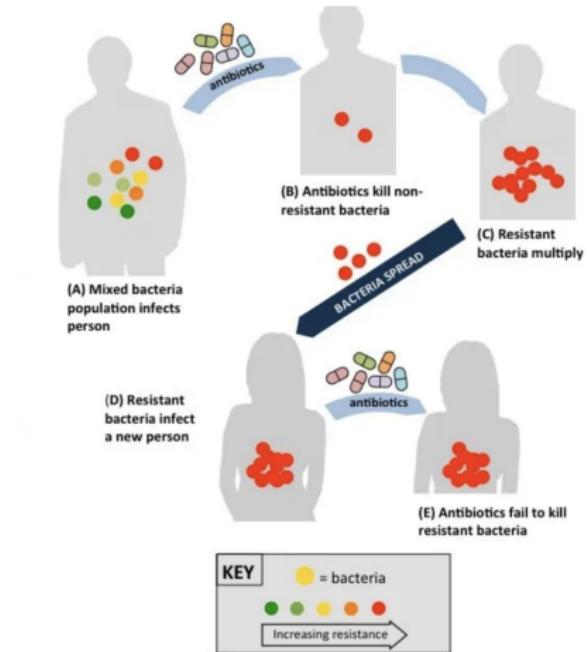
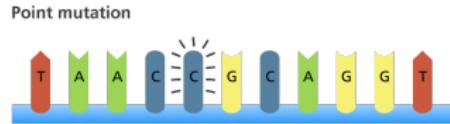
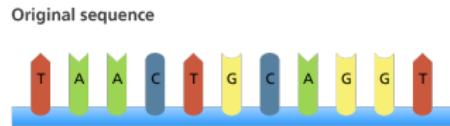


# RECONSTRUCTING TRANSMISSION WITH GENOMIC DATA: INTRODUCTION

Caroline Colijn  
Jessica Stockdale

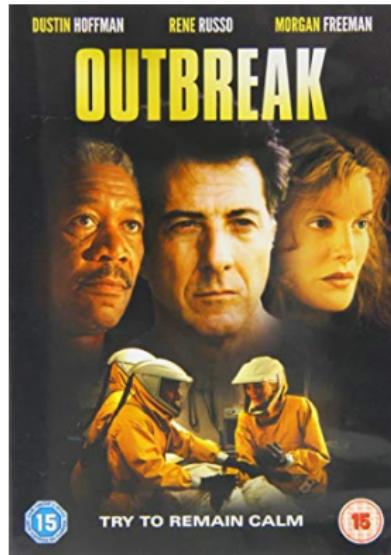
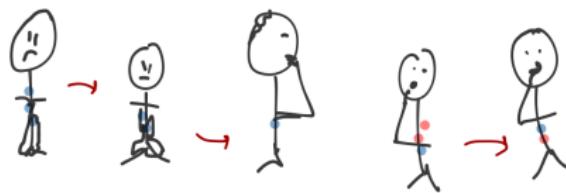
# INFECTIONS ARE EVOLVING

Take a look at the DNA of your favourite bacteria:

