# Persistent Homology in Genomics
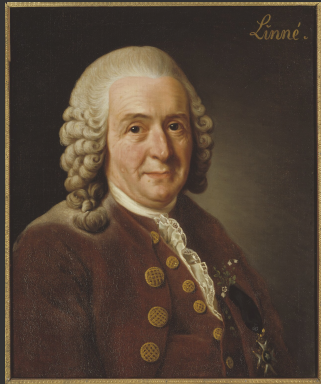
Plan: ① Review of Evolutionary Biology

② Appearance of non-trivial Topology
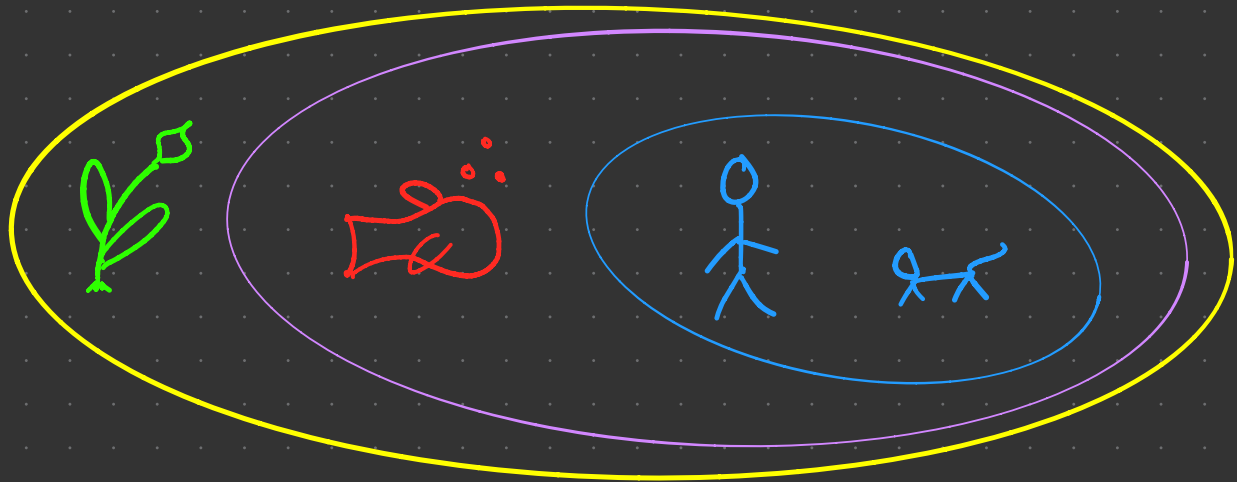   & how Persistent Homology can help

③ Example: Influenza

References: [1] "Topology of viral evolution" (Chen, Carlsson, Rabadan)

