# Seasonal influenza viruses: Limited predictability of evolution & Inference of reassortment networks

Pierre Barrat-Charlaix

Team of Richard Neher
Biozentrum, University of Basel
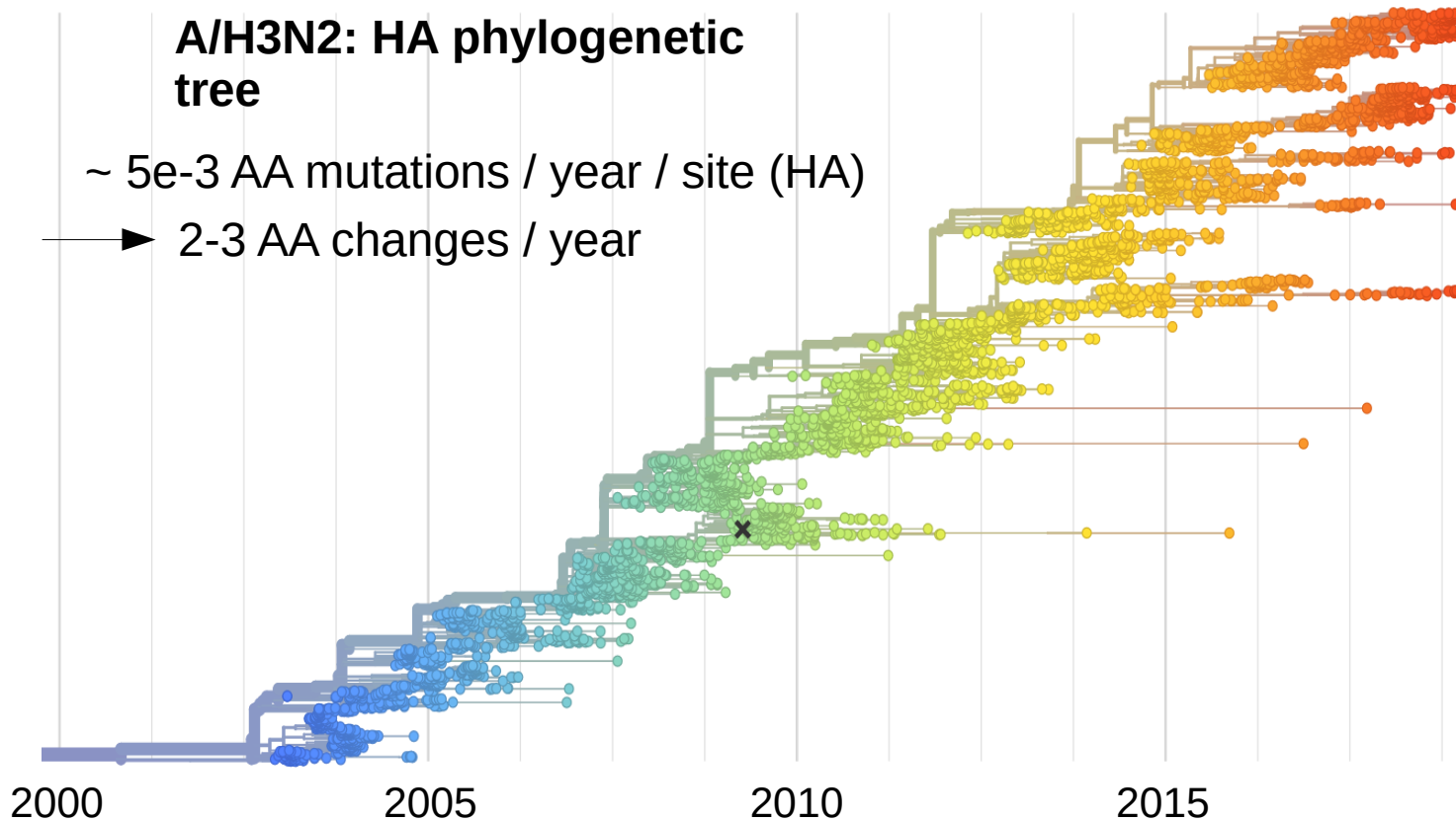
# Human seasonal influenza virus

~ hundreds of million cases per year ──────► 5-10 % of humans

In constant evolution (especially surface proteins HA & NA)

**A/H3N2: HA phylogenetic tree**

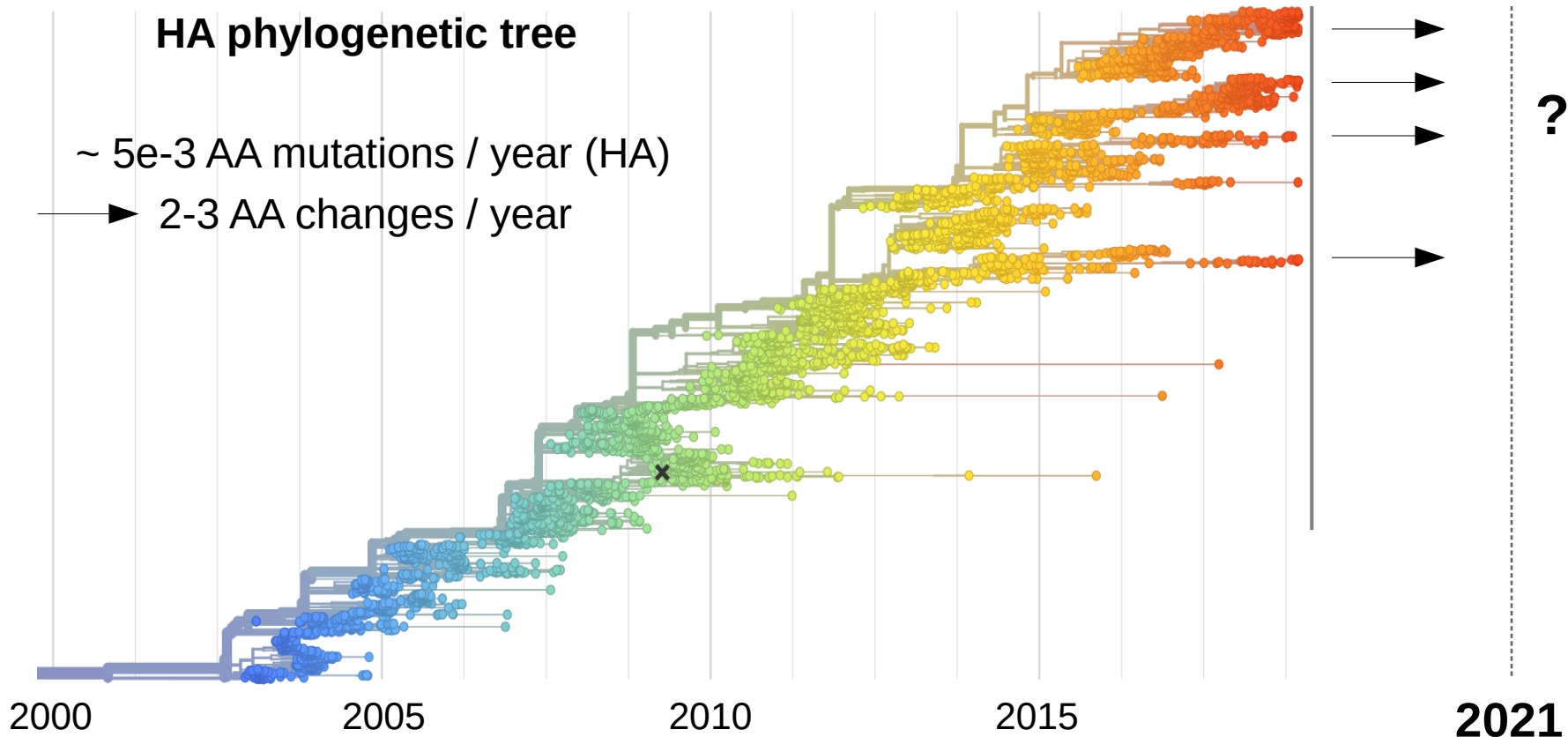~ 5e-3 AA mutations / year / site (HA)

──────► 2-3 AA changes / year

Variability in the present population

2000                    2005                    2010                    2015
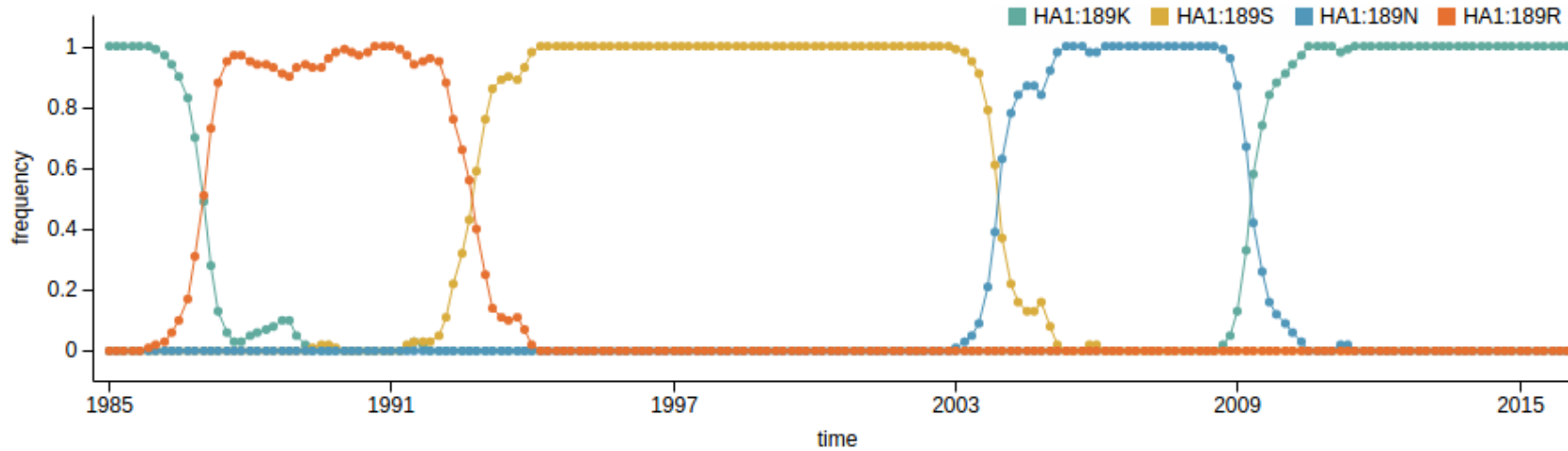
# Human seasonal influenza virus

## Can we understand/predict its evolution?

Which present clade will take over ?



**HA phylogenetic tree**

~ 5e-3 AA mutations / year (HA)

→ 2-3 AA changes / year

?

2000        2005        2010        2015        **2021**

# Selection in viral proteins



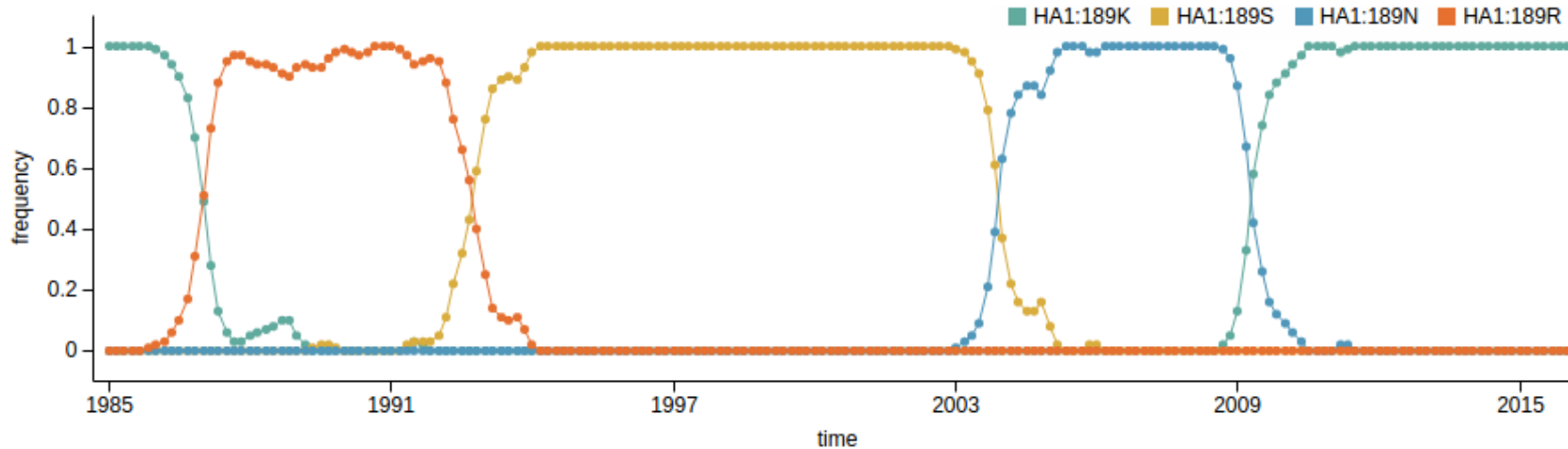Frequency of amino acid mutations at position HA1:189

Selective pressure to avoid **human immunity** → Adaptive mutations with a **fitness advantage**

# Selection in viral proteins

Frequency of amino acid mutations at position HA1:189
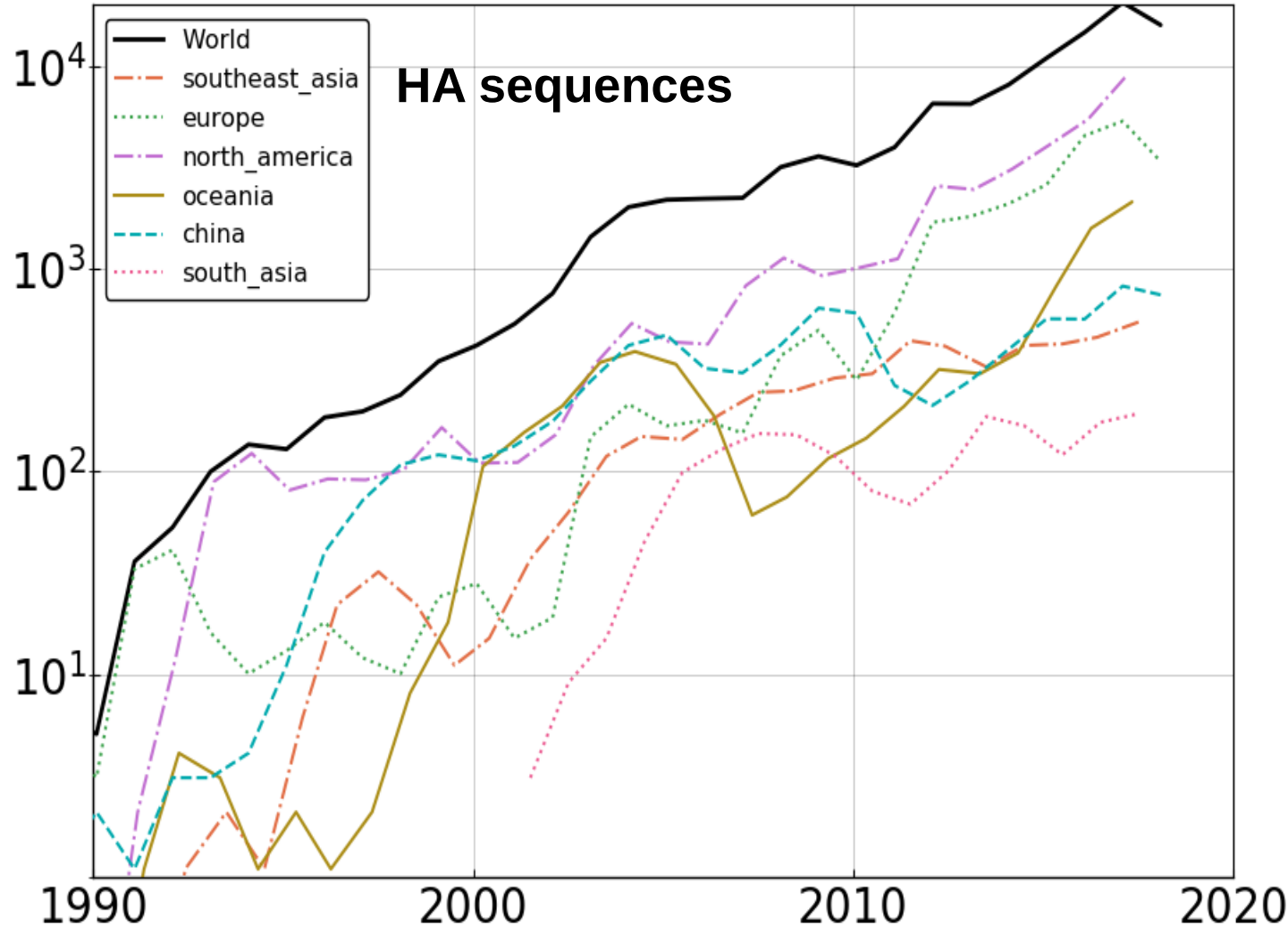


Selective pressure to avoid **human immunity** → Adaptive mutations with a **fitness advantage**

**Questions**

- Is this **way of viewing the data** correct ?
- Can we use it to **predict future evolution** ?

