



Clinical and Epidemiological Virology,  
Rega Institute, Department of Microbiology  
and Immunology  
KU Leuven, Belgium.



## Introduction to molecular epidemiology and infectious disease phylodynamics

Philippe Lemey<sup>1</sup> and Marc Suchard<sup>2</sup>

1.Regia Institute, Department of Microbiology and Immunology, K.U. Leuven, Belgium.

2.Departments of Biomathematics and Human Genetics, David Geffen School of Medicine at UCLA. Department of Biostatistics, UCLA School of Public Health

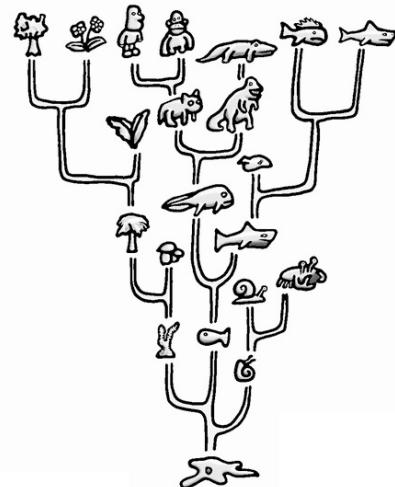
SISMID, July 10-12, 2019

### This course (SISMID module 05)

- Wednesday, July 10
    - Introduction
    - Alignment, substitution models and phylogenetic inference
  - Thursday, July 11
    - Phylogenetic inference practical
    - Bayesian phylogenetics
    - Molecular clocks and model testing
    - BEAST practical
  - Friday, July 12
    - Viral epidemiology and the coalescent
    - BEAST practical
    - Phylogeography
    - BEAST practical
  - Bonus
    - Phylo-Alignment
    - Recombination
    - Robust counting
    - Antigenic cartography
    - Phylo-Factor analysis
- (We are here to cater for your needs!)*

# Molecular evolution and phylogenetics

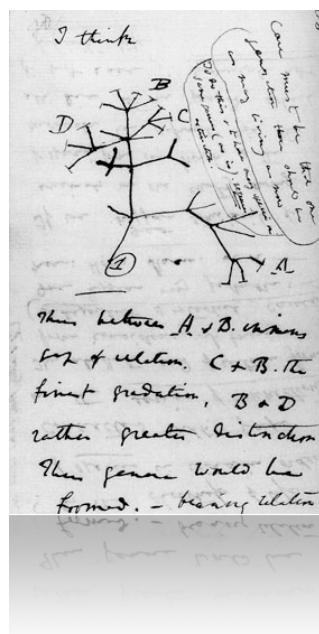
- biological **sequences** (DNA, RNA, protein) contain information about the processes and events that formed them
  - this information is often **scrambled, fragmentary, hidden, or lost** completely
  - our aim is to use **mathematical models** to recover and decipher this information
  - The central concept is a **phylogeny**: a diagram depicting the ancestral relationships among characters or genetic sequences



HIV-1 (UK) ATC---TGCTAAAGCATATGACACAGAGGTACA**TAATGTT**  
HIV-1 (USA) ATC**GGA**TGCTAGAGCTTATGATACAGAGGTACA---TGT

# Phylogenetics

- ## • Darwin, 1837



- Haeckel, 1866



# Fitch & Margoliash 1967

- molecular change in cytochrome c gene

## Construction of Phylogenetic Trees

A method based on mutation distances as estimated from cytochrome *c* sequences is of general applicability.

Walter M. Fitch and Emanuel Margoliash

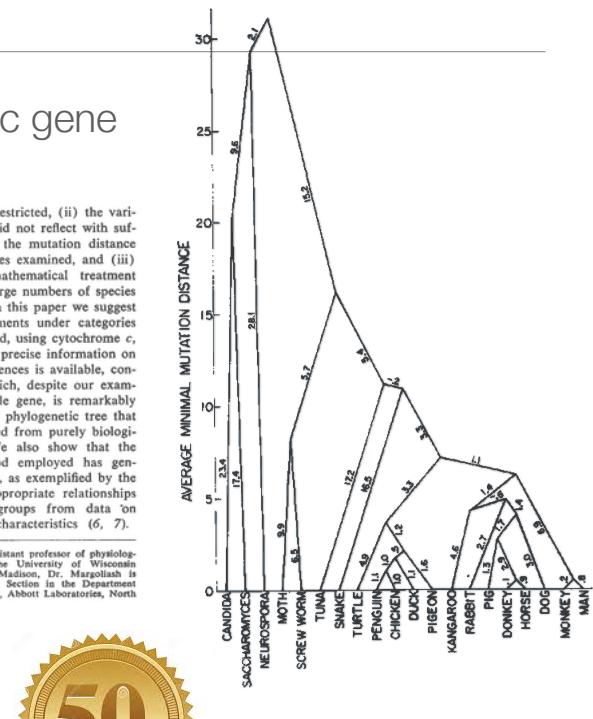
Biochemists have attempted to use quantitative estimates of variance between substances obtained from different species to construct phylogenetic trees. Examples of this approach include studies of the degree of interspecific hybridization of DNA (1), the degree of cross reactivity of antisera to purified proteins (2), the number of differences in the peptides from enzymic digests of purified homologous proteins, both as estimated by paper electrophoresis-chromatography or column chromatography and as estimated from the amino acid compositions of the proteins (3), and the number of amino acid replacements between homologous proteins whose complete primary structures had been determined (4). These methods have not been completely satisfactory because (i) the portion of the genome examined

was often very restricted, (ii) the variable measured did not reflect with sufficient accuracy the mutation distance between the genes examined, and (iii) no adequate mathematical treatment for data from large numbers of species was available. In this paper we suggest several improvements under categories (ii) and (iii) and, using cytochrome *c*, for which much precise information on amino acid sequences is available, construct a tree which, despite our examining but a single gene, is remarkably like the classical phylogenetic tree that has been obtained from purely biological data (5). We also show that the analytical method employed has general applicability, as exemplified by the derivation of appropriate relationships among ethnic groups from data on their physical characteristics (6, 7).

Dr. Fitch is an assistant professor of physiological chemistry at the University of Wisconsin Medical School in Madison. Dr. Margoliash is head of the Protein Section in the Department of Molecular Biology, Abbott Laboratories, North Chicago, Illinois.

20 JANUARY 1967

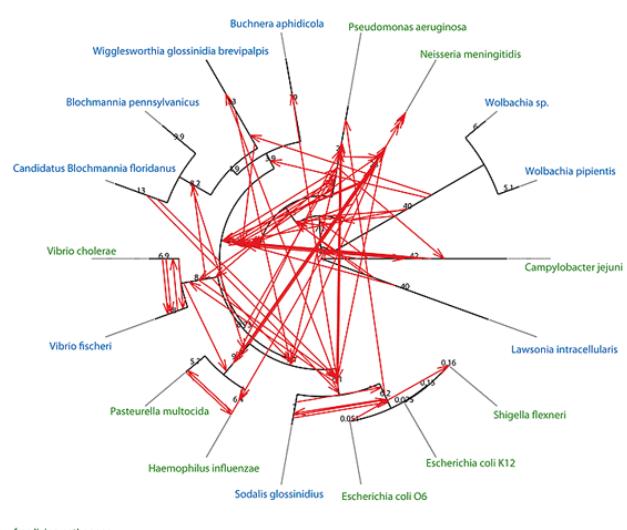
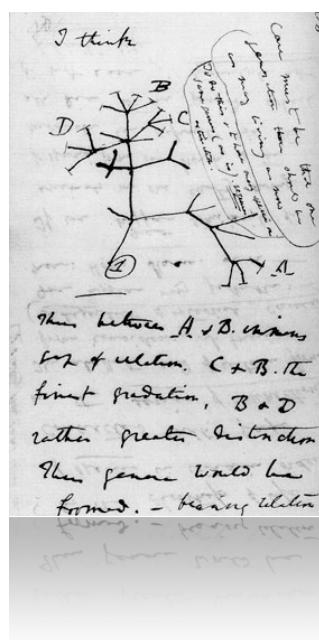
Science. 1967 Jan 20;155(3760):279-84. Construction of phylogenetic trees. Fitch WM, Margoliash E.



5

# Phylogenetics

- Darwin, 1837



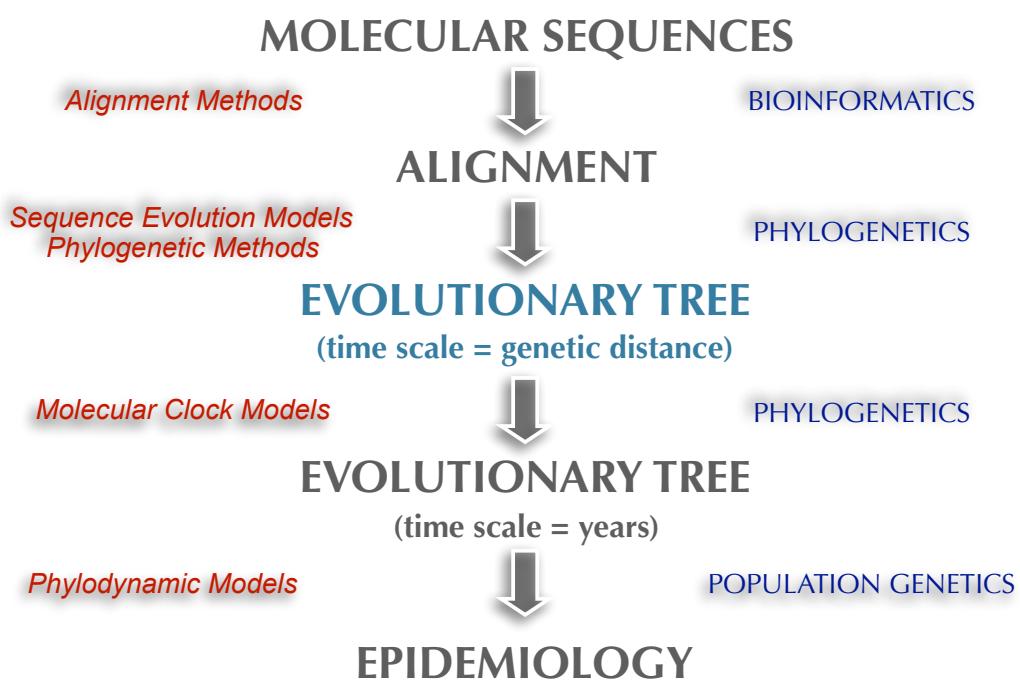
Genome distance tree with lateral gene transfers

# Information in (viral) molecular sequences

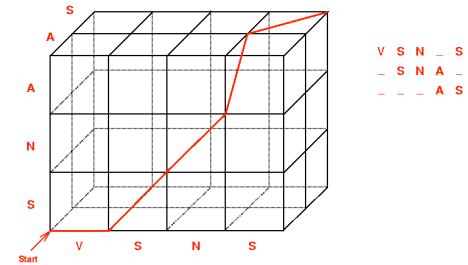
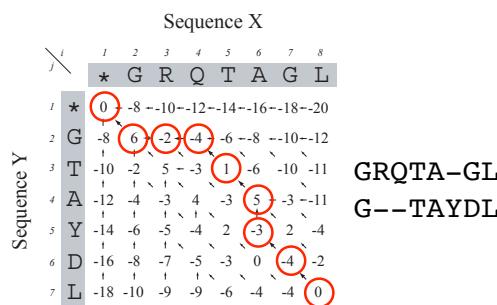
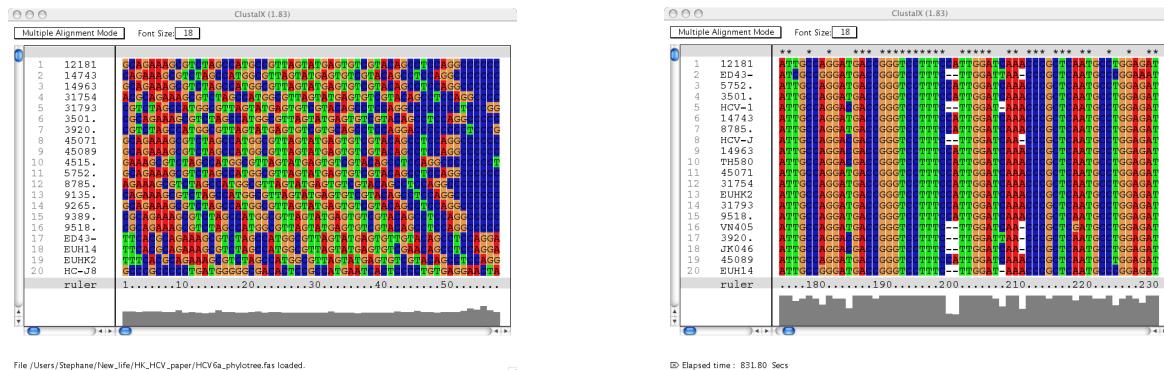
- Genetic distances among strains
- Phylogeny
  - subtyping/classification
  - identification of transmission clusters
  - association with risk factors / traits
  - forensics
- Dates of historical events
- Evolutionary processes
  - recombination
  - natural selection
- Epidemiological processes
  - transmission rates
  - movement among locations
- Phenotypic trait evolution?

```
HIV-1 (UK)  ATC---TGCTAAAGCAATATGACACAGAGGTACATAATGTTT  
HIV-1 (USA)  ATCGGATGCTAGAGCTATGATACAGAGGTACA---TGTTT
```