[2] "TDA for Genomics and Evolution" (Rabadan, Blumberg)
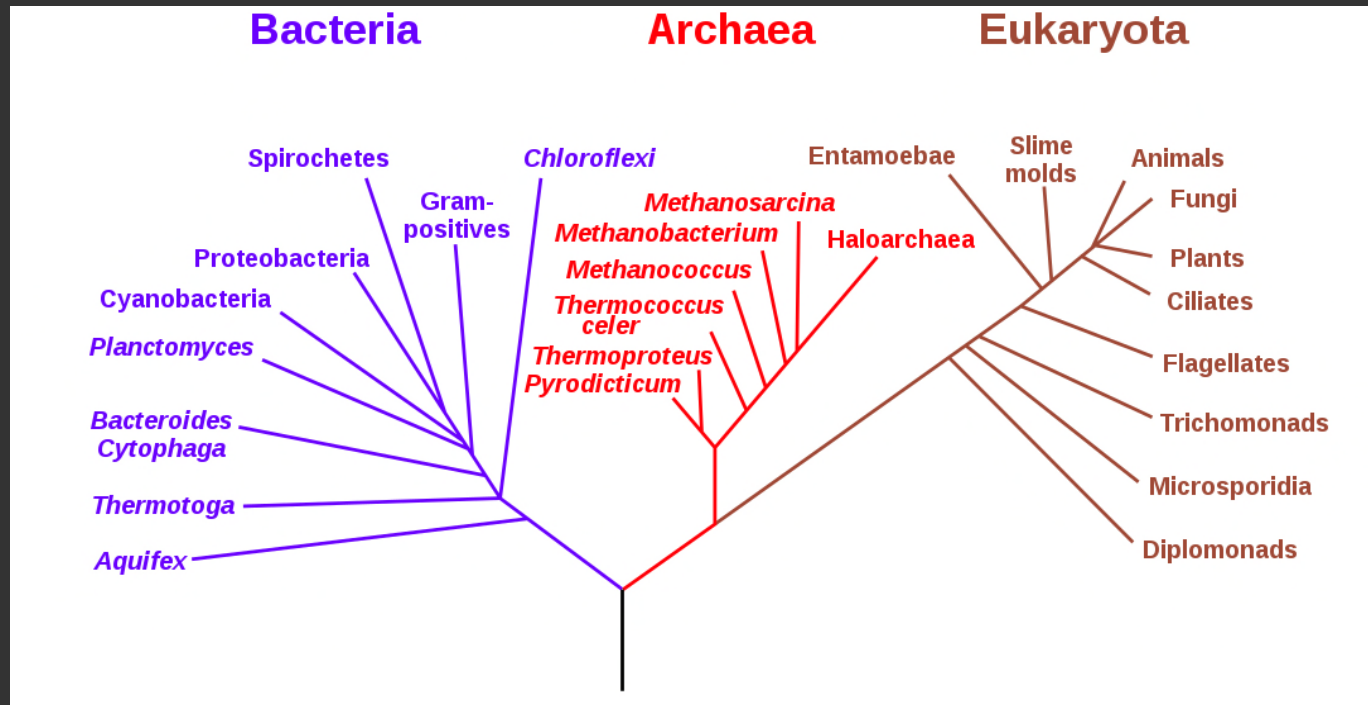
# ① Review of Evolutionary Biology

Carl von Linné 1755: Attempt to classify all living organisms on earth

# CAROLI LINNÆI    REGNUM ANIMALE.

## I. QUADRUPEDIA.
*Corpus hirsutum. Pedes quatuor. Femina vivipara, lactifera.*

| | |
|---|---|
| Homo. | Nosce te ipsum. |
| Simia. | |
| Bradypus. | |
| Ursus. | |
| Leo. | |
| Tigris. | |
| Felis. | |
| Mustela. | |
| Didelphis. | |
| Lutra. | |
| Odobenus. | |
| Phoca. | |
| Hyæna. | |
| Canis. | |
| Meles. | |
| Talpa. | |
| Erinaceus. | |
| Vespertilio. | |
| Hystrix. | |
| Sciurus. | |
| Castor. | |
| Mus. | |
| Lepus. | |
| Sorex. | |
| Equus. | |
| Hippopotamus. | |
| Elephas. | |
| Sus. | |
| Camelus. | |
| Cervus. | |
| Capra. | |
| Ovis. | |
| Bos. | |
| Ordines. | Genera. |

## II. AVES.
*Corpus plumosum. Alæ duæ. Pedes duo. Rostrum osseum. Femina ovipara.*

| | |
|---|---|
| Psittacus. | |
| Strix. | |
| Falco. | |
| Psittacus. | |
| Corvus. | |
| Cuculus. | |
| Picus. | |
| Certhia. | |
| Sitta. | |
| Upupa. | |
| Ispida. | |
| Grus. | |
| Ciconia. | |
| Ardea. | |
| Platalea. | |
| Pelecanus. | |
| Cygnus. | |
| Anas. | |
| Anser. | |
| Mergus. | |
| Graculus. | |
| Colymbus. | |
| Larus. | |
| Hæmatopus. | |
| Charadrius. | |
| Vanellus. | |
| Tringa. | |
| Numenius. | |
| Fulica. | |
| Struthio. | |
| Casuarius. | |
| Otis. | |
| Pavo. | |
| Meleagris. | |
| Gallina. | |
| Tetrao. | |
| Columba. | |
| Turdus. | |
| Sturnus. | |
| Alauda. | |
| Motacilla. | |
| Luscinia. | |
| Parus. | |
| Hirundo. | |
| Loxia. | |
| Ampelis. | |
| Fringilla. | |

## III. AMPHIBIA.
*Corpus nudum, vel squamosum. Dentes molares nulli: reliqui semper. Pinnæ nullæ.*

| | |
|---|---|
| Testudo. | |
| Rana. | |
| Lacerta. | |
| Anguis. | |

### PARADOXA.

HYDRA corpore anguino, pedibus duobus, collis totidem, & totidem capitibus, alarum expers, asservatur Hamburgi, similisque in Apocalypsi & Dan. descripta Apocal. cap. 12 & 13.