# Influenza pandemic: a retrospective view



# of seqs per year

**HA sequences**

Legend:
- World
- southeast_asia
- europe
- north_america
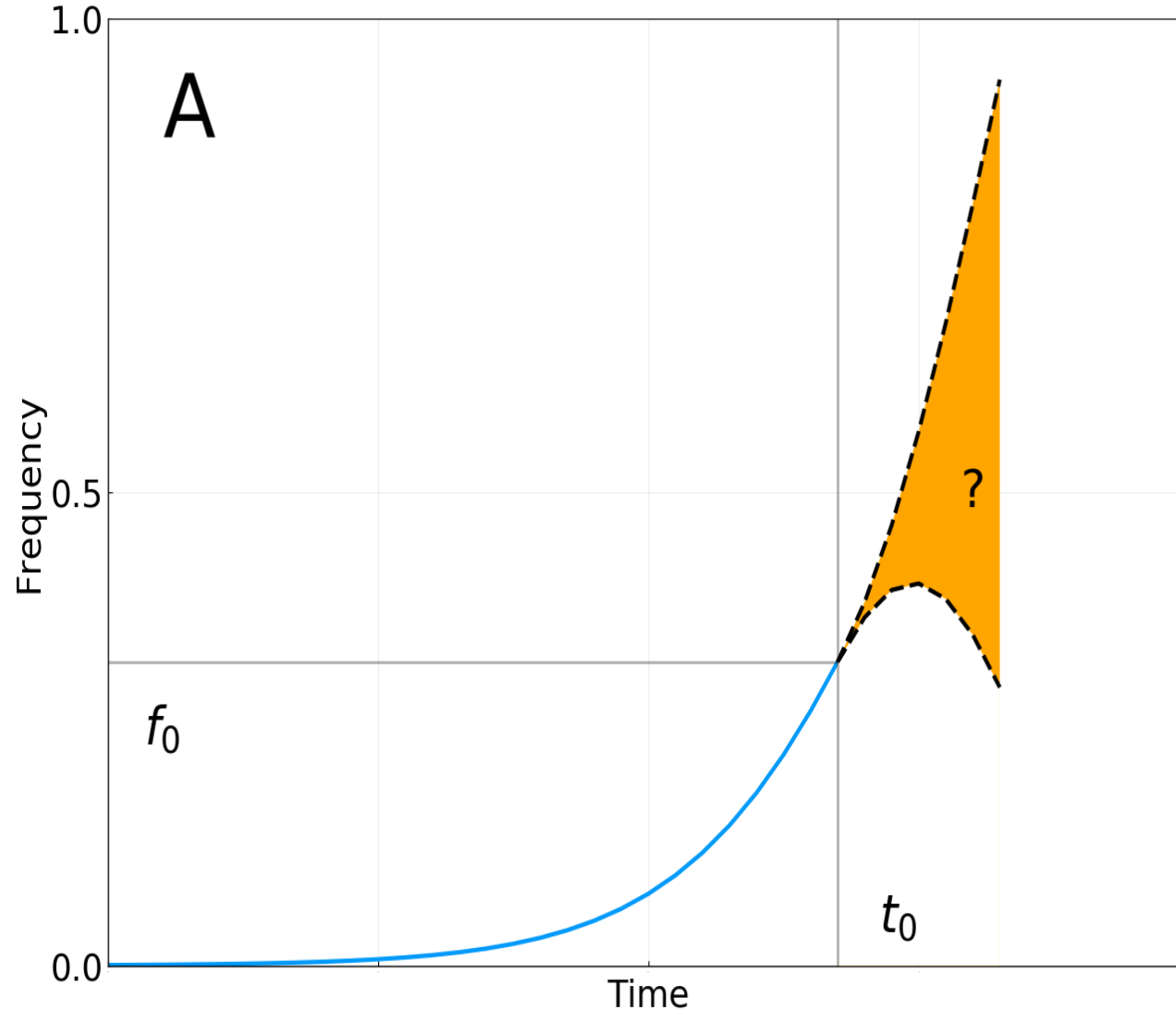- oceania
- china
- south_asia

**What can we learn from these past sequences?**

Time-binning of the past sequences by 1 month intervals

→ **Snapshots of the population**
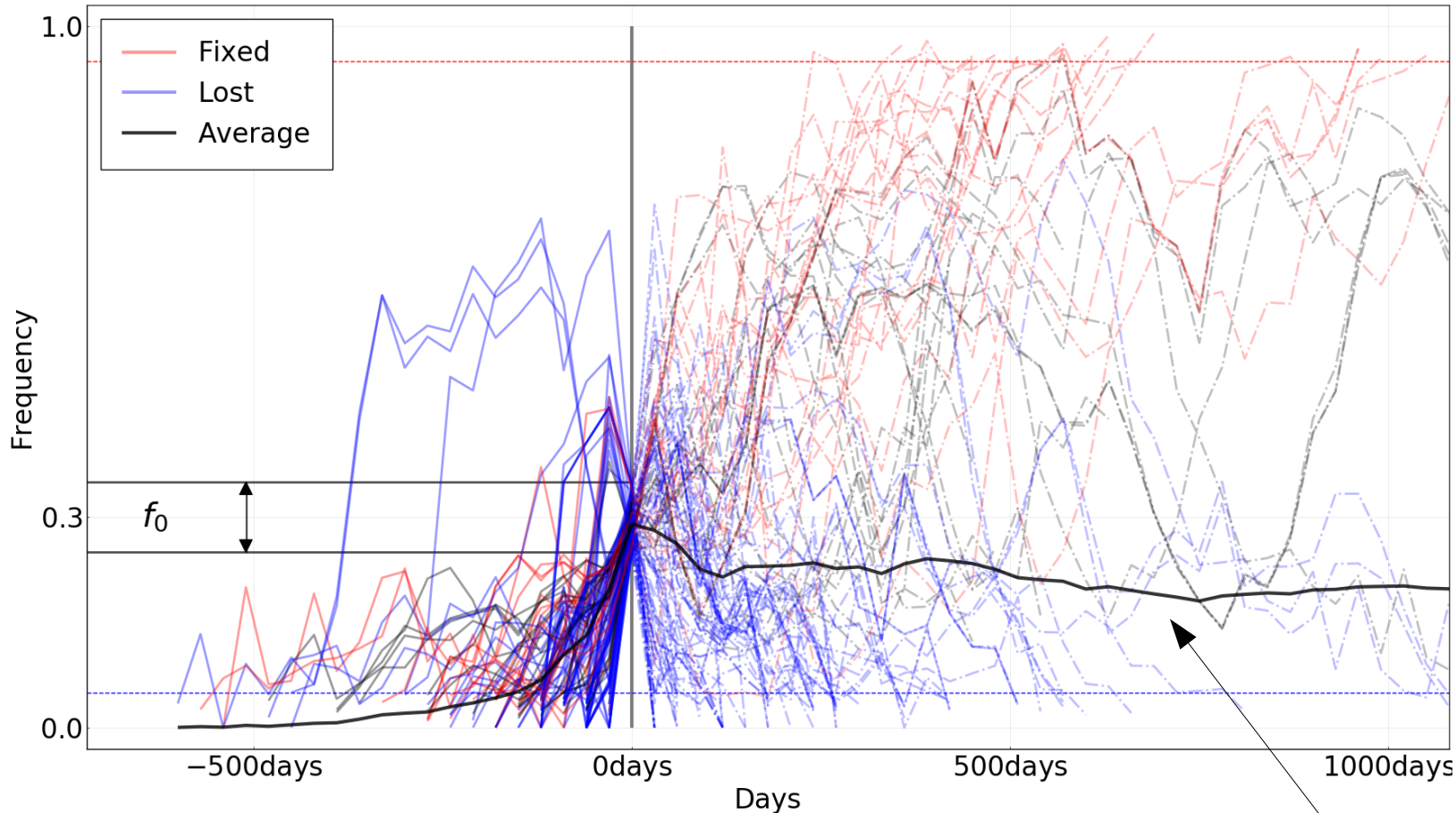
→ **Frequency trajectories**

# Short term prediction



Frequency trajectories of amino acid mutations

Frequency distribution at t0+dt?

Statistics from 460 **rising** frequency trajectories from year 2000

# Short term prediction
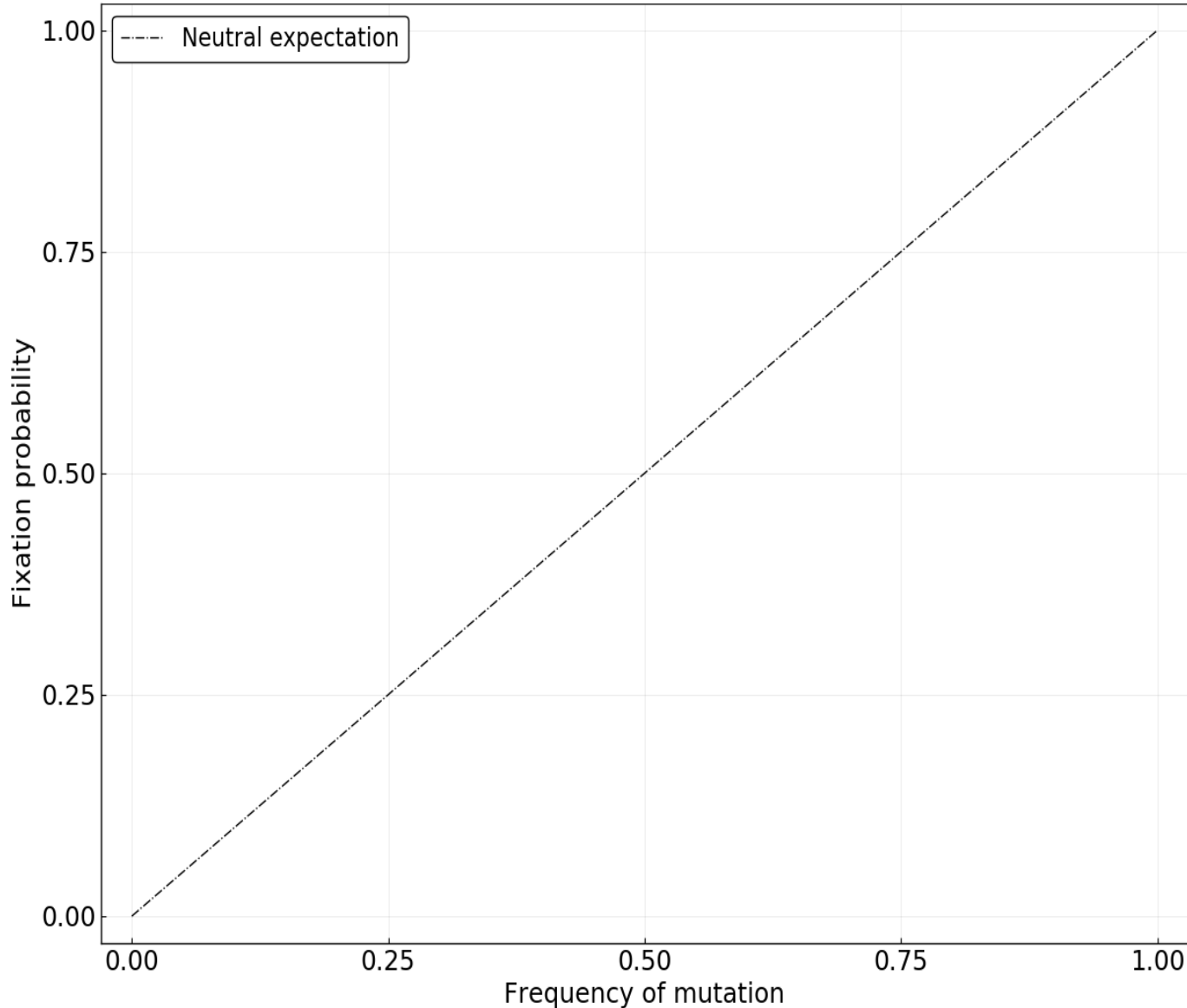
Influenza h3n2, HA protein



**Mutations:**
- Absent in the past
- Seen around f0=30%

Average trajectory is **flat** after the conditioning

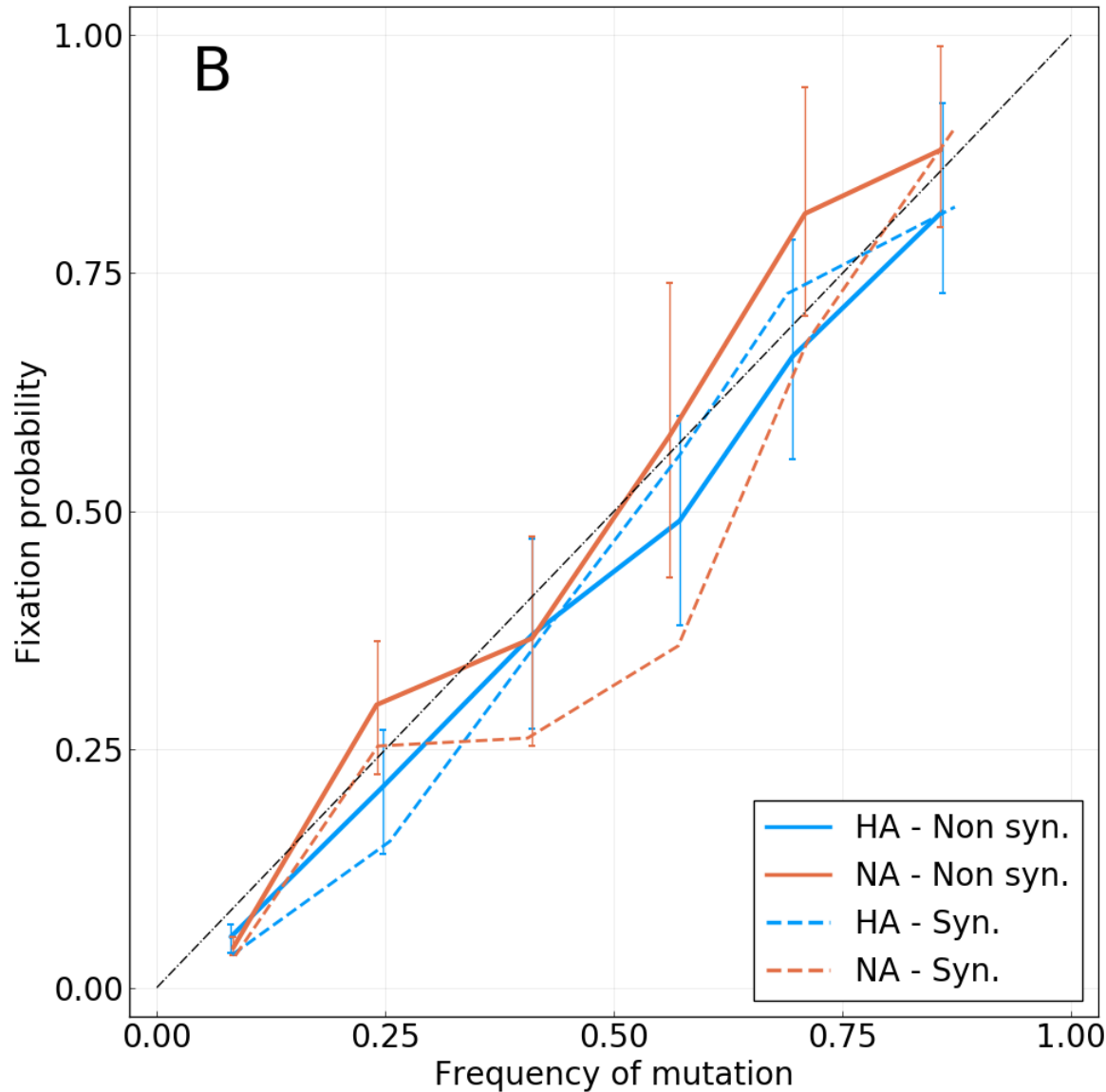**No inertia**

# Fixation probability

**Neutral evolution, *e.g.* Wright-Fisher model**

No selective advantage

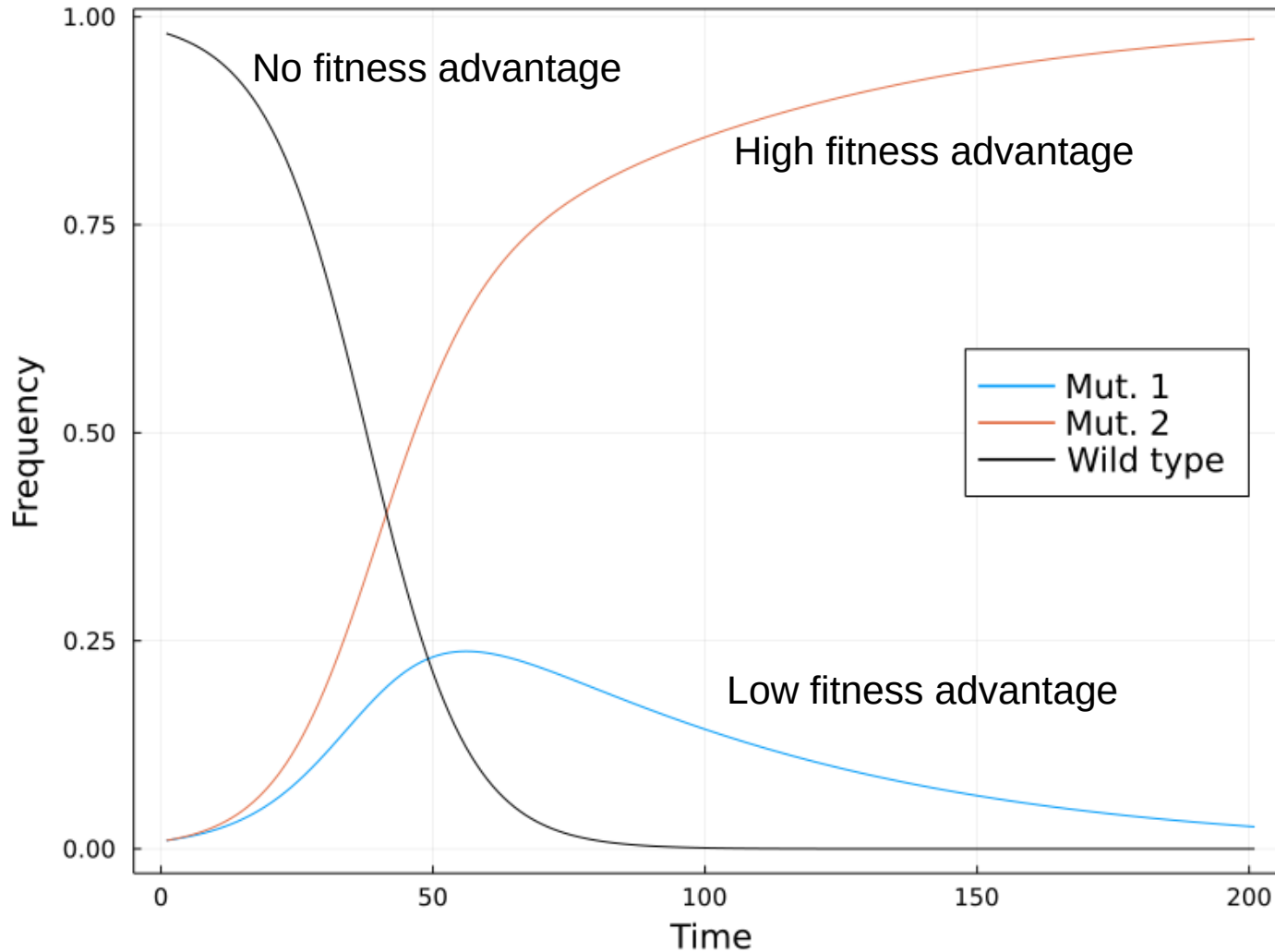Probability of fixation is equal to frequency in the population