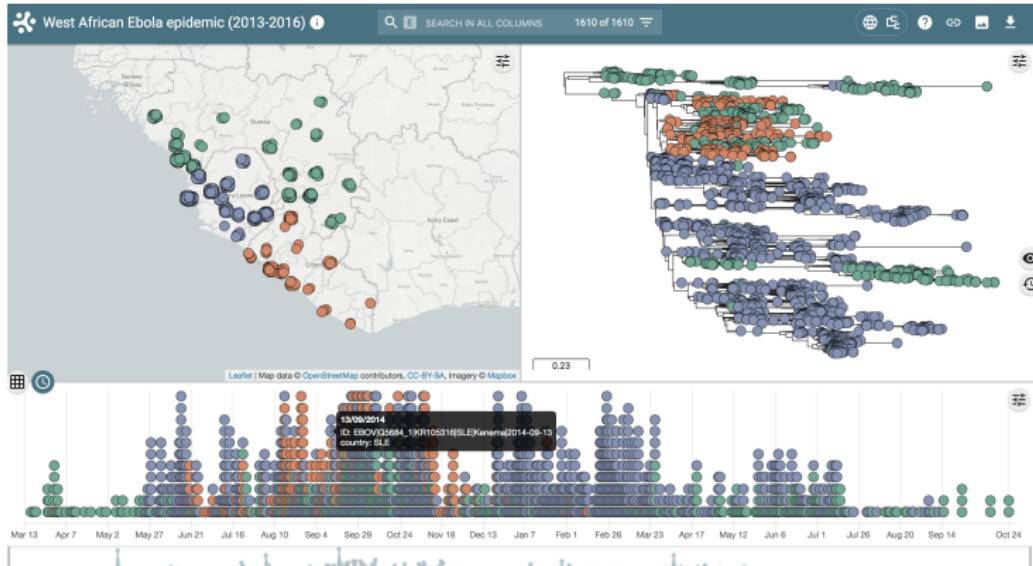


Vivian Chou, Harvard

# IT IS HARD TO CONTROL OUTBREAKS



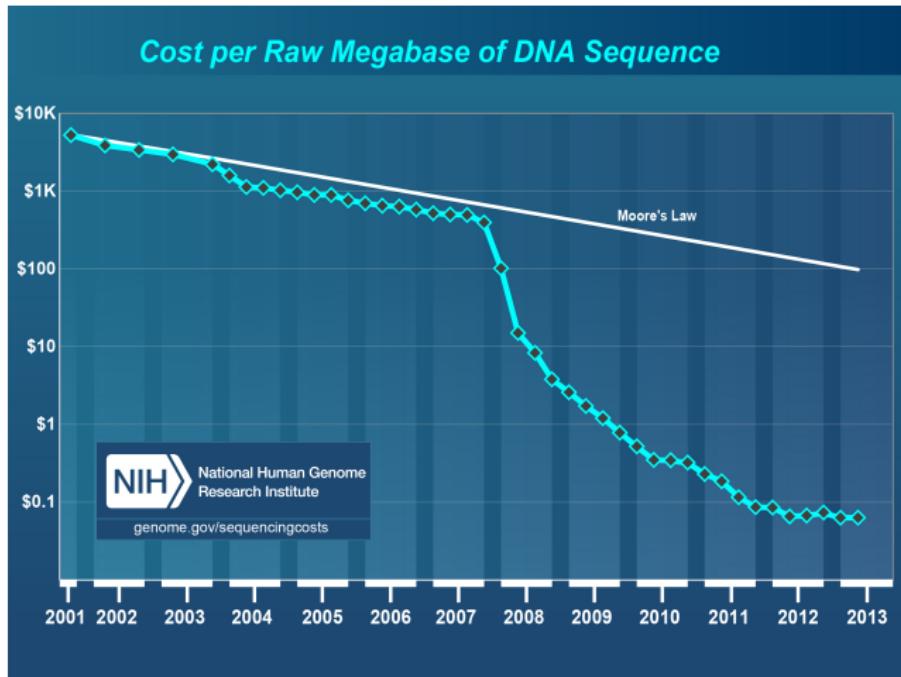
# NEW DATA - A BIG OPPORTUNITY



microreact.org, by D. Aanensen and colleagues

<https://microreact.org/project/west-african-ebola-epidemic?tt=rc>

# SEQUENCING IS LESS EXPENSIVE NOW THAN EVER

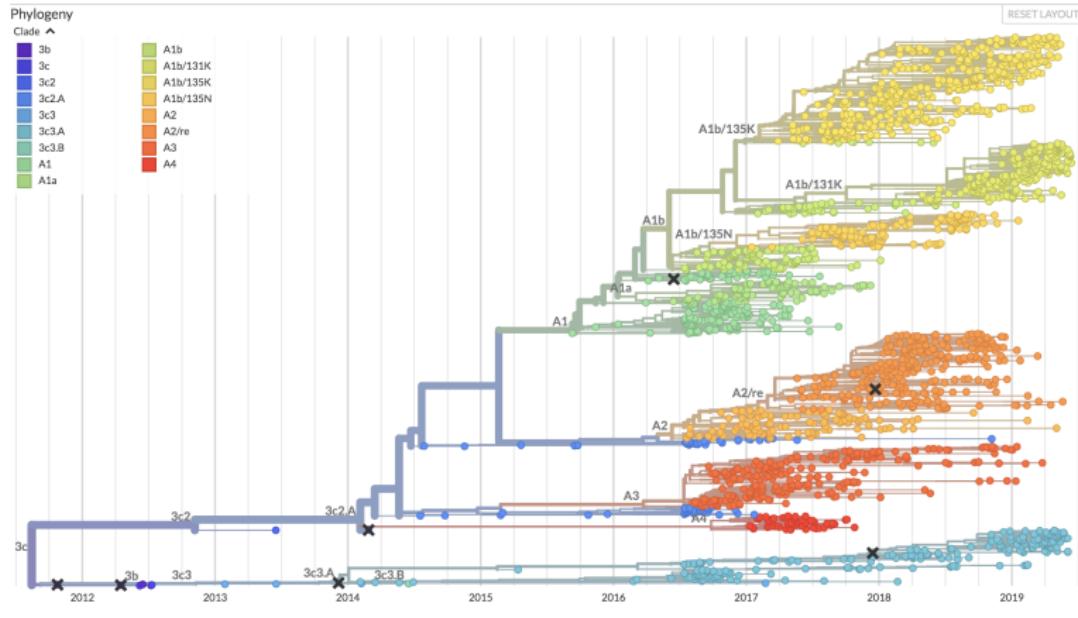


We can read the DNA of viruses and bacteria that cause infection, and see how they are changing as they spread.