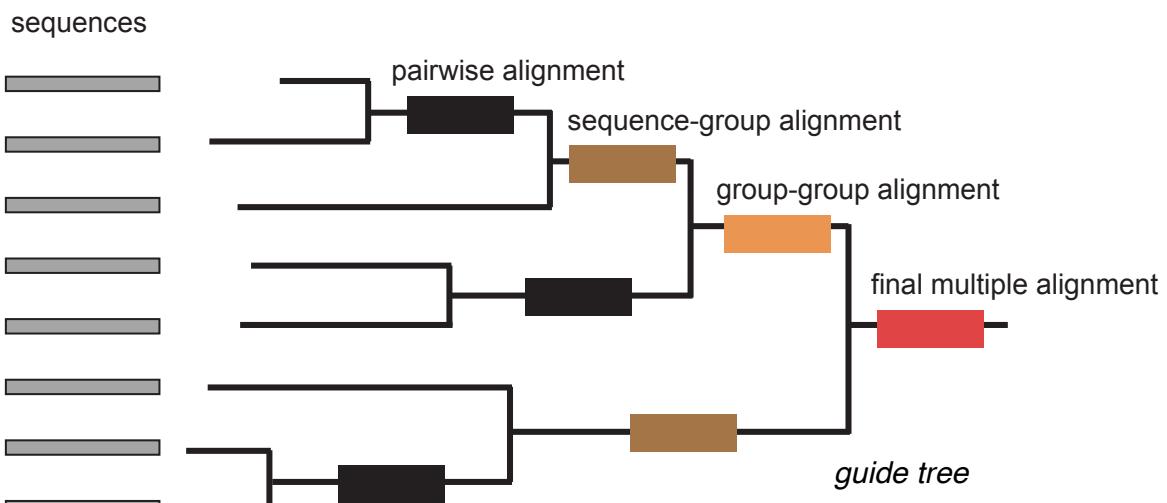
## Our goal



# Sequence alignment



# Progressive alignment



# Genetic distances

SIVcpz HIV-1	ATGGGTGCGA ATGGGTGCGA	GAGCGTCAGT GAGCGTCAGT	TCTAACAGGG ATTAAAGCAGG	GGAAATTAG GGAGAATTAG	ATCGCTGGGA ATCGATGGGA
SIVcpz HIV-1	AAAAGTTCCG AAAAATTCCG	CTTAGGCCCG TTAAGGCCAG	GGGAAAGAAA GGGGAAAGAA	AAGATATATG AAAATATAAA	ATGAAACATT TTAAAACATA
SIVcpz HIV-1	TAGTATGGGC TAGTATGGGC	AAGCAGGGAG AAGCAGGGAG	CTGGAAAGAT CTAGAACGAT	TCGCATGTGA TCGCAGTTAA	CCCGGGCTA TCCTGGCCTG
SIVcpz HIV-1	ATGGAAAGTA TTAGAAACAT	AGGAAGGGATG CAGAAGGCTG	TACTAAATTG TAGACAAATA	TTACAACAAT CTGGGACAGC	TAGAGCCAGC TACAACCATC
SIVcpz HIV-1	TCTCAAAACA CCTTCAGACA	GGCTCAGAAG GGATCAGAAG	GACTGCGGTC AACTTAGATC	CTTGTAAAC ATTATATAAT	ACTCTGGCAG ACAGTAGCAA
SIVcpz HIV-1	TACTGTGGTG CCCTCTATTG	CATACATAGT TGTGCATCAA	GACATCACTG AGGATAGAGA	TAGAAGACAC TAAAAGACAC	ACAGAAAGCT CAAGGAAGCT
SIVcpz HIV-1	CTAGAACAGC TTAGACAAGA	TAAAGCGGCA TAGAG--GAA	TCATGGAGAA -----GAGCA	CAACAGAGCA AAACAAAAGT	AAACTGAAAG AA---GAAA
SIVcpz HIV-1	TAACTCAGGA AACACAGCA	AGCCGTGAAG AGC-----AG	GGGGAGCCAG CAGCTGACA-	TCAAGGCCT -CAGGACAC-	AGTGCCTCTG AG---CAGC--
SIVcpz HIV-1	CTGGCATTAG CAGG--TCAG	TGGAAATTAC CCAAAATTAC			

chimpanzee SIV vs HIV-1 envelope gene

# Not all mutations are equally likely

- some point substitutions are more likely to occur than others:  
transitions are more likely than transversions

- transitions:

purine↔purine or  
pyrimidine↔pyrimidine

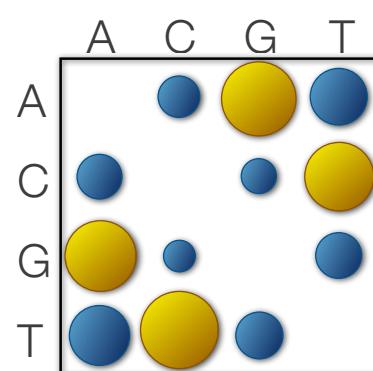


- transversions:



● Transversions

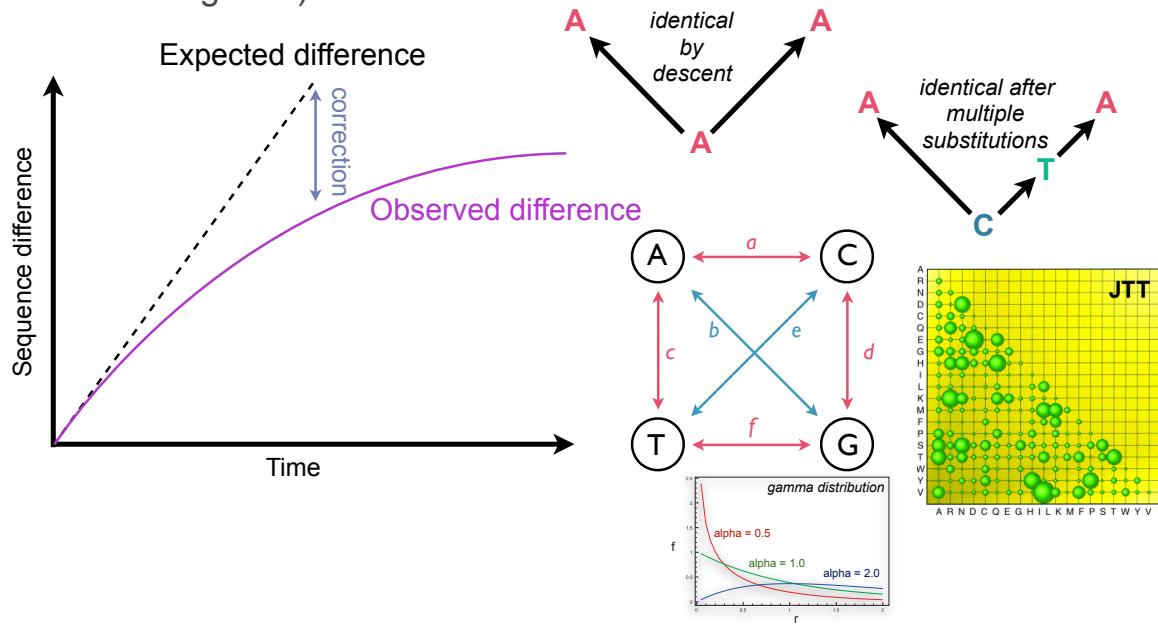
● Transitions



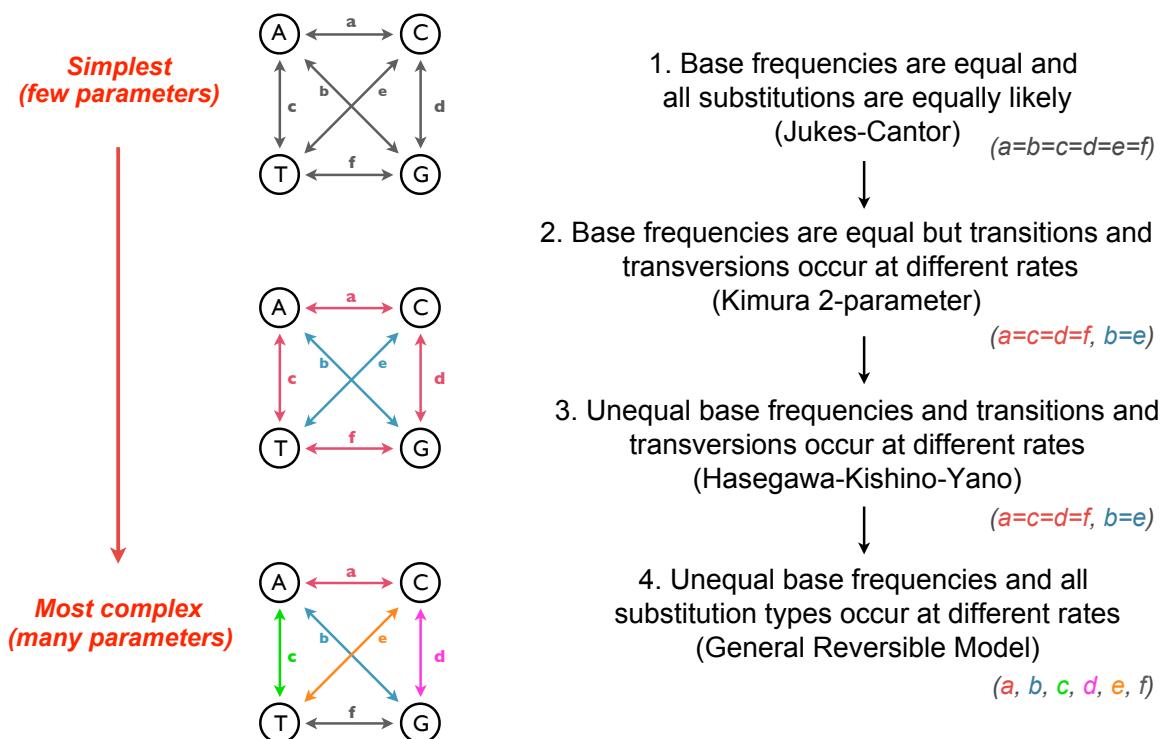
Unambiguous changes on most parsimonious tree of Ciliate SSUrDNA

# Substitution models

- During evolution, ‘multiple hits’ can occur at a single position: the evolutionary distance is almost always larger than the dissimilarity (% nt or aa divergence)



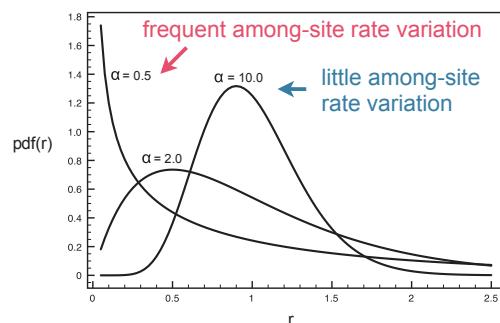
## Nucleotide substitution models



# Does this matter?

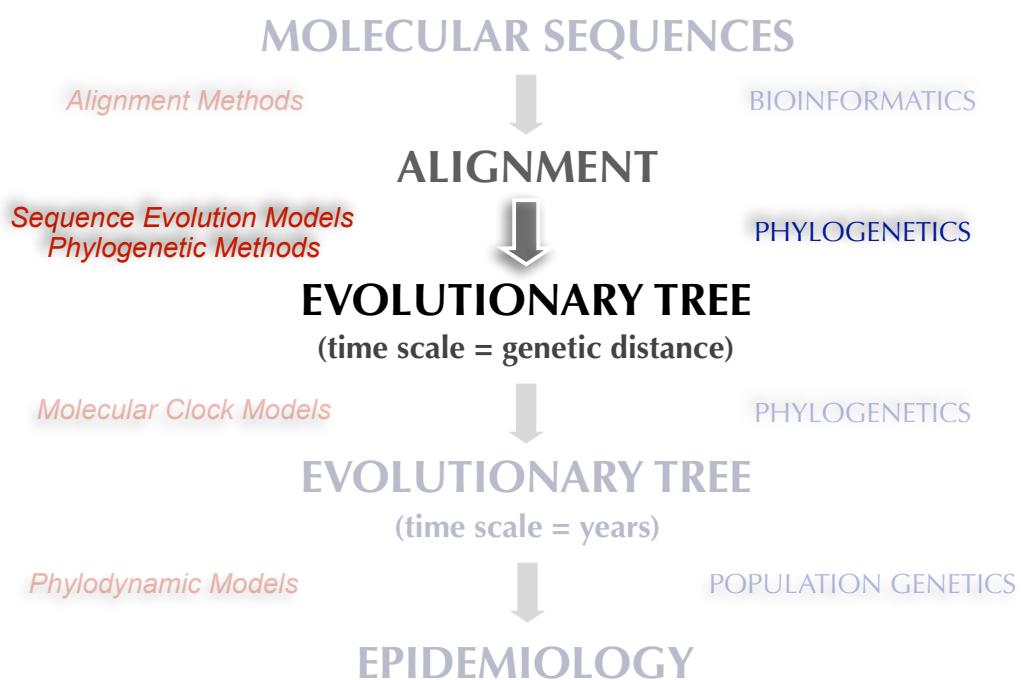
Estimated genetic distances between SIVcpz and HIVlai,  
under different substitution models:

Observed % mismatches	= 0.406
JC (Jukes-Cantor)	= 0.586
HKY (Hasegawa-Kishino-Yano)	= 0.611
GTR (General Time Reversible)	= 0.620
GTR + gamma	= 1.017