RANA-PISCIS seu RANA IN PISCEM METAMORPHOSIS.

MONOCEROS Pisces, corpore aequino, pedibus binis...

PELECANUS rostro vulnus infligens femori suo...

SATYRUS caudatus, hirsutus, barbatus, humanum...

BOROMETZ sive AGNUS SCYTHICUS...

DRACO corpore anguino, duobus pedibus, duabus alis...

AUTOMA MORTIS Horologii minimi...

## IV. PISCES.
*Corpus apodum, pinnis veris instructum, nudum, vel squamosum.*

| | |
|---|---|
| Thrichechus. | |
| Catodon. | |
| Monodon. | |
| Balæna. | |
| Delphinus. | |
| Raja. | |
| Squalus. | |
| Acipenser. | |
| Petromyzon. | |
| Lophius. | |
| Cyclopterus. | |
| Ostracion. | |
| Balistes. | |
| Gasterosteus. | |
| Zeus. | |
| Cottus. | |
| Trigla. | |
| Trachinus. | |
| Perca. | |
| Sparus. | |
| Labrus. | |
| Mugil. | |
| Scomber. | |
| Xiphias. | |
| Gobius. | |
| Gymnotus. | |
| Muræna. | |
| Blennus. | |
| Gadus. | |
| Pleuronectes. | |
| Ammodytes. | |
| Coryphæna. | |
| Echeneis. | |
| Esox. | |
| Salmo. | |
| Osmerus. | |
| Coregonus. | |
| Clupea. | |
| Cyprinus. | |
| Cobitus. | |
| Syngnathus. | |

## V. INSECTA.
*Corpus crusta ossea cutis loco tectum. Caput antennis instructum.*

| | |
|---|---|
| Blatta. | |
| Dytiscus. | |
| Meloë. | |
| Forficula. | |
| Notopeda. | |
| Mordella. | |
| Curculio. | |
| Buceros. | |
| Lucanus. | |
| Scarabæus. | |
| Dermestes. | |
| Cassida. | |
| Chrysomela. | |
| Coccinella. | |
| Gyrinus. | |
| Necydalis. | |
| Attalabus. | |
| Cantharis. | |
| Carabus. | |
| Cicindela. | |
| Lepura. | |
| Cerambyx. | |
| Buprestis. | |
| Papilio. | |
| Libellula. | |
| Ephemera. | |
| Hemerobius. | |
| Panorpa. | |
| Raphidia. | |
| Apis. | |
| Ichneumon. | |
| Musca. | |
| Gryllus. | |
| Lampyris. | |
| Formica. | |
| Cimex. | |
| Notonecta. | |
| Nepa. | |
| Scorpio. | |
| Pediculus. | |
| Pulex. | |
| Monoculus. | |
| Acarus. | |
| Araneus. | |
| Cancer. | |
| Oniscus. | |
| Scolopendra. | |

## VI. VERMES.
*Corporis Musculi ab una parte basi cuidam solidæ affixi.*

| | |
|---|---|
| Gordius. | |
| Tænia. | |
| Lumbricus. | |
| Hirudo. | |
| Limax. | |
| Cochlea. | |
| Nautilus. | |
| Cypræa. | |
| Haliotis. | |
| Patella. | |
| Dentalium. | |
| Concha. | |
| Lepas. | |
| Tethys. | |
| Echinus. | |
| Asterias. | |
| Medusa. | |
| Sepia. | |
| Microcosmus. | |

Charles Darwin 1859: This clustering also reflects the __origin__ of these species.



Source: Wiki

"Phylogenetic Tree"

Since the 1970s : Availability of _molecular data_
for determining evolutionary history

# How does this work?

All information about an organism is stored in a DNA or RNA molecule

DNA
```
... A T G C T G C C A ...
    | | | | | | | | |
... T G C A T G T C A ...
```

RNA ... G U C A U A A G ...

**A**denin
**C**ytosin
**G**uanin
**T**hymin

---

**U**racil

# Protein - Biosynthesis

Proteins

Source: Wiki

Ribosome

MESLVP...