# BUGS ACQUIRE SMALL GENETIC VARIATION AS THEY SPREAD

## Real-time tracking of influenza A/H3N2 evolution

Showing 2169 of 2169 genomes sampled between Oct 2011 and Jun 2019 and comprising 17 clade memberships, 10 regions, 118 countries and 43 authors.



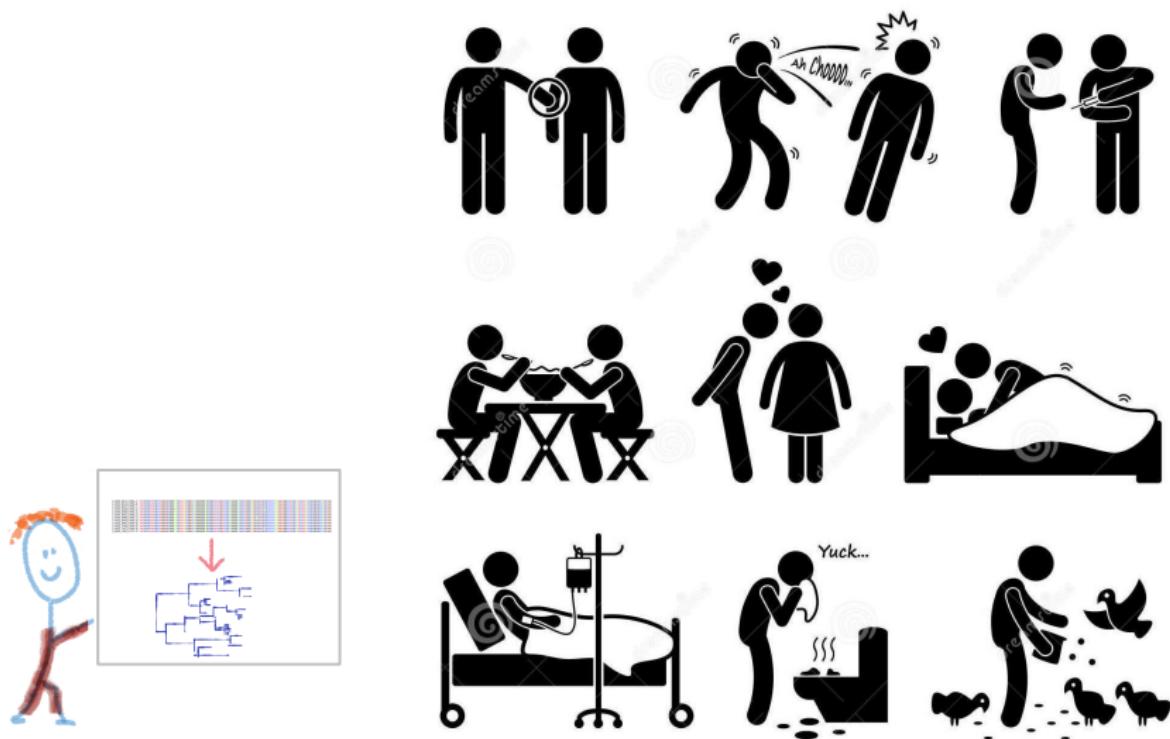
Visualization of flu evolution. Image: [nextflu.org](http://nextflu.org)

# BUGS EVEN VARY FROM PERSON TO PERSON

REMEMBER THE OLD GAME OF “TELEPHONE”?