# Fixation probability



For rising trajectories

**No signs of selection !**

The **rise in frequency** of a mutation does **not** inform us about its future **fixation**

# Is this expected? Clonal interference



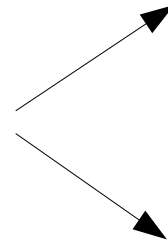Adaptive mutations appearing on different individuals

**Competition**

Figure labels:
- No fitness advantage
- High fitness advantage
- Low fitness advantage
- Frequency (y-axis)
- Time (x-axis)

Legend:
- Mut. 1
- Mut. 2
- Wild type

# Genetic linkage: toy model

Simple fitness lanscape $\qquad f(\vec{s}) = \sum_{i=1}^{L} h_i s_i$

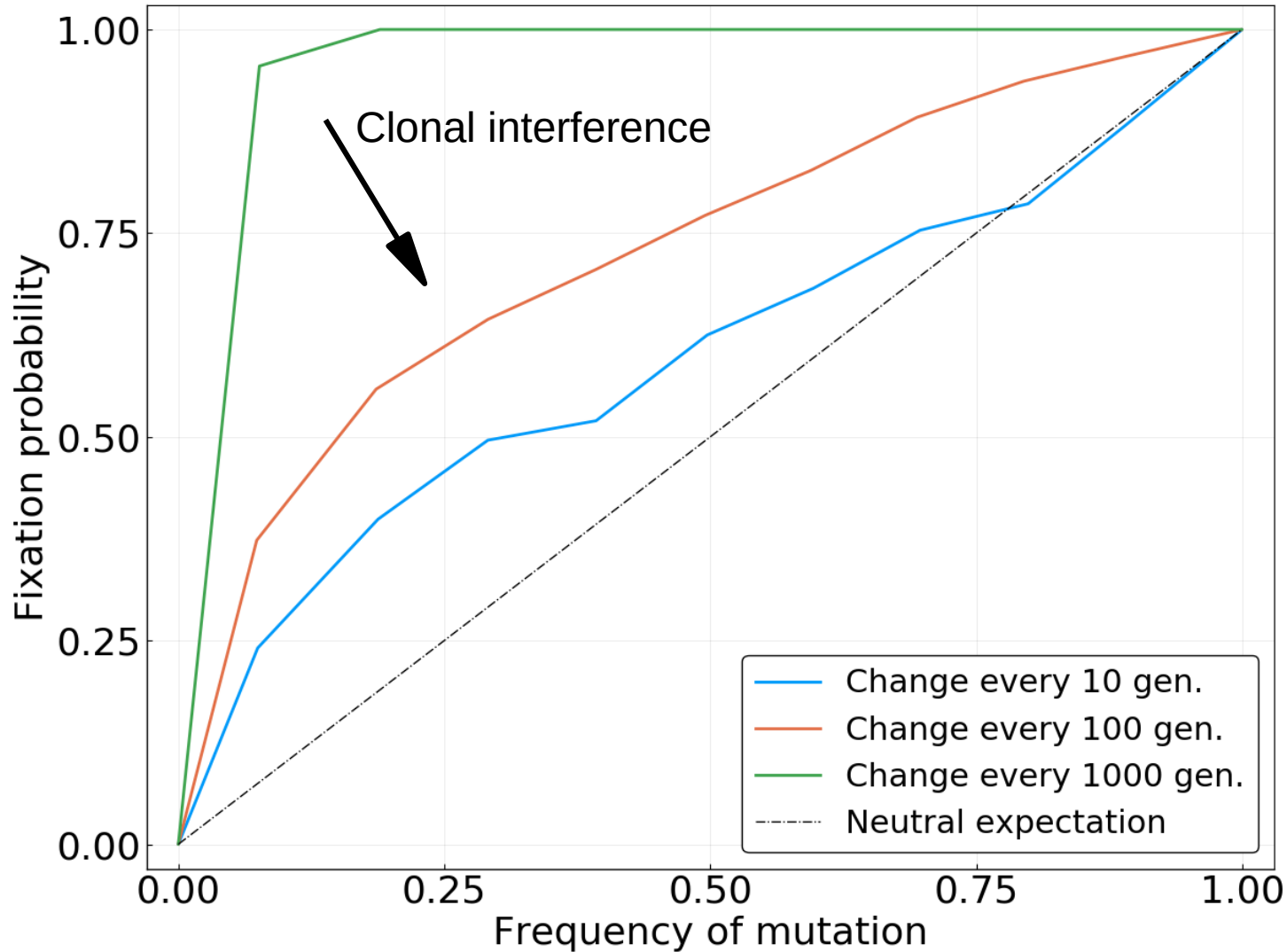Change the fitness landscape periodically

Slow rate of change
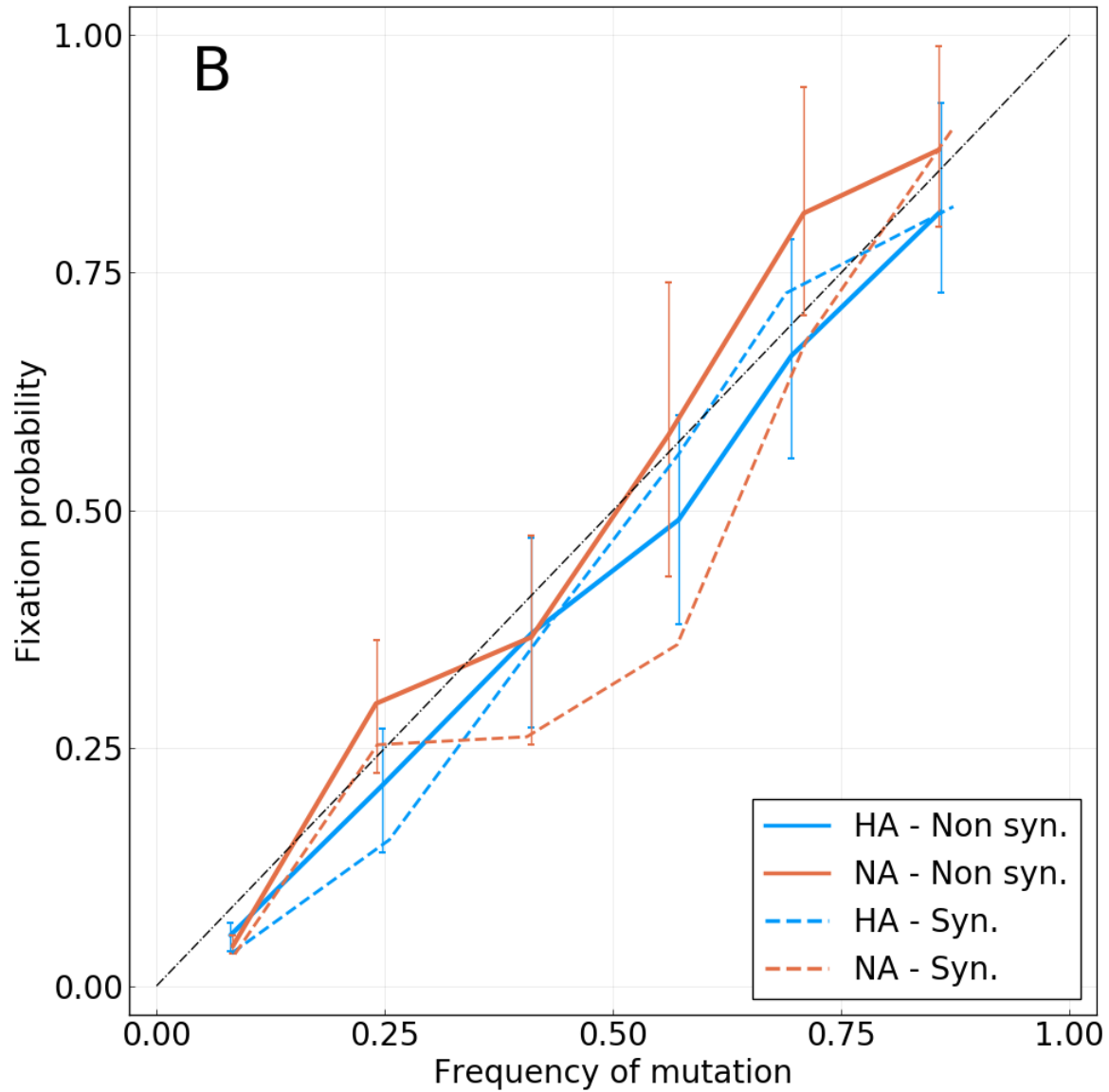**Clean sweeps**

High rate of change
**Clonal interference**

# Genetic linkage: toy model

Sweep time ~400 generations

(vs ~3 years for flu)

**It's hard to mimic neutrality!**



Clonal interference

# Fixation probability

# Predictors of fixation?

# Predictors of fixation?

## Epitope positions

- Targeted by human immune system
- Expected to be under strong selection
- Used in models of selection in influenza

But often ascertained *post-hoc*

**Shih *et. al.*** 2007
**Koel *et. al.*** 2013
**Luksza & Lässig** 2014
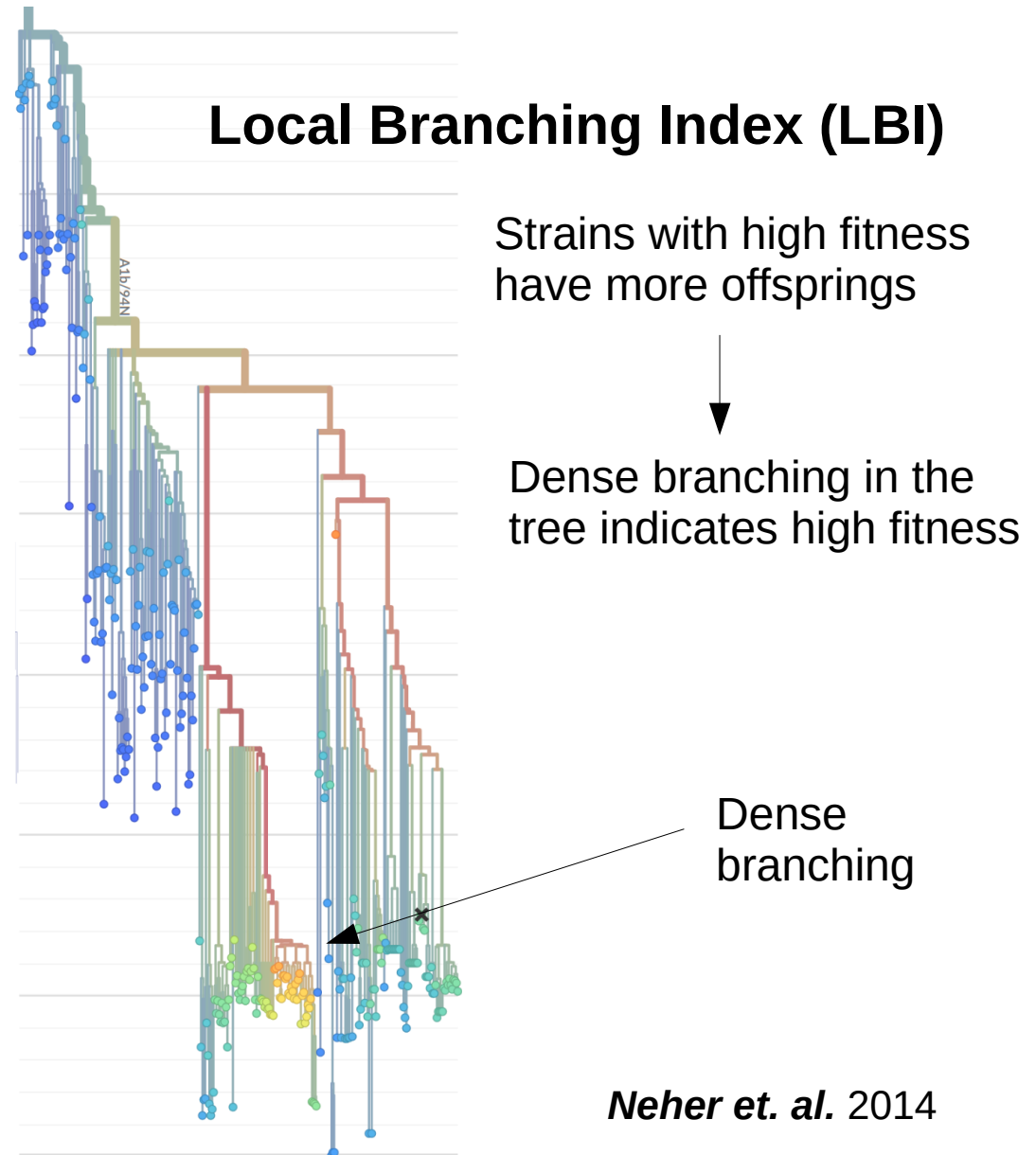
# Predictors of fixation?

## Epitope positions

- Targeted by human immune system
- Expected to be under strong selection
- Used in models of selection in influenza

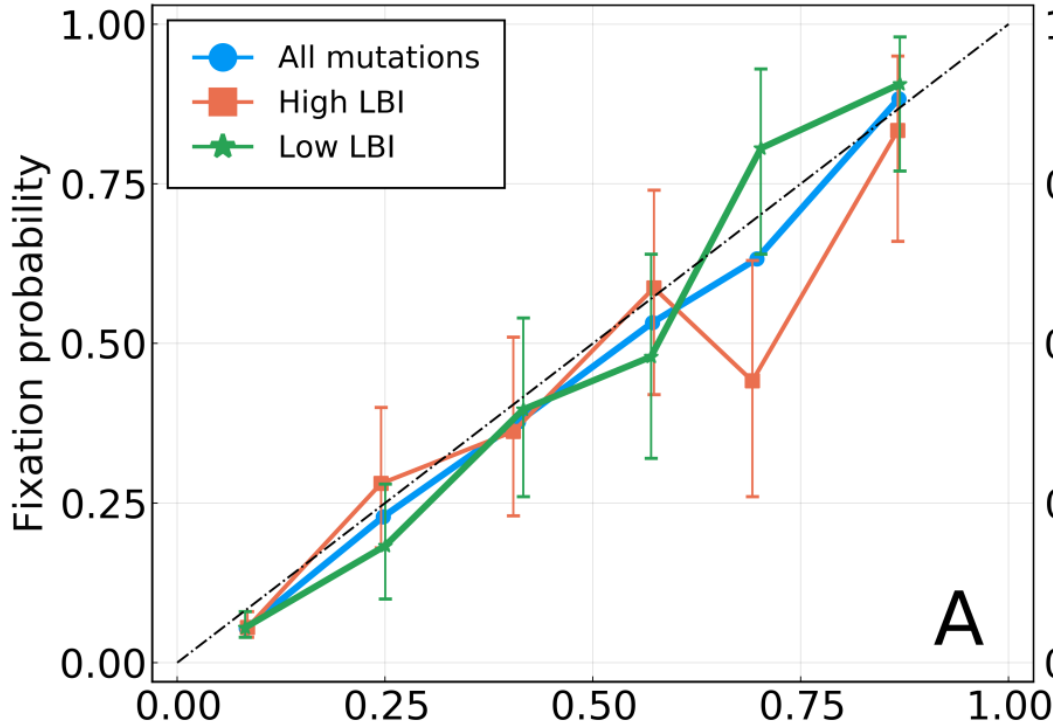But often ascertained *post-hoc*

**Shih** *et. al.* 2007
**Koel** *et. al.* 2013
**Luksza & Lässig** 2014

## Local Branching Index (LBI)

Strains with high fitness have more offsprings

↓

Dense branching in the tree indicates high fitness

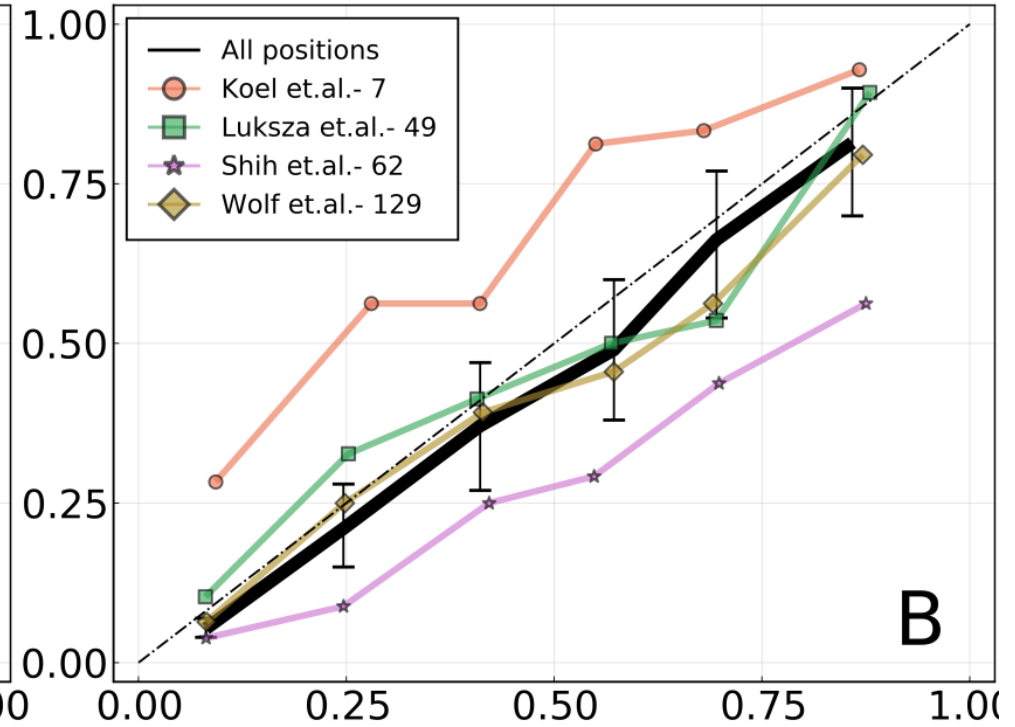Dense branching

*Neher et. al.* 2014

# Fixation probability: for specific mutations?



**Local Branching Index (LBI) Measure of fitness**

**Epitope positions Potential targets of ABs**

➤ **Current models do not predict fixation!**

# Summary

A/H3N2 influenza is under adaptive selection, but …

- Predictibility of frequency trajectories is low
- Fixation probability is equal to present frequency

