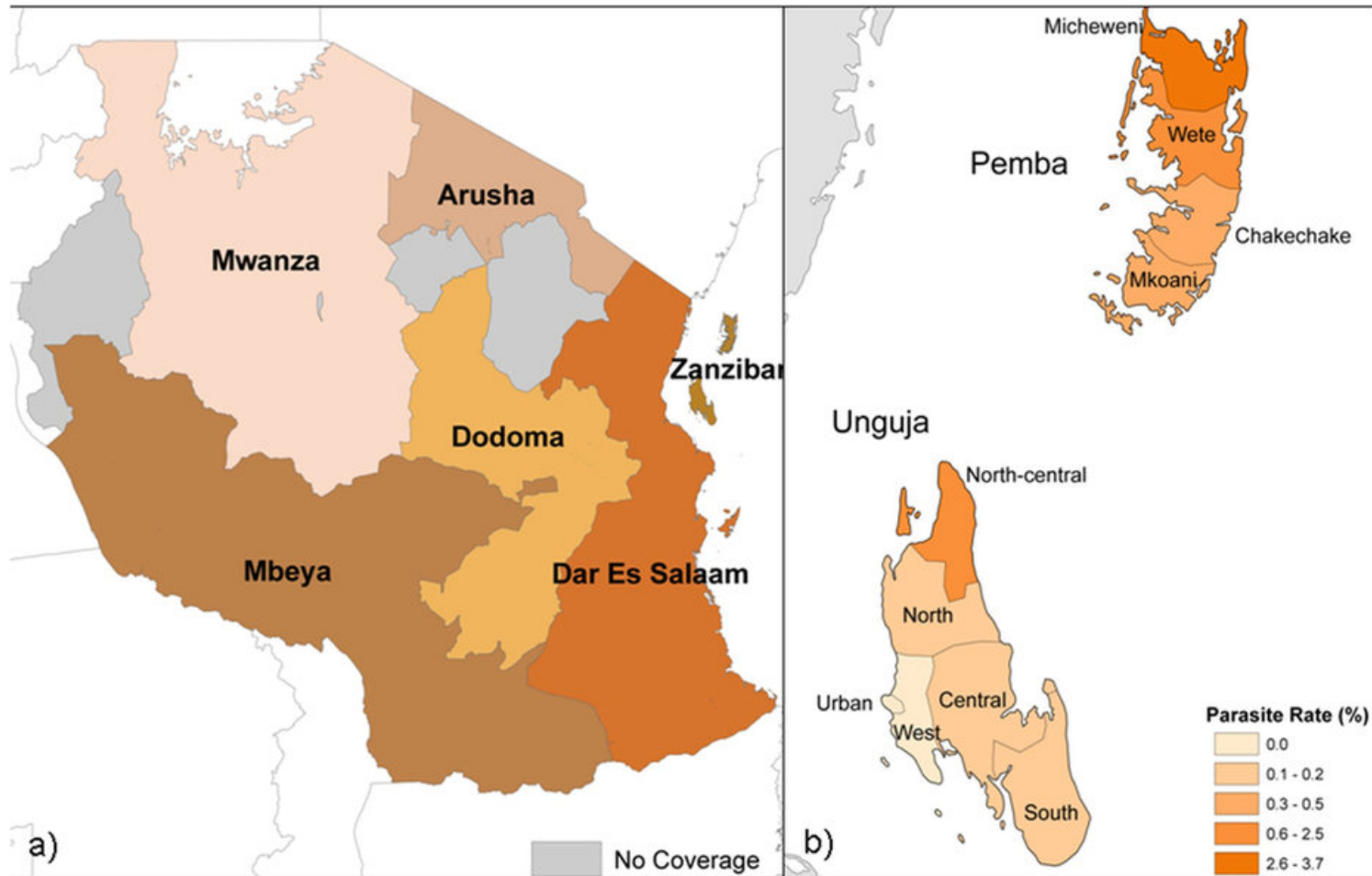
Also questions re privacy



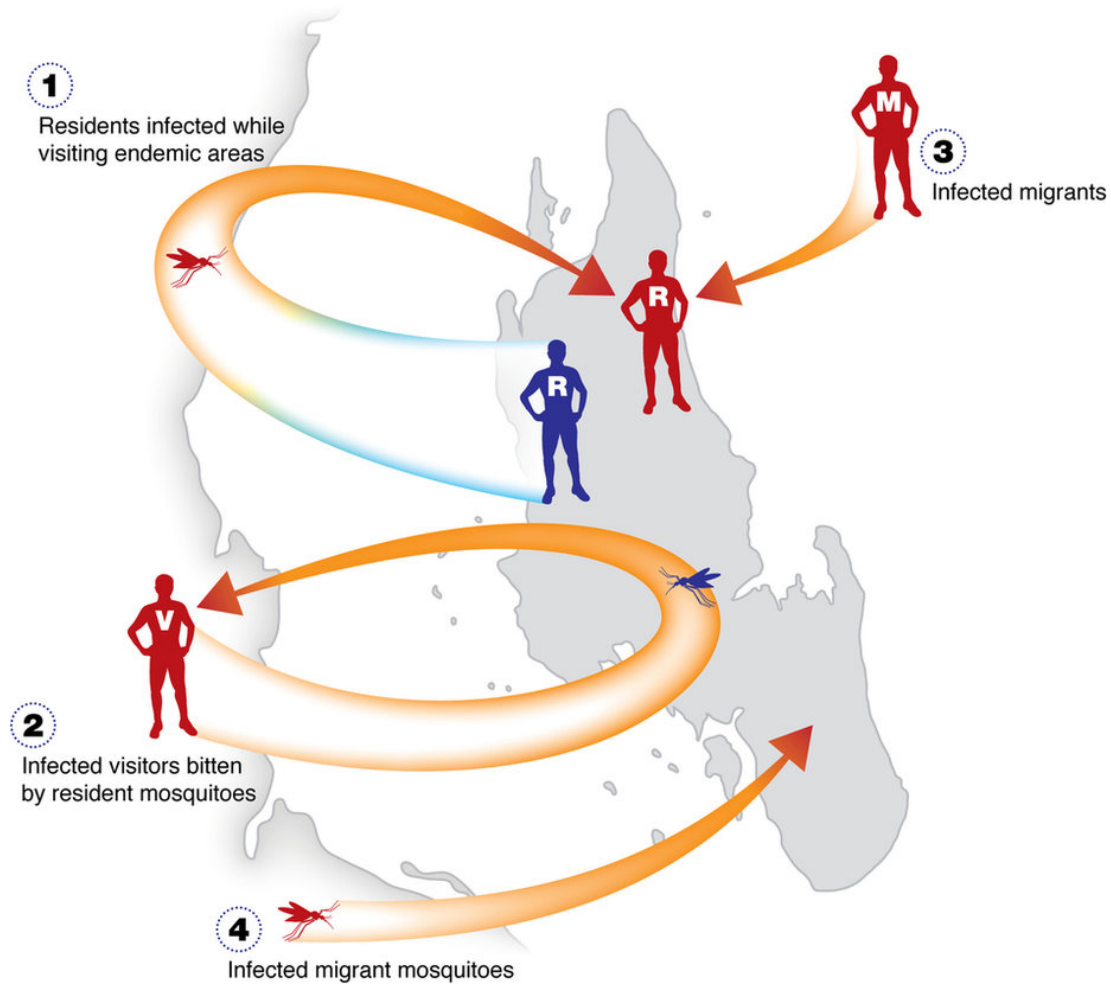
Location information using Call Detail Records  
(dotted line = real trajectory, solid line = route based on CDR (2 SMSs and one call))

## e.g. malaria in Zanzibar

Malaria infections in Zanzibar largely result from malaria imported from parts of mainland Tanzania where malaria is endemic and subsequent transmission.