# CAN WE USE SEQUENCES TO UNDERSTAND TRANSMISSION?



Can we learn about who infected whom and when from sequence data?

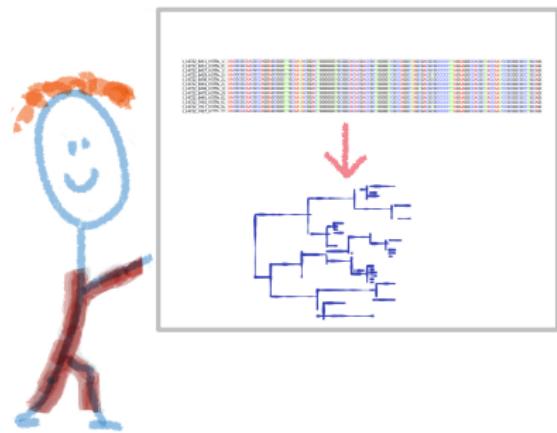
# WHY DO WE WANT TO KNOW WHO INFECTED WHOM?

- We don't really need to know whether Alice infected Bob, Bob infected Eve, or Eve infected Jackson
- But we do want to know when, where and how transmission takes place
- Drilling into the details of transmission can help understand this
- And it can help to know when infection did not occur

# DATA ON GENETIC VARIATION CAN HELP!

Data about this variation:

- help with who infected whom - improve outbreak control
- bigger scale: help with choosing best vaccines, best antibiotics

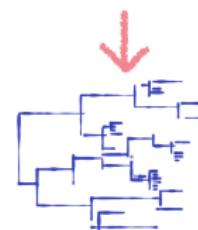
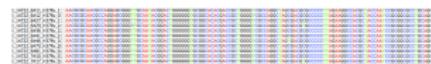


# CHALLENGES IN GENOMIC EPIDEMIOLOGY AND BEYOND

## Outbreak questions:

- How fast do we have to find cases?
- How do we find missing cases?
- What are early signs of a big outbreak?

← GAP →



## RELATEDNESS IS KEY

The answers to our questions aren't just in the data, but in the connections among data points.

## The sequence

AACCATAGGT

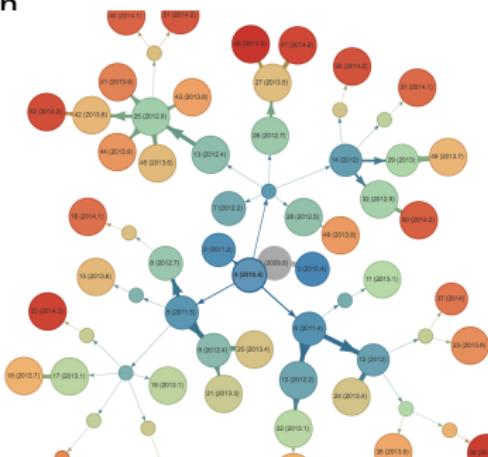
doesn't mean much for transmission on its own

But with two:

AACCATAGGT

GACCATAGGT

we know we have two very similar things.

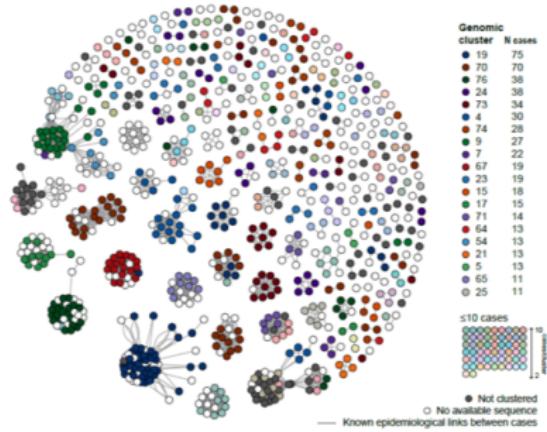


# THREE ROLES FOR SEQUENCES IN PUBLIC HEALTH

- ① Infer global routes of movement
  - ▶ nextstrain.org
  - ▶ Challenges from different sampling and sequencing in different places
- ② Infer population dynamics back through time
  - ▶ Field of phylodynamics - see Julia Palacios' module!
  - ▶ Limited in terms of direct public health action; more ecological
- ③ Analysis of transmission in localised outbreaks
  - ▶ Potentially directly useful for public health in the short term