⟶ « Apparent neutrality »
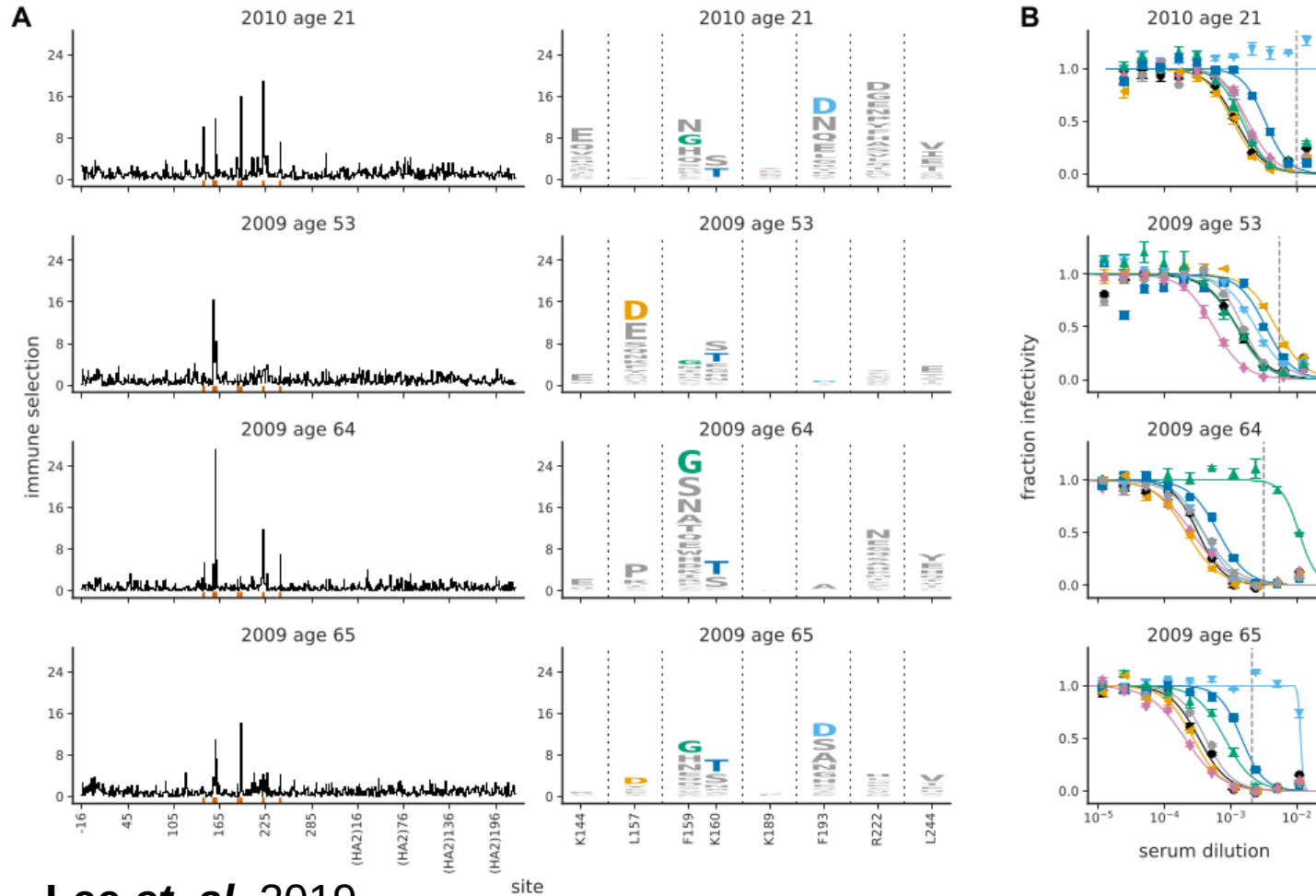
- Hard to find predictors of fixation / fitness

➔ LBI

➔ epitopes

⟶ **Influenza does not behave like models suggest !**

# Open question: what could explain these results?

# Open question: what could explain these results?

**Diversity in human immune response** →

- Adaptive mutations only allow escape to a fraction of the population

- Fitness advantage expires before fixation



Lee *et. al.* 2019

Can this result in **apparent neutrality** ?

# Epidemiological considerations

**Influenza is a seasonal virus**

⟶ In temperate regions : exponential increase (winter) followed by bottleneck
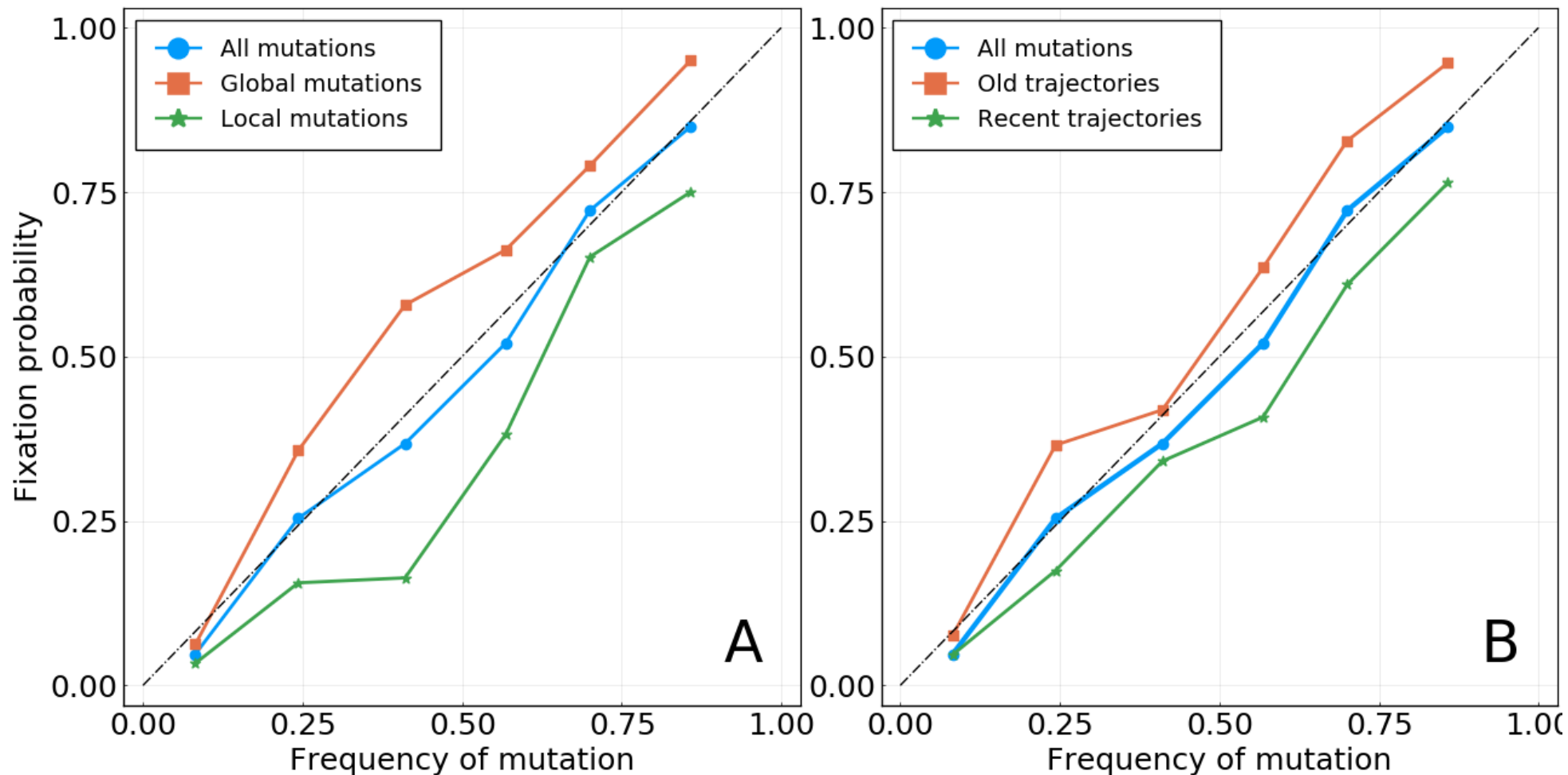
**Geographical structure**

⟶ Frequency of variants varies in different regions

⟶ **What does this mean for frequency trajectories ?**

**To be investigated...**

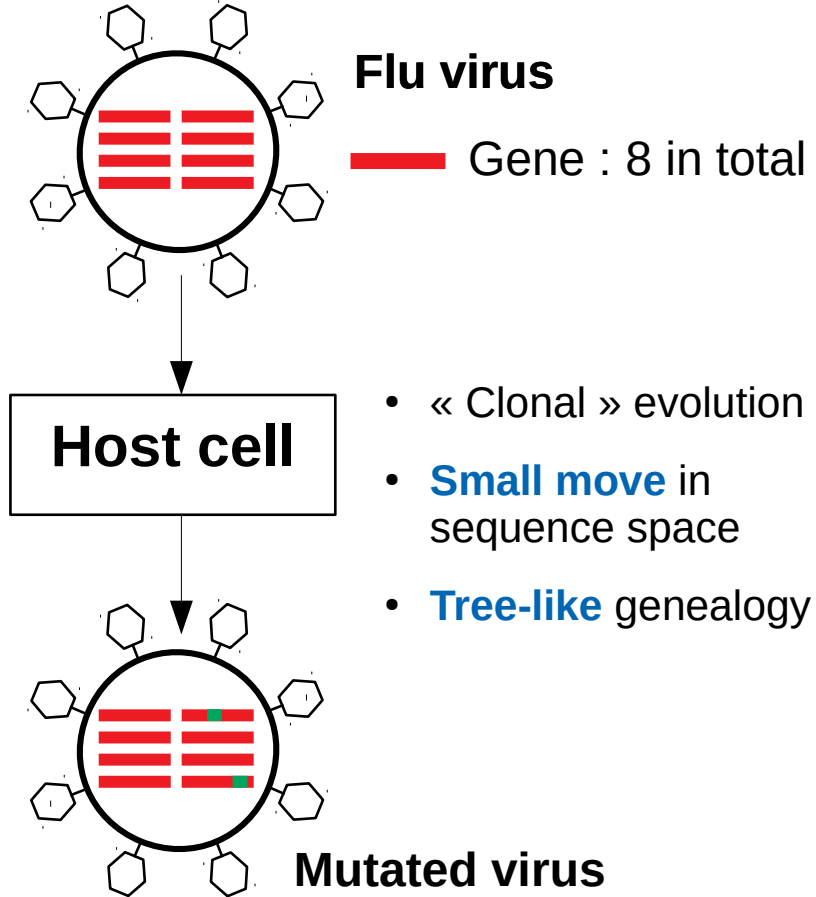# Epidemiological considerations

**For A/H3N2 - HA**

# Seasonal influenza viruses:
## Limited predictability of evolution
## & Inference of reassortment networks

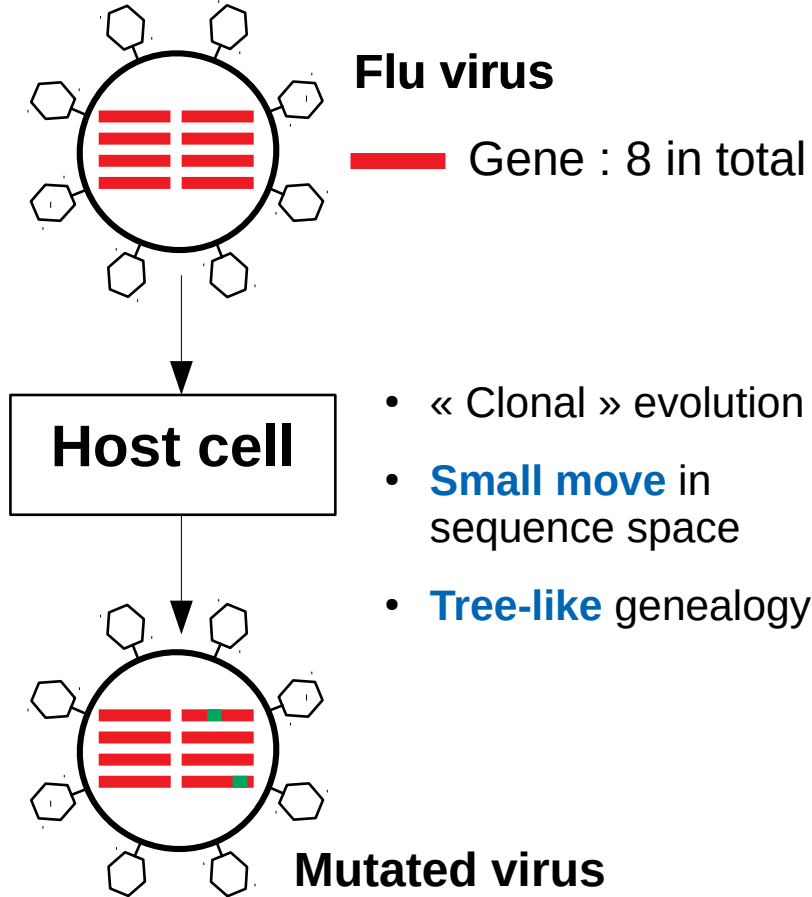Pierre Barrat-Charlaix
Biozentrum, University of Basel

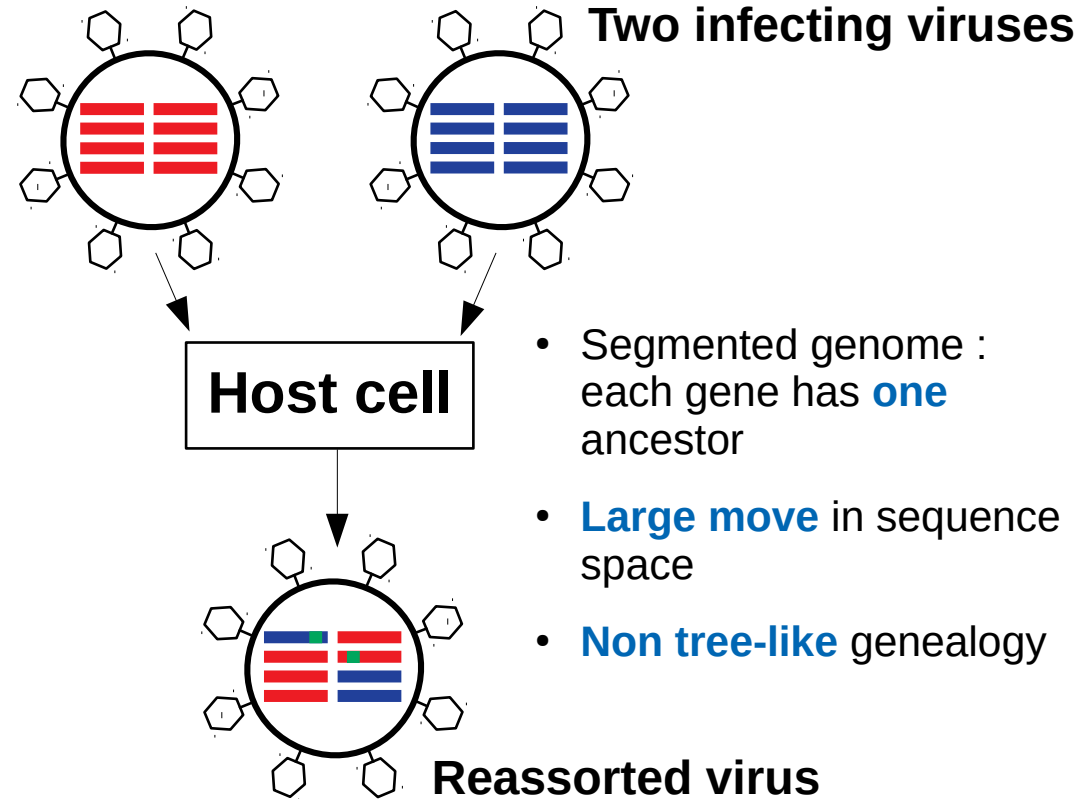# Evolution of influenza: Mutations and reassortment

## Mutation



**Flu virus**

━━━ Gene : 8 in total

**Host cell**

- « Clonal » evolution
- **Small move** in sequence space
- **Tree-like** genealogy

**Mutated virus**

# Evolution of influenza: Mutations and reassortment

## Mutation



**Flu virus**

▬ Gene : 8 in total

**Host cell**

- « Clonal » evolution
- **Small move** in sequence space
- **Tree-like** genealogy

**Mutated virus**

## Reassortment

**Two infecting viruses**

**Host cell**

- Segmented genome : each gene has **one** ancestor
- **Large move** in sequence space
- **Non tree-like** genealogy

**Reassorted virus**

# Reassortment in influenza

- Combines strains from **different subtypes**, or from **human/animal** hosts.

- Origin of many **pandemics**
  - Asian flu – 1957
  - Hong Kong flu – 1968
  - H1N1 pandemic – 2009

- Also happens at "smaller" scale: within a subtype.