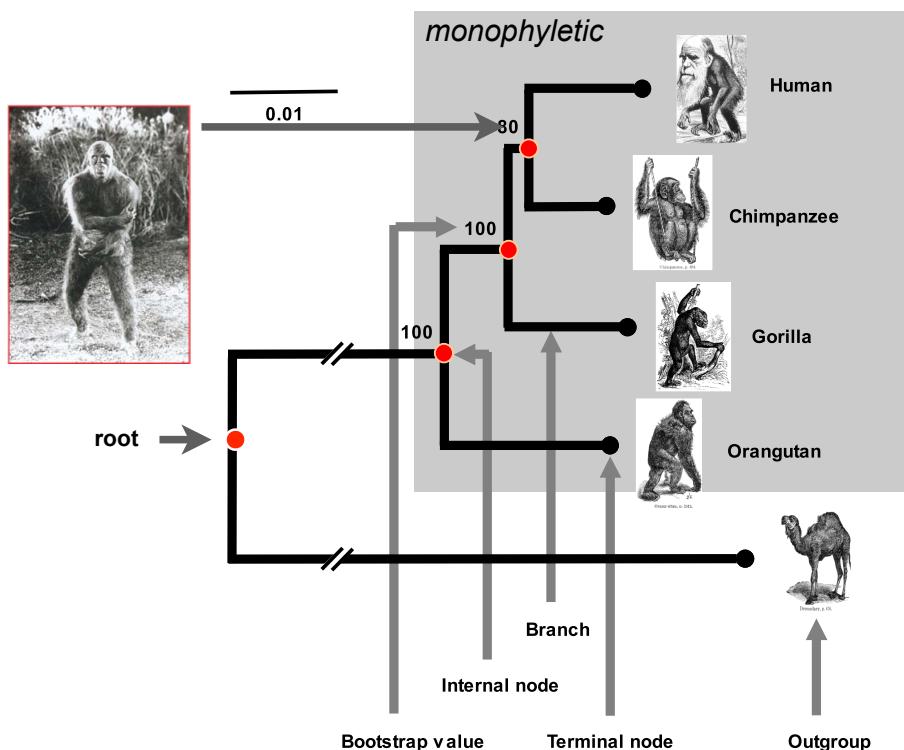


Gene	$\alpha$
Prolactin	1.37
Albumin	1.05
C-myc	0.47
Cytochrome $\beta$ (mtDNA)	0.44
Insulin	0.40
D-loop (mtDNA)	0.17
12S rRNA (mtDNA)	0.16

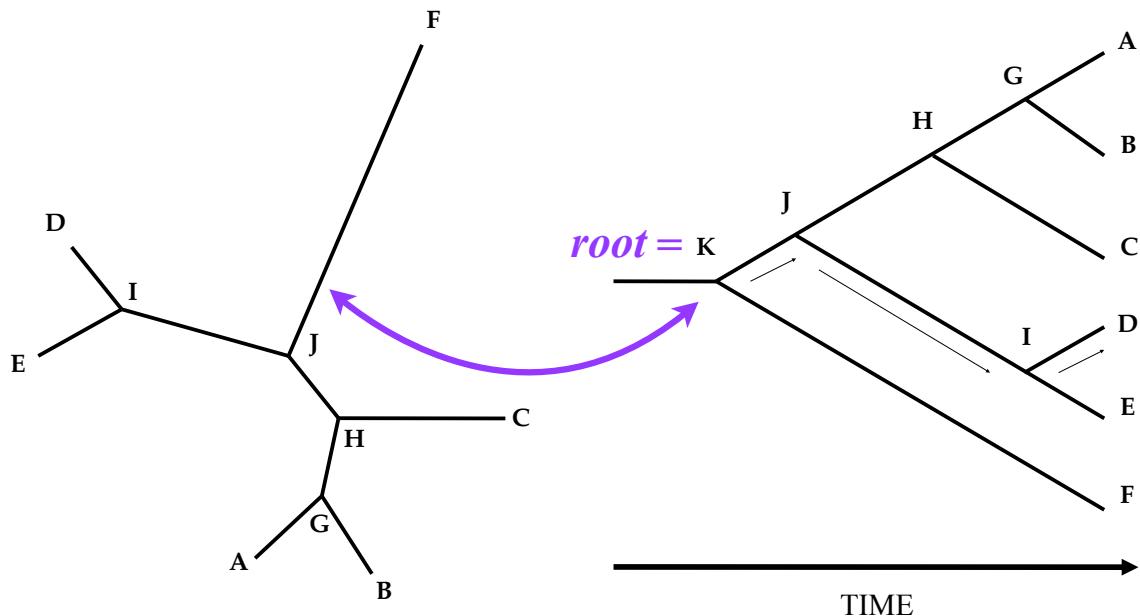
# Phylogenetic reconstruction



# What is a tree?



## Tree terminology: unrooted and rooted



# Phylogenetic reconstruction

- **CLUSTERING APPROACHES:** These begin with a genetic distance between each pair of sequences. A ‘clustering algorithm’ then transforms the genetic distances into a tree.
  - e.g. UPGMA, Neighbour-Joining
  - Simple, faster.
  - No measure of how good the estimated tree is (non-statistical)
- **OPTIMALITY METHODS:** These define a score for each possible tree. ‘Search algorithms’ are then used to find the tree with the highest score.
  - e.g. Parsimony, Maximum Likelihood (& Bayesian Inference)
  - More complex, slower. Search may not locate the ‘best’ tree.
  - Quality of each tree can be directly compared (statistical)

# Phylogenetic reconstruction

- For  $n$  taxa, there are:

$$(2n-3)!/[(2^{n-2})^*(n-2)!]$$

rooted, binary trees

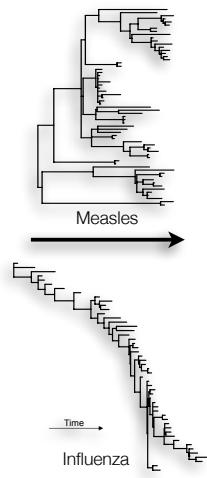
# taxa	# trees	
4	15	enumerable by hand
5	105	enumerable by hand on a rainy day
6	945	enumerable by computer
7	10395	still searchable very quickly on computer
8	135135	a bit more than the number of hairs on your head
9	2027025	population of Glasgow
10	34459425	≈ upper limit for exhaustive searching; about the number of possible combinations of numbers in the National Lottery
20	$8.20 \times 10^{21}$	≈ upper limit for branch-and-bound searching
48	$3.21 \times 10^{70}$	≈ the number of particles in the universe
136	$2.11 \times 10^{267}$	=number of trees to choose from in the “Out of Africa” data (Vigilant et al., 1991)

# Phylogenetic inference: books

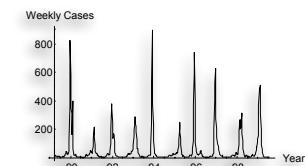
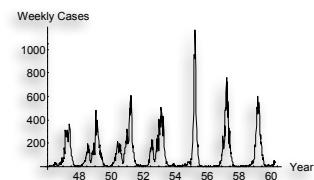


## Phyldynamics™

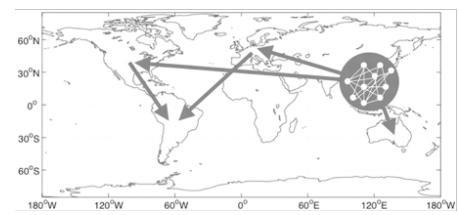
GENETIC DIVERSITY  
(phylogenetics &  
molecular evolution)



EPIDEMIC DYNAMICS  
(mathematical epidemiology)



NATURAL SELECTION  
(population genetics &  
immunology)



## Unifying principle

“ Rapidly evolving pathogens are unique in that their ecological and evolutionary dynamics occur on the same timescale and can therefore potentially interact. ”

Pybus & Rambaut (2009) Nat. Rev. Genetics 10:540-50

## Fundamental Phylodynamic Questions

---

- How genetically diverse is a pathogen population?
- How do pathogen genomes change through time?
- How does pathogen genetic diversity vary through time and space?
- What are the effects of pathogen genetic diversity on virulence, transmissibility, resistance to treatment, etc.

## Specific questions

---

- When did a epidemic start?
- Where did it come from?
- How fast is it transmitting?
- In what direction is it spreading?
- Are hosts X, Y & Z epidemiologically linked?
- Of how many strains is the epidemic composed?
- Are strains associated with particular transmission routes?
- What adaptations has it accrued?

## Fundamental Phylodynamic Questions

---

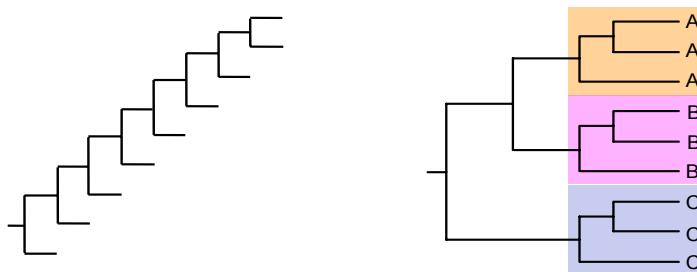
- How genetically diverse is a pathogen population?
- How do pathogen genomes change through time?
- How does pathogen genetic diversity vary through time and space?
- What processes and/or events determine these changes?
- What are the effects of pathogen genetic diversity on virulence, transmissibility, resistance to treatment, etc.

## Measuring sequence diversity

- Not as straightforward as you might think...
- Are your pathogen sequences all sampled at the same time?

If sequences not sampled over time it's difficult to separate the effects of diversity and divergence on genetic diversity.
- Are you measuring sample diversity or population diversity?

The former is simply a summary of your data, the latter is an inference about the population you have sampled. Sequences should be sampled randomly to estimate the latter.

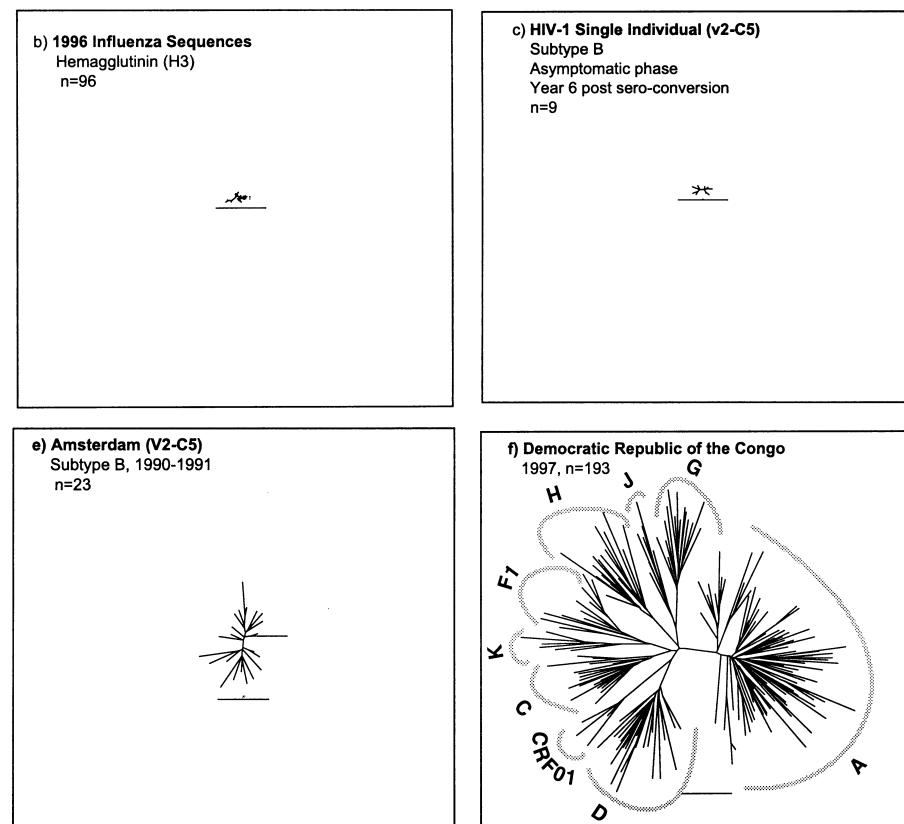


## Measuring sequence diversity

- Are you studying an inter-host or intra-host population?

For the former, each sequence represents a different infection.  
For the latter, each sequence represents a different virion within an infected individual. The measure of diversity must be interpreted accordingly.
- How do we deal with intra-host diversity when studying the inter-host level?
- Intra-host diversity is low for most acute infections (e.g. influenza) but can be high for chronic infections (e.g. HIV).

## Example: diversity of HIV-1 versus influenza



Scale bar represents a genetic distance of 0.1 substitutions per site.