# TRANSMISSION IN LOCALISED OUTBREAKS

- The first step is to broadly find out what you have.
- Clustering: put similar sequences into groups



Tracking the COVID-19 pandemic in Australia using genomics

● Torsten Seemann, Courtney Lane, Norelle Sherry, Sebastian Duchene, Anders Goncalves da Silva, Leon Caly, Michelle Sait, Susan A Ballard, Kristy Horan, Mark B Schultz, Tuyet Hoang, Marion Easton, Sally Dougall, Tim Stinear, Julian Druce, mike Catton, Brett Sutton, Annaliese van Diemen, Charles Alpren, Deborah Williamson, Benjamin P Howden

doi: <https://doi.org/10.1101/2020.05.12.20099929>

# HOW TO PUT SEQUENCES INTO GROUPS: CLUSTERING

A cluster: more similar to each other than to objects outside of the group'

- SNP: single nucleotide polymorphism. A change from (eg) A to C at a single site.
- Simplest clustering method: Sequences from Bob and Alice are placed in the same cluster if they differ by  $k$  SNPs or fewer
- More sophisticated clustering methods:
  - ▶ Account for time, rate variation, selection: eg our own 'transcluster' (Stimson et al MBE 2019 Beyond the SNP threshold")
  - ▶ Phylogenetic methods: cluster picker, cluster picker II, cluster matcher (developed for HIV primarily)
  - ▶ Integrated methods – combine diagnosis, genotyping, and genetic similarity (Poon et al Lancet HIV)

# SEQUENCE CLUSTERS IN PUBLIC HEALTH 1

Exclude suspected transmission events:

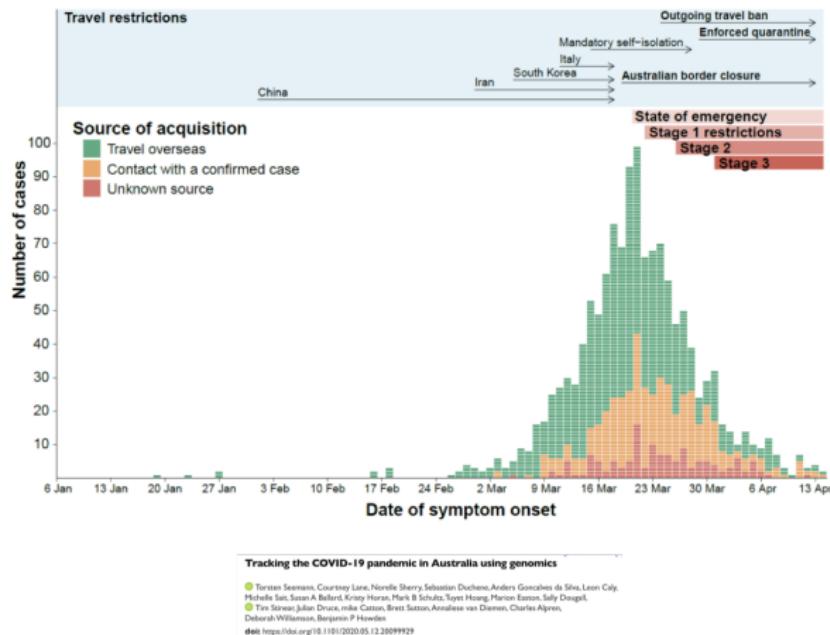
- ① For example: 2 cases that are epidemiologically linked
- ② If their viral (or bacterial) sequences are very different, there was likely no direct transmission
- ③ The apparent link was spurious ("false positive" epi link)

## SEQUENCE CLUSTERS IN PUBLIC HEALTH 2

Identify sources of infection that did not have epidemiological links

- ① Sequences firmly place Bob in a cluster with Alice and Eve (for example, sequences are 1 SNP away)
- ② Bob has no epidemiological links to Alice, Eve or their contacts
- ③ This can help identify previously unknown exposures: true links were missing ("false negative" epi link)