Amino Acids

AGCUAGCUAGUCUAGUCAGUCAGCAGUC

# Simplest mode of reproduction: <u>Clones</u>

( Not a realistic
depiction! )



... AAA**A** ...

Random
point mutations

... AAA**G** ...

Evolation
by natural
selection

Conversely, given

* A number of species

* together with their genomes AAA, AAG, AGG, GGG

* and a way to <u>compare them</u>

AAA ——————— ...

AAG ——————— ... . . . What is the tree?

AGG ——————— ...

GGG ——————— ...

## Comparison

Given an alphabet $\mathcal{A}$ (here $\#\mathcal{A} = 4$) we can define the **Hamming-distance** on $\mathcal{A}^n$, $n \in \mathbb{N}$, as

$$d : \mathcal{A}^n \times \mathcal{A}^n \longrightarrow \mathbb{R} \qquad \left( \begin{array}{l} x = AAA \\ y = A\Lambda G \end{array} \longrightarrow d(x,y) = 1 \right)$$

$$(x \ , \ y) \longmapsto d(x,y)$$

$$= \#\{ i \in \{1, \dots, n\} \mid x_i \neq y_i \}$$

$\implies$ For a number of genomes of given length, we get a finite metric space $(X, d_X)$.

**Def.** A <u>tree</u> is a finite, weighted and connected graph without loops and s.t. all vertices have either degree 1 or $\geq$ 3.

A ——3——•——2——•——3—— E
         |1              |3
         B              D

with F (1) and E (3) at top right node, C (4) at bottom.

$\Longrightarrow$ We get a finite set $\{A, \ldots, F\}$ with a "tree metric"

**Question:** Is there an isometry (at least approximately)

$$(X, d_{Ham}) \xrightarrow{\sim} (X, d_{Tree}) \text{ ?}$$

**Facts:**
* Not every metric arises from a tree metric

* There are (heuristic?) algorithms to construct the tree, if possible

$$\longrightarrow \text{ "Neighbor - Joining"}$$

$$|d_{Ham}(x,y) - d_{Tree}(x,y)| \leq \frac{1}{2} \min d_{Tree}(x,y)$$

② <u>Appearance</u> of <u>Topology</u>

* Implication of the "tree-paradigm" is that species have to be reproductively <u>isolated</u>

* <u>But</u> there are counterexamples: (Already known to Darwin)



Fertile (!) hybrids of plants

(E.g. some orchids)

Many more examples of "horizontal" evolution or "reticulate events"

* "Conjugation" of bacteria & archaea



* Symbionts



Eucaryot - ancestor          Bacterium - ancestor

Eucaryotic cell with mitochondria

# Tree of Life ?

One more thing...



Bacteria    Archaea    Eukaryota

Viruses

# Reproduction of viruses

Viruses need a __host__ to replicate



\* If genome is built into gametes, it is inherited

$\Longrightarrow$ 5-8% of human genome

Also for bacteria / archaea – "_Bacteriophages_"



Source: Wiki

( Possibly responsible for the 2011 "EHEC" outbreak)

There is also mixing of genetic material between viruses

"Coinfection"  "Recombination"



( see Example )

# Instead of a tree, we need a network



Phylogenetic Tree

Phylogenetic Network

# In principle, it is possible to construct the network



**Parental Strains**

Gene A
Gene B
Gene C

Gene A
Gene B
Gene C

Gene A
Gene B
Gene C

Gene A
Gene B
Gene C

**Reassortant Strain**

Source: [2].

However, this seems to be extremely difficult in real life, due to * computational complexity * biological interpretation issues



← "Split - Network"

Hard to get relevant information

# How Persistent Homology is able to help

Philosophy : * Evolution happens in a complicated, high-dimensional space ($?$) and we went to know its structure.

* How is this structure related to relevant information about evolution, e.g. type, scale & statistics of reticulate events.

Hint: Reticulate events prevent contractibility of the network.

dim 0
dim 1
dim 2
filtration

Source: [2]

# Simplest example

Suppose
* $A = \{0, 1\}$
* Every site only mutates once ("infinite-sites assumption")

For a realistic example:

# Theorem (Chan, Carlsson, Rabadan)

$(X, d_X)$ a tree-like finite metric space, $\varepsilon \geq 0$.