Korber et al. 2001. *British Medical Bulletin* 58:19-42

## Phylogenetic Patterns

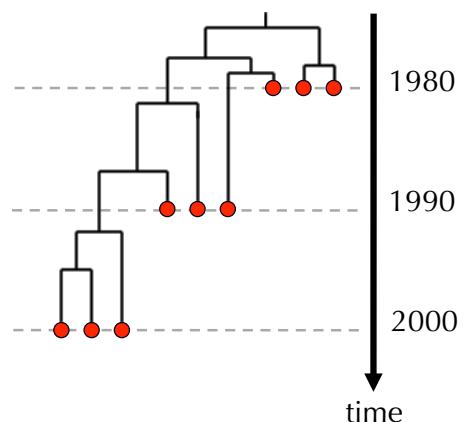
	Continual Immune Selection	Weak/No Immune Selection	
Idealised Phylogeny Shapes		Population dynamics	Spatial dynamics
		Population growth	Strong spatial structure
		Population decline	Weak spatial structure
<b>Examples</b>	Human influenza A within-host HIV	among-host HIV among-host HCV	Measles Rabies, Dengue

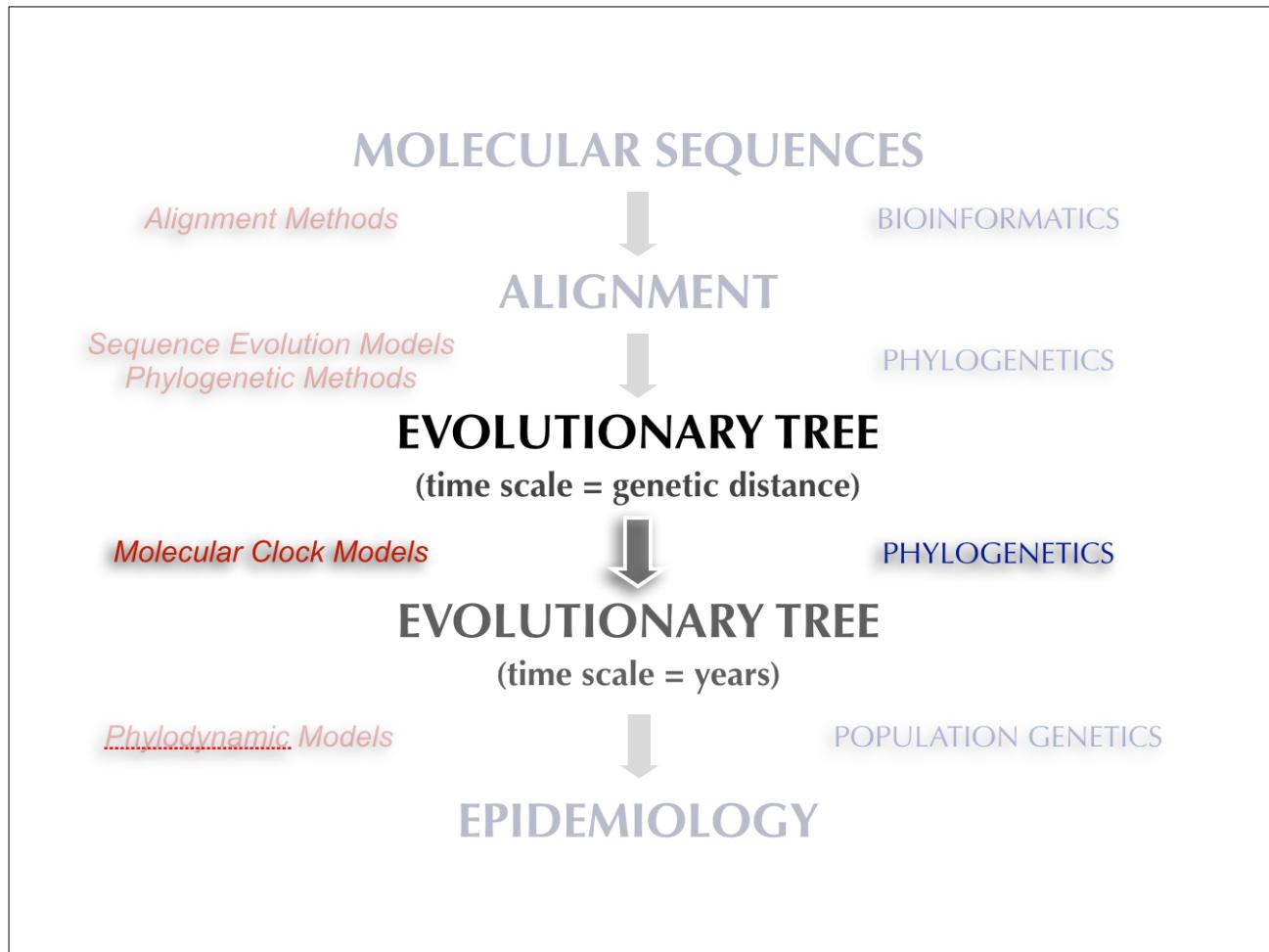
## Fundamental Phylodynamic Questions

- How genetically diverse is a pathogen population?
- How do pathogen genomes change through time?
- How does pathogen genetic diversity vary through time and space?
- What processes and/or events determine these changes?
- What are the effects of pathogen genetic diversity on virulence, transmissibility, resistance to treatment, etc.

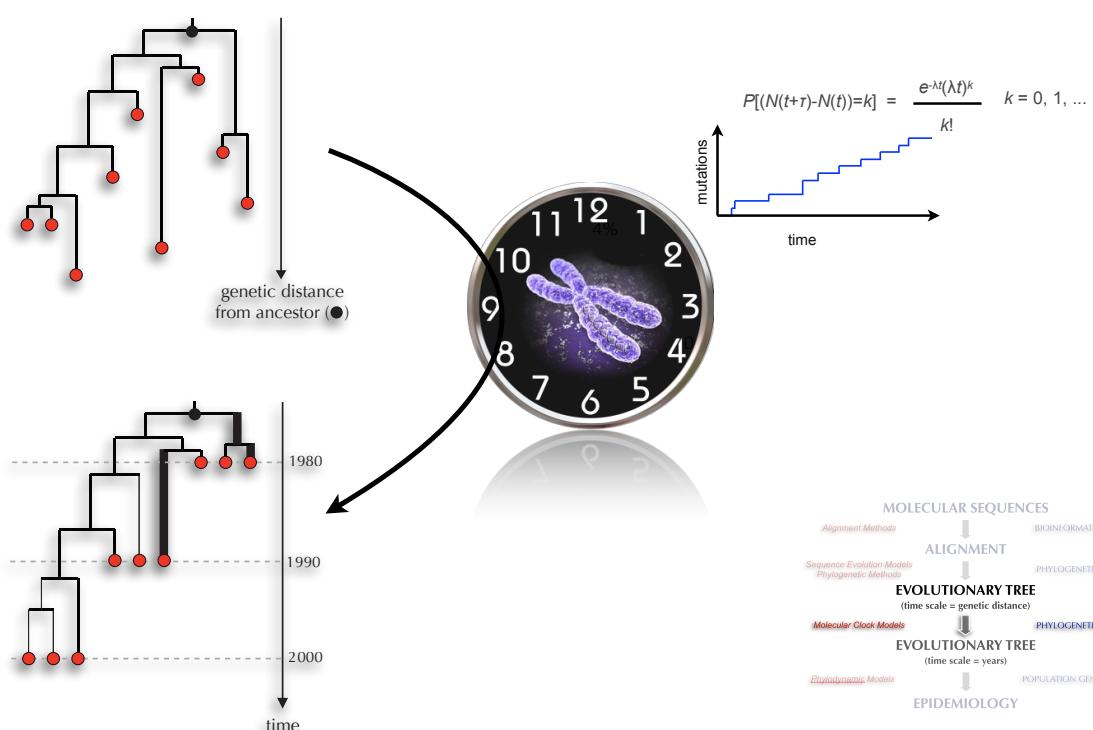
## 'Phylodynamic' Data

- Pathogen genomes are sampled at different points in time and from different locations.
- Hence transmission history is estimated on a real time-scale (e.g. years).

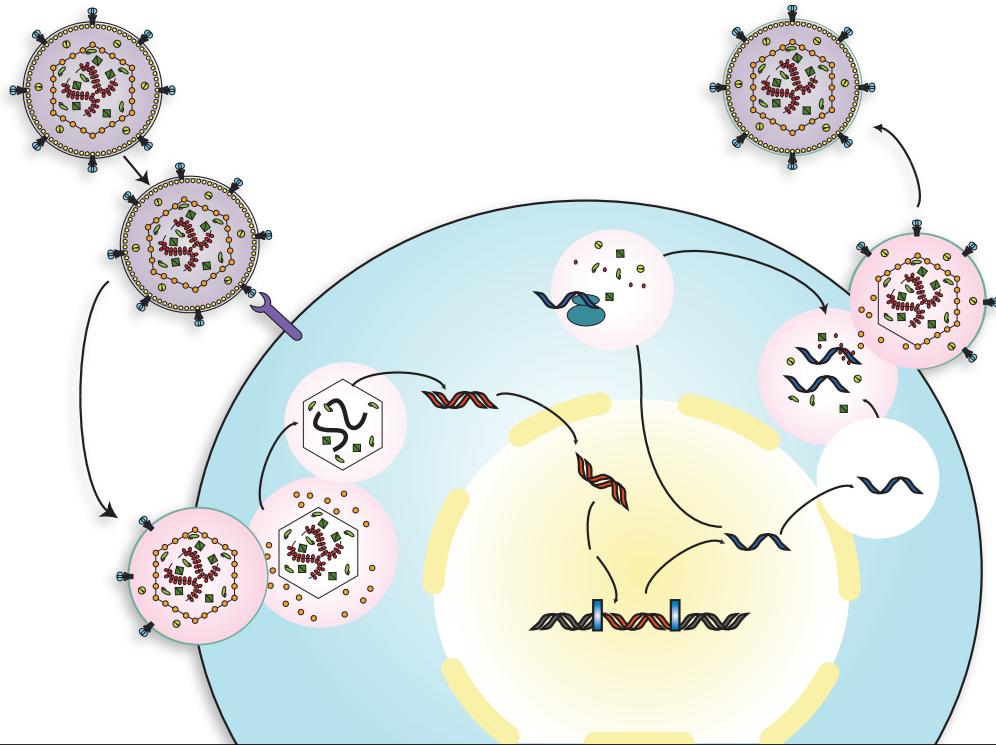




# Molecular clocks



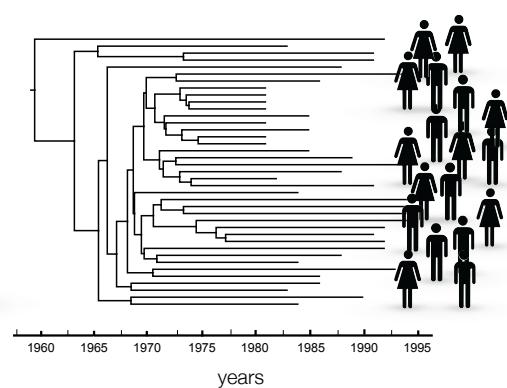
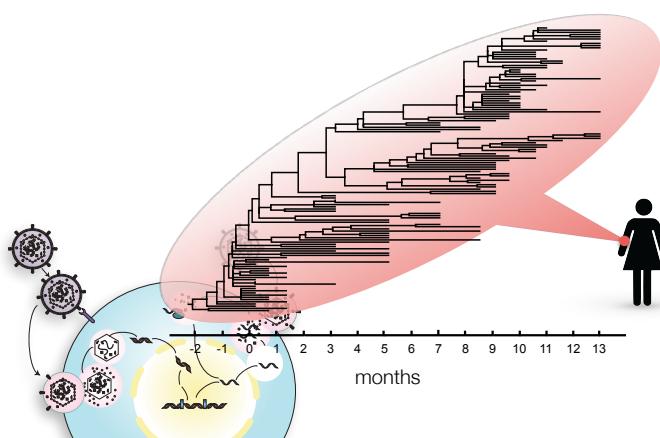
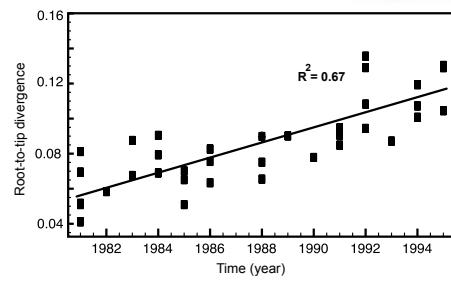
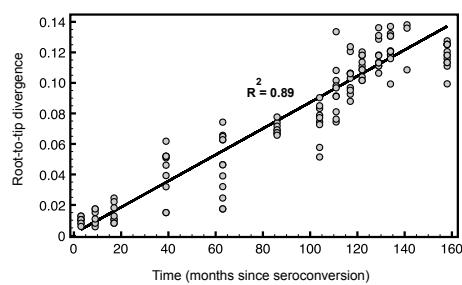
# HIV: the ultimate evolver