# EXAMPLE: COVID-19 GENOMIC EPIDEMIOLOGY IN AUSTRALIA



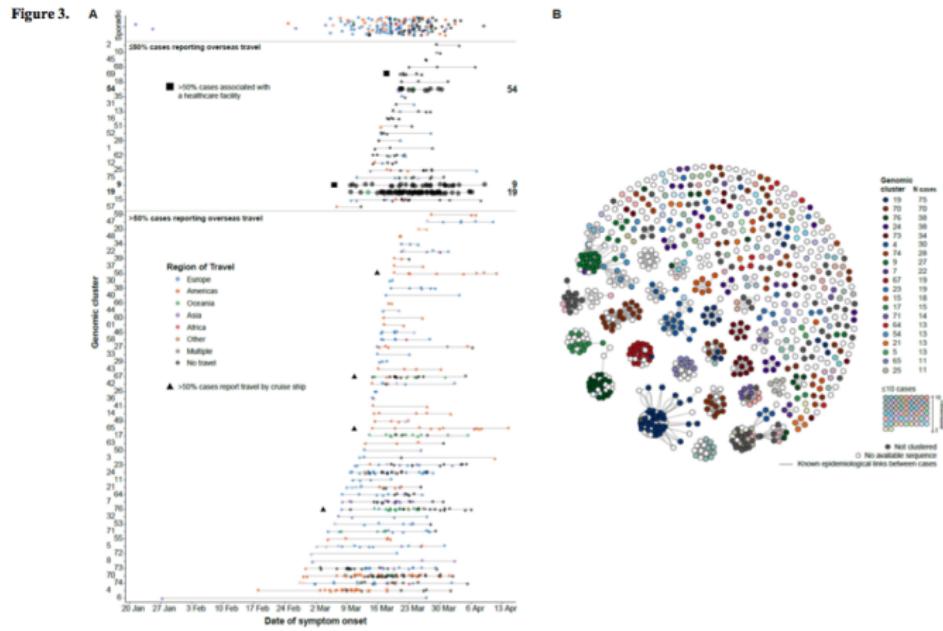
# SARS-CoV-2 DATA FROM AUSTRALIA

## Brief data description

- 1388 lab-confirmed cases in Victoria
- 62% travellers
- 27% known contacts
- 10% unknown source of exposure
- 1242 sequenced
- 1085 passed quality control
- Maximum 15 SNPs compared to Wuhan 1

# CLUSTERS IN AUSTRALIAN (VICTORIA) SARS-CoV-2 DATA

- The authors used ClusterPicker to divide the genomes into clusters
- 737 of 1085 were in any cluster
- 76 clusters: median size 5, median duration 13 days This suggests good control (and could provide a serial interval estimate too!)
- 34 clusters were entirely overseas travellers
- 34 were mixed, typically the first case was a traveller
- 81 sequences with unknown exposure (from the epi point of view) land in 24 clusters
- This suggests what the exposure was



## SPECIAL CLUSTERS

- Epidemiological (epi) clusters: groups of cases thought to be linked together on the basis of epidemiological (not genome) data, eg where people live, timing, health care, suspected exposure
- Genomic clusters: similar genomes grouped together on the basis of (sort of) genetic distance
- Four distinct epidemiological clusters one genomic cluster -find links they didnt know about
- One big epi cluster separated into 4 distinct genomic clusters: exclude links they thought they knew about

**Data required for this nice work:** SARS-CoV-2 sequences, suspected exposure times and sources from epidemiology.

## EARLY WORK: TRANSMISSION WITH SEQUENCES

- Cottam et al: Integrating genetic and epidemiological data to determine transmission pathways of foot-and-mouth disease virus. Proc. Biol. Sci. 2008.
  - ▶ Phylogenetic trees constrain transmission events
- Jombart T, Eggo RM, Dodd PJ, Balloux F (2011) Reconstructing disease outbreaks from genetic data: a graph approach. Heredity 106: 383390
  - ▶ Parsimonious: transmission tree with fewest mutations

## EARLY WORK CONTINUED

- 2012 Ypma RJ, Bataille AM, Stegeman A, Koch G, Wallinga J, et al. Unravelling transmission trees of infectious diseases by combining genetic and epidemiological data. Proc Biol Sci 279:
- 2012 Morelli MJ, Thebaud G, Chadoeuf J, King DP, Haydon DT, et al. A Bayesian inference framework to reconstruct transmission trees using epidemiological and genetic data. PLoS Comput Biol 8
  - ▶ Both use a unified likelihood of genetic and epidemiological data - needs full sampling
- 2014 Outbreaker: Jombart T, Cori A, Didelot Z, Cauchemez A, Fraser C, Ferguson N, Bayesian Reconstruction of Disease Outbreaks by Combining Epidemiologic and Genomic Data, PLOS CB

