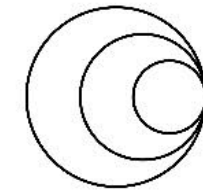**Viral Genomics and Bioinformatics (Latin America and the Caribbean)**

**10–14 October 2022 - Virtual course**

wellcome
connecting
science

# Introduction to multiple sequence alignments

**Carolina Torres**

*Universidad de Buenos Aires, Facultad de Farmacia y Bioquímica, Instituto de Investigaciones en Bacteriología y Virología Molecular (IBaViM), Buenos Aires, Argentina.*

*Consejo Nacional de Investigaciones Científicas y Técnicas (CONICET), Buenos Aires, Argentina.*

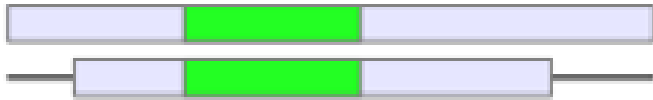caro.torr@gmail.com | ctorres@ffyb.uba.ar

@CaroTorr

# Goals for sampling and sequencing

- **Genomic surveillance and molecular epidemiology:**
  frecuency of lineages, variants or mutations.

- **Origin and evolution:** introductions, diversification pattern.

- **Outbreaks and transmission chains:** common sources of infection.

- **Evolutionary dynamics:** ancestral ages, rates of evolution, viral demography, dispersion rates and patterns, ancestral locations, predictors.
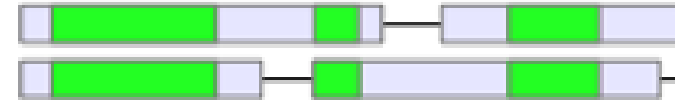
# Alignments

✓ **An arrangement of DNA, RNA or amino acid sequences in which homologous sites are in the same position (they are "aligned").**
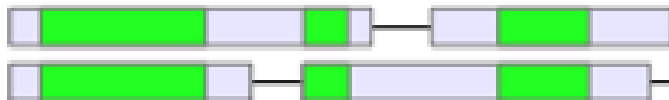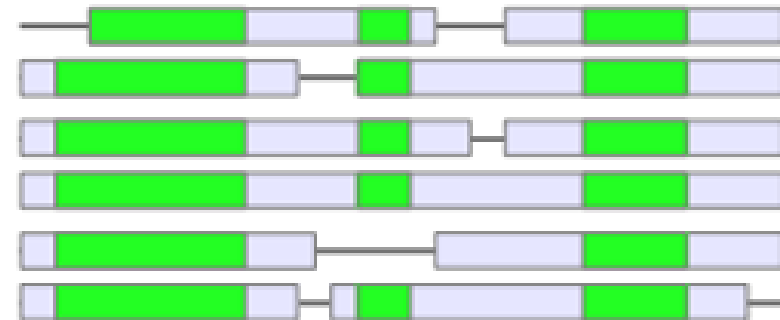
- **Local:** Align one or more stretches of similarity

- **Global:** Align sequences end-to-end
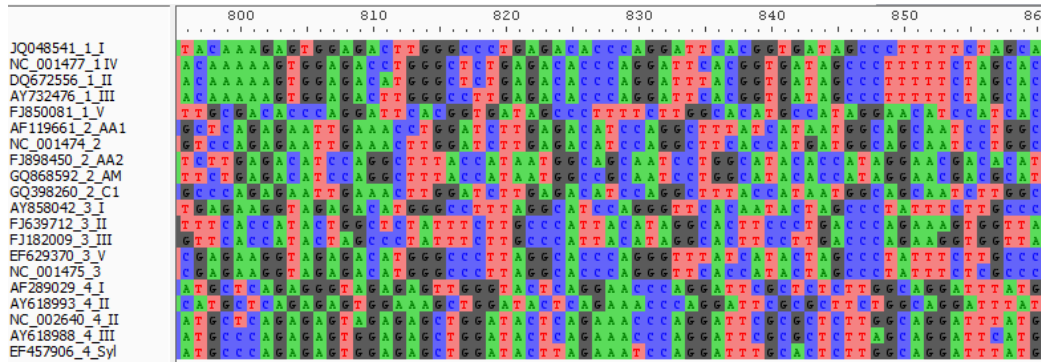
- **Pairwise:** Align two sequences

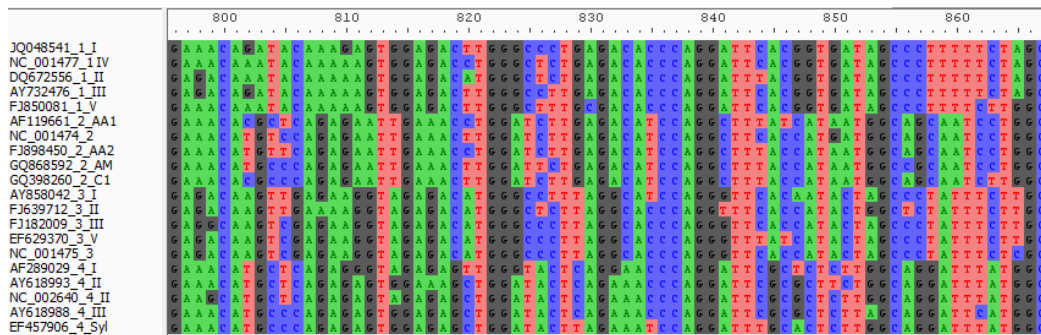- **Multiple:** Align more than two sequences

# Multiple sequence alignments (MSA)

✓ **Positional homology**


*Non-aligned sequences*


*Aligned sequences*

- The nts or AAs found at a position (column) in the sequences are considered descendants from a common ancestral site (homologous).