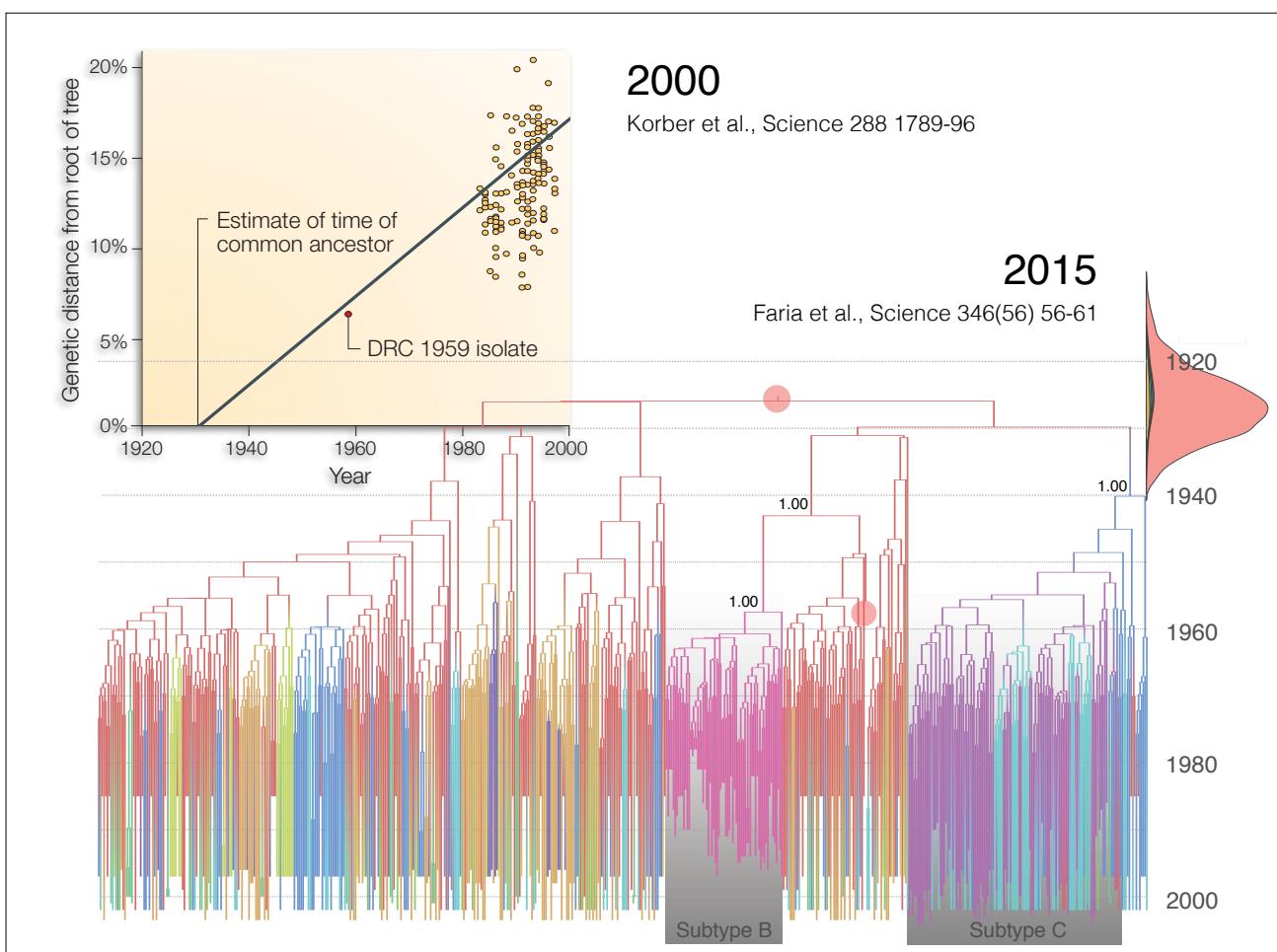
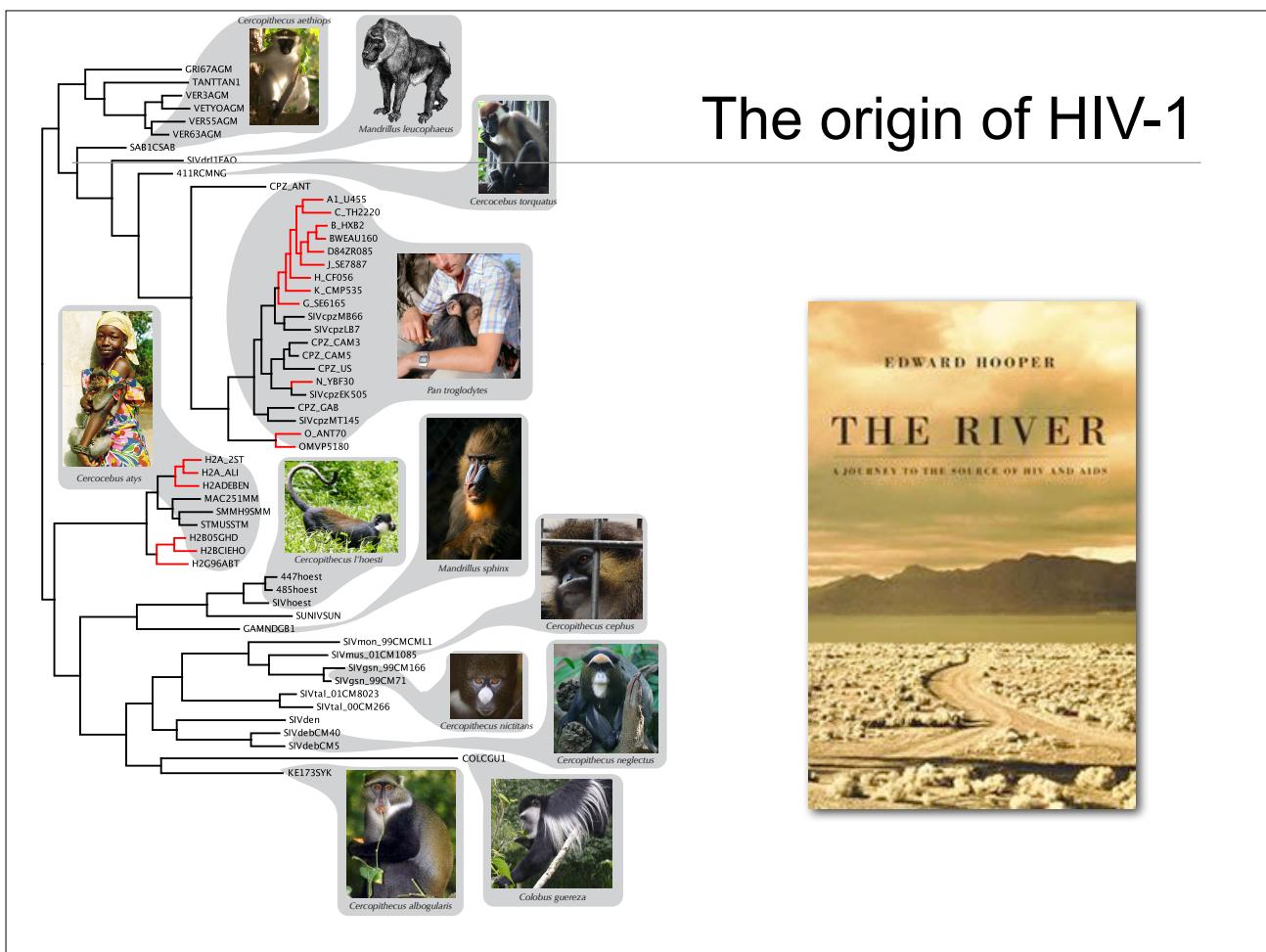
## measurable evolution of HIV-1

Continuous immune selection	Weak/no immune selection
Population dynamics	Spatial dynamics
Population growth	Strong spatial selection
Population decline	Weak spatial selection
Population stable	None
Mixed phylogeny shapes	None

Examples: Human influenza A (continuous), Human HIV-1 (weak/no), Monkey rotavirus (mixed).

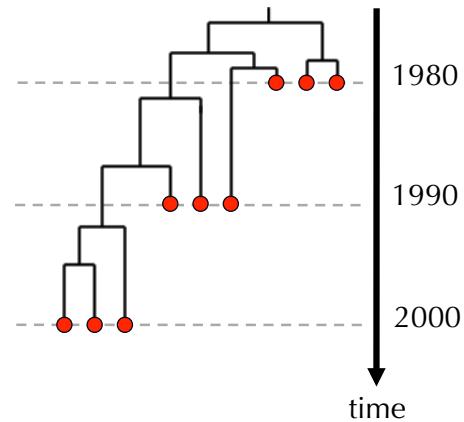


Lemey et al 2006 AIDS Rev

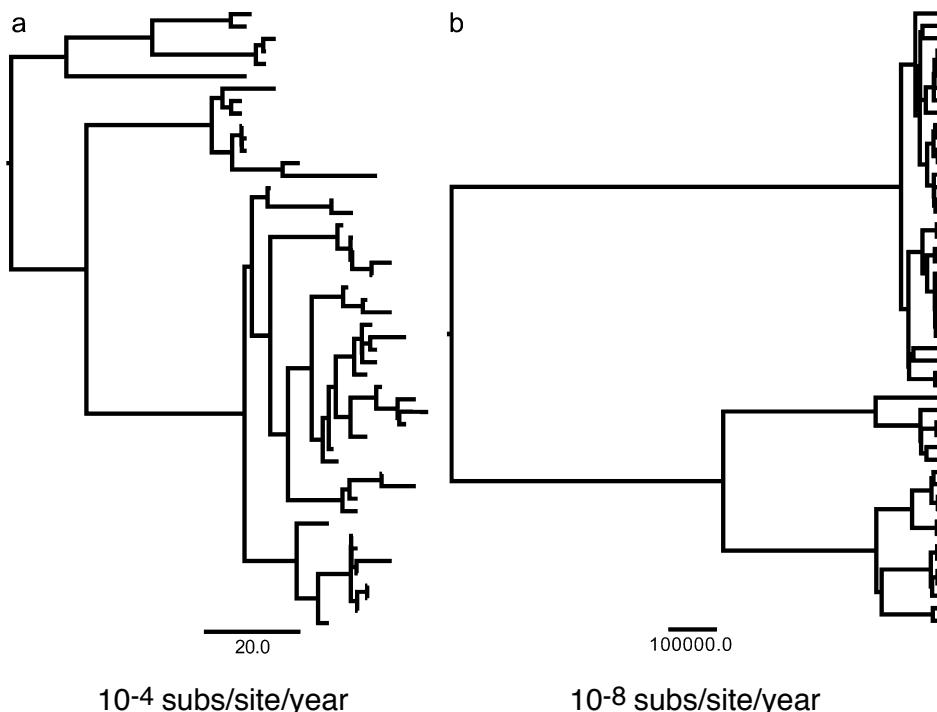


## 'Phylodynamic' Data

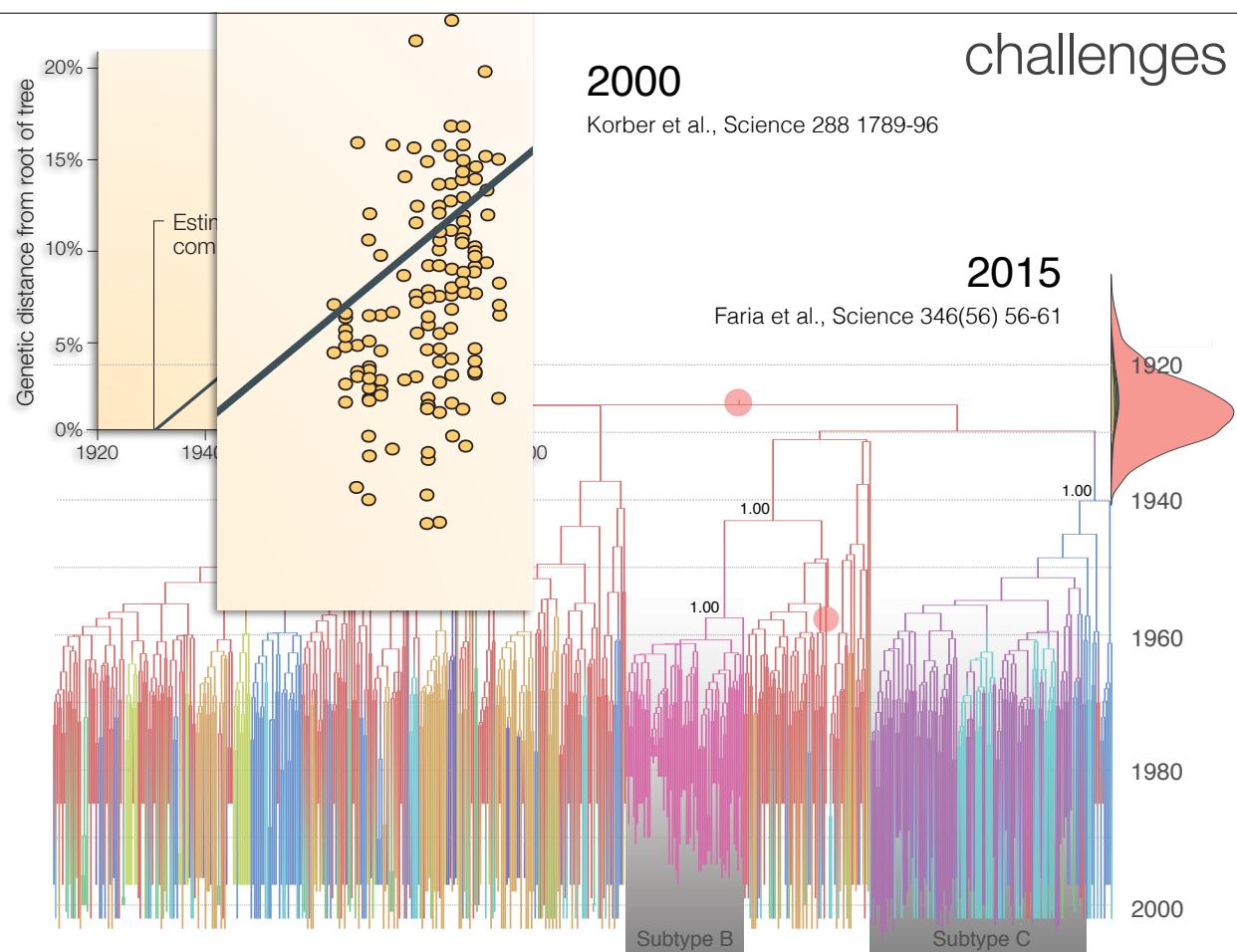
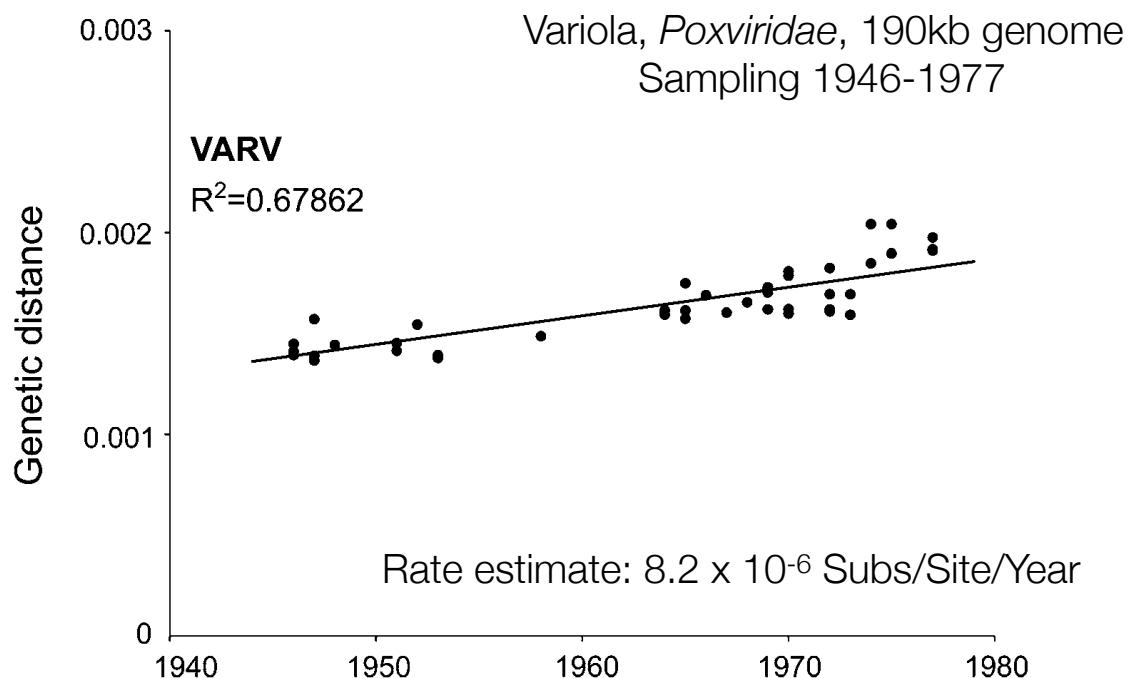
- Pathogen genomes are sampled at different points in time and from different locations.
- Hence transmission history is estimated on a real time-scale (e.g. years).
- The ability to genetically distinguish sequences sampled at different times depends on:
  - (i) the rate of evolution of the gene/genome that is obtained
  - (ii) the length of time between samples
  - (iii) the sequence length of the gene/genome that is obtained