From Le Menach et al. (2011) Travel risk, malaria importation and malaria transmission in Zanzibar  
*Scientific Reports* 1:93



Mobile phone usage data used to refine models of importation of *Plasmodium falciparum* malaria to Unguja – help quantify risk from infectious visitors and returnees.

Most malaria on Unguja is from returning islanders, rather than visitors to the island.

Improved malaria control measures based as a result contribute towards elimination of malaria

Big data have the potential to:

- Augment existing surveillance systems
- **Provide an early warning of a disease outbreak**
- Provide a basis for research on health and related matters

# Early warning of disease outbreaks through tracking consumer behaviour

Simple idea = human behaviour as reflected in our data shadows can be used as an early warning of disease outbreak and to track the spread of illness

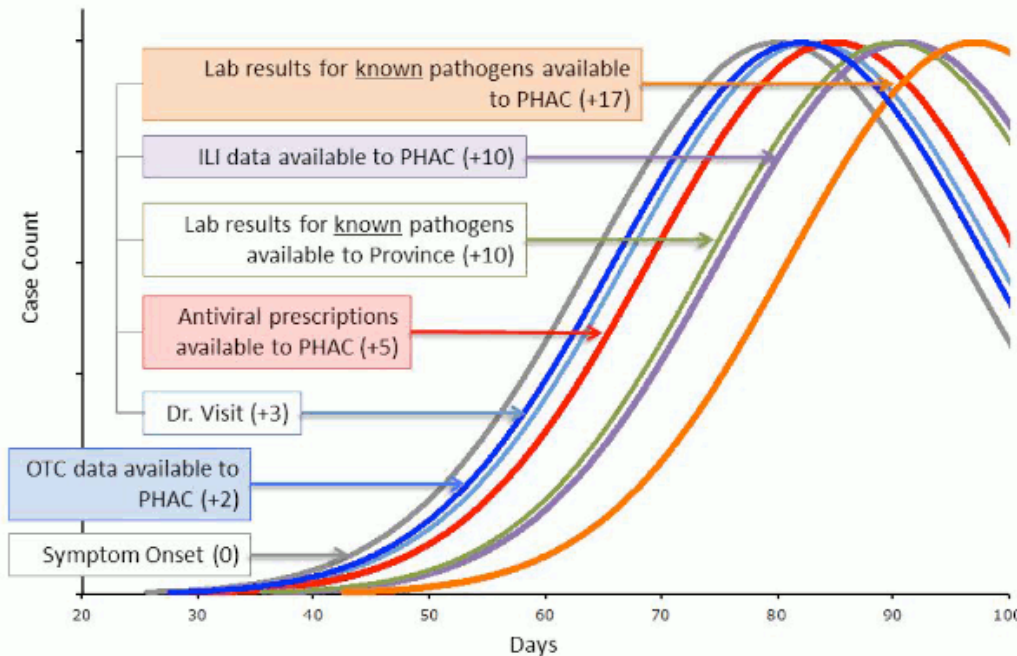


Figure depicts the timelines of pharmacy, clinical and laboratory data relative to the estimated onset of illness and the availability of the information for the purpose of respiratory surveillance.

From Muchaal et al. (2015)

*Evaluation of a national pharmacy-based syndromic surveillance system* CCDR online 41 (9)

Data on the purchasing of drugs from pharmacies available several days/weeks before hospital data on positive disease tests available

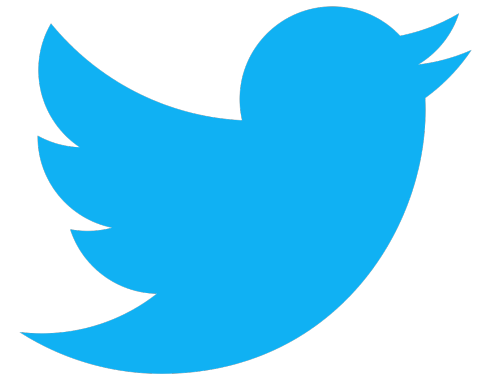
Simple idea – but fraught with difficulties (e.g. agreement over definition of terms such as fever, malaria, dengue, uncertainties in self-diagnosing etc)

Big data have the potential to:

- Augment existing surveillance systems
- Provide an early warning of a disease outbreak
- **Provide a basis for research on health and related matters**



## *Infodemiology* in an age of Twitter



Social media platforms (SMPs) mean that the public is no longer a passive recipient of information

Public now play a larger role in knowledge translation, including information generation, filtering and amplification

Public health authorities can use information from SMPs to monitor public perceptions of and responses to health risks, and the effectiveness/penetration of health campaigns

E.g, public health researchers use Twitter – established in 2006 – to interact with the public and to mine the platform for data

## Taxonomy of use of twitter generated data in health articles, 2010-2015\*

Taxonomy	Description	Articles, No. (%)	Examples
<b>Use of Twitter data</b>			
Content analysis	Assessment of body of tweets for themes in relation to a specific subject	77 (56)	Smoking, diabetes, obesity, concussion
Sentiment analysis	Assessment of body of tweets for positive or negative discussion of a specific subject	21 (15)	Schizophrenia, vaccination, trans health
Image analysis	Assessment of images within body of tweets for themes in relation to a specific subject	1 (1)	#thinspo
Surveillance	Monitoring of Twitter traffic for mentions of a particular topic above the normal background level of discussion	36 (26)	Influenza, Ebola, adverse drug reactions
Prediction	Using Twitter to estimate prevalence of disease or behavior	7 (5)	Heart disease mortality, influenza infection, Affordable Care Act enrollment, asthma emergency department visits
Engagement	Assessing impact of discussion on Twitter by analyzing presence of an account, number of retweets, favorites, followers, etc.	19 (14)	Nutrition public health marketing campaign, social media impact of local health departments, social media adoption by pharmaceutical companies
Network analysis	Assessing the relationship and interactions between Twitter users about a certain topic	5 (4)	Communities of cancer patients, sharing of health information by health organizations

\*See Sinneberg et al. (2017) Twitter as a tool for health research: A systematic review. *AJPH* 107, e1-8.



<https://birdiq.net/twitter-search>

Application Programme Interface (API)s such as BirdIQ provide a means of downloading information on tweets that can then be analysed (location, content etc).

BirdIQ allows information to be exported as an EXCEL file


NCapture exports social media information to Nvivo

<http://www.qsrinternational.com/support/faqs/what-is-ncapture>

Other online tools – such as the One Million Tweet Map – maps the tweets for a particular hashtag or keyword



# The one million tweet map

powered by 

Tweets since page load  12 points

9 9 9 9 9 9 7

Legend

-  tweet cluster
-  latest tweet

Filters

cluster view  heatmap view

Q dengue

# hashtags filter

5 most popular hashtags

dengue (1) [essalud](#) (1) [muertos](#) (1) [plura](#) (1)

last 4 hours

If you want to start with a blank map, click on reset map button (reload the page to see all tweets again).

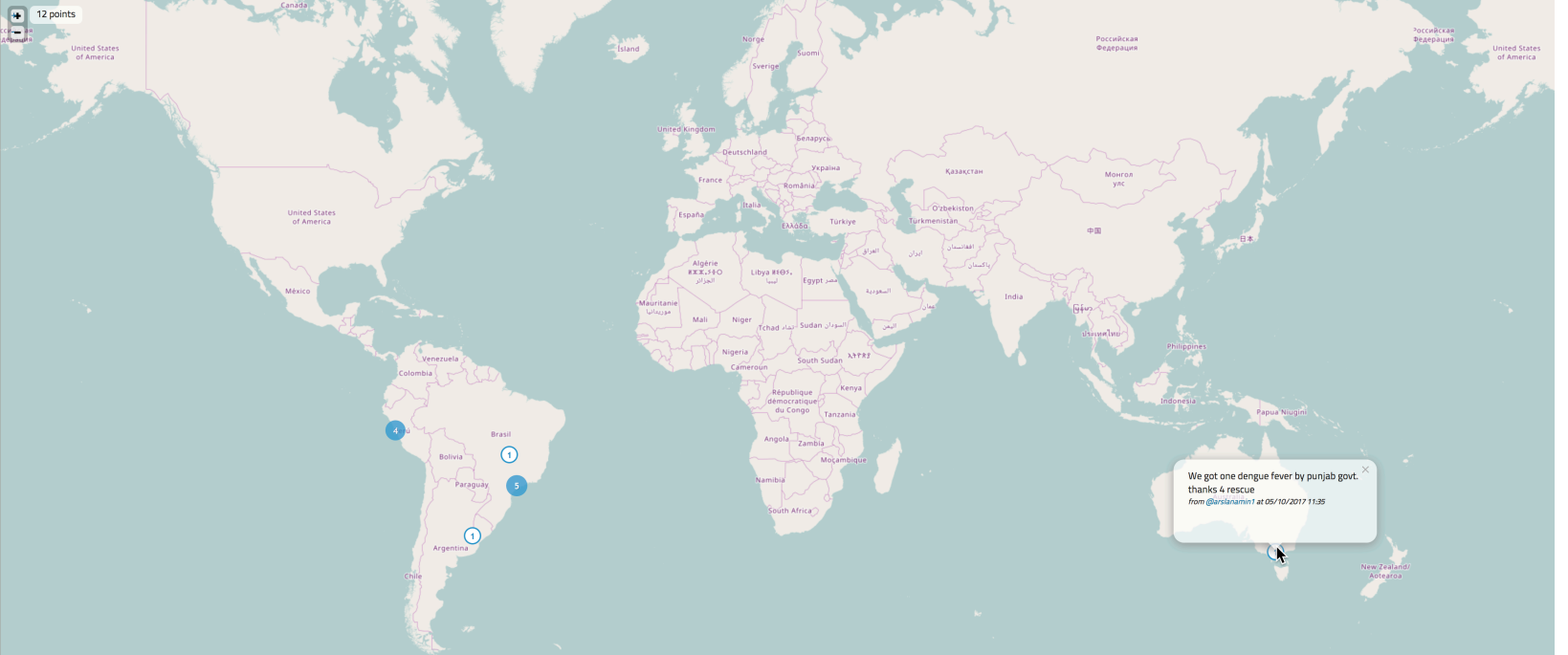
[Reset map](#)

Maptimize

A geographical clustering engine for online maps to display and analyse big geolocalized data.

[Get more info](#)

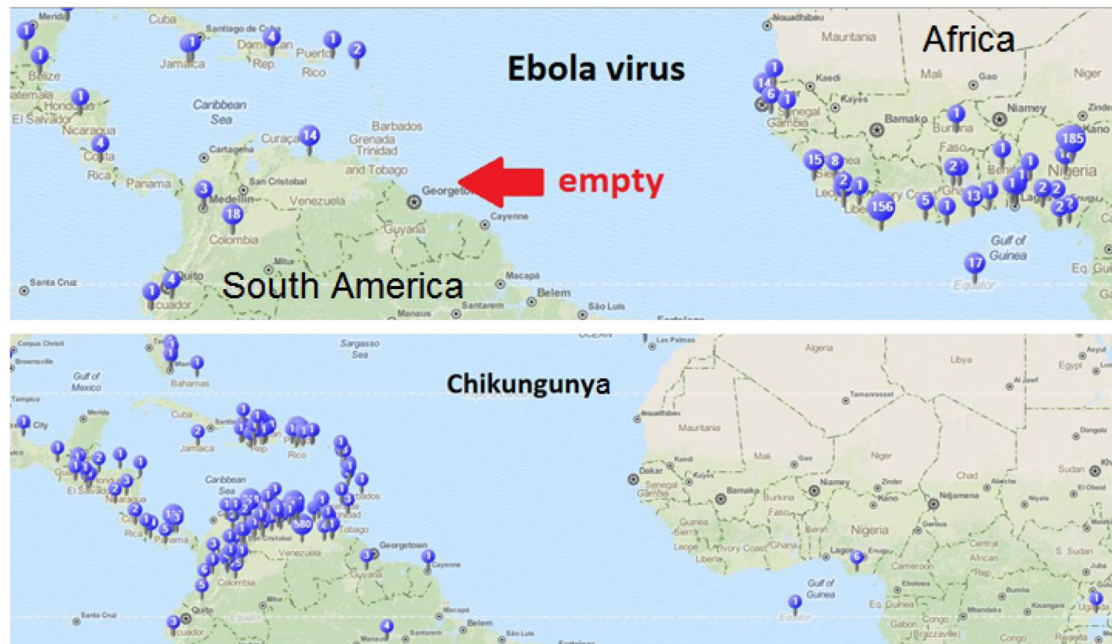
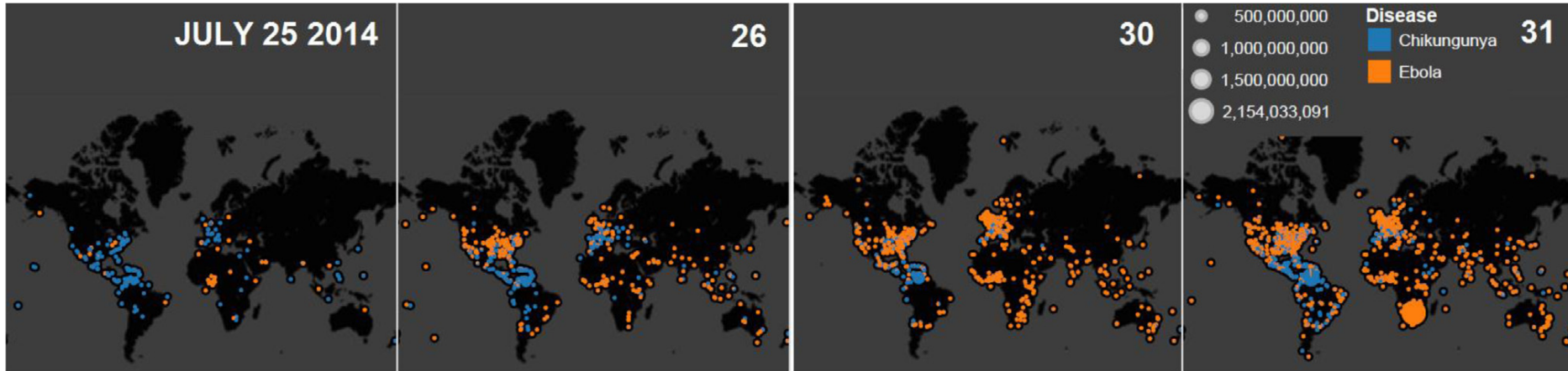
[Tweet](#)



We got one dengue fever by punjab govt. thanks 4 rescue from @arslanamin1 at 05/10/2017 11:35