# OUTBREAKER AND SARS

The data:

- Sequences
- Times of sample collection
- Number of cases

Additional inputs (derived from data):

- Generation time: time from infection to symptom onset
- Time from infection to sample collection

# THE MATH BEHIND OUTBREAKER

Outbreaker is a Bayesian MCMC method, and it uses *augmentation*.

“Augmented” quantities for each case  $i$ :

- The infector for each case, or most recent sampled infector
- The number of unknown intermediates between  $i$  and  $i$ 's ancestor
- The date when  $i$  was infected

The likelihood relies on these augmented data.

# OUTBREAKER'S LIKELIHOOD IN BRIEF

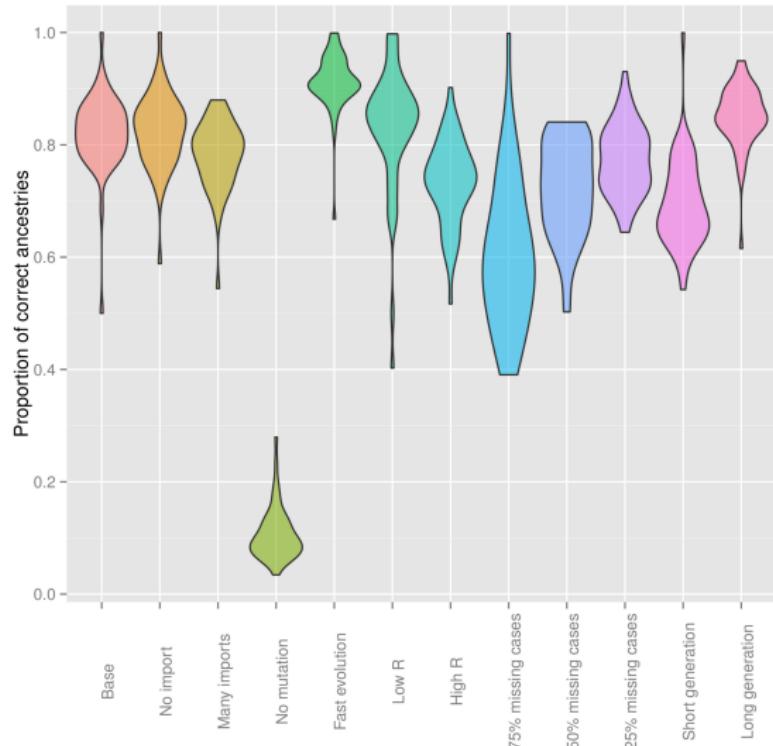
Likelihood for case  $i \propto$

$$\begin{aligned} & P(\text{sequence} \mid \text{ancestor's sequence, intermediates, evolution model}) \times \\ & P(\text{sampled time} \mid \text{infection time}) \times \\ & P(\text{infection time} \mid \text{parent's infection time, intermediates}) \end{aligned}$$

## A 60-second MCMC briefing:

- Propose augmented data: transmission tree, intermediates, times of infection
- Compute the above likelihood
- Accept or reject (in a principled way, depending on the proposal method and the likelihood)
- Run for a long time. Collect “posterior samples” of these quantities

# OUTBREAKER RESULTS: SIMULATED DATA



<https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1003457>

# APPLICATION: SARS (THE FIRST SARS!)

- 13 SARS “genomes” – but note: sequencing was *very* different then
- Ruan et al, *Lancet*, Volume 361, Issue 9371, 24 May 2003

<https://www.sciencedirect.com/science/article/pii/S0140673603134149>

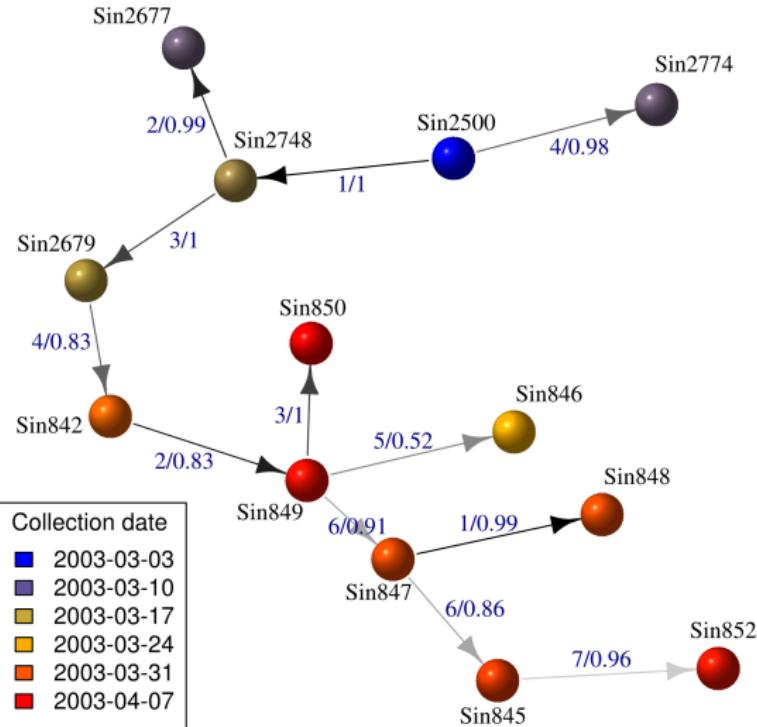
- Liu et al *PLOS Medicine* <https://doi.org/10.1371/journal.pmed.0020043>

- Generation time:  
gamma with mean 8.4 days

- Time to sampling: same
- Lots more variation  
than TB or SARS-CoV-2

Index case	Singapore cases			Overseas cases						Variations frequency	ORF/protein	AA change (SNP2000→ variation)
	Primary SN2500	Secondary SN2774	Secondary SN2670	Canada URBAN	Hong Kong OUMW1	S China B01	N China B03	N China B04	T			
29711	29706	29705	29712	29727	29727	29429	29430	24774	T	2	Gtrta (Met>p65)	Silent
8559	C	C	C	T	C	T	T	T	T	2	Gtrta (Asp1)	A→Y
8559	T	T	T	T	T	T	C	T	T	2	Gtrta (Asp1)	Silent
8572	G	G	G	G	G	G	T	G	G	2	Gtrta (Asp1)	Y→G
3440	T	T	T	T	T	T	C	C	C	5	Gtrta (Asp1)	Y→K
3479	T	T	T	T	T	T	C	C	T	2	Gtrta (Asp1)	W→K
8560	G	G	G	G	G	G	C	C	T	2	Gtrta (Asp1)	W→Y
17564	I	E	T	T	T	T	G	G	G	6	Gtrta (Ile>Phe, deleted)	D→E
19064	A	A	A	A	A	G	A	G	A	2	Gtrta (Asp1)	Silent
19084	T	T	T	C	C	C	C	C	C	4	Gtrta (Asp1)	I→T
19038	A	A	A	A	A	A	A	G	G	4	Gtrta (Asp1)	Silent
19038	G	G	G	G	G	G	A	G	G	2	Spine phosphorylation	O→O
22222	I	T	T	T	T	T	C	C	C	5	Spine phosphorylation	I→T
27243	C	C	C	C	C	C	T	T	N	3	Putative protein	P→A
27827	T	T	T	T	T	T	C	C	C	6	Noncoding	Q→P
29279	A	A	A	A	A	A	A	C	A	2	Nucleosidepal	Q→P

# SARS (1) RESULTS FROM OUTBREAKER



## DISCUSSION

- You'll hear more about the advent of sequencing - much more data now!
- Outbreaker is now in the broader package 'outbreaker2'
- there was more in these papers that I did not describe
- Both these methods have limitations:
  - ▶ they do not capture the shared ancestry of the pathogens
  - ▶ for that, we need phylogenetic trees!
  - ▶ they do not accommodate variation within hosts

# WHAT'S NEXT?

- Reconstructing transmission trees!
- Introduction to genomics for genomic epidemiology
- Non-phylogenetic outbreak reconstructions in outbreaker
- Phylogenetic trees: theory and practice
- TransPhylo: genomic epi with trees
- Research frontends: bringing in more data
- Research frontends: SARS-CoV-2 and COVID-19
- Discussion