



# Βιοπληροφορική

## Ενότητα 3<sup>η</sup>: Πολλαπλή ευθυγράμμιση

Σ. Γκέλης  
Τμήμα Βιολογίας



Ευρωπαϊκή Ένωση  
Ευρωπαϊκό Κοινωνικό Ταμείο



ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ ΚΑΙ ΘΡΗΣΚΕΥΜΑΤΩΝ  
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης



ΕΥΡΩΠΑΪΚΟ ΚΟΙΝΩΝΙΚΟ ΤΑΜΕΙΟ

# Άδειες Χρήσης

- Το παρόν εκπαιδευτικό υλικό υπόκειται σε άδειες χρήσης Creative Commons.
- Για εκπαιδευτικό υλικό, όπως εικόνες, που υπόκειται σε άλλου τύπου άδειας χρήσης, η άδεια χρήσης αναφέρεται ρητώς.



# Χρηματοδότηση

- Το παρόν εκπαιδευτικό υλικό έχει αναπτυχθεί στα πλαίσια του εκπαιδευτικού έργου του διδάσκοντα.
- Το έργο «Ανοικτά Ακαδημαϊκά Μαθήματα στο Αριστοτέλειο Πανεπιστήμιο Θεσσαλονίκης» έχει χρηματοδοτήσει μόνο τη αναδιαμόρφωση του εκπαιδευτικού υλικού.
- Το έργο υλοποιείται στο πλαίσιο του Επιχειρησιακού Προγράμματος «Εκπαίδευση και Δια Βίου Μάθηση» και συγχρηματοδοτείται από την Ευρωπαϊκή Ένωση (Ευρωπαϊκό Κοινωνικό Ταμείο) και από εθνικούς πόρους.



# Περιεχόμενα ενότητας

- Ευθυγράμμιση ακολουθιών
- Αλγόριθμοι
- Δομή προγραμμάτων
- Ευθυγράμμιση (ή στοίχιση)
- Πολλαπλή ευθυγράμμιση (π.ε.)
- Οπτικοποιώντας μια ευθυγράμμιση



# Ευθυγράμμιση ακολουθιών: τι;

σύγκριση ακολουθιών (DNA/RNA - πρωτεΐνες)



που έχουμε εισαγωγή ή διαγραφή σε κάθε ακολουθία

καλύτερη ευθυγράμμιση:

- ✓ μέγιστος αριθμός ταύτισης
- ✓ ελάχιστος αριθμός κενών & μη όμοιων



# Ευθυγράμμιση ακολουθιών: γιατί;

ομοιότητα (similarity) μεταξύ των ακολουθιών



μάλλον υπάρχει και κάποιας μορφής ομολογία (homology)

**ομοιότητα (similarity):** μοιάζουν με βάση κάποιο χαρακτηριστικό τους, π.χ. στην ακολουθία

**ομολογία (homology):** δύο ή περισσότερες ακολουθίες έχουν ένα κοινό πρόγονο (εξελικτική ιστορία)

[βασική υπόθεση: τα γονίδια και οι πρωτεΐνες σχετίζονται εξελικτικά]



# Ευθυγράμμιση ακολουθιών: πώς; (1/2)

```
chite ---ADKPKRPLSAYMLWLNLSARES IKRENPDFK-VTEVAKKGGELWRGLKD
wheat --DPNKPKRAPSAFFVFMGEFREEFKQKNPKNKSVAAVGKAAGERWKSLSLSE
trybr  KKDSNAPKRAMTSFMFFSSDFRS----KHS DLS-IVEMSKAAGA AWKELGP
mouse  -----KPKRPRSA YNIYVSESFQ-----EAKDDS-AQGK LKLVNEAWKNLSP
      ***. ::: :. . . . : . . . * . * : *

chite  AATAKQNYIRALQ EYERNGG-
wheat  ANK LKGEYNKAI AAYNKGESA
trybr  AEKDKERYKREM-----
mouse  AKDDRIRYDNEMKSWEEQMAE
      * : .* . :
```

γίνεται και «με το χέρι»;



# Ευθυγράμμιση ακολουθιών: πώς; (2/2)

γιατί χρειάζεται η πληροφορική;

αποδοτικότερη διαχείριση & αποθήκευση της πληροφορίας

πολλές ακολουθίες [συνεχώς αλλάζουν & προστίθενται καινούργιες]

μείωση χρόνου υπολογισμών

υπολογισμός στατιστικών, πχ. e-value

διάθεση αποτελεσμάτων & βιβλιογραφίας στο διαδίκτυο





---

**...ας πάμε λίγο πίσω...**  
**[ιντερμέδιο πληροφορικής]**



# Αλγόριθμοι (1/2)

συγκεκριμένη ακολουθία βημάτων που οδηγεί στη λύση ενός προβλήματος, πχ. ευθυγράμμιση ακολουθιών πρέπει να:

1. βρίσκει την σωστή λύση σε αποδεκτά δεδομένα
2. υπολογίζει την λύση σε λογικό χρόνο

**χαρακτηριστικά αλγόριθμοι** [στην πληροφορική]:

1. πολυπλοκότητα (χρονική, χωρική)
2. ολοκλήρωση (completeness)
3. βέλτιστος (optimal)



# Αλγόριθμοι (2/2)

## Χαρακτηριστικά αλγορίθμου:

1. πολυπλοκότητα – συμβολίζεται με  $O()$   
χωρική: μνήμη που χρειάζεται  
χρονική: πόσος χρόνος χρειάζεται να “τρέχει”  
διαφοροποίηση ανάμεσα στην μέση πολυπλοκότητα

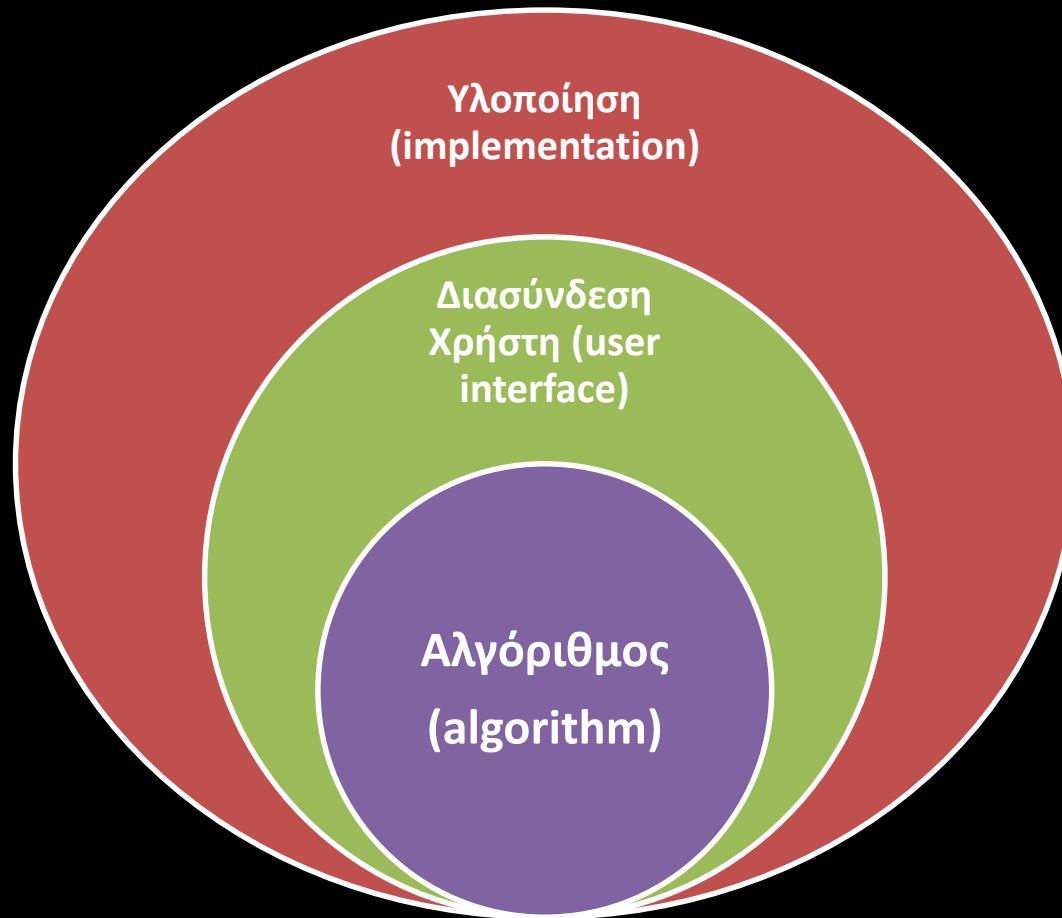
και

στη χειρότερη περίπτωση

1. ολοκλήρωση  
αν υπάρχει έστω μία λύση, εγγυάται ό,τι θα την βρει
1. βέλτιστος  
εγγυάται ό,τι θα βρει την καλύτερη δυνατή λύση



# Δομή προγραμμάτων (1/2)



# Δομή προγραμμάτων (2/2)

## αλγόριθμος:

BLAST

CLUSTALW2

κλπ

## υλοποίηση:

σε τι γλώσσα προγραμματισμού υλοποιήθηκε  
ποιές δομές δεδομένων χρησιμοποιήθηκαν  
τι δομές προγραμματισμού χρησιμοποιήθηκαν

## διασύνδεση χρήστη

γραφικά (Graphical User Interface - GUI)

γραμμή εντολών (Command Line Interface - CLI)



# 1. ...επιστροφή στην ευθεία...



# Ευθυγράμμιση (ή στοίχιση) (1/5)

ένα παράδειγμα



# Ευθυγράμμιση (ή στοίχιση) (2/5)

Τύποι ευθυγράμμισης:

- δυαδική (pairwise): δύο ακολουθίες μεταξύ τους
- πολλαπλή (multiple): πολλές ακολουθίες μαζί





# Ευθυγράμμιση (ή στοίχιση) (3/5)

Τύποι ευθυγράμμισης:

- καθολική (global - σε όλο το μήκος)
- τοπική (local - σε υποπεριοχές)



# Ευθυγράμμιση (ή στοίχιση) (4/5)

Αλγόριθμοι δυαδικής ευθυγράμμισης:

- Needleman-Wunsch, 1970 καθολική
- Smith -Waterman, 1981 τοπική
- FASTA (Pearson 1990) τοπική
- BLAST (Basic Local Alignment Search Tool), Altschul *et al.* 1990, 1997 τοπική



# Ευθυγράμμιση (ή στοίχιση) (5/5)

Διαφορές αλγορίθμων ευθυγράμμισης:

- διαφέρουν ως προς την πολυπλοκότητά τους
  - μικρή χωρική και χρονική πολυπλοκότητα, π.χ. BLAST
  - μεγάλη χωρική και χρονική πολυπλοκότητα, π.χ. Smith-Waterman
- κάποιοι δεν είναι βέλτιστοι, π.χ.
  - BLAST → μη βέλτιστος αλγόριθμος
  - Smith – Waterman → βέλτιστος αλγόριθμος



# Πολλαπλή ευθυγράμμιση (1/29)

```
Taxon 1  ---ADKPKRPLSAYMLWLNSARES IKRENPDFK-VTEVAKKGGELWRGLKD
Taxon 2  --DPNKPKRAPSAFFVFMGEFREEFKQKNPKNKSVAAVGKAAGERWKSLSLSE
Taxon 3  KKDSNAPKRAMTSFMFFSSDFRS----KHS DLS-IVEMSKAAGA AWKELGP
Taxon 4  -----KPKRPRSAYNIYVSESFQ----EAKDDS-AQGKIKLVNEAWKNLSP
          ***. ::: :. . . . : . . . * . * : *
```

```
Taxon 1  AATAKQNYIRALQEYERNGG-
Taxon 2  ANKCLKGEYINKAIAAYNKGESA
Taxon 3  AEKDKEFYKREM-----
Taxon 4  AKDDRIFDYDNEMKSWEEQMAE
          * : . * . :
```

- **δομικά κριτήρια**

- τα στοιχεία τοποθετούνται έτσι ώστε αυτά που παίζουν παρόμοιο ρόλο να είναι στην ίδια στήλη.

- **εξελικτικά κριτήρια**

- τα στοιχεία τοποθετούνται έτσι ώστε αυτά που έχουν τον ίδιο πρόγονο να βρίσκονται στην ίδια στήλη

- **κριτήρια ομοιότητας**

- όσο περισσότερα παρόμοια στοιχεία γίνεται να βρίσκονται στην ίδια στήλη



# Πολλαπλή ευθυγράμμιση (2/29)

- ευθυγράμμιση περισσότερων των δύο ακολουθιών
- ιδανική προσέγγιση: ευθυγράμμιση όλων των ακολουθιών μεταξύ τους ταυτόχρονα
- γρήγορη προσέγγιση: σταδιακή δυαδική ευθυγράμμιση  
(progressive pairwise alignment)



# Πολλαπλή ευθυγράμμιση (3/29)

Τι περιέχει μια π.ε. ;

- ό,τι θέλεις εσύ να περιέχει
- μια π.ε. μπορεί να ιδωθεί ως:
  - μια εξελικτική καταγραφή
  - μια σύνοψη οικογένειας πρωτεϊνών
  - μια συλλογή πειραμάτων που έκανε για σένα η φύση



# Πολλαπλή ευθυγράμμιση (4/29)

και σε τι  
χρησιμεύει;

## Main applications of multiple sequence alignments

εφαρμογή	Procedure
δημιουργία οικογένειας πρωτεϊνών	A good multiple alignment can help convincing you that an uncharacterized sequence is really a member of a protein family.
δημιουργία φυλογενετικών δέντρων – εύρεση εξελικτικών σχέσεων	If you carefully chose the sequences to include in your multiple alignment, you can reconstruct the history of these proteins.
αναγνώριση συντηρητικών περιοχών	By discovering very conserved positions you can identify a region that is characteristic of a function (in proteins or in nucleic acid sequences).
σύγκριση ακολουθιών μεταξύ διαφορετικών οργανισμών	It is possible to turn a multiple sequence alignment into a profile that describes a protein family or a protein domain. You can use this profile to scan databases for new members of the family.
εύρεση ρυθμιστικών στοιχείων	You can turn a DNA multiple alignment of a binding site into a weight matrix and scan other DNA sequences for potential similar binding sites.
πρόβλεψη δευτεροταγούς – τριτοταγούς δομής πρωτεΐνης	A good multiple alignment can give you an almost perfect prediction of your protein secondary structure for proteins or RNA. Sometimes it can also help building a 3-D model.

παράδειγμα



# Πολλαπλή ευθυγράμμιση (5/29)

Και που είναι το πρόβλημα;

Biology:

What is **a** good alignment?

Computation:

What is **the** good alignment?





# Πολλαπλή ευθυγράμμιση (6/29)

υπάρχουν πολλές (πόσες;) πιθανές ευθυγραμμίσεις και τι κάνουμε;  
διευθέτηση δύο ή περισσότερων αλληλουχιών (νουκλεοτιδικών ή πρωτεϊνικών σε ένα πλέγμα (μήτρα)

```
ELNGSLTLIMELMDMSMYDYI ---  
--NKSLTLIMHLMDINLYEYM ---  
E-EENAVLVLEFLRSDLAAVIRDG  
S-QKTVYMIFFEYADNDLSGLL ---  
V--FDVYMA MEYIENDVKNWI ---  
V-SSSLYLVFEYMDHDLVG-L ---
```

Στοιχεία (νουκλεοτίδια, αμινοξέα) της ίδιας σειράς προέρχονται από το ίδιο βιολογικό μακρομόριο (πρωτεΐνη ή νουκλεϊκό οξύ). Τα στοιχεία διευθετούνται με τη σειρά που εμφανίζονται στο μακρομόριο:

- από το N στο C άκρο στις πρωτεΐνες  
από το 5' στο 3' στα νουκλεϊκά οξέα



# Πολλαπλή ευθυγράμμιση (7/29)

κάθε κελί περιλαμβάνει ένα μόνο στοιχείο [είτε ένα στοιχείο είτε ένα κενό (gap)]



τα στοιχεία της ίδιας στήλης είναι

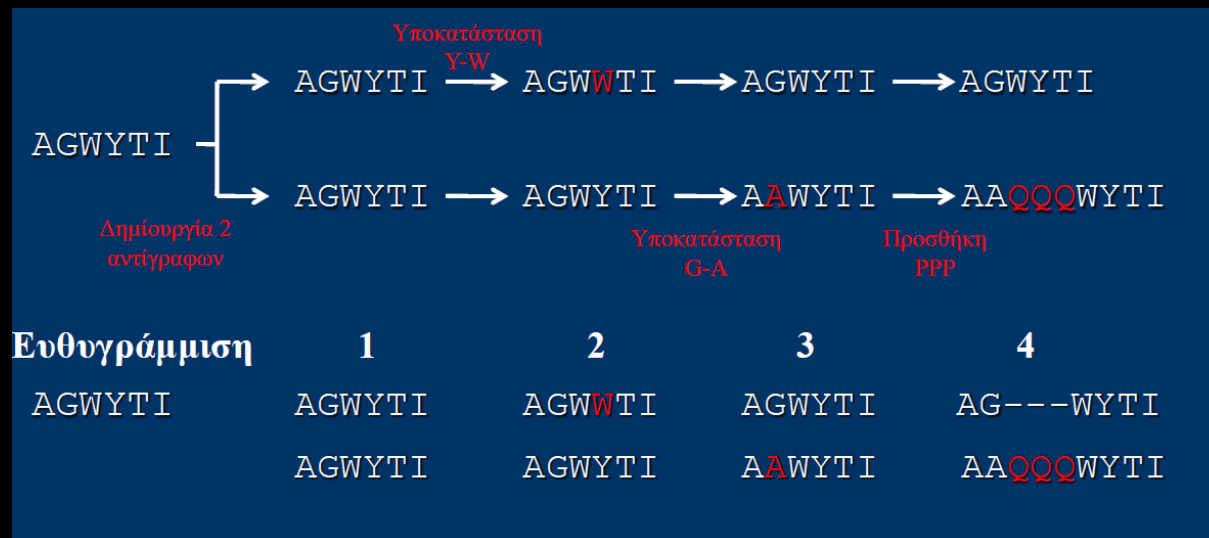
- δομικά ισοδύναμα
- εξελικτικά ισοδύναμα (ομόλογα)



# Πολλαπλή ευθυγράμμιση (8/29)

εξελικτική ισοδυναμία=ομολογία

αναφερόμενοι στην ίδια στήλη, η ιστορία κάθε στοιχείου θα πρέπει να αναζητηθεί στο αντίστοιχο στοιχείο της προγονικής αλληλουχίας, όπου κάθε αλλαγή οφείλεται σε σημειακές αλλαγές



# Πολλαπλή ευθυγράμμιση (9/29)

ένα ζεύγος αλληλουχιών μπορεί να ευθυγραμμιστεί γράφοντας την μία αλληλουχία κάτω από την άλλη με τέτοιο τρόπο ώστε να μεγιστοποιηθεί ο αριθμός των νολεοτιδίων που ταιριάζουν, βάζοντας κενά (gaps) στην μια ή στην άλλη αλληλουχία όταν απαιτείται

```
AF486227      TACGAAAACACCACCCAATCCTAAGAA
AF486228      TACGAAAACACGACCCAATCCTAAAAA
AF486223      TACGAAAACACCACCCTATCCTAAAAA
```

η ευθυγράμμιση γίνεται συνήθως με ειδικά υπολογιστικά πακέτα, που χρησιμοποιούν συγκεκριμένους αλγόριθμους. Οι περισσότεροι αλγόριθμοι αρχίζουν συγκρίνοντας την ομοιότητα των αλληλουχιών ανά ζεύγη, και ευθυγραμμίζουν πρώτα τις δύο αλληλουχίες με τη μεγαλύτερη ομοιότητα. Οι άλλες αλληλουχίες, βάσει της σειράς ομοιότητας, προστίθενται σταδιακά.



# Πολλαπλή ευθυγράμμιση (10/29)

Όταν σε μια ομάδα αλληλουχιών έχουν προστεθεί κάποια κενά, τότε το τελικό alignment συχνά βελτιώνεται από τον ίδιο τον ερευνητή με manual editing. Η απόκτηση μιας καλής ευθυγράμμισης είναι ίσως το πιο σημαντικό βήμα ώστε να εκτιμήσουμε ένα σωστό φυλογενετικό δέντρο

```
AF486227      TACGAA--AACACCACC---CAATCCTAAGAA
AF486228      TACGAA--AACACGACCGGGCAATCCTAAAAA
AF486223      TACGAATTAACACCACCGGGCTATCCTAAAAA
```

Είναι αναγκαίο να ορίσουμε τον αριθμό των gaps ώστε το τελικό αποτέλεσμα να έχει βιολογική υπόσταση

Για το λόγο αυτό χρησιμοποιείται ένα σύστημα σκοραρίσματος όπου τα ταιριάσματα παίρνουν ένα θετικό βαθμό και τα κενά ένα αρνητικό, που είναι γνωστό ως «gap penalty»



# Πολλαπλή ευθυγράμμιση (11/29)

υπάρχουν πολλές (πόσες;) πιθανές ευθυγραμμίσεις

για παράδειγμα συγκρίνετε

```
-GCGC-ATGGATTGAGCGA  
TGCGCCATTGAT-GACC-A
```

με

```
-----GCGCATGGATTGAGCGA  
TGCGCC----ATTGATGACCA--
```

ποια είναι η καλύτερη;

**παράδειγμα 2**



# Πολλαπλή ευθυγράμμιση (12/29)

και που είναι (πάλι) το πρόβλημα;

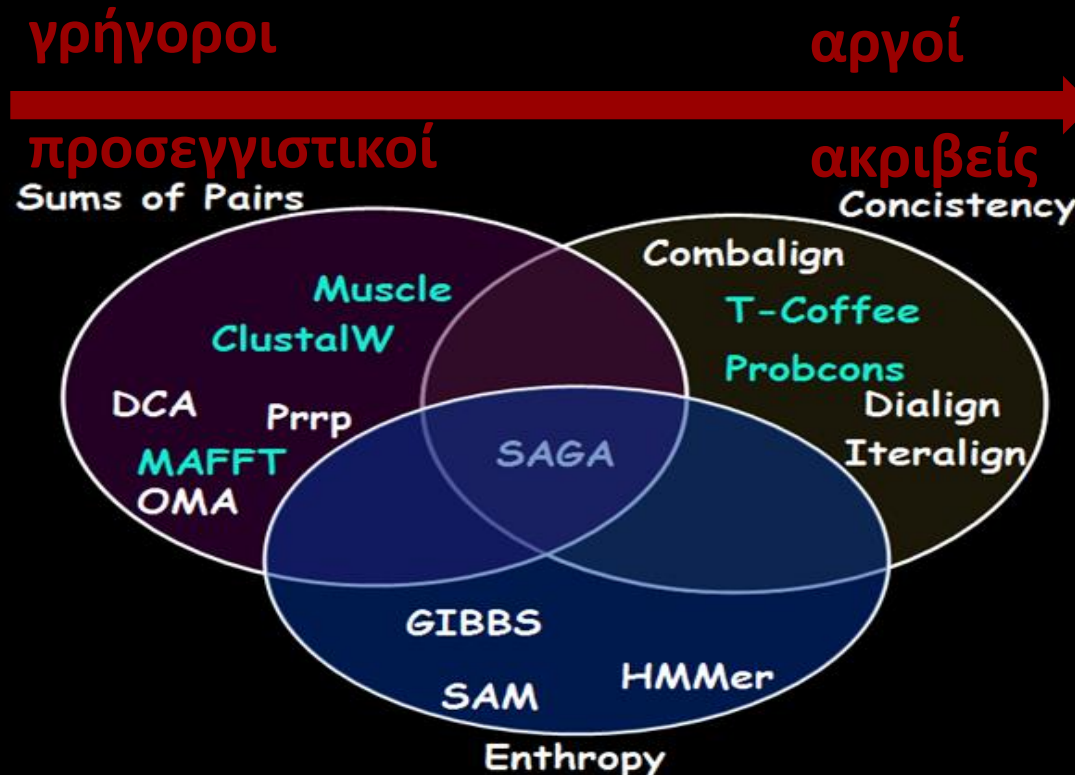
- ο αριθμός των πιθανών συνδυασμών ευθυγράμμισης αυξάνει ραγδαία σε σχέση με των αριθμό των ακολουθιών
- χρονοβόρα διαδικασία
- συχνά δύσκολο να πούμε ποια ευθυγράμμιση είναι η καλύτερη



# Πολλαπλή ευθυγράμμιση (13/29)

και ποιες είναι οι λύσεις;

μια κάποια ταξινόμηση των αλγορίθμων





# Πολλαπλή ευθυγράμμιση (14/29)

## τύποι αλγόριθμων π.ε

- σταδιακοί [progressive]: **ClustalW**
- επαναληπτικοί [iterative]: **Muscle**
- consistency based: **T-Coffee** και **Probcons**



# Πολλαπλή ευθυγράμμιση (15/29)

---

Πολλοί αλγόριθμοι online:

<http://www.ebi.ac.uk/Tools/msa/>



# Πολλαπλή ευθυγράμμιση (16/29)

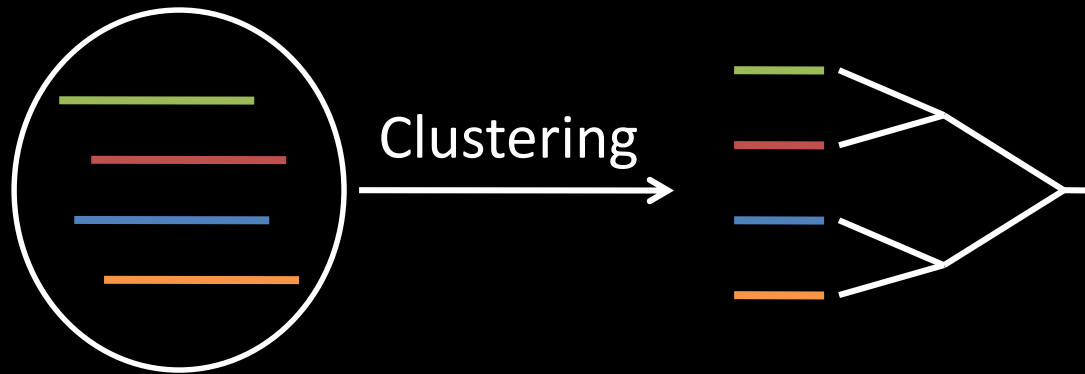
## ClustalW2

- Global multiple sequence alignment
- δυνατότητα παραγωγής φυλογενετικών δέντρων
- ευθυγράμμιση είναι γρήγορη (progressive alignment)
- δυνατότητα επεξεργασίας με JalView της ευθυγράμμισης
- μέχρι 500 ακολουθίες ή 10MB αρχεία
- <http://www.ebi.ac.uk/Tools/clustalw2>



# Πολλαπλή ευθυγράμμιση (17/29)

πως δουλεύει το ClustalW2;



1. Εκτελεί δυαδική ευθυγράμμιση για όλα τα πιθανά ζεύγη των ακολουθιών – Οι ομοιότητες και τα αποτελέσματα των ευθυγραμμίσεων αποθηκεύονται σε ένα πίνακα



# Πολλαπλή ευθυγράμμιση (18/29)

πως δουλεύει το ClustalW2;

2. Μετατρέπει τον πίνακα με τις ομοιότητες σε εξελικτική απόσταση μεταξύ των ακολουθιών
3. Φτιάχνει ένα βοηθητικό δέντρο (guide tree) που δείχνει πως θα συνδυαστούν οι προηγούμενες δυαδικές ευθυγραμμίσεις
4. Κάνει σταδιακή ευθυγράμμιση των ακολουθιών προσθέτοντας σταδιακά κάθε ακολουθία με βάση την σειρά που προέκυψε από το 2ο βήμα



# Πολλαπλή ευθυγράμμιση (19/29)

The screenshot displays the ClustalW2 web interface on the EMBL-EBI website. The page is titled "ClustalW2 - Multiple Sequence Alignment" and provides instructions for using the tool. It is divided into four main steps:

- STEP 1 - Enter your input sequences:** A text area for entering or pasting a set of DNA or protein sequences. A dropdown menu is set to "DNA". Below the text area is a "Choose File" button.
- STEP 2 - Set your Pairwise Alignment Options:** Includes an "Alignment Type" selector (Slow/Fast) and "Fast Pairwise Alignment Options" for parameters like KTUP, WINDOW LENGTH, SCORE TYPE, TOPDIAG, and PAIRGAP.
- STEP 3 - Set your Multiple Sequence Alignment Options:** Includes options for DNA Weight Matrix, GAP OPEN, GAP EXTENSION, GAP DISTANCES, NO END GAPS, ITERATION, NUMITER, and CLUSTERING. There are also "OUTPUT Options" for FORMAT and ORDER.
- STEP 4 - Submit your job:** A checkbox for "Be notified by email" and a "Submit" button.

Alignment: χρήση γρήγορου αλγόριθμου (σταδιακή δυαδική ευθυγράμμιση) ή όχι

Iteration: ο αριθμός των επαναλήψεων μέχρι να τερματίσει ο αλγόριθμος και να μας δώσει αποτελέσματα

Fast Pairwise Alignment: επιλογές που σχετίζονται με γρήγορη ευθυγράμμιση KTUP, Window, Score, Topdiag, Pairgap

Gaps: βαθμοί σε διάφορους τύπους κενών

<http://www.ebi.ac.uk/Tools/msa/clustalw2/help/index.html#matrix>

Εικόνα 1: το user interface του ClustalW2.



# Πολλαπλή ευθυγράμμιση (20/29)

## παράδειγμα - πρωτεΐνη fosb

### εισαγωγή σε Fasta μορφοποίηση:

#### >FOSB\_MOUSE Protein fosB

```
MFQAFPGDYD SGSRCS SSPS AESQYLSSVD SFGSPPTAAA SQECAGLGEM
PGSFVPTVTA ITTSQDLQWL VQPTLISSMA QSQGQPLASQ PPAVDPYDMP
GTSYSTPGLS AYSTGGASGS GGPSTSTTTS GPVSARPARARPRRPREETL
TPEEEKRRV RRERNKLA AAKCRNRRRELT DRLQAETDQL EEEKAELESE
IAELQKEKER LEFVLVAHKP GCKIPYEEGP GPGPLAEVRD LPGSTSAKED
GFGWLLPPPP PPPLPFQSSR DAPPNLTASL FTHSEVQVLG DPFVNVPSY
TSSFVLTCP E VSAFAGAQR T SGSEQPSDPL NSPSLLAL
```

#### >FOSB\_HUMAN Protein fosB

```
MFQAFPGDYD SGSRCS SSPS AESQYLSSVD SFGSPPTAAA SQECAGLGEM
PGSFVPTVTA ITTSQDLQWL VQPTLISSMA QSQGQPLASQ PPVDPYDMP
GTSYSTPGMS GYSSGGASGS GGPSTSGTTS GPGPARPARARPRRPREETL
TPEEEKRRV RRERNKLA AAKCRNRRRELT DRLQAETDQL EEEKAELESE
IAELQKEKER LEFVLVAHKP GCKIPYEEGP GPGPLAEVRD LPGSAPAKED
GFSWLLPPPP PPPLPFQTSQ DAPPNLTASL FTHSEVQVLG DPFVNVNPSY
TSSFVLTCP E VSAFAGAQR T SGSDQPSDPL NSPSLLAL
```



# Πολλαπλή ευθυγράμμιση (21/29)

The screenshot shows the EMBL-EBI ClustalW2 Results page. The page title is "ClustalW2 Results" and it includes navigation tabs for "Alignments", "Result Summary", "Guide Tree", "Submission Details", and "Submit Another Job". Under the "Alignments" tab, there are two buttons: "Download Alignment File" and "Show Colors", with the latter being circled in blue. The main content area displays a multiple sequence alignment of FOSB\_MOUSE and FOSB\_HUMAN sequences. The alignment is shown in a text-based format with sequence identifiers and positions on the right. The sequences are as follows:

```
FOSB_MOUSE      MFQAFPGDYDSGSRCS SSPSAESQYLSSVDSFGSPPTAAASQECAGLGEMPGSFVPTVTA  60
FOSB_HUMAN      MFQAFPGDYDSGