

# Fast Computation of Genome Distances

---

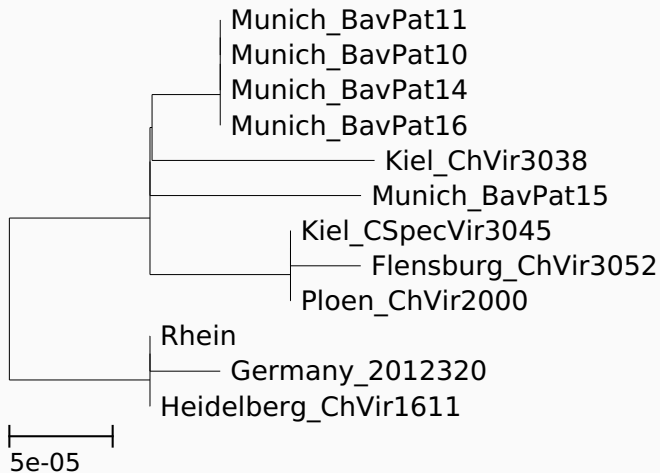
Fabian Klötzl & Bernhard Haubold

Aquavit 2020

MPI für Evolutionsbiologie

Slides: [kloetzl.info](mailto:kloetzl.info)

# Sars-CoV-2



# Multiple Sequence Alignment

A AACGTTGTGCA  
B AATG-TGAGC-  
C -ACGTTGTG--  
D AAC-TTGTGA-

$L$ : Length of genome sequence

$N$ : Number of sequences

Multiple Sequence Alignment:  $N \cdot L$

Computation of optimal MSA: NP-Complete,  
think:  $O(2^{NL})$

## Example

$L = 11$  bp,  $N = 4$

$NL = 44$ ,  $2^{NL} = 1.7 \cdot 10^{13}$

# Multiple Sequence Alignment

A AACGTTGTGCA  
B AATG-TGAGC-  
C -ACGTTGTG--  
D AAC-TTGTGA-

$L$ : Length of genome sequence

$N$ : Number of sequences

Multiple Sequence Alignment:  $N \cdot L$

Computation of optimal MSA: NP-Complete,  
think:  $O(2^{NL})$

## Example

$L = 11$  bp,  $N = 4$

$NL = 44$ ,  $2^{NL} = 1.7 \cdot 10^{13}$

## Sars-CoV-2

$L \approx 30\,000$  bp

$N = 12$

$NL = 3.6 \cdot 10^5$

$2^{NL} = 6 \cdot 10^{108370}$

# Multiple Sequence Alignment

A AACGTTGTGCA  
B AATG-TGAGC-  
C -ACGTTGTG--  
D AAC-TTGTGA-

$L$ : Length of genome sequence

$N$ : Number of sequences

Multiple Sequence Alignment:  $N \cdot L$

Computation of optimal MSA: NP-Complete,  
think:  $O(2^{NL})$

## Example

$L = 11$  bp,  $N = 4$

$NL = 44$ ,  $2^{NL} = 1.7 \cdot 10^{13}$

## Sars-CoV-2

$L \approx 30\,000$  bp

$N = 12$

$NL = 3.6 \cdot 10^5$

$2^{NL} = 6 \cdot 10^{108370}$

## Ensembl *E. coli*

$L \approx 5$  Mb

$N = 2681$

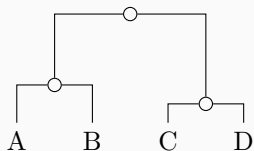
$NL = 1.3 \cdot 10^{10}$

$2^{NL} \equiv \text{Inf}$

# Faster Phylogeny Reconstruction

>A  
AACGTTGTGCA  
>B  
AATGTGAGC  
>C  
ACGTTGTG  
>D  
AACTTGTGA

$$\Rightarrow \begin{pmatrix} 0 & 0.1 & 0.25 & 0.3 \\ 0.1 & 0 & 0.3 & 0.3 \\ 0.25 & 0.3 & 0 & 0.05 \\ 0.3 & 0.3 & 0.05 & 0 \end{pmatrix} \Rightarrow$$



$O(N \cdot L)$

$\Rightarrow$   
?

$O(N^2)$

$\Rightarrow$   
 $O(N^3)$

$O(N)$

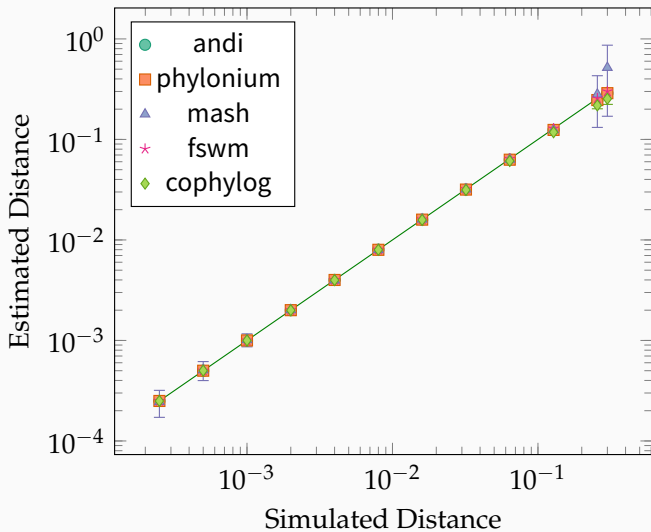
## Pairwise Comparison

A     AACGTTGTGCA  
B     AATG TGAGC

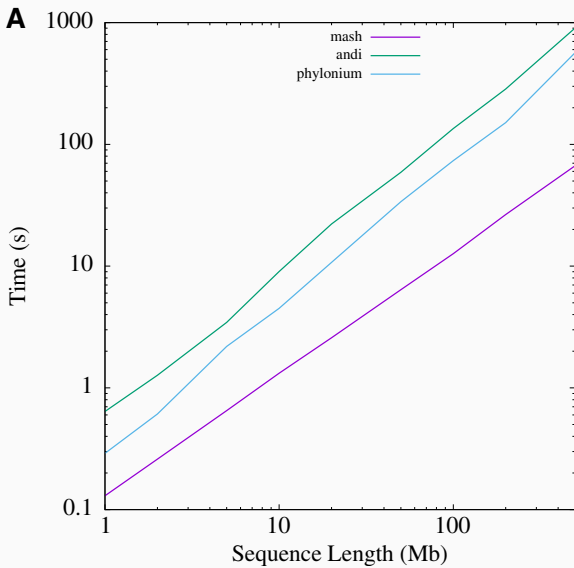
- Align long exact matches
- Disregard indels
- Don't focus on every substitution
- Count substitutions, quickly

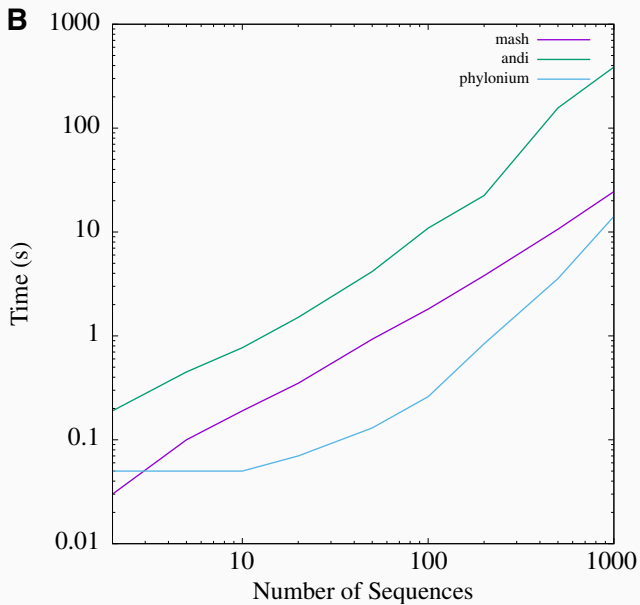
$O(L_1 + L_2)$  comparison

# Accuracy









## A simplified MSA

A     AACGTTGTGCA  
B     AATG TGAGC  
C     ACGTTGTG  
D     AAC TTGTGA

- Align long exact matches against a single reference
- Disregard indels
- Don't focus on every substitution
- Count substitutions, quickly

$O(N^2L)$  comparison



0.02

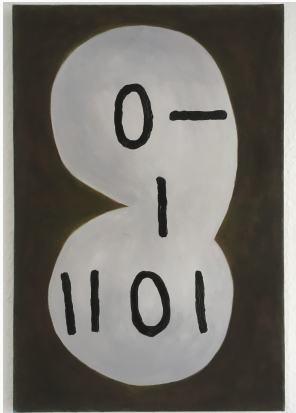
$L \approx 5 \text{ Mb}$

$N = 2681$

13 min

14 GB

# Phylonium



Art: Christian Rothmaler

Paper:

- [doi.org/10.1093/bioinformatics/btz903](https://doi.org/10.1093/bioinformatics/btz903)

Get it:

- `apt install phylonium`
- `brew install science/phylonium`
- `aura -A phylonium`
- [github.com/evolbioinf/phylonium](https://github.com/evolbioinf/phylonium)