<http://onemilliontweetmap.com/>

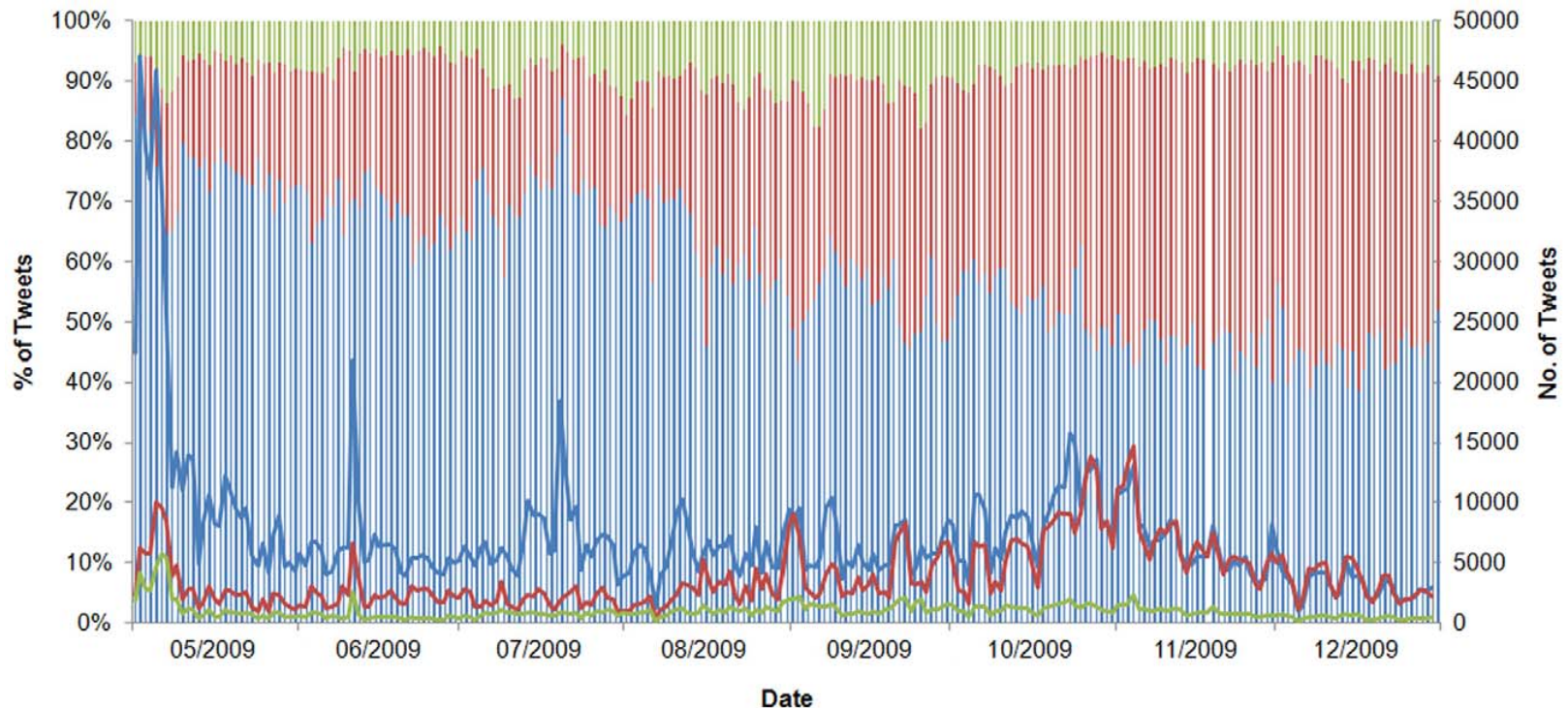
# e.g. Ebola Virus Disease outbreak, 2014



Daily geographic spread of tweets mentioning Ebola Virus Disease (EVD).

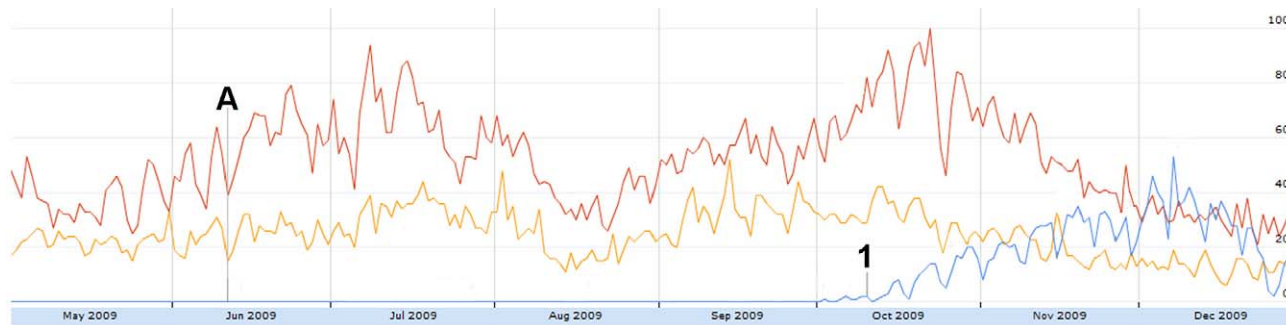
Following July 29 2014 health advisory announcement by CNC, US – awareness of EVD spread, but possibly at expense of attention to Chikungunya outbreak in Caribbean everywhere but in Caribbean