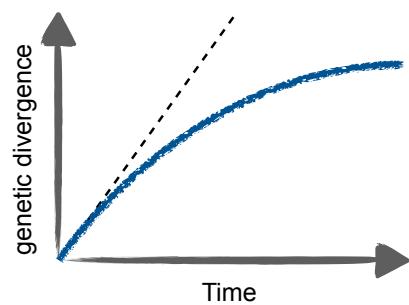
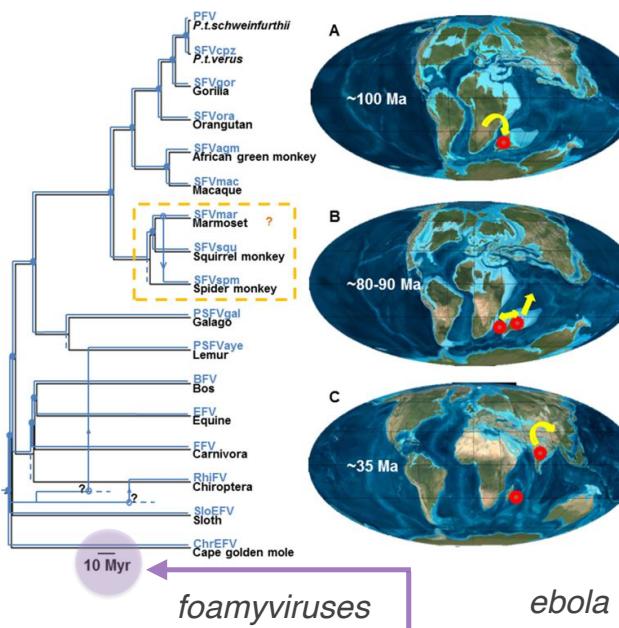
## 'Phylodynamic' Data



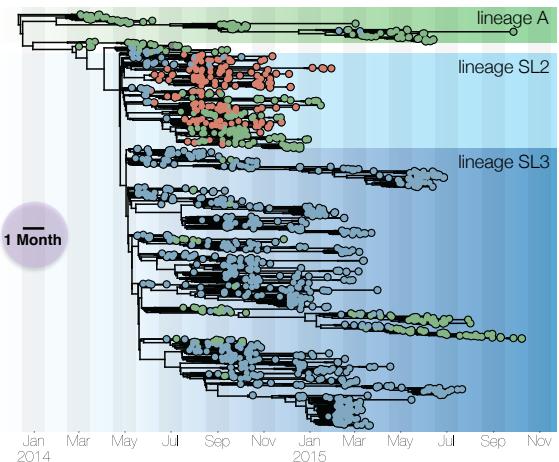
## A DNA virus (smallpox)



Katzourakis et al., Retrovirology, 2014.

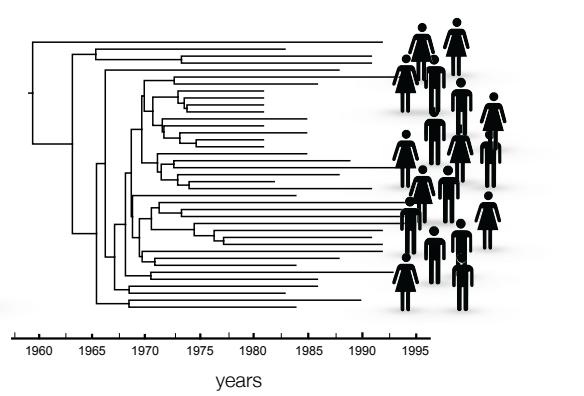
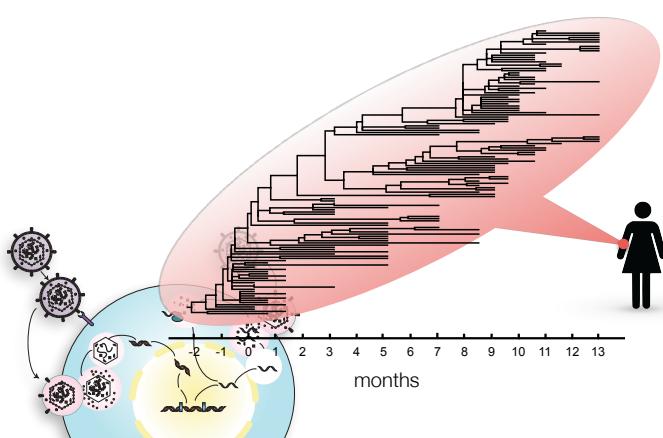
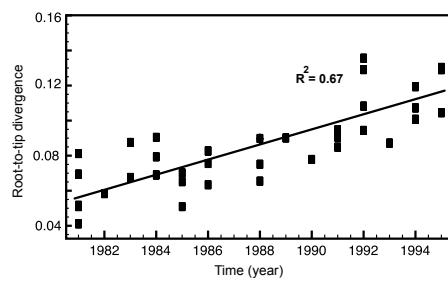
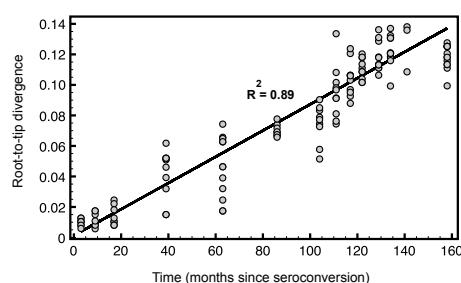


Dudas et al., Nature, 2017.



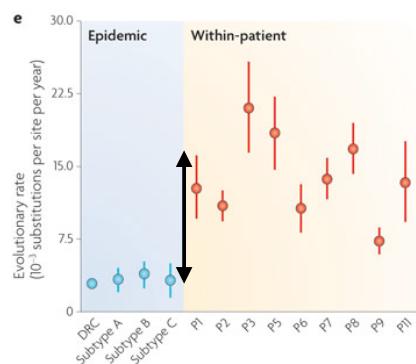
challenges

opportunities



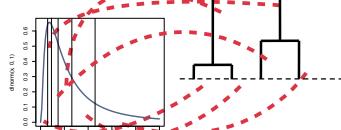
Lemey et al 2006 AIDS Rev

# opportunities

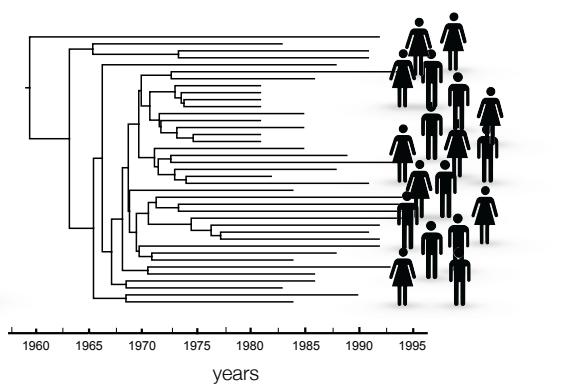
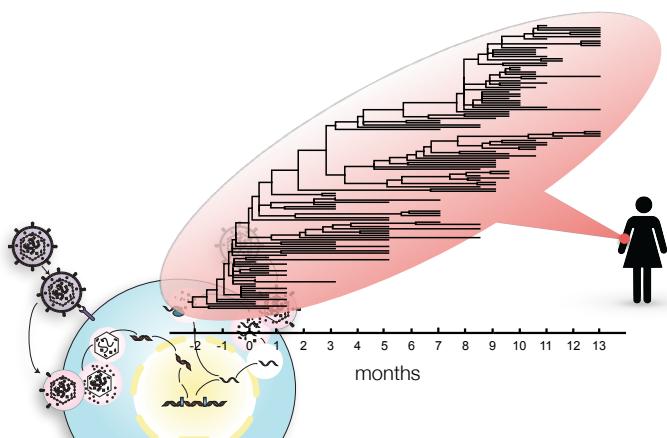


Vrancken et al., PLoS Comp Bio, 2014

Mixed effects model:  
 $\log \mu_i = \theta_i + \beta X_i$

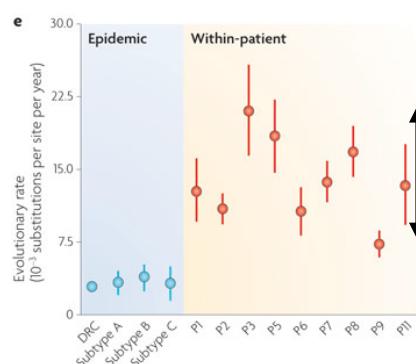


Pybus and Rambaut, NGR, 2009



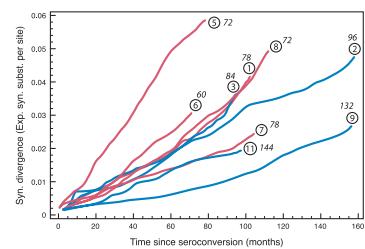
Lemey et al 2006 AIDS Rev

# opportunities



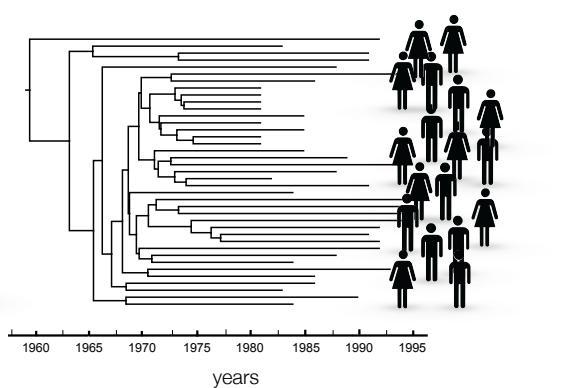
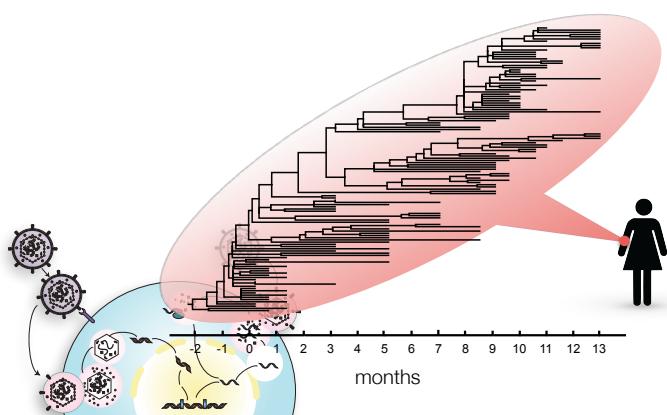
Edo-Matas et al., Mol Biol Evol, 2011

$$\log \theta_i = \beta_0 + \delta_{LTNP} \beta_{LTNP} LTNP_i + \delta_{\Delta 32} \beta_{\Delta 32} \Delta 32_i + \varepsilon_i$$



Pybus and Rambaut, NGR, 2009

Lemey et al., PLoS Comp Bio, 2007



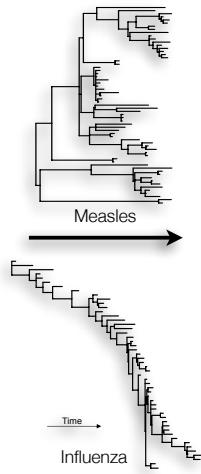
Lemey et al 2006 AIDS Rev

## Fundamental Phylodynamic Questions

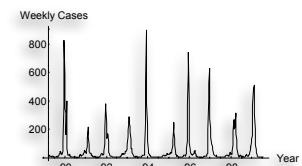
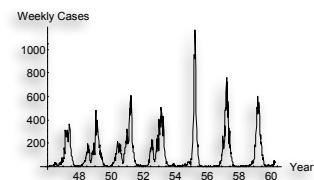
- How genetically diverse is a pathogen population?
- How do pathogen genomes change through time?
- How does pathogen genetic diversity vary through time and space?
- What processes and/or events determine these changes?
- What are the effects of pathogen genetic diversity on virulence, transmissibility, resistance to treatment, etc.

## Phylodynamics™

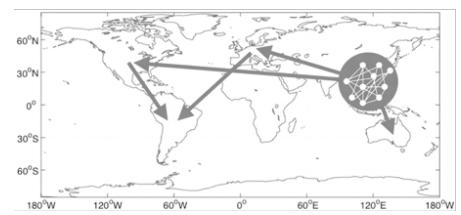
GENETIC DIVERSITY  
(phylogenetics &  
molecular evolution)



EPIDEMIC DYNAMICS  
(mathematical epidemiology)

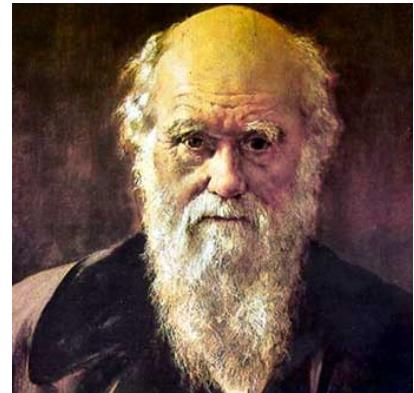


NATURAL SELECTION  
(population genetics &  
immunology)



# Evolutionary processes: natural selection

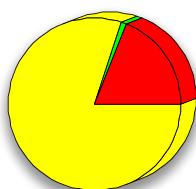
- “the preservation of favourable variations and the rejection of injurious variations, i call natural selection. variations neither useful nor injurious would not be affected by natural selection, and would be left a fluctuating element”
    - darwin, the origin of species



## Evolutionary processes: natural selection

most fixed mutations  
are neutral

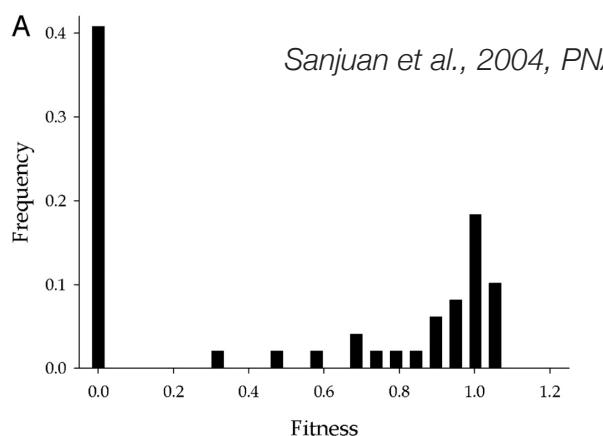
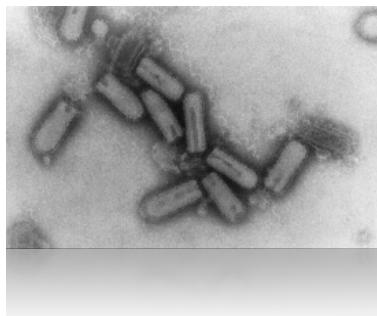
# neutralist model motoo kimura



  $s > 0$   
  $s \approx 0$   
  $s < 0$

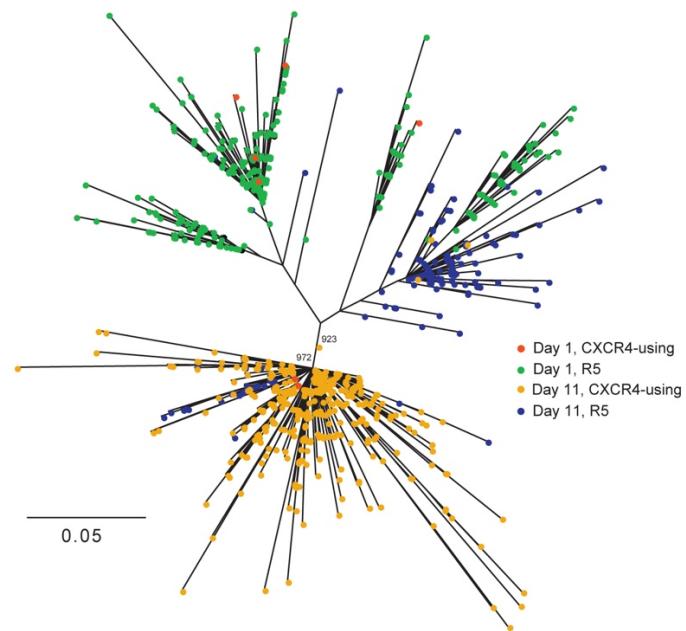
most fixed mutations  
are advantageous

## selectionist model john gillespie



# Evolutionary processes: natural selection

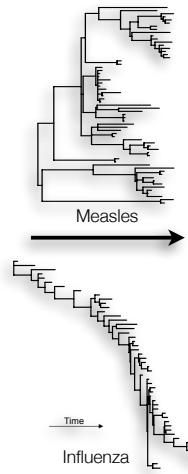
- Immune escape (antibodies\*, T-cells\*, innate immune responses)
- Antiviral drug resistance
- Vaccine escape mutations
- Cell & tissue tropism
- Inter-host viral transmission (i.e. for viral emergence)



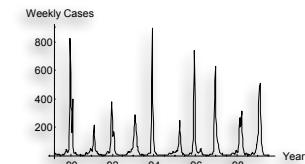
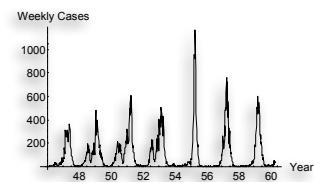
→ **module 15:** Pathogen evolution, selection and immunology

## Phylodynamics™

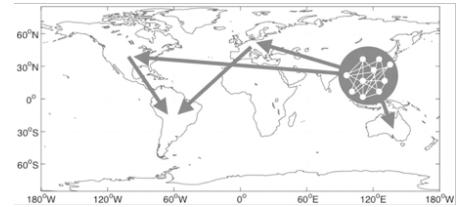
GENETIC DIVERSITY  
(phylogenetics &  
molecular evolution)



EPIDEMIC DYNAMICS  
(mathematical epidemiology)



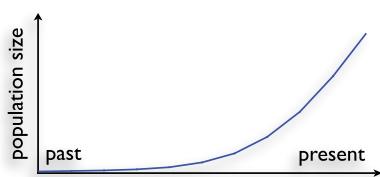
NATURAL SELECTION  
(population genetics &  
immunology)



# Phylodynamic Patterns

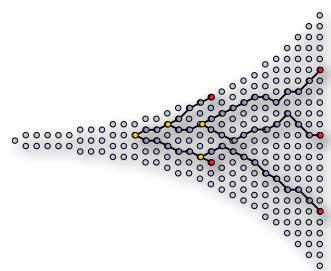
	Continual Immune Selection	Weak/No Immune Selection	
Idealised Phylogeny Shapes		Population dynamics	Spatial dynamics
Examples	Human influenza A within-host HIV	among-host HIV among-host HCV	Measles Rabies, Dengue

## Demography and coalescent theory



- The rate at which lineages ‘coalesce’ depends on population size and population structure.

Kingman JFC (1982) *Journal of Applied Probability* 19A:27–43

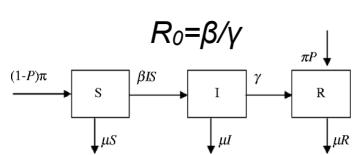


- Population dynamics can be reconstructed using parametric or flexible nonparametric models (the ‘skyline or skyride plot’ method)

Pybus et al. (2000) *Genetics* 155:1429-37

Drummond, Rambaut, Shapiro & Pybus (2005) *Mol Biol Evol* 22:1185-92

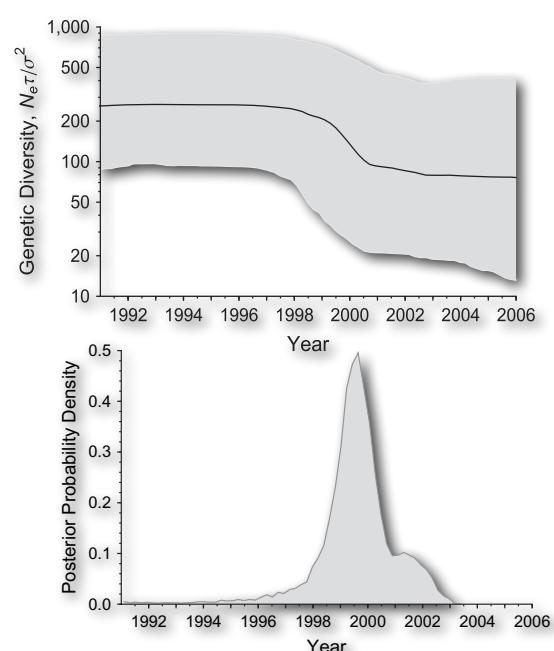
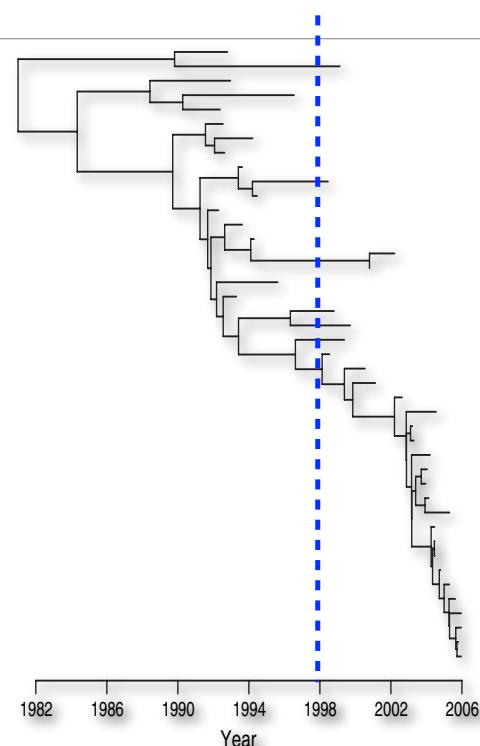
Minin, Bloomquist and Suchard (2008) *Mol Biol Evol* 25:1459-71



- Birth-death models can also be used as the tree-generative model and just like coalescent models they can be parametrized in terms of compartmental epidemic models.