- How often does it happen?
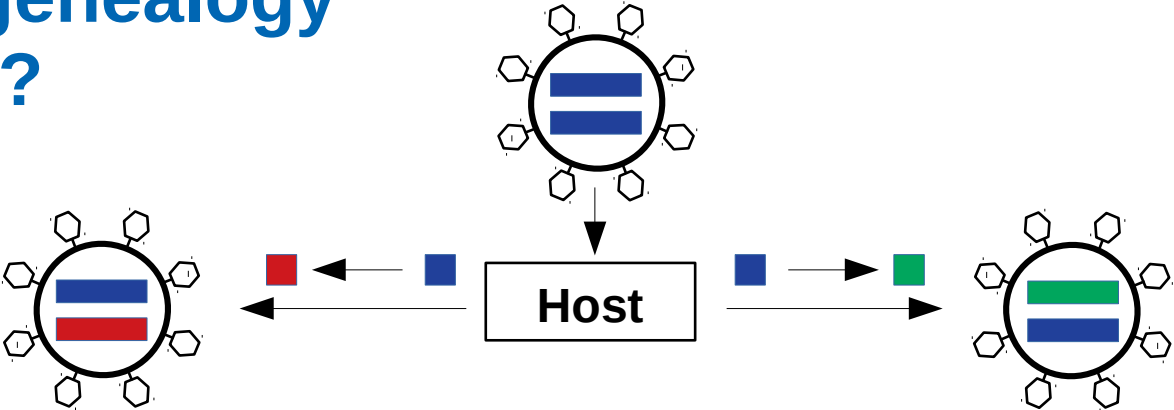- Contribution to immune escape and adaptation?

**?**

**Reassortment**



**Two infecting viruses**

**Host cell**

**Reassorted virus**

- Segmented genome : each gene has **one** ancestor

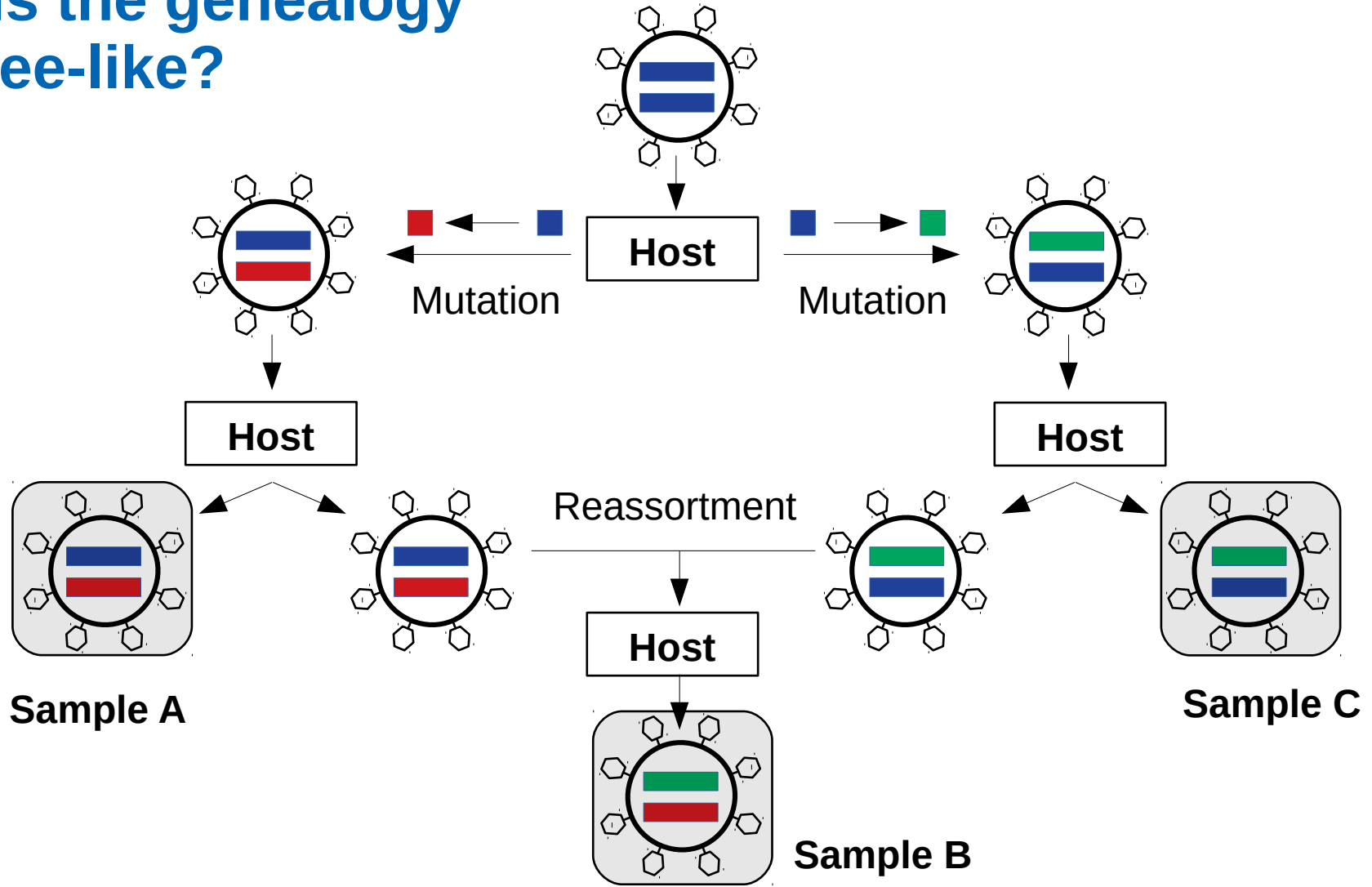- **Large move** in sequence space

- **Non tree-like** genealogy

**Reassortments are hard to infer from sequences!**

# Why is the genealogy not tree-like?

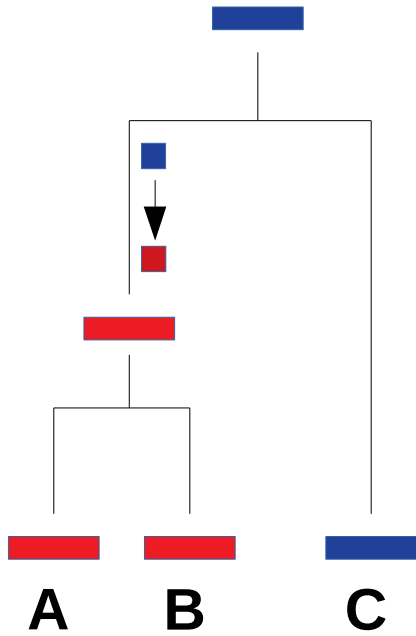# Why is the genealogy not tree-like?
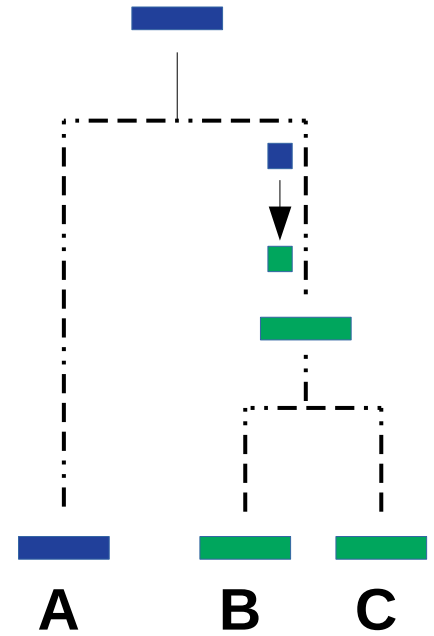
# Ancestral Reassortment Graph

**Observed sequences**

Reconstructed segment trees

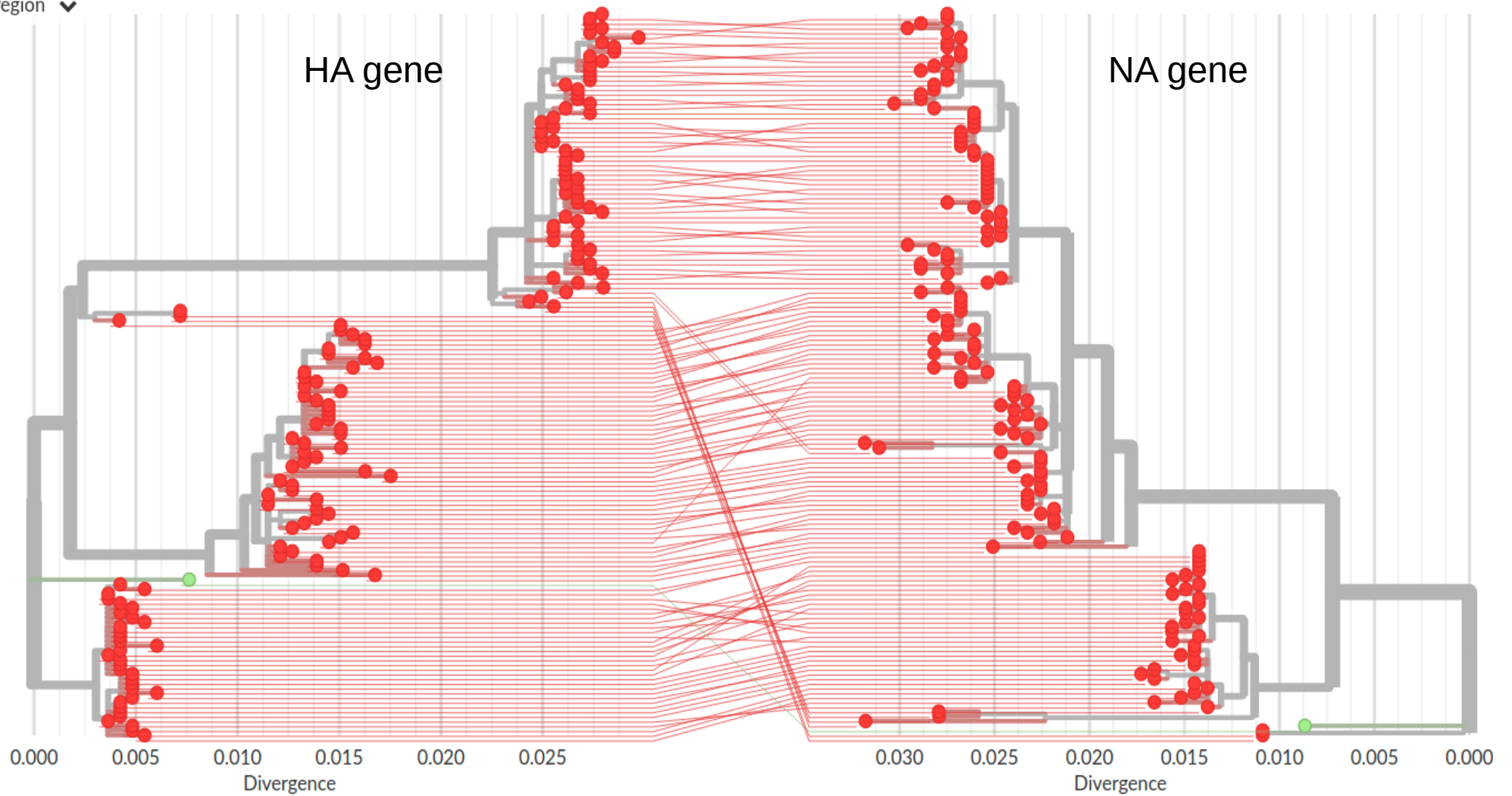Topological differences

Genealogy of first gene

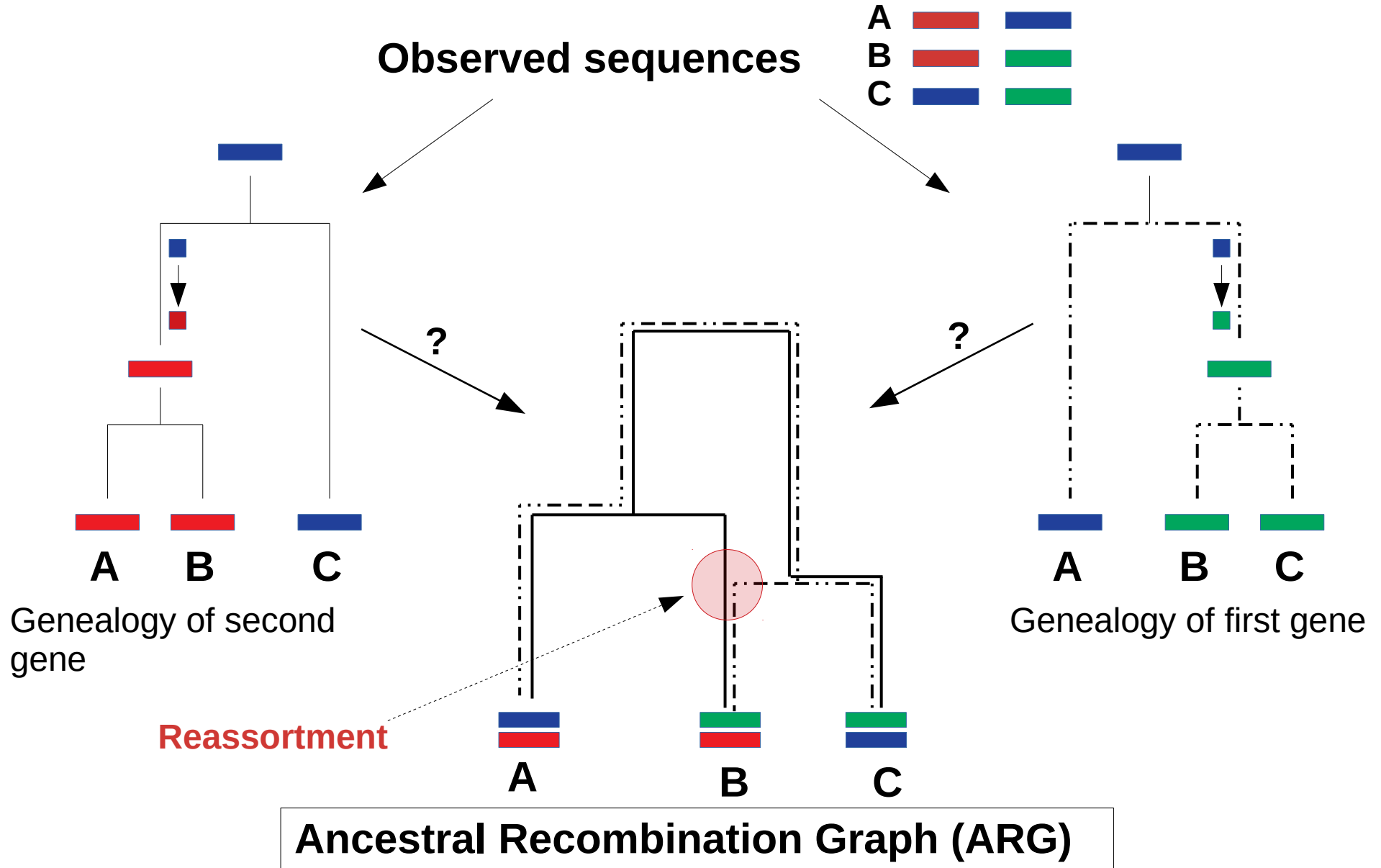Genealogy of second gene

# Example of flu trees

# Why is the genealogy not tree-like?

**Observed sequences**

Genealogy of second gene

Genealogy of first gene

**Reassortment**

**Ancestral Recombination Graph (ARG)**

# Inferring reassortments / Reconstructing the ARG

**Existing methods**

- Manual inspection of trees
  (*e.g.* **[Holmes et. al. 2005], [Boni et. al. 2010]**)

- Methods based on genetic distance **[Rabadan et. al. 2008]**

- Trees + mutation methods [**Villa & Lässig 2017**]

- Tree topology based methods **[Nagarajan & Kingsford 2011]**

→ Finds a subset of reassortment events

- Maximum likelihood methods **[Müller et. al. 2020]** → Accurate but slow
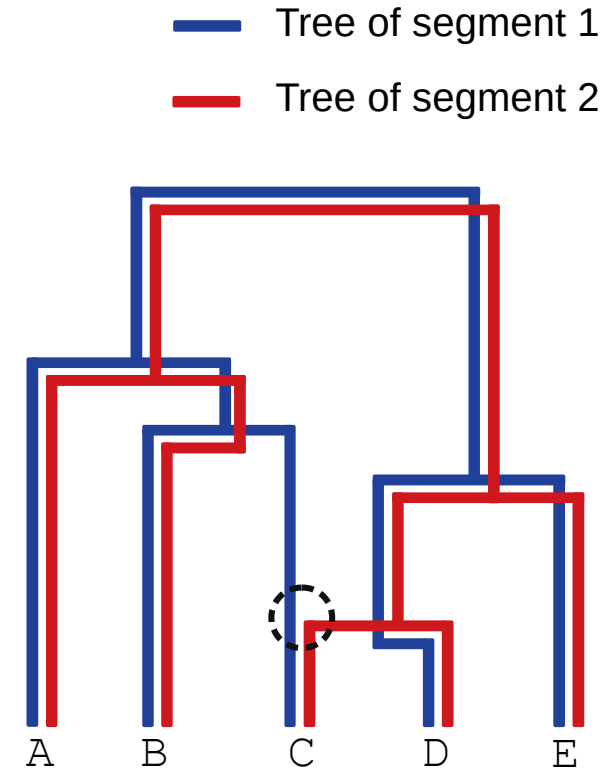
→ **No "reference" method**