- All methods attempt to maximize matches/similarities and minimize mismatches/differences between sequences.

- Different programs search for the "best" alignment through different methods.

# Multiple sequence alignments (MSA)

✓ **Positional homology**



- The goal of alignment methods is to maximise a score based on:
  - rewarding matches (+ score).
  - penalising mismatches/rare substitutions (- score).
  - penalising gaps or indels (- score), which requires a gap opening and extension penalty scheme (these values are chosen arbitrarily).

- Gaps are introduced in the sequences or at the ends of the alignment.

- In coding regions, gaps are usually introduced in triplets.

# Programs to build alignments

- **Clustal W/X: Progressive alignment.**

- **Muscle: Iterative  method.**

- **MAFFT: multiple methods.**

- **T-Coffee:** consistency-based alignment (capable of combining a collection of multiple/pairwise, global/local alignments into one).

- **ProbCons:** combination of probabilistic modeling and consistency-based alignment techniques (protein sequences).

- **Probalign:** combines amino acid posterior probability estimation using partition function methods and computation of maximal expected accuracy alignment.

# Methods to build alignments

- The **progressive algorithm** consists of three main stages:

  **(i) All pairs of sequences are aligned** separately (**pairwise alignments**) in order to calculate a "distance" matrix (this is done using dynamic programming);

  **(ii) A guide tree** is built from the distance matrix (using a clustering algorithm, such as Neighbor Joining);

  **(iii) The guide tree is used to cluster and align the sequences progressively according to the branching order** (starting with the two closest sequences and ending with the most distant).



(a) input

(b) pairwise distances

(c) Guide tree

(d) Progressive alignment

✓ Gaps are inserted at identical positions in all sequences of a cluster and are preserved.

- **Programs:** Clustal W/X

# Methods to build alignments

- The **iterative methods** work similarly to progressive methods, but they repeatedly realign the initial sequences and add new sequences to the MSA.

- **MUSCLE:** there are three main stages:

  **(i) Draft progressive** to build a MSA (from a distance matrix with approximate values)**;**

  **(ii) Improved progressive:** a new distance matrix is created from the first MSA and sequences realigned to reflect new guide tree;

  **(iii) Refinement:** The tree is split into 2 subtrees, profiles are built and aligned, different bipartitions are tried until convergence is reached.

# Methods to build alignments

- **MAFFT** offers a range of multiple alignment strategies:



(a) FFT–NS–1, FFT–NS–2 — Progressive methods

Distance matrix based on the number of shared 6-tuples
→ Constructing guide tree
→ Progressive alignment → FFT-NS-1 / NW-NS-1
→ Re-constructing guide tree
→ Re-alignment → FFT-NS-2 / NW-NS-2

(b) FFT–NS–i, NW–NS–i — Iterative refinement method

Distance matrix based on the number of shared 6-tuples
→ Constructing guide tree
→ Progressive alignment → FFT-NS-1 / NW-NS-1
→ Re-constructing guide tree
→ Re-alignment → FFT-NS-2 / NW-NS-2
→ Iterative refinement (WSP) → FFT-NS-i / NW-NS-i

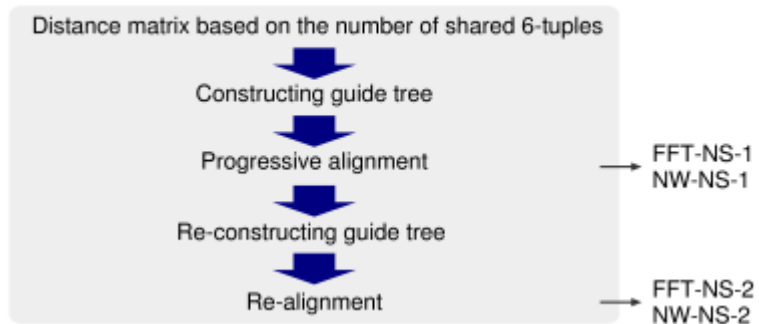(c) L–INS–i, E–INS–i, G–INS–i — Iterative refinement methods using WSP and consistency scores

Distance matrix based on all-pairwise aignments
→ Constructing guide tree
→ Progressive alignment → G-INS-1 / L-INS-1
→ Iterative refinement (WSP + COFFEE) → G-INS-i / L-INS-i

In general, there is a tradeoff between speed and accuracy. The order of speed is a > b > c, whereas the order of accuracy is a < b < c.

https://mafft.cbrc.jp/alignment/software/algorithms/algorithms.html

# Editing alignments

- **It is a good practice to always visually inspect the alignment** to check the position of gaps, outlier sequences and poorly-aligned regions.