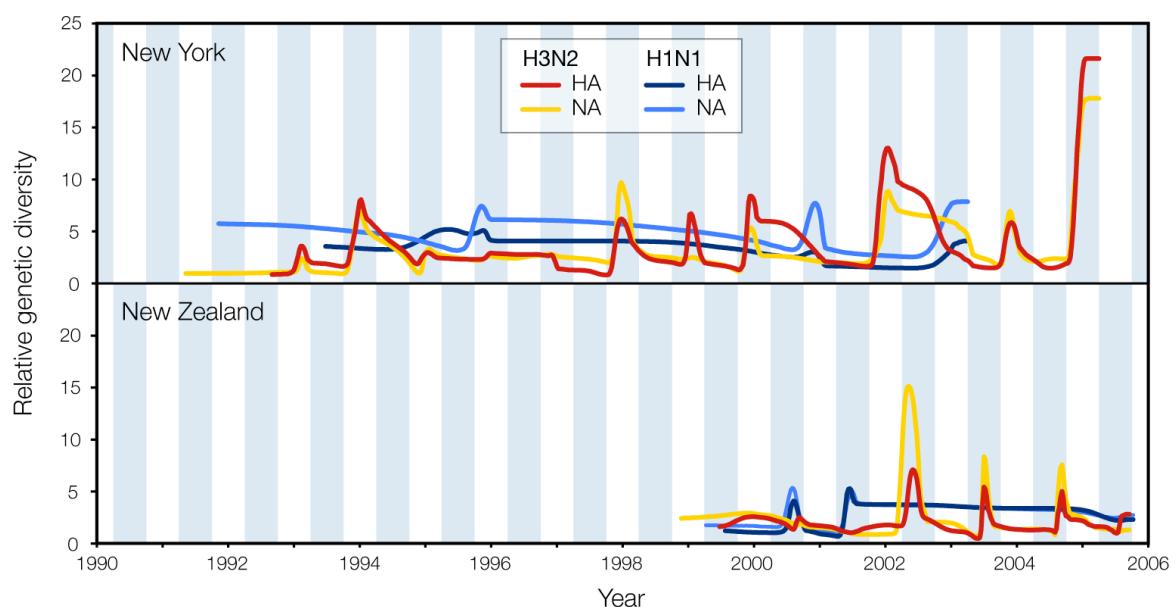
Stadler et al. (2012) *MBE* 29:347-357

# HBV Vaccination in Amsterdam



van Ballegooijen et al. 2009. Am. J. Epidemiol. 170:1455-63

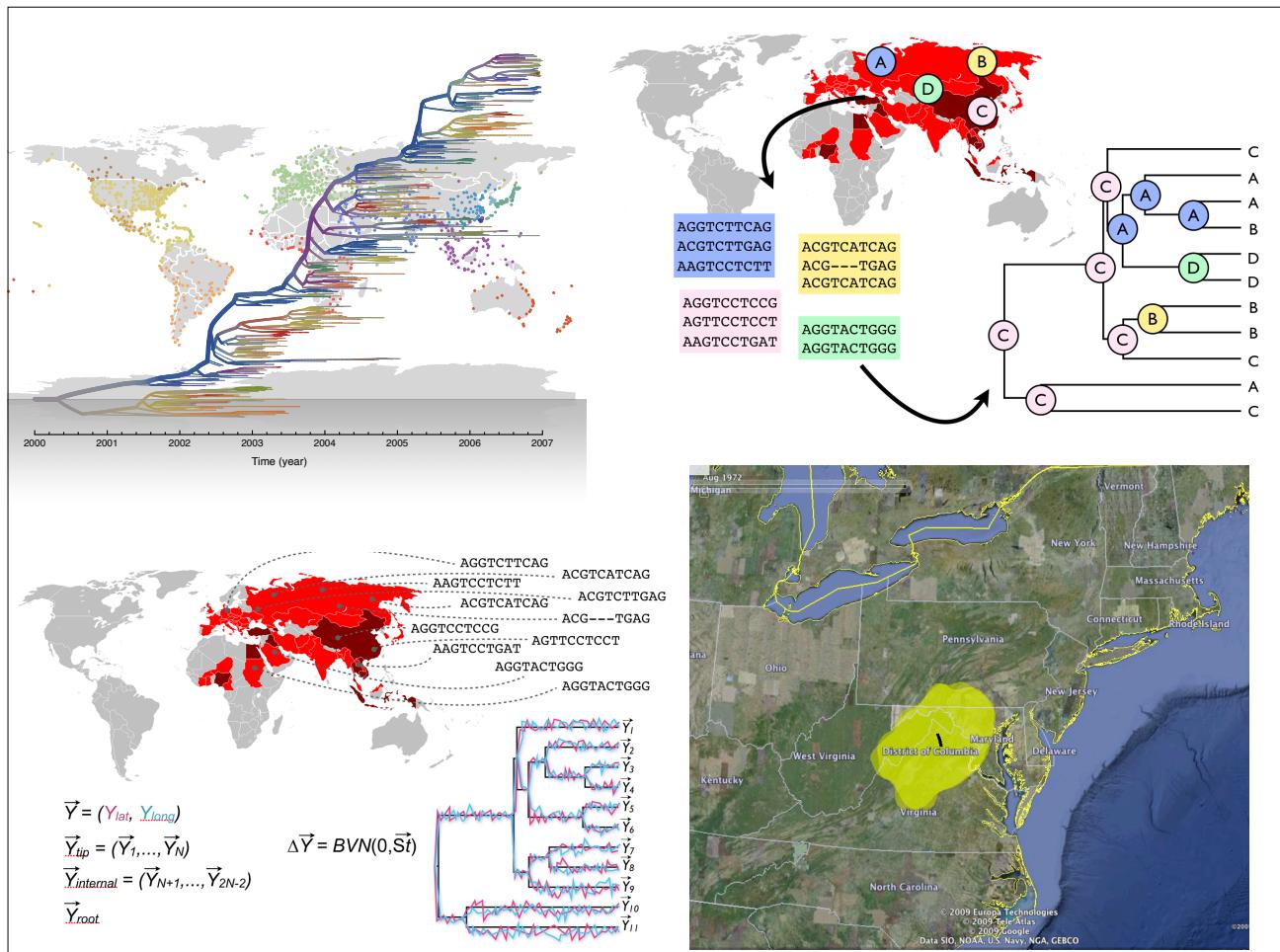
# Influenza H3N2 epidemic dynamics



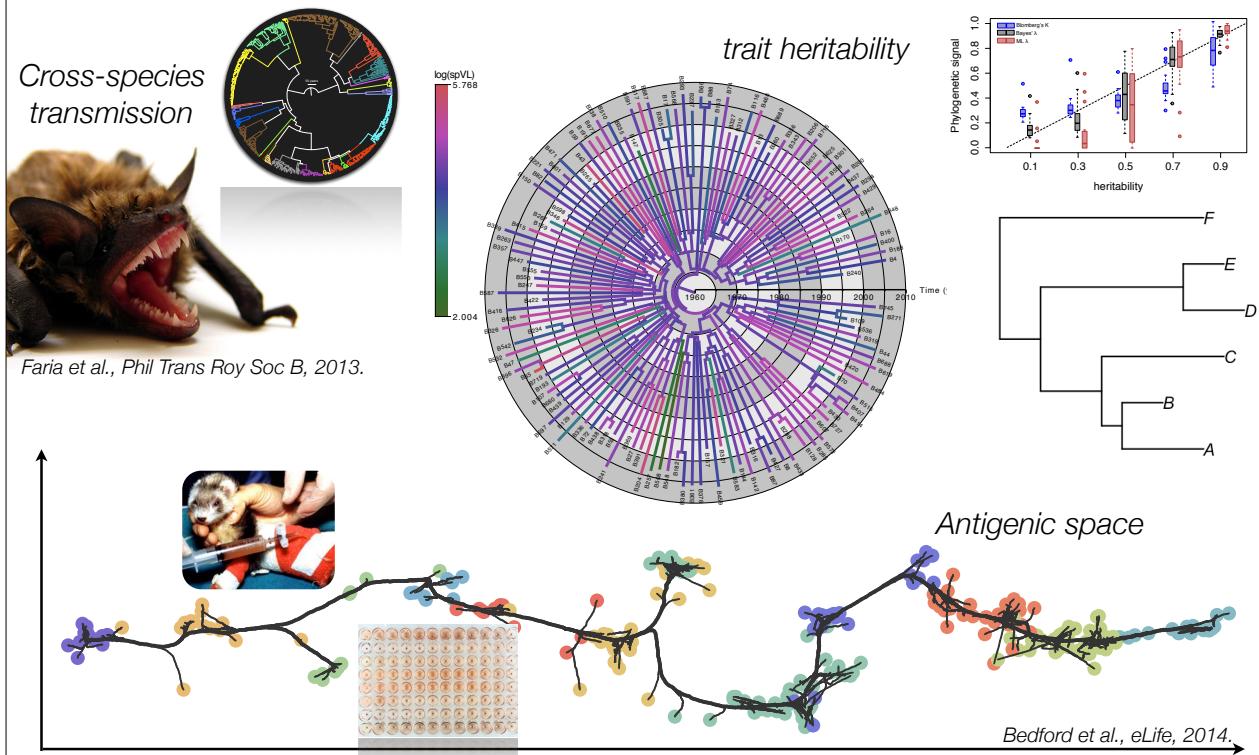
Rambaut et al. 2009. Nature

# PhyloGEOdynamic Patterns

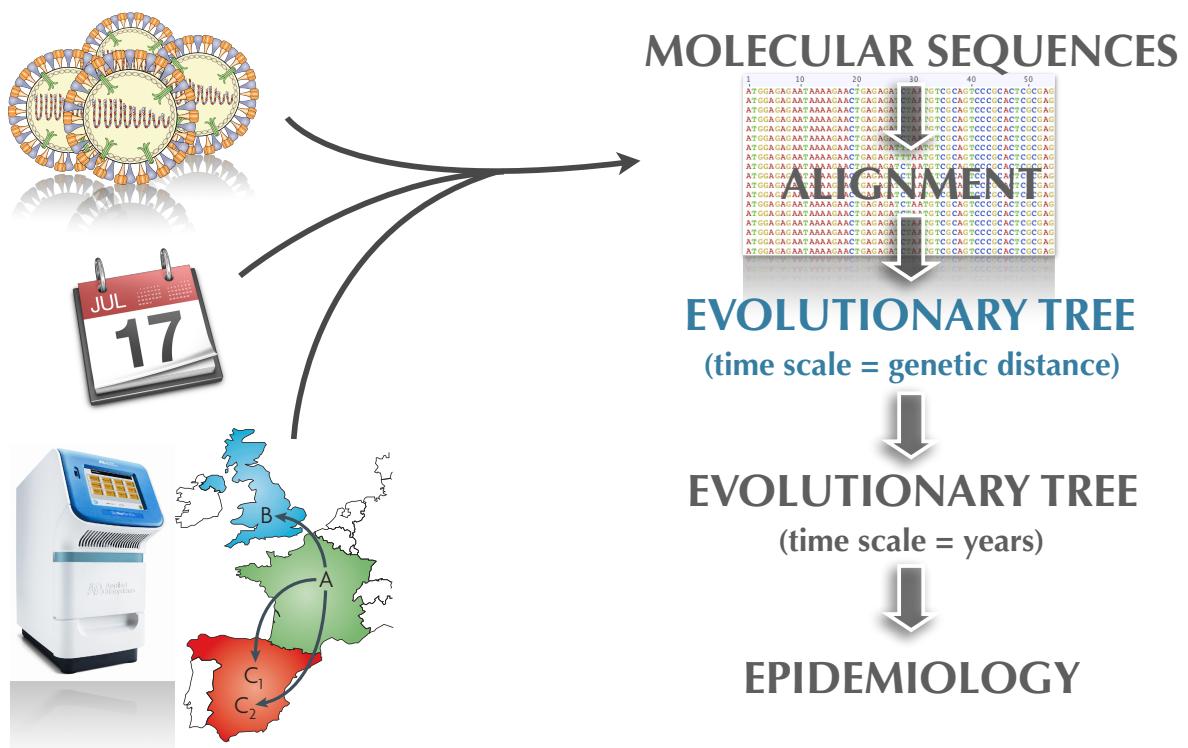
	Continual Immune Selection	Weak/No Immune Selection	
Idealised Phylogeny Shapes		Population dynamics	Spatial dynamics
		Population growth	Strong spatial structure
Population decline			Weak spatial structure
Examples	Human influenza A within-host HIV	among-host HIV among-host HCV	Measles Rabies, Dengue



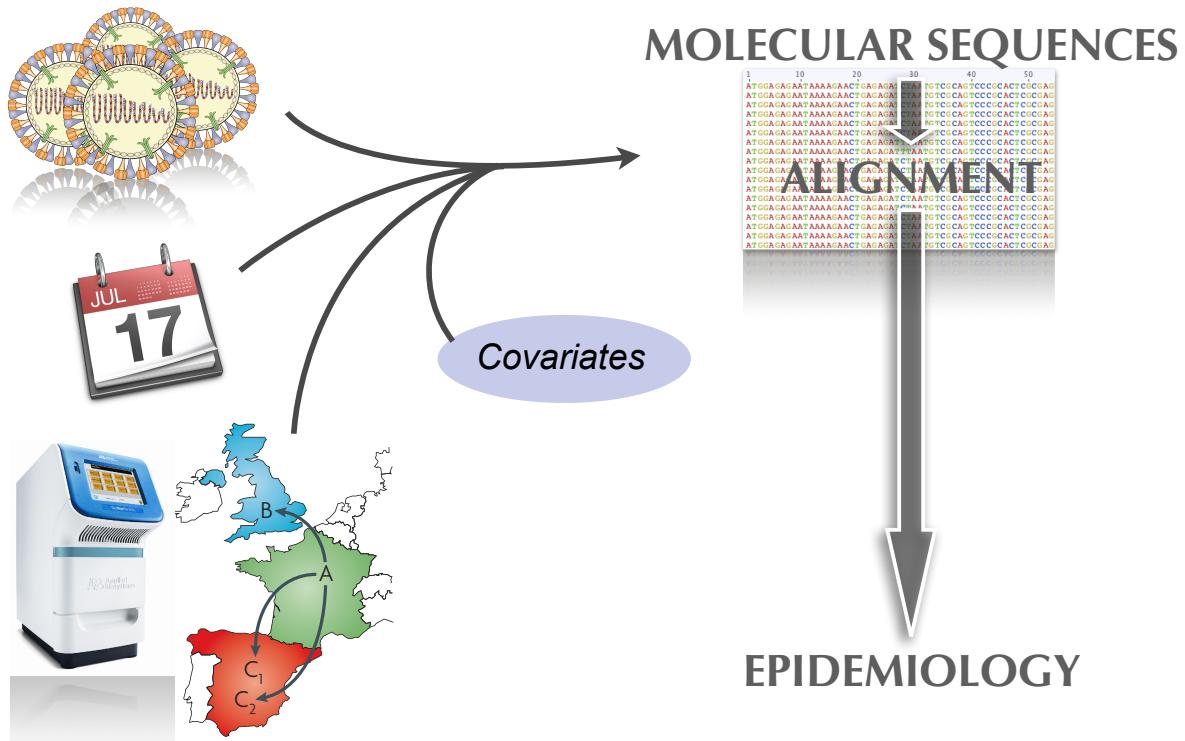
# Trait evolution and the comparative approach



# Bayesian Evolutionary Analysis Sampling Trees



# Bayesian Evolutionary Analysis Sampling Trees



# Bayesian Evolutionary Analysis Sampling Trees

