# Incremental algorithms for large homologous gene families

Jean-François Dufayard

Manolo Gouy

François Rechenmann
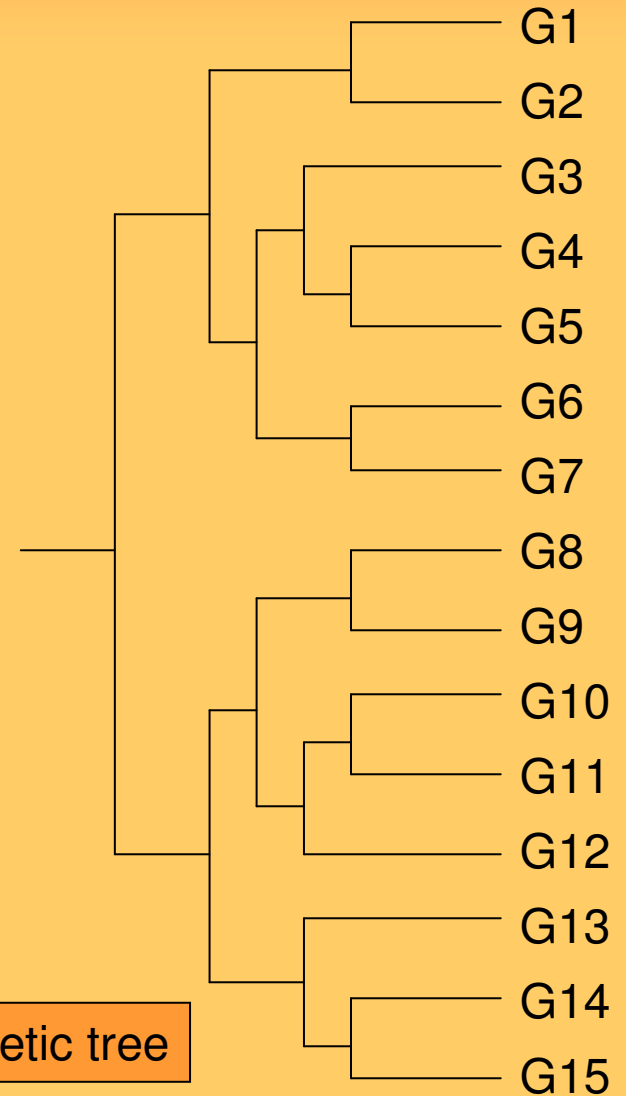
# Problematic (1/3)

**Phylogenic analysis**

```
L-QGNDLLEDSPYEPVNSRLSDIFRAVPFISDVFQQVEHI-SKGNNCLDAAKACNLDDTC
L-QGNDLLEDSPYEPVNSRLSDIFRAVPFISDVFQQVEHI-SKGNNCLDAAKACNLDDTC
LTEGEEFYEASPYEPVTSRLSDIFRLASIFSGTGADP-VVSAKSNHCLDAAKACNLNDNC
LTEGEEFYEASPYEPVTSRLSDIFRLASIFSGTGADP-VVSAKSNHCLDAAKACNLNDNC
L-QGNDLLEDSPYEPVNSRLSDIFRLAPIVS-----VEPVLSKGNNCLDAAKACNLNDTC
LAEGEEFYEASPYEPITSRLSDIFRLASIFSGM--DP-ATNSKSNHCLDAAKACNLNDNC
LTEGEEFYEASPYEPVTSRLSDIFRLASIFSGTGTDP-AVSTKSNHCLDAAKACNLNDNC
L-QGNDLLEDSPYEPVNSRLSDIFRAVPFIS-----VEHI-SKGNNCLDAAKACNLDDTC
L-QGNDLLEDSPYEPVNSRLSDIFRVVPFIS-----VEHI-PKGNNCLDAAKACNLDDIC
L-QGNDLLEDSPYEPVNSRLSDIFRVVPFISDVFQQVEHI-PKGNNCLDAAKACNLDDIC
L-QGNDLLEDSPYEPVNSRLSDIFRAVPFIS-----VEHI-SKGNNCLDAAKACNLDDTC
LTEGEEFYEASPYEPVTSRLSDIFRLASIFSGTGTDP-AVSTKSNHCLDAAKACNLNDNC
LMEGMNVLESSPYEPFIRGF-DYVRLASITAGSENEVTQV----NRCLDAAKACNVDEMC
---GEEFYEASPYEPVTSRLSDIFRLASIFSGTGADP-VVSAESNHCLDAAKACNLNDNC
---G--------------------------TGADP-VVSAESNHCLDAAKACNLNDNC
```
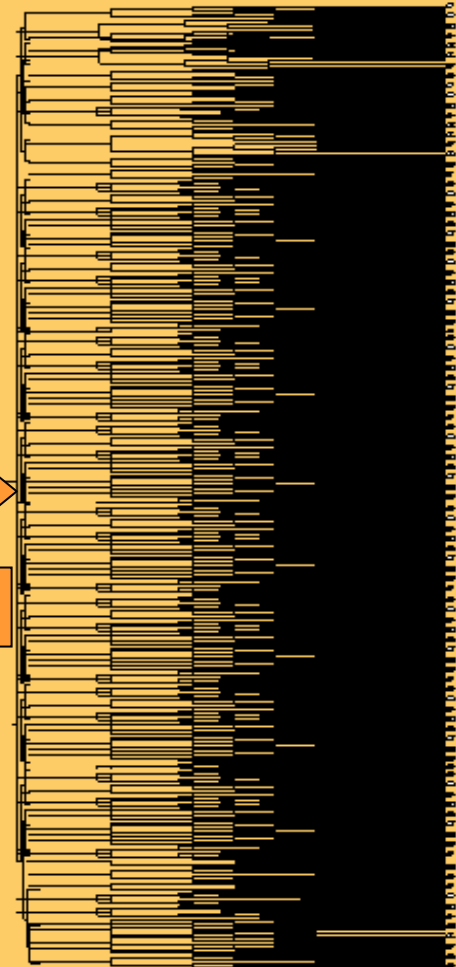
**Multiple alignment**