- **Manual edition may be necessary:**

  - Trim the ends or specific positions that cannot be aligned unambiguously (Gblocks may be useful).

  - Realign blocks.

  - Use biological knowledge to improve the alignment.

- **Program:** Aliview, BioEdit, Seaview, MEGA, UGENE, others.

# Editing alignments

- **It is a good practice to always visually inspect the alignment** to check the position of gaps, outlier sequences and poorly-aligned regions.
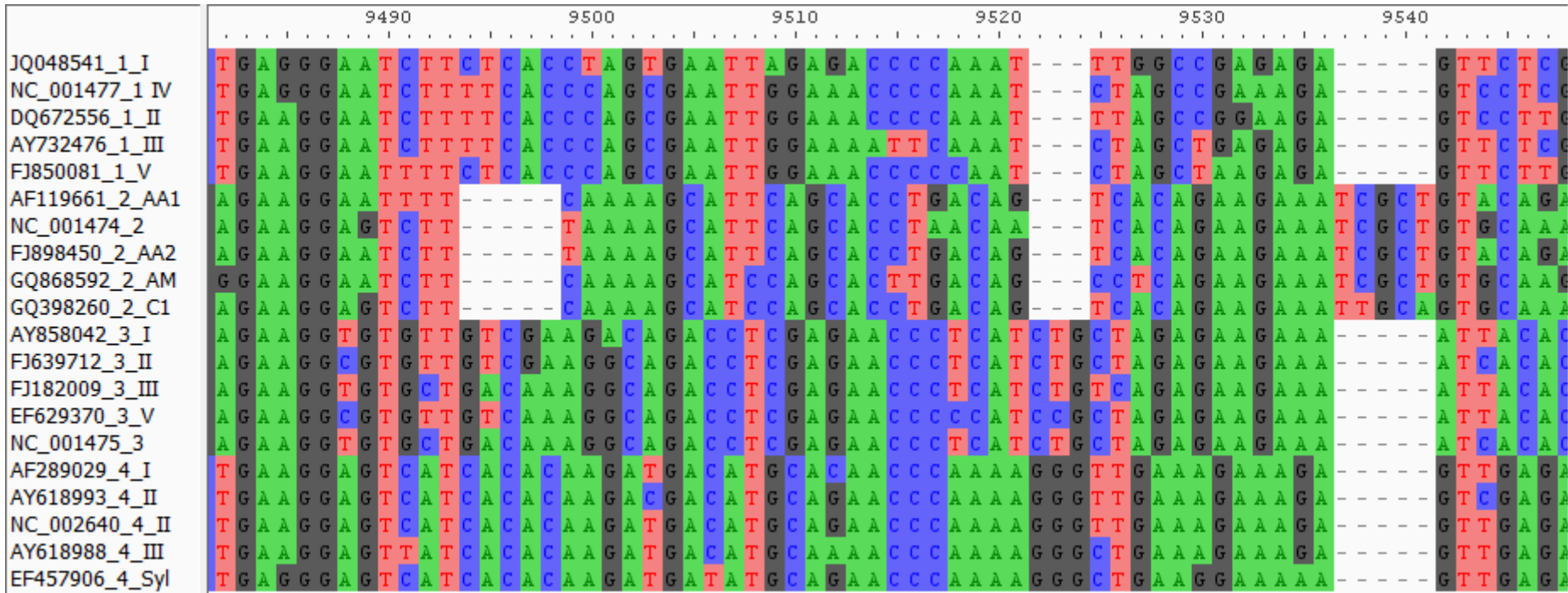


- **The final decision on what to include or exclude is yours.** ➡️ **Sensitivity analysis (check the impact of your decisions)**

# References

- The Phylogenetic Handbook: A Practical Approach to DNA and Protein Phylogeny. Marco Salemi, Anne-Mieke Vandamme (Eds). Cambridge University Press. (2009).

- Katoh, Standley 2013. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Molecular Biology and Evolution* 30:772-780.

- Clustal W and Clustal X version 2.0. Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H, Valentin F, Wallace IM, Wilm A, Lopez R, Thompson JD, Gibson TJ, Higgins DG. Bioinformatics. 2007 Nov1;23(21):2947-8.

- Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. Nucleic Acids Res. 2004 Mar 19;32(5):1792-7.

- BLAST® Help. Bethesda (MD): National Center for Biotechnology Information (US); 2008 (http://www.ncbi.nlm.nih.gov/books/NBK1762/).

- Talavera, G., and Castresana, J. (2007). Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments. Systematic Biology 56, 564-577. http://molevol.cmima.csic.es/castresana/Gblocks.html