We want something that is

- **Fast** : can be easily applied to new sequences
- Finds **all reassortments**, and not only large obvious ones
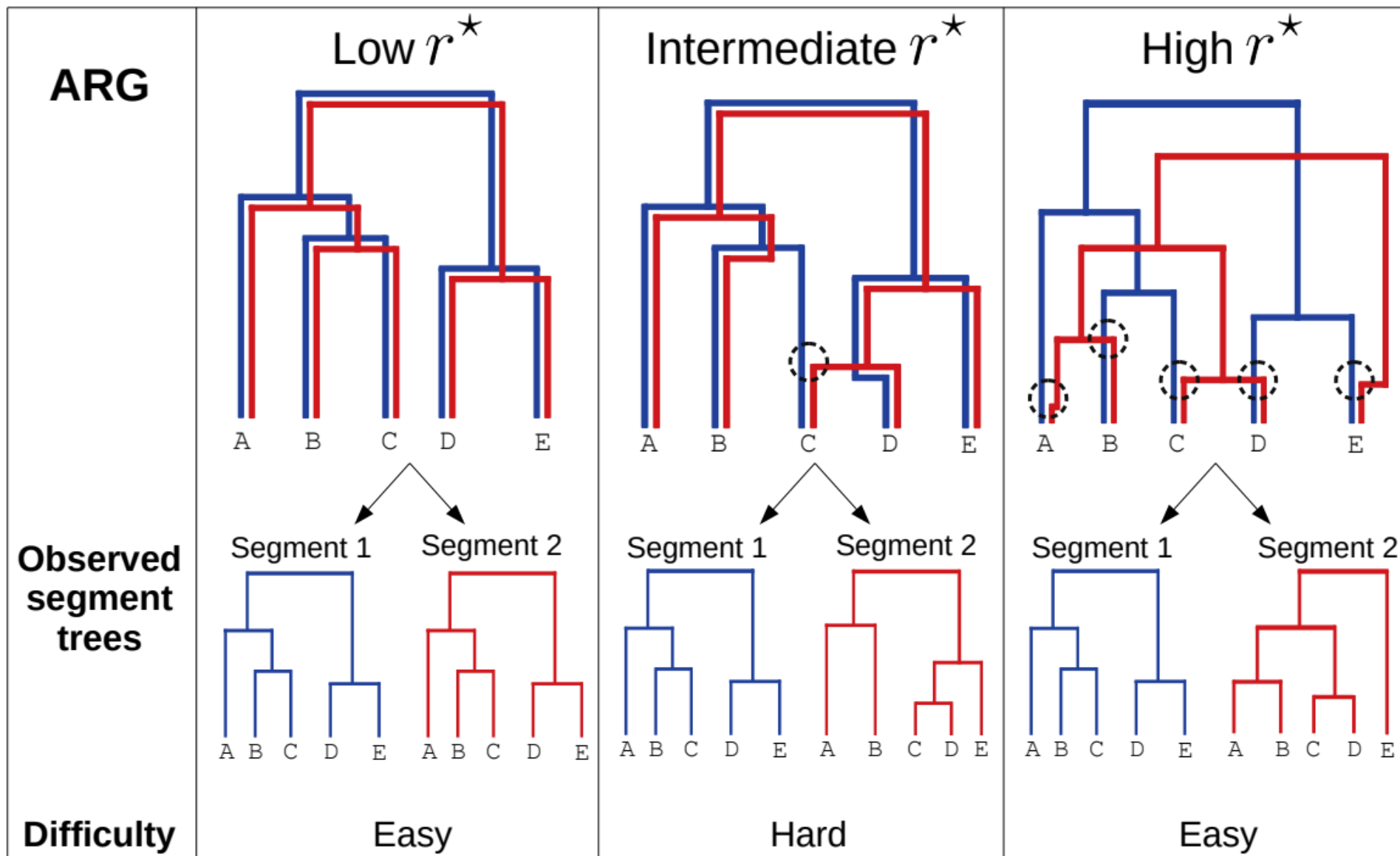- Works for the 2-genes case (simplicity)

# Inferring the ARG: the Treeknit method

We want something that is

- **Fast** : can be easily applied to new sequences
- Finds **all reassortments**, and not only large obvious ones
- Works for the 2-genes case (simplicity)

Main idea :

- The ARG is a **collage of gene trees**
- We can **infer each tree** from sequences (iqtree, RaxML, …)
- **Topological differences** between these trees are **due to reassortment**

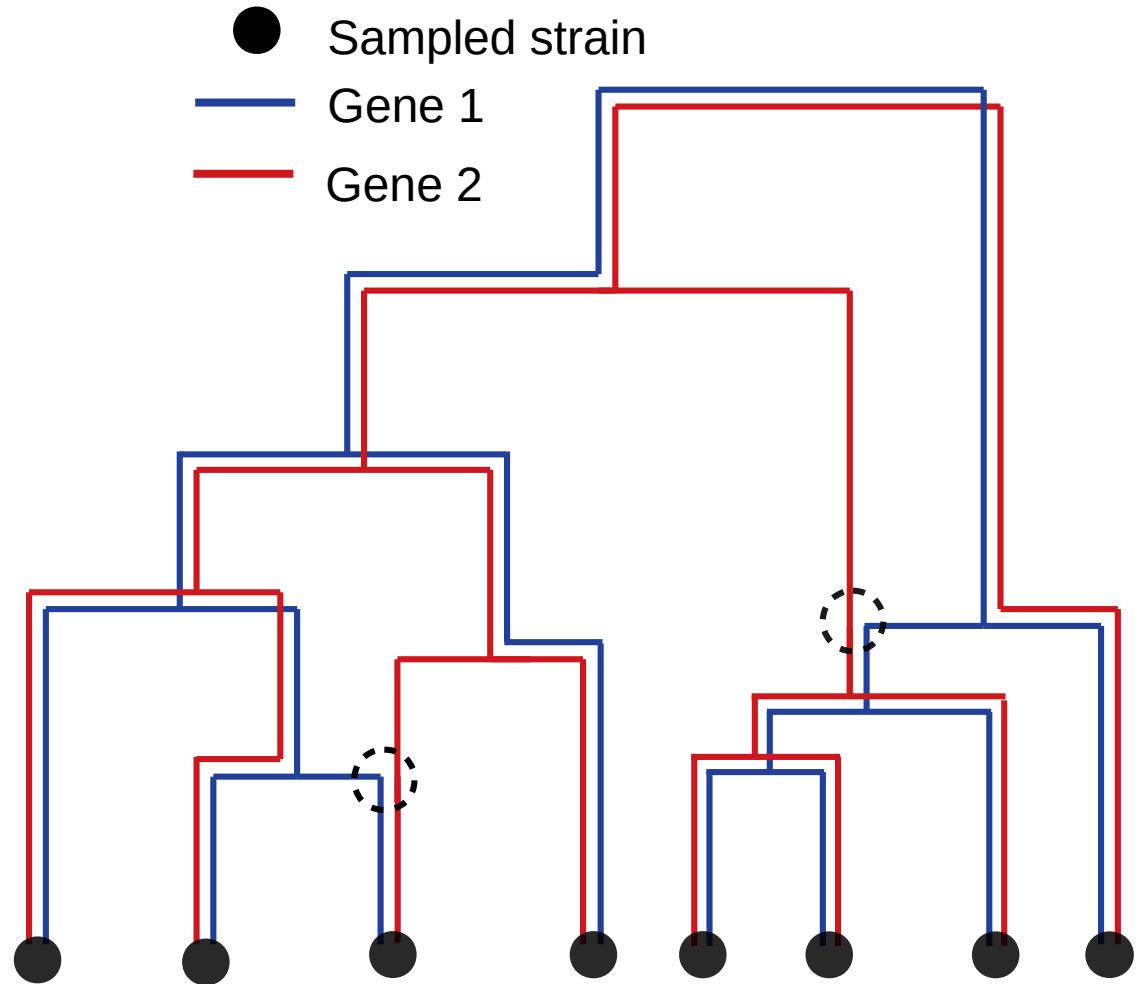⟶ **Method based on topological differences between trees**

— Tree of segment 1
— Tree of segment 2

A  B  C  D  E

# Inferring the ARG

# Maximally compatible clades (MCCs)

The ARG is a **collage of gene trees**
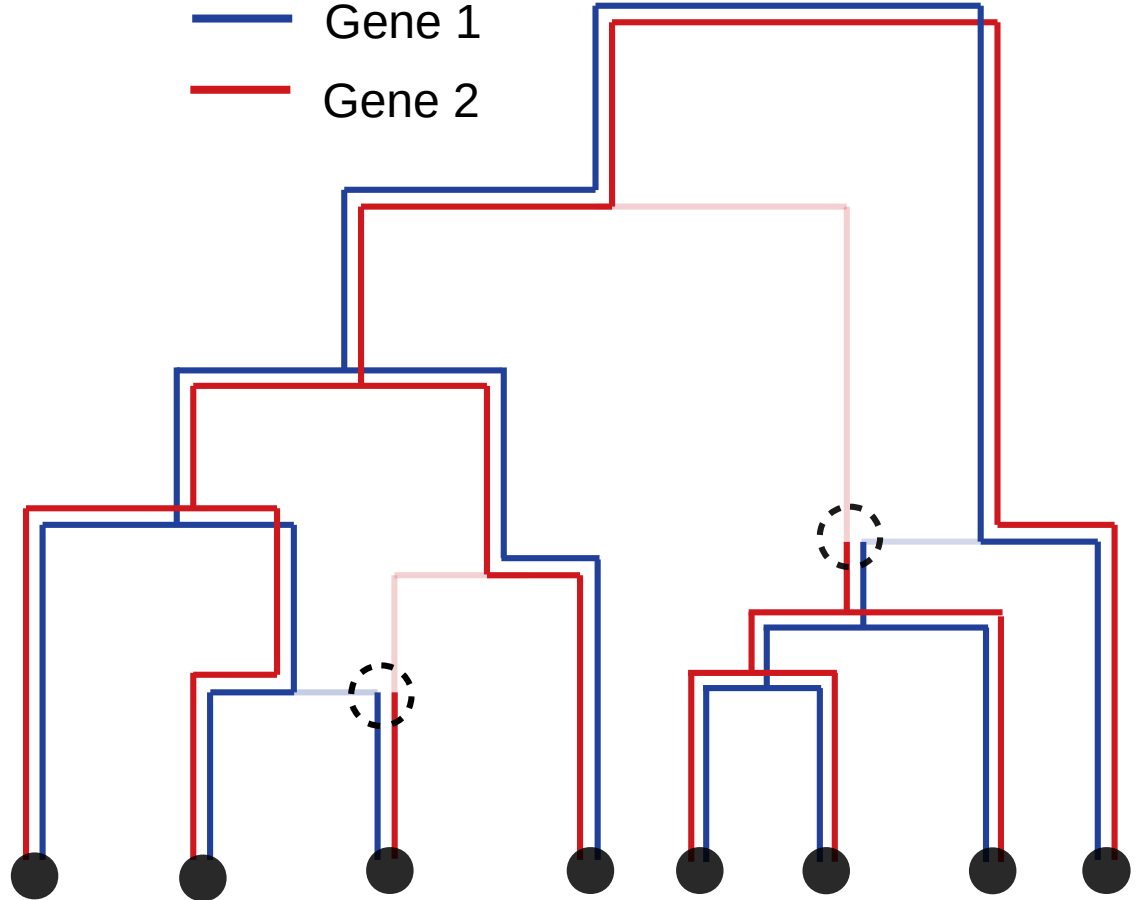
# Maximally compatible clades (MCCs)

The ARG is a **collage of gene trees**

Restricting to branches that
**belong to both trees**

↓

**Maximally compatible clades**

● Sampled strain
— Gene 1
— Gene 2

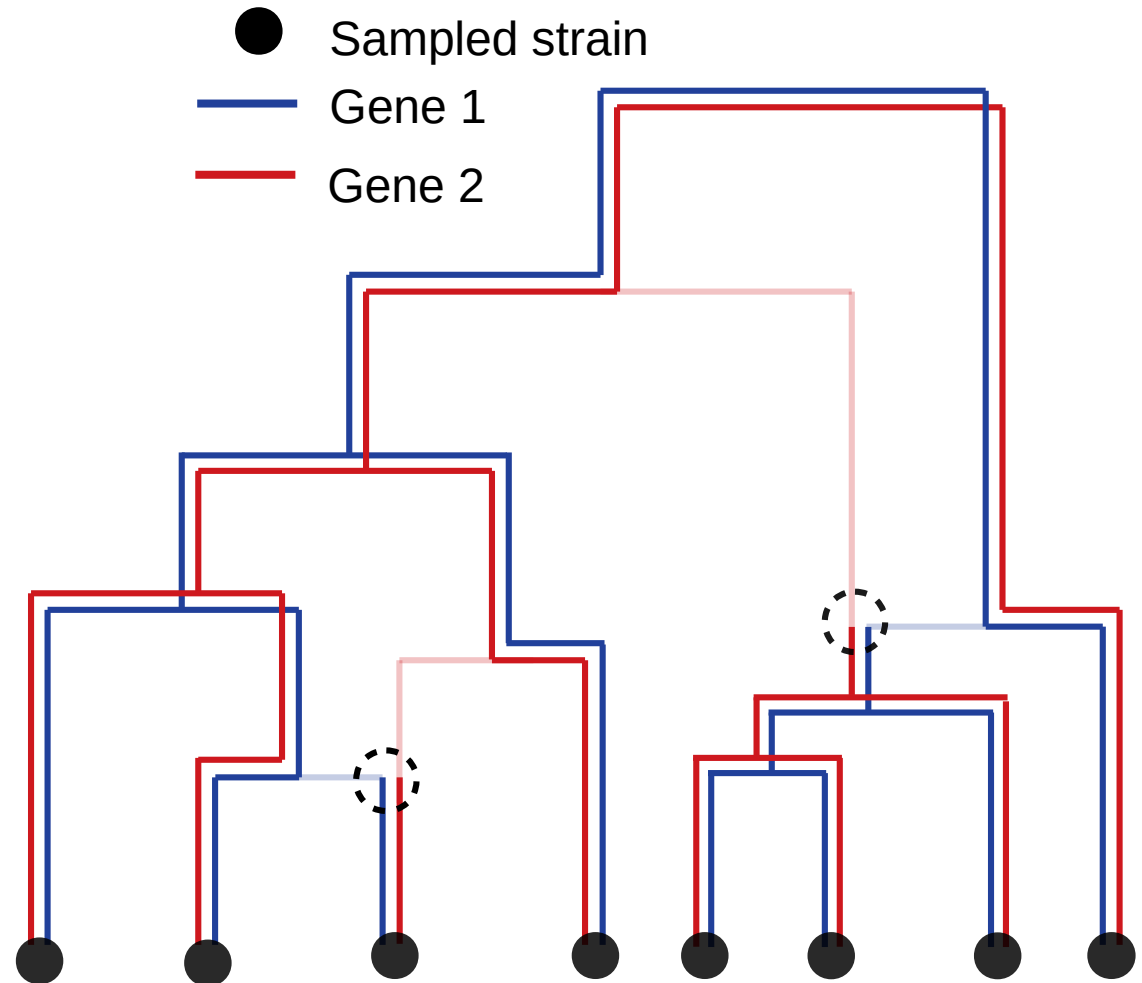# Maximally compatible clades (MCCs)

The ARG is a **collage of gene trees**

Restricting to branches that
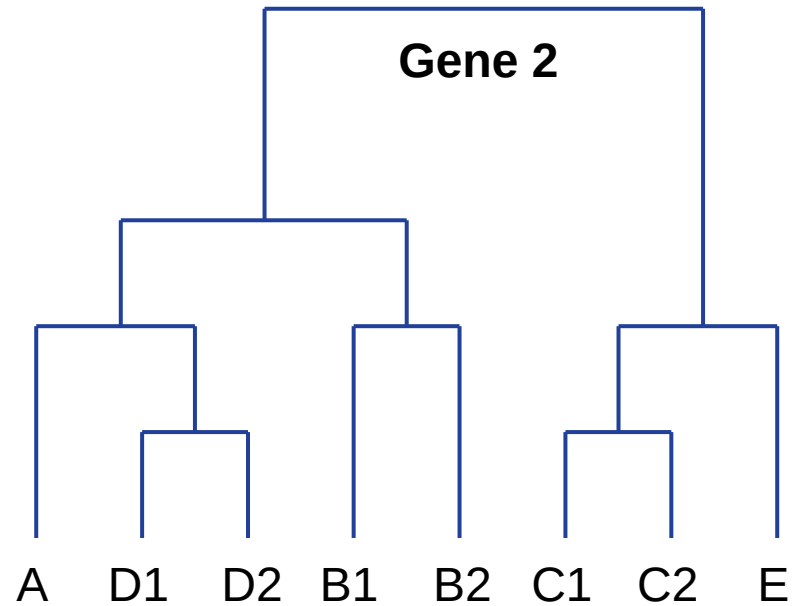**belong to both trees**