## e.g. H1N1 pandemic, 2009

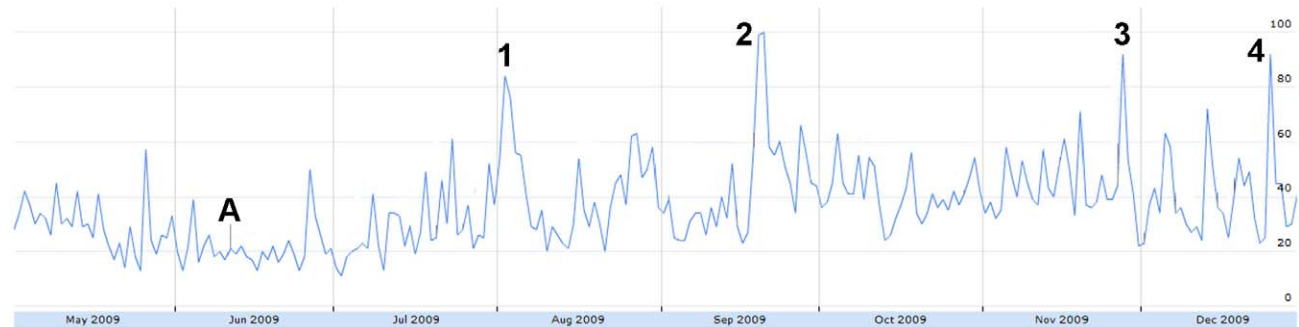


Tweets containing H1N1, swine flu, or both from May to December 2009. Lines = absolute number. Bars = relative percentage. Blue = “swine flu” or swineflu. Red = H1N1. Green = (“swine flu” or swineflu) AND H1N1

From: Chew & Eysenbach (2010) Pandemics in the age of twitter: content analysis of tweets during the 2009 H1N1 outbreak *PLoS ONE* 5, e14118



Relative proportion of tweets sharing personal experiences, 1 May-31 December, 2009. Figure was scaled to the highest peak on Oct. 20. Red = indirect (family/friend) experience. Yellow = personal/direct experience. Blue = vaccination experience. **A** = June 11: WHO pandemic level 6 announcement. **1** = Oct 6: H1N1 vaccinations arrive in the US.



Relative proportion of tweets expressing misinformation, 1 May-31 December, 2009. Figure was scaled to the highest peak on Sept. 20. **A** = June 11: WHO pandemic level 6 announcement. **1** = Aug 2: CBS reports on parental concerns about H1N1. **2** = Sept 18–21: Ten swine flu lies told by the mainstream media. **3** = Nov 27: WHO and drug companies in collusion. **4** = Dec 25: Carbon monoxide poisoning can create same symptoms as H1N1

## Quick summary

(Re)-emergence & spread of (existing and) new infectious diseases occur at a high rate

*facilitated by human activity, including climate change impacts*

Rate means that conventional surveillance methods – where they exist – are not fit for purpose

Big data offer opportunities to health researchers/practitioners

But not straightforward, as big data

*are not the same as total data*

*may not be representative*

*require new forms of analysis – current statistical techniques not designed to handle the variety, volume and velocity of big data?*