Then $H_i \left( VR_\varepsilon (X, d_X) \right) = \{0\}$ for $i \geq 0$.

_____

$\implies$ Topological obstruction : $L^\infty$- norm on bars.

( Identification of "noise" by usual stability
results  —  GH - metric bounds  bottleneck - metric )

# Example: In/luenza

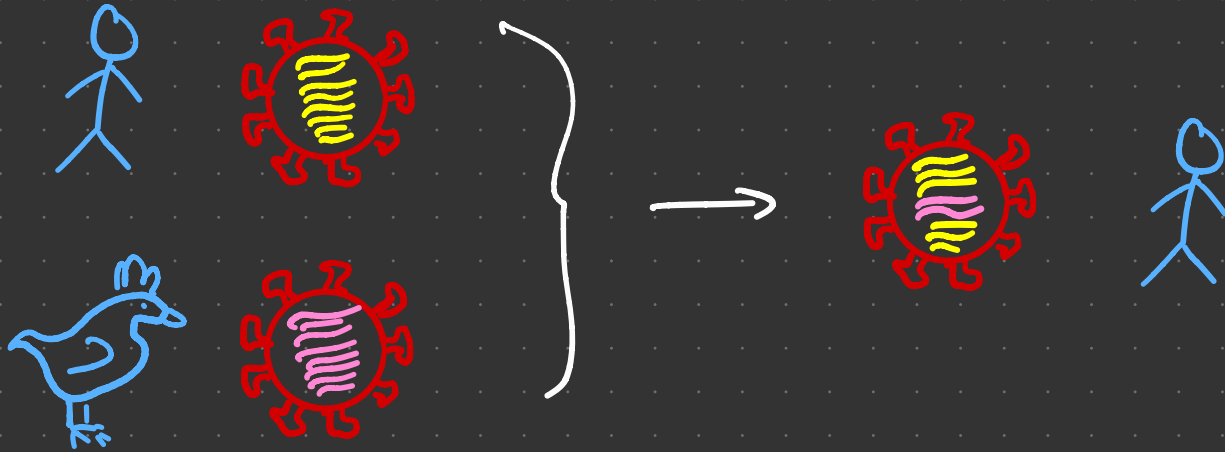* RNA virus
* 8 "segments"
* mutates very fast



Source: [2]



Source: Wiki

In principle two modes of evolution:

* vertically ~ $10^{-3}$ substitutions / nucleotide · year
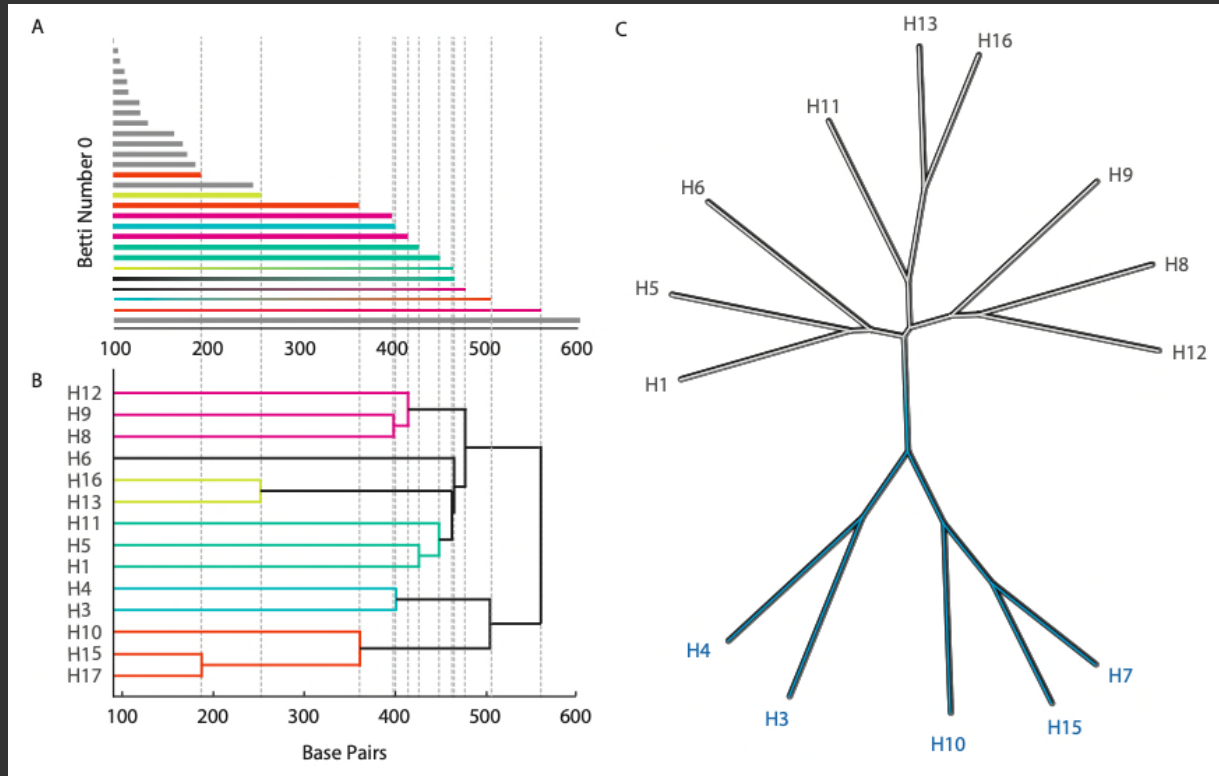* horizontally, through coinfections

Important feature: "Reassortment"

# Questions concerning vaccines etc.

* Reassortment hotspots?
* Location?
* Rate?
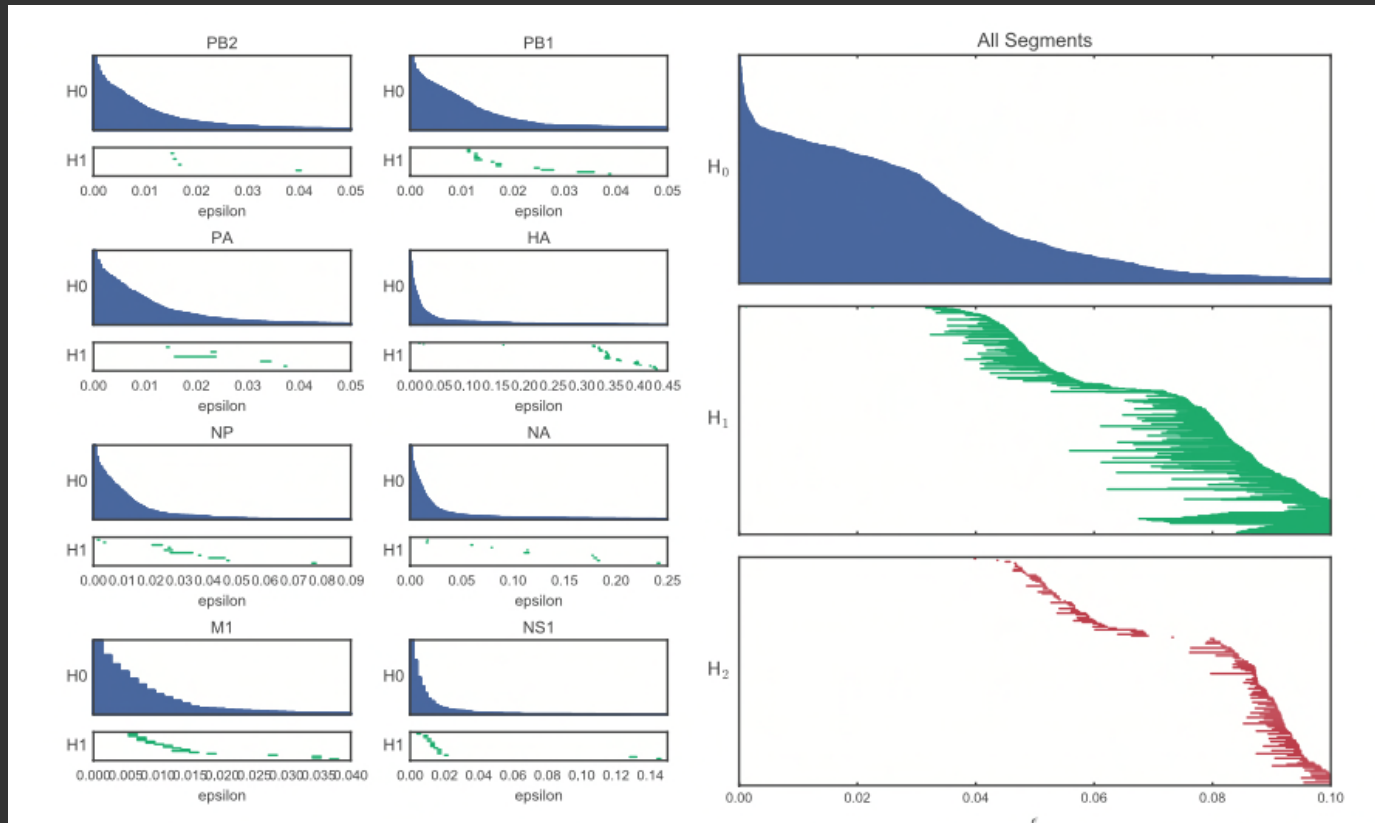* History?

$\Longrightarrow$ Use TDA

# For __one__ segment:



Source: [2]

Only relevant homology in dim. 0

→ no recombination

For the whole genome:
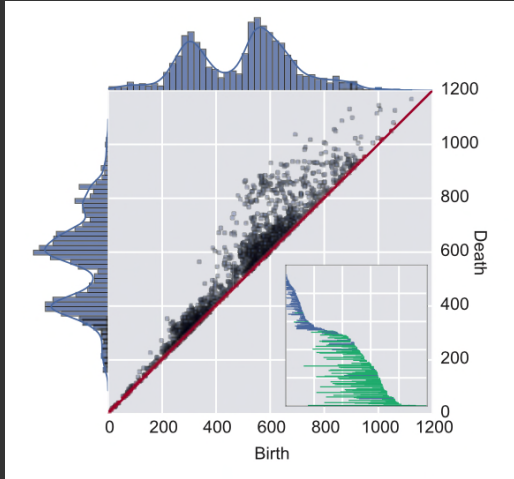
Barcodes provide useful information:

   * cycles correspond to reticulate events
     of various types

   * Statistics of cycles reveal reassortment
     hotspot preserving vital functions

   * Reassortment location and rates
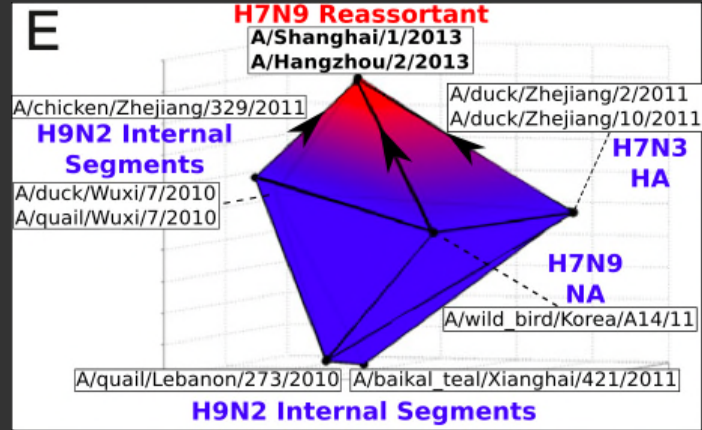

In particular, there is interesting information

... in the persistence diagram



Source: [2]

Infra - & Inter - subtype reassortments

... in higher dim. cycles



Source: [2]

Triple reassortment lead to 2013 H7N9 avian flu outbreak in China.

| Persistent Homology | Viral Evolution |
| --- | --- |
| Filtration value ε | Genetic distance (evolutionary) scale |
| 0-dimensional Betti number at filtration value ε | Number of clusters at scale ε |
| Generators of 0-dimensional homology | A representative element of the cluster |
| Hierarchical relationship among generators of 0-0-dimensional homology | Hierarchical clustering |
| 1-dimensional Betti number | Number of irreducible recombination/reassortment events |
| Generators of 1-dimensional homology | Recombinant/reassortant events |
| Generators of 2-dimensional homology | Complex horizontal genomic exchange |
| Number of higher dimensional generators in time frame | Lower bound on recombination/reassortment rate |
| Non-zero high dimensional homology (topological obstruction to phylogeny) | No phylogenetic representation |

# Summary

* Persistent homology provides a new way to think about evolution beyond trees & networks

* Barcodes and their statistics yield insights into evolutionary history on earth and pathogens

* Possible applications for vaccines, antibiotics and _pandemics_

# Thank you!