**Phylogenetic tree**

G1
G2
G3
G4
G5
G6
G7
G8
G9
G10
G11
G12
G13
G14
G15

# Problematic (2/3)



New sequence

Incremental algorithms

# Problematic (3/3)

Main goal: add sequences in a phylogenetic tree and its alignment
- ➢ Avoid redondant calculations
- ➢ Preserve the quality of the tree and the alignment

Application: *European Small Subunit Ribosomal RNA database*
- ➢ Family of 35 000 sequences (~1500 nucleotids per sequence)
- ➢ A tree and an alignment is computed for 10 000 sequences, **manually**.
- ➢ New sequences frequently added to the database.

HELiX
BIOINFORMATICS

UMR 5558
Biométrie, biologie
évolutive.

INRIA
RHÔNE-ALPES

# Adding a new sequence, main problems

- How to find the location of the new sequence in the alignment ?

- How to find the location of the new sequence in the tree ?

- How to insert the new sequence in the alignment, knowing its location ?

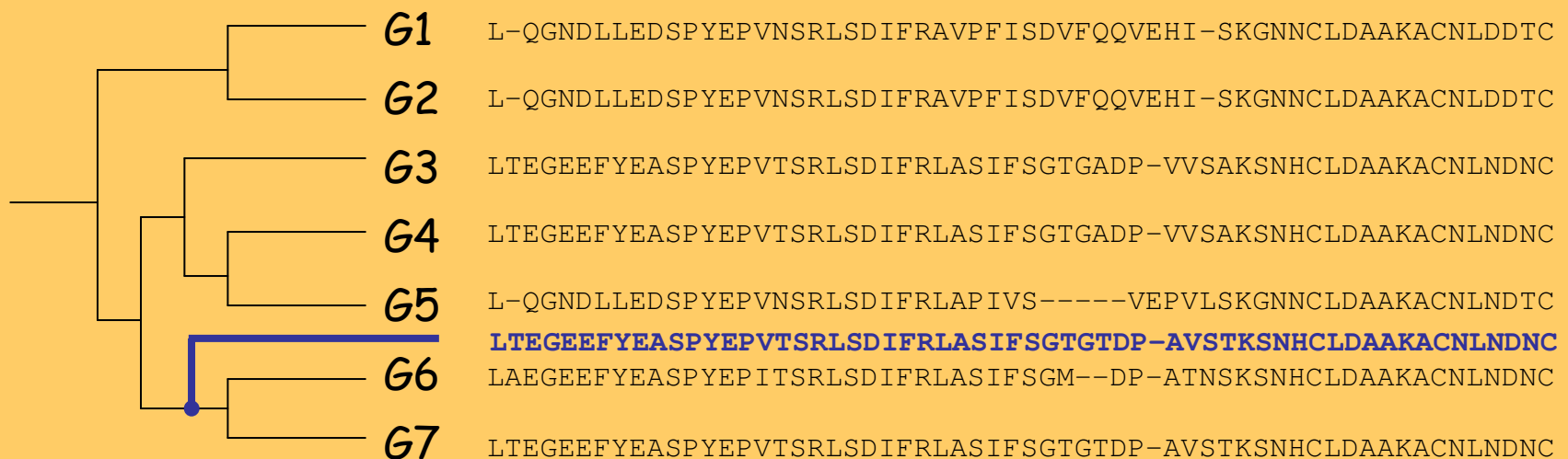- How to insert the new sequence in the tree, knowing its location ?

H E L i X
BIOINFORMATICS

UMR 5558
Biométrie, biologie
évolutive.

I N R I A
RHÔNE-ALPES

# Location of the new sequence (1/2)

- How to find the location of the new sequence in the alignment ?

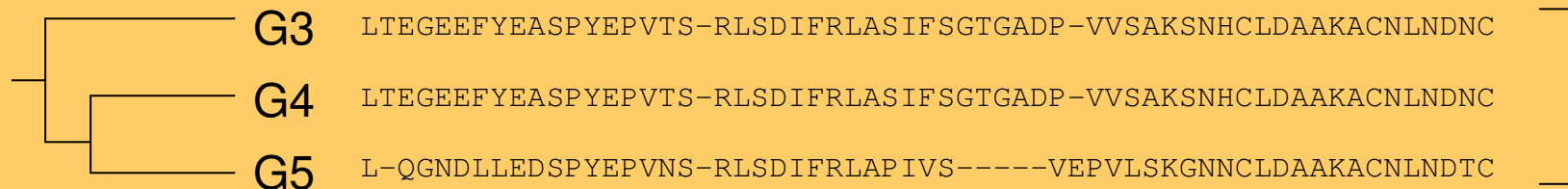  $\Updownarrow$

- How to find the location of the new sequence in the tree ?

| | |
|---|---|
| G1 | L-QGNDLLEDSPYEPVNSRLSDIFRAVPFISDVFQQVEHI-SKGNNCLDAAKACNLDDTC |
| G2 | L-QGNDLLEDSPYEPVNSRLSDIFRAVPFISDVFQQVEHI-SKGNNCLDAAKACNLDDTC |
| G3 | LTEGEEFYEASPYEPVTSRLSDIFRLASIFSGTGADP-VVSAKSNHCLDAAKACNLNDNC |
| G4 | LTEGEEFYEASPYEPVTSRLSDIFRLASIFSGTGADP-VVSAKSNHCLDAAKACNLNDNC |
| G5 | L-QGNDLLEDSPYEPVNSRLSDIFRLAPIVS-----VEPVLSKGNNCLDAAKACNLNDTC |
| | **LTEGEEFYEASPYEPVTSRLSDIFRLASIFSGTGTDP-AVSTKSNHCLDAAKACNLNDNC** |
| G6 | LAEGEEFYEASPYEPITSRLSDIFRLASIFSGM--DP-ATNSKSNHCLDAAKACNLNDNC |
| G7 | LTEGEEFYEASPYEPVTSRLSDIFRLASIFSGTGTDP-AVSTKSNHCLDAAKACNLNDNC |

If S and the outgroup are grouped:

- •Remove independant gaps in the founded block, and in S.

- •Align S and the founded block.

- •Restart the research in the founded block.

Stop when the research doesn't bring an amelioration.

```
G3  LTEGEEFYEASPYEPVTS-RLSDIFRLASIFSGTGADP-VVSAKSNHCLDAAKACNLNDNC

G4  LTEGEEFYEASPYEPVTS-RLSDIFRLASIFSGTGADP-VVSAKSNHCLDAAKACNLNDNC

G5  L-QGNDLLEDSPYEPVNS-RLSDIFRLAPIVS-----VEPVLSKGNNCLDAAKACNLNDTC
```

```
S   LTEGEEFYEA-------SSRLSDIFRLASIFSGTGTDP-AVSTKSNHCLDAAKACNLNDNC
```

HELiX
BIOINFORMATICS

UMR 5558
Biométrie, biologie
évolutive.

INRIA
RHÔNE-ALPES

# Adding the new sequence to the alignment

• How to insert the new sequence in the alignment, knowing its location ?

$\Updownarrow$

• What would be the alignment, if entirely recomputed with S, using the *progressive multiple alignment* method ?

➢ Recompute block alignments from current node to the root.

```
G1   LQ-GNDLLEDSPYEPVNS-RLSDIFRAVPFISDVFQQVEHIS-KGNNCLDAAKACNLDDTC
G2   LQ-GNDLLEDSPYEPVNS-RLSDIFRAVPFISDVFQQVEHIS-KGNNCLDAAKACNLDDTC
G3   LTEGEEFYEASPYEPVTS-RLSDIFRLASIFSGTGADP-VVSAKSNHCLDAAKACNLNDNC
G4   LTEGEEFYEASPYEPVTS-RLSDIFRLASIFSGTGADP-VVSAKSNHCLDAAKACNLNDNC
G5   L-QGNDLLEDSPYEPVNS-RLSDIFRLAPIVS-----VEPVLSKGNNCLDAAKACNLNDTC
S    LTEGEEFYEA-------SSRLSDIFRLASIFSGTGTDP-AVSTKSNHCLDAAKACNLNDNC
G6   LAEGEEFYEASPYEPITSRLSDIFRLASIFSGM---DP-ATNSKSNHCLDAAKACNLNDNC
G7   LTEGEEFYEASPYEPVTSRLSDIFRLASIFSGTGT-DP-AVSTKSNHCLDAAKACNLNDNC
```

# Prototype

# Discussion / Improvements

- Placement method: greedy algorithm. Errors are definitive, and generate other errors.
  - ➢ Local recomputing algorithm currently in test.

- Alignment of large blocks of sequences (several thousands) too complex.
  - ➢ Find a small set of representative sequences is a difficult problematic.

- Simplify the algorithm for simple cases.
  - ➢ Identify simple cases.
  - ➢ Preserve the alignment quality.