↓

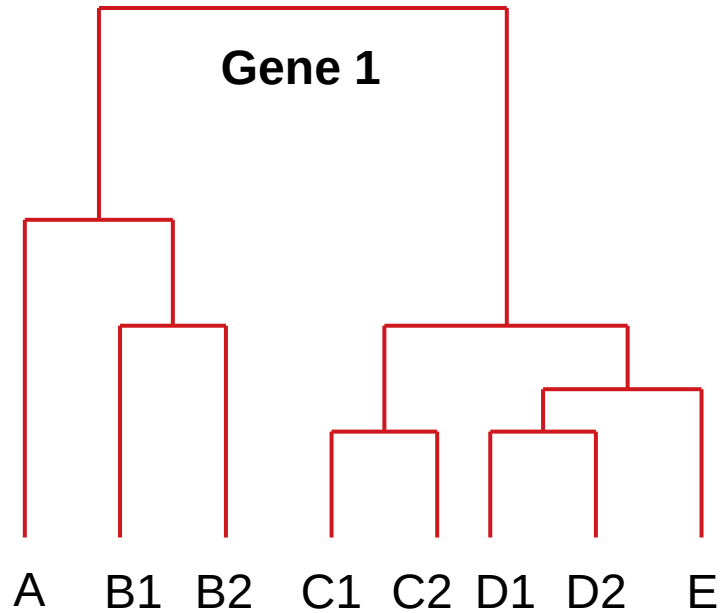**Maximally compatible clades**

- The **root of an MCC** is either
  - A reassortment
  - The root of both trees

- If **both trees** and **all MCCs** are known, then the **ARG** is known

● Sampled strain
— Gene 1
— Gene 2

# Inferring the ARG ⟶ Inferring MCCs

**First step**: naive estimation of MCCs
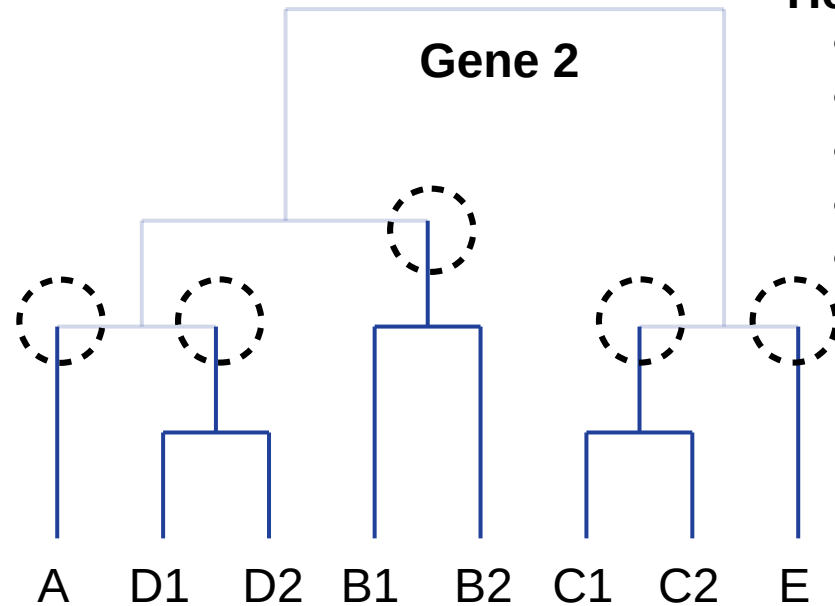
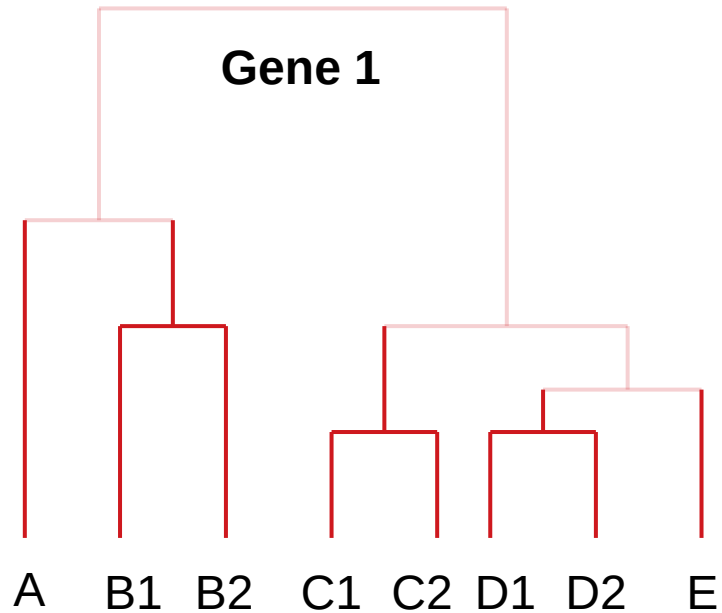⟶ Take clades that have exactly matching topologies

# Inferring the ARG ⟶ Inferring MCCs

**First step**: naive estimation of MCCs

⟶ Take clades that have exactly matching topologies



**Gene 1**

A  B1  B2  C1  C2  D1  D2  E

**Gene 2**

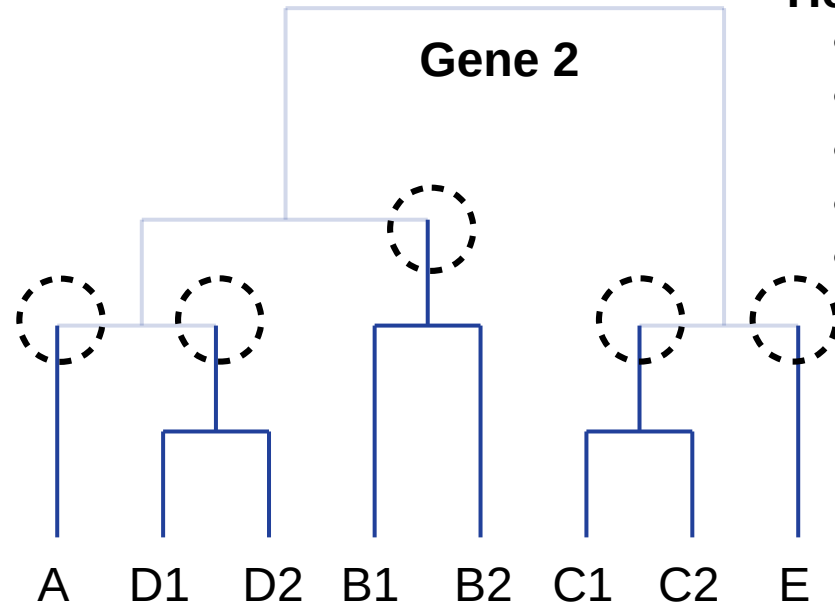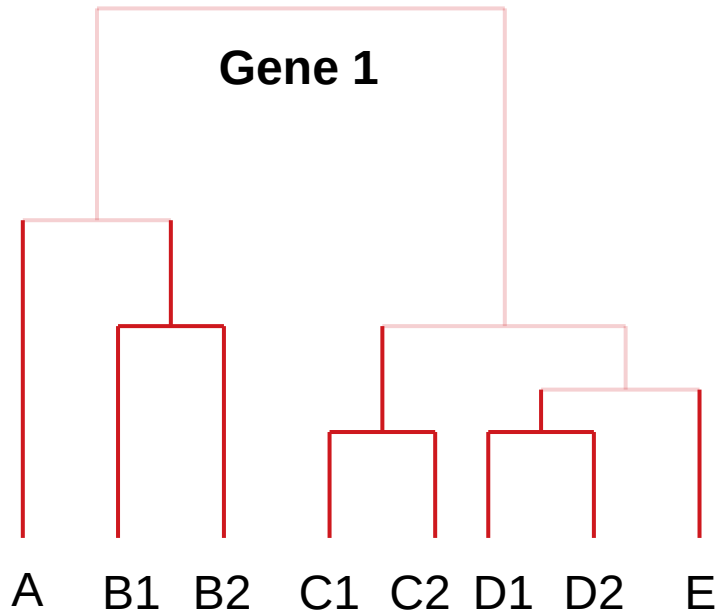A  D1  D2  B1  B2  C1  C2  E

**Here** : 5 naive MCCs
- A
- B1, B2
- C1, C2
- D1, D2
- E

5 reassortments !

# Inferring the ARG ⟶ Inferring MCCs

**First step**: naive estimation of MCCs

⟶ Take clades that have exactly matching topologies



**Gene 1**

A  B1  B2  C1  C2  D1  D2  E

**Gene 2**

A  D1  D2  B1  B2  C1  C2  E

**Here** : 5 naive MCCs
- A
- B1, B2
- C1, C2
- D1, D2
- E

↓

5 reassortments !

**Naive estimation** :

Finds too many MCCs ⟶ Too many reassortments

Conservative approach ⟶ Does not overextend MCCs

# Inferring MCCs



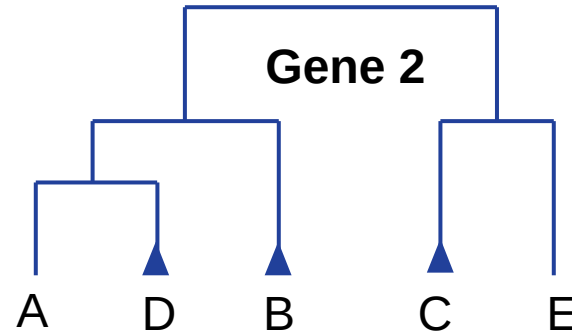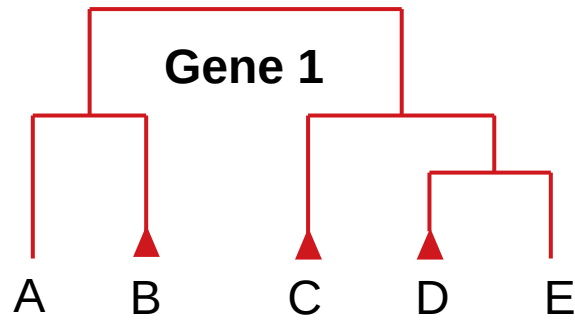**Gene 1**

A    B    C    D    E

**Gene 2**

A    D    B    C    E

**Second step:** "reduce" to naive MCCs

- (B1, B2) ⟶ B
- (C1, C2) ⟶ C
- (D1, D2) ⟶ D

# Inferring MCCs: Parsimonious approach



**First step:** "reduce" to naive MCCs

- (B1, B2) $\longrightarrow$ B
- (C1, C2) $\longrightarrow$ C
- (D1, D2) $\longrightarrow$ D

**By eye:**
D is the reassorted clade.
How can we **formalize** this?

Surrounding of each leaf: **clade** defined by parent:

- A $\longrightarrow$ (A,B) / (A,D)
- B $\longrightarrow$ (A,B) / (A,D,B)
- C $\longrightarrow$ (C,D,E) / (C,E)

- D $\longrightarrow$ (D,E) / (A,D)
- E $\longrightarrow$ (D,E) / (C,E)

$\longrightarrow$ **5 incompatibilities**

# Inferring MCCs: Parsimonious approach



**First step:** "reduce" to naive MCCs

- (B1, B2) ⟶ B
- (C1, C2) ⟶ C
- (D1, D2) ⟶ D

**By eye:**
D is the reassorted clade.
How can we **formalize** this?

Surrounding of each leaf: **clade** defined by parent:

- A ⟶ (A,B) / (A,D)
- B ⟶ (A,B) / (A,D,B)
- C ⟶ (C,D,E) / (C,E)

- D ⟶ (D,E) / (A,D)
- E ⟶ (D,E) / (C,E)

⟶ **5 incompatibilities**

Hypothesis: D is a reassortant ⟶ Remove it from the trees

- A ⟶ (A,B) / (A,B)
- B ⟶ (A,B) / (A,B)
- C ⟶ (C,E) / (C,E)

- ~~D ⟶ (D,E) / (A,D)~~
- E ⟶ (C,E) / (C,E)

⟶ **0 incompatibilities**
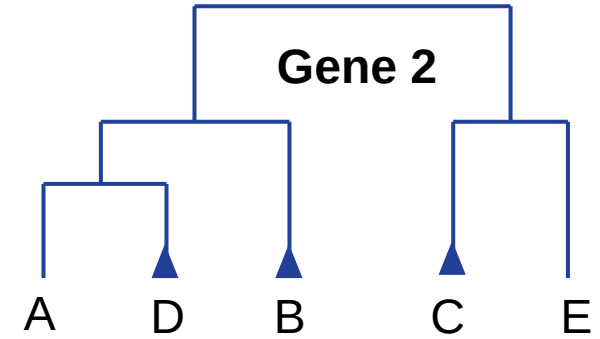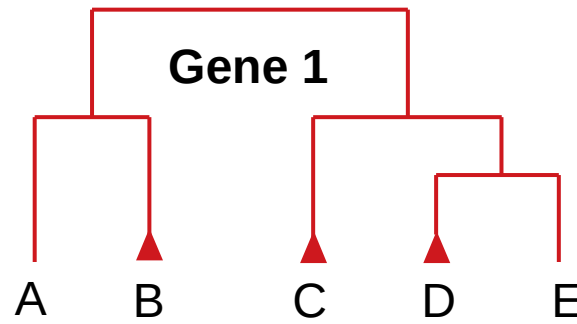
**0 remaining reassortments!**

# Inferring MCCs

For each leaf **n**

$$\longrightarrow \sigma_n \begin{cases} \text{1 if we } \textbf{remove } \textbf{\textit{n}} \\ \text{0 otherwise} \end{cases}$$



**Gene 1**

A    B    C    D    E

**Gene 2**

A    D    B    C    E

$$\longrightarrow \vec{\sigma} = (\sigma_1 \ldots \sigma_L) \quad : \text{"configuration" vector}$$

$$\longrightarrow \Delta(n, \vec{\sigma}) \begin{cases} \text{1 if } \textbf{incompatibility} \text{ above } \textbf{\textit{n}} \\ \text{0 otherwise} \end{cases}$$

# Inferring MCCs

For each leaf **n**

$$\longrightarrow \sigma_n \begin{cases} \text{1 if we \textbf{remove} } n \\ \text{0 otherwise} \end{cases}$$



**Gene 1**

A   B   C   D   E

**Gene 2**

A   D   B   C   E

$$\longrightarrow \vec{\sigma} = (\sigma_1 \ldots \sigma_L) \quad : \text{"configuration" vector}$$

$$\longrightarrow \Delta(n, \vec{\sigma}) \begin{cases} \text{1 if \textbf{incompatibility} above } n \\ \text{0 otherwise} \end{cases}$$

**# of incompatibilties**

**# of removed leaves**

$$\text{Minimize} \quad N_\gamma(\vec{\sigma}) = \underbrace{\sum_{n \in leaves} \Delta(n, \vec{\sigma})\sigma_n} + \gamma \underbrace{(L - |\vec{\sigma}|)}$$

**(Simulated annealing)**

**Minimize incompatibilities with a minimal number of reassortments**

# Inferring MCCs: summary

**Given two trees:**

Compute naive MCCs

Iterate

Minimize $N_\gamma(\vec{\sigma})$
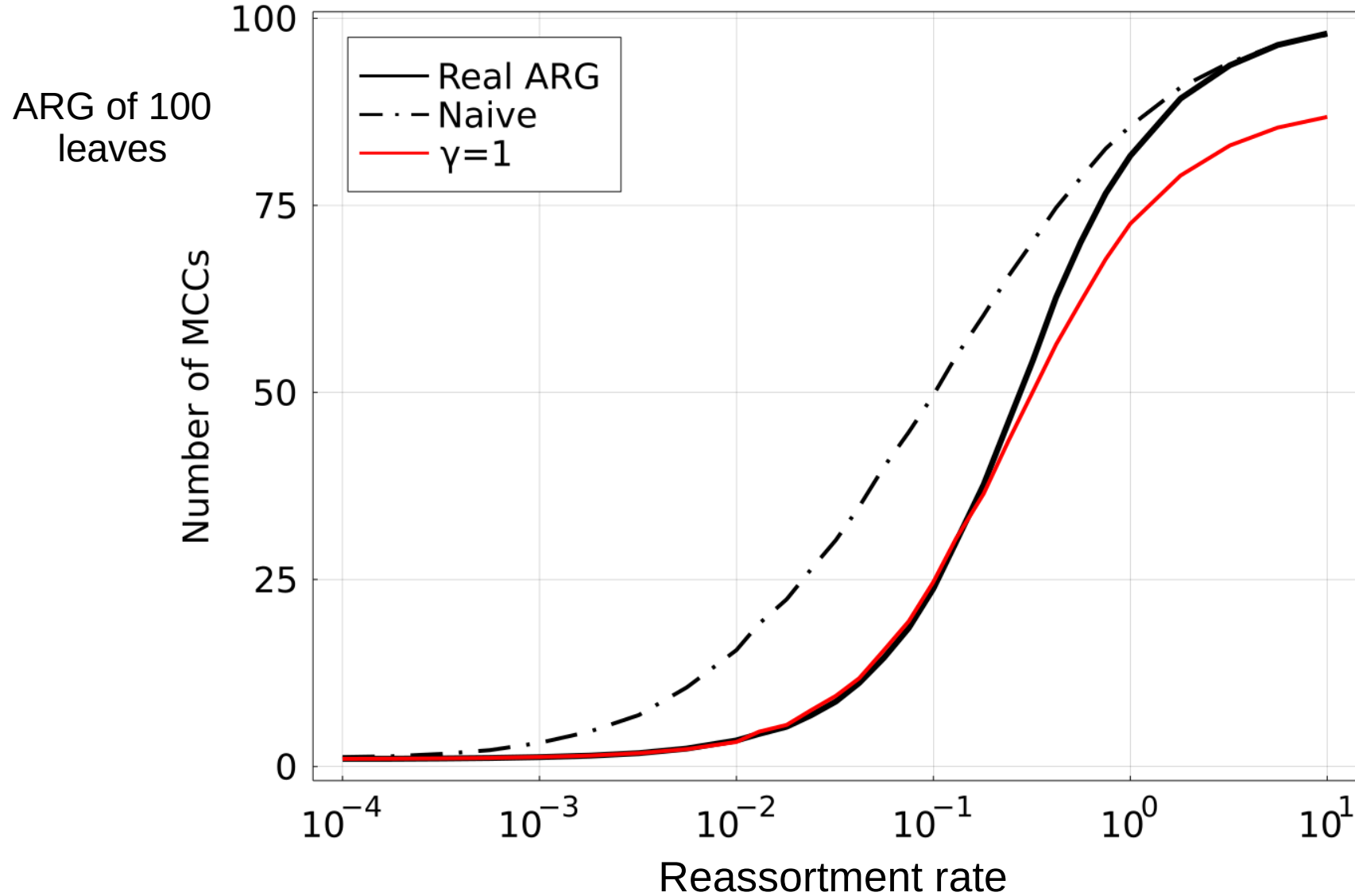
Remove all leaves such that $\sigma_n = 0$

Stop if only one naive MCC is found: **trees match perfectly**

# Interpretation of gamma

$$N_\gamma(\vec{\sigma}) = \sum_{n \in leaves} \Delta(n, \vec{\sigma})\sigma_n + \gamma(L - |\vec{\sigma}|)$$

- $\gamma \longrightarrow \infty$  **Infinite cost** for removing leaves  $\longrightarrow$  **Naive approach**
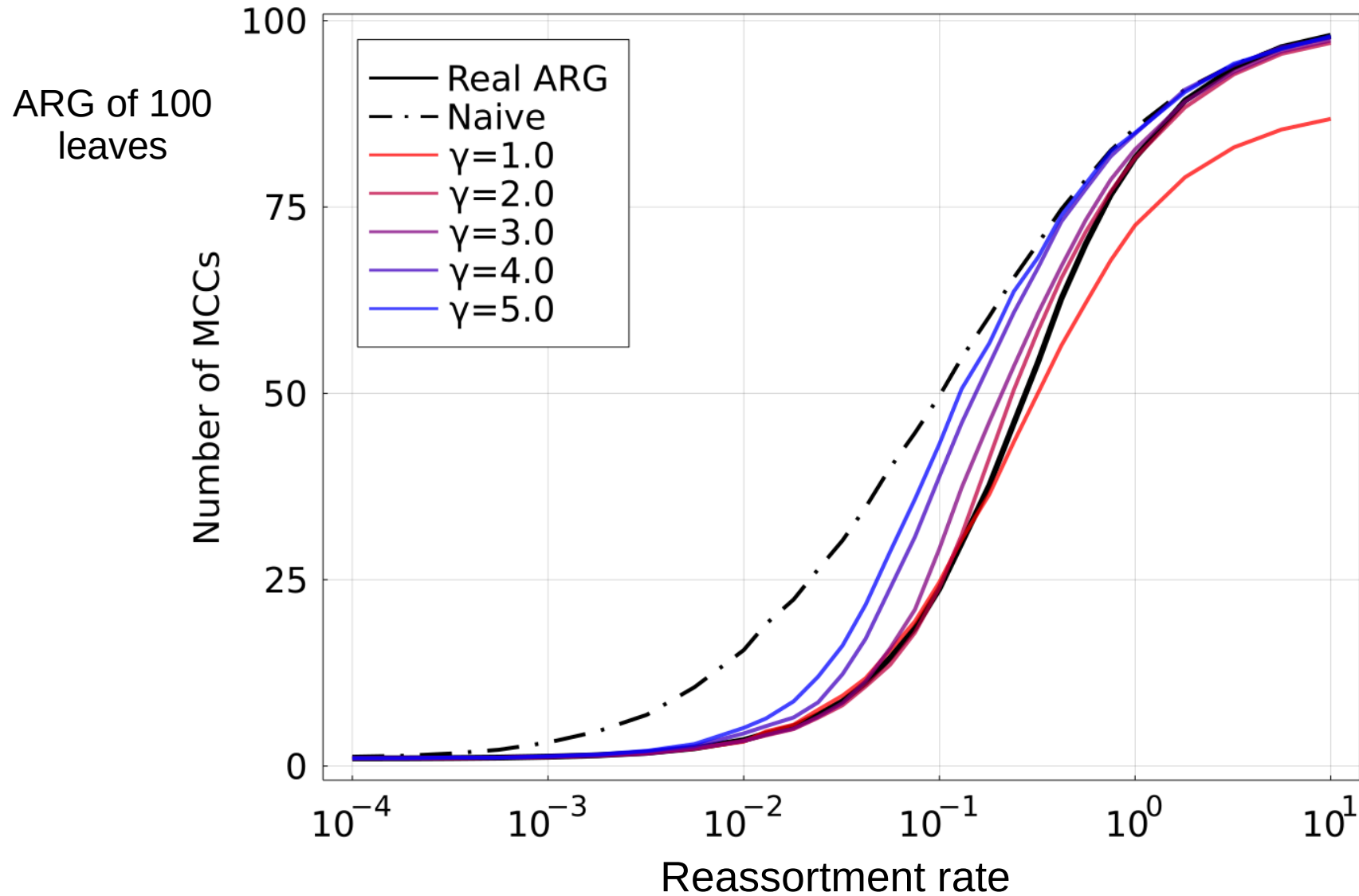
# Interpretation of gamma

$$N_\gamma(\vec{\sigma}) = \sum_{n \in leaves} \Delta(n, \vec{\sigma})\sigma_n + \gamma(L - |\vec{\sigma}|)$$

- $\gamma \to \infty$     **Infinite cost** for removing leaves   ⟶   **Naive approach**

---

- $\gamma = 1$     $N(\vec{\sigma})$ = # **incompatibilities** + # **removed leaves**

                             ↑                  ↑

Reassortments w. naive approach     Enforced reassortments

$N(\vec{\sigma})$ = **Total number of reassortments**   ⟶   **"Parsimonious" approach**

---

- Intermediate $\gamma$   ⟶   **Interpolate** between **naive** and **"parsimonious"**

# Interpretation of gamma



ARG of 100 leaves

# Interpretation of gamma



ARG of 100 leaves

# Evaluating the method: Simulated data

**Simulate an ARG**: coalescent process with reassortment rate $\rho$ $\longrightarrow$ **Apply the method**

**How can we evaluate the inference of MCCs?**

**Example**: for strains (A,B,C,D,E)
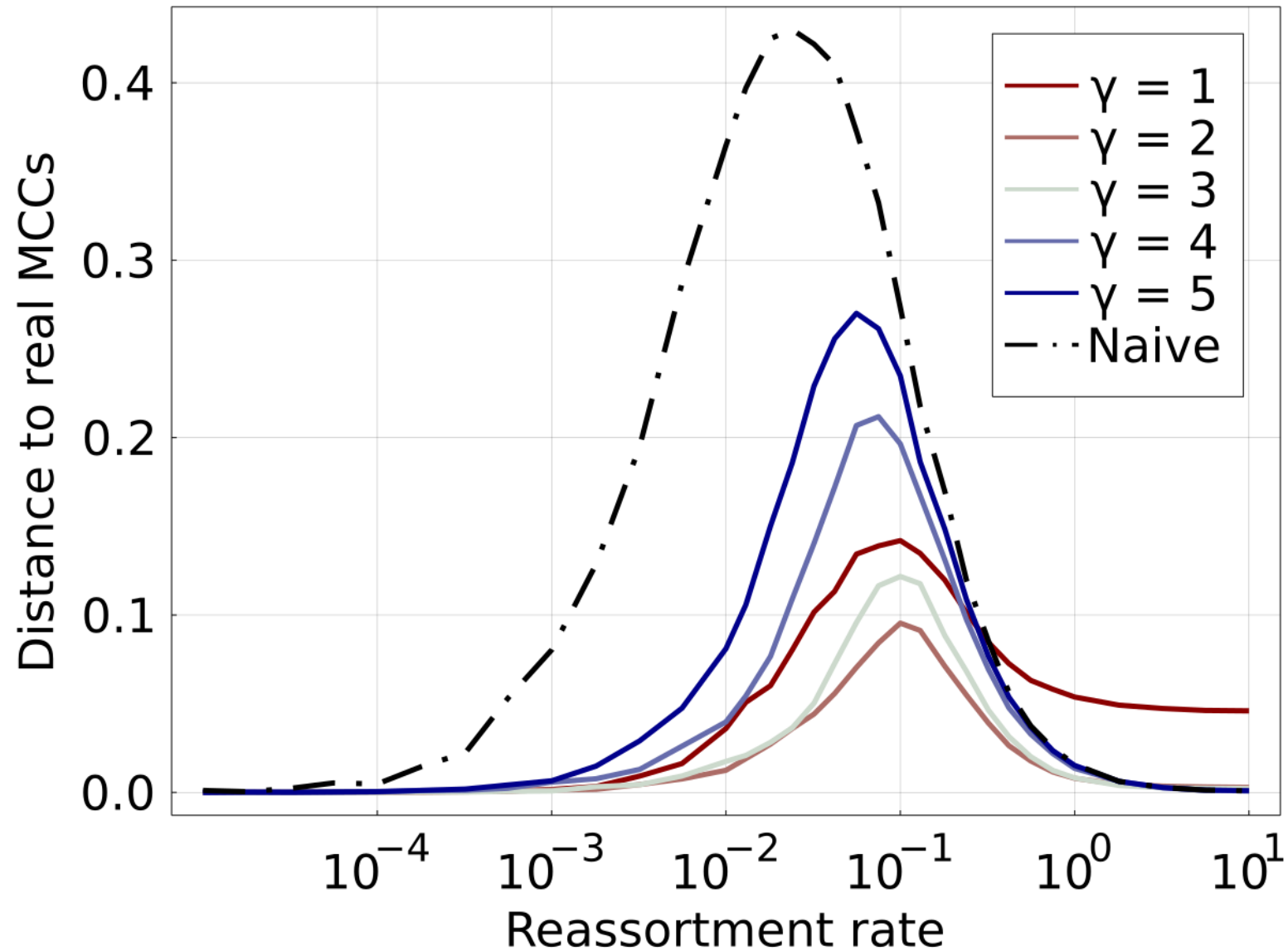- Real MCCs: (A,B,C), (D,E)
- Inferred: (A,B), (C), (D,E)

$\longrightarrow$ Defines a **partition of strains**

Using the **Variation of Information (VI)**: distance between partitions of a set

**[Meilă 2007]**

# Choosing gamma

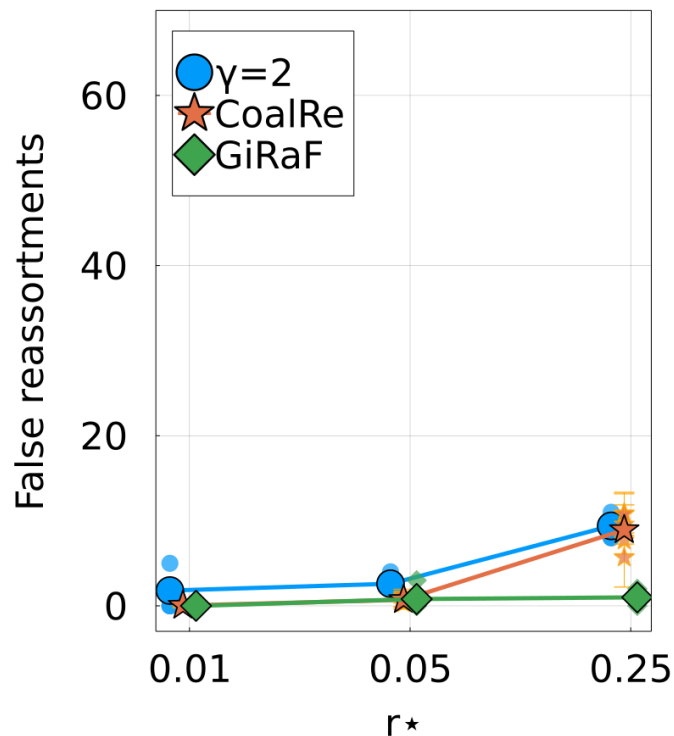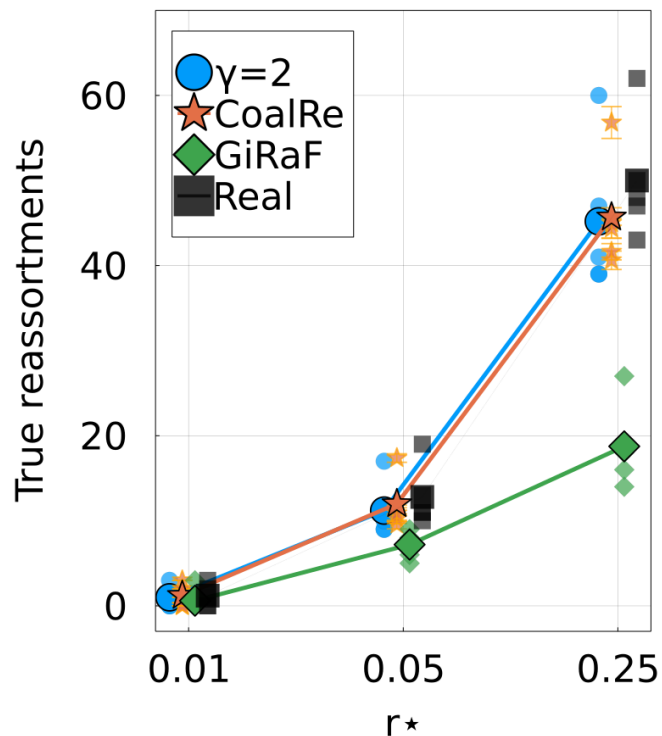**Distance: Variation of Information (VI)**



$$\gamma = 2$$

# Comparison w. other methods

CoalRe: ML based [Müller et. al. 2020]

GiRaF: topology based [Nagarajan & Kingsford 2011]



Runtime

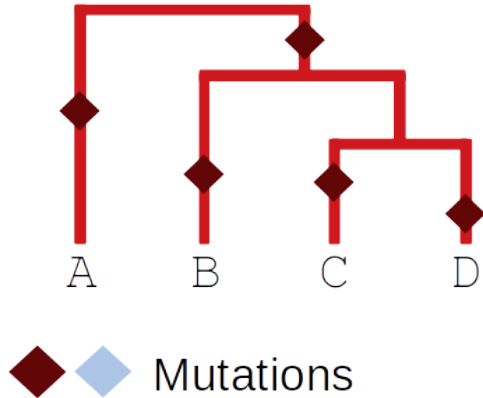| | CoalRe | GiRaF | Treeknit |
|---|---|---|---|
| Inferring trees | | 20min | 30s |
| Inferring the ARG | ~hours | 40s | 40ms |

for 100 leaves

# Application: better resolved trees



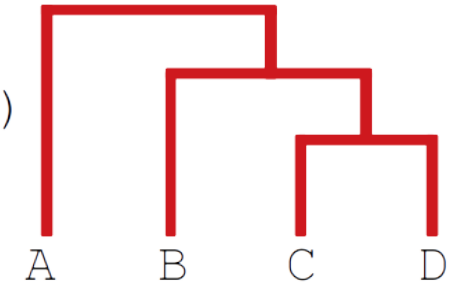Using information from both segments

**Real trees** | **Observed trees** | **Resolved trees**

Mutations
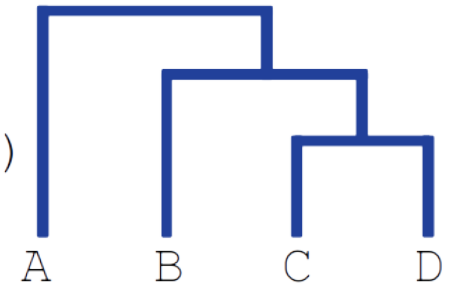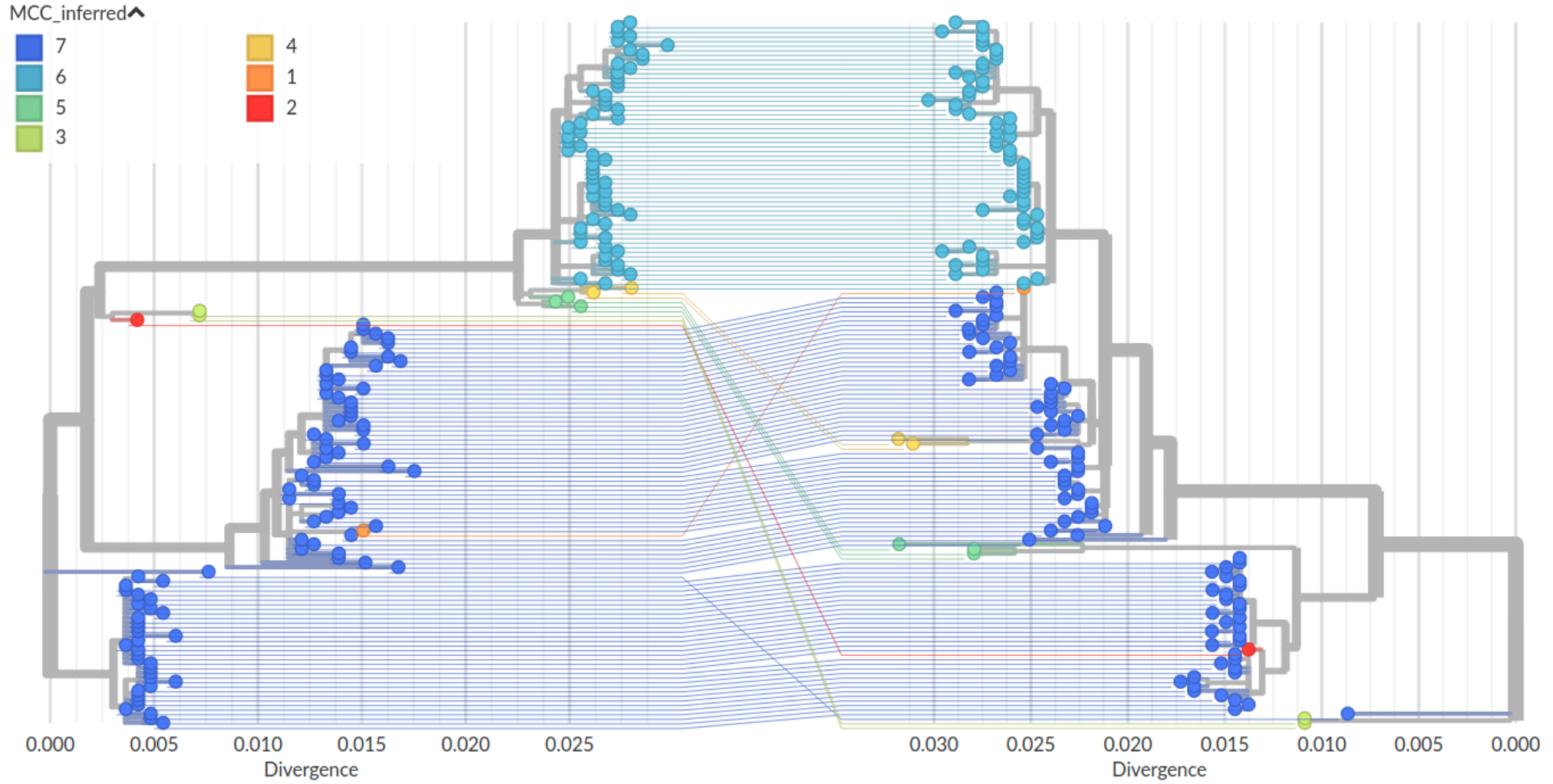
Split (CD)

Split (BCD)

# Application: better resolved trees

# Application: disentangling tanglegrams



HA gene

NA gene

150 sequences
from New York
**[Holmes et. al. 2005]**

**Without the knowledge of reassortments: hard problem**

# Application: disentangling tanglegrams



**With the knowledge of reassortments: easy**

# Summary

Available at **github.com/PierreBarrat/TreeKnit**

## Results

- **Treeknit:** Heuristic to infer ARGs from two trees

- Depends on **one parameter**, **interpolating** between naive and parsimonious inference

- Very **fast** runtime

- **Good performance** on **simulated data** for all reassortment rates

## Applications

- Resolve trees

- Visualisation: disentangle tanglegrams

- Knowledge of the ARG $\longrightarrow$ Effect of reassortment on influenza evolution

# Thank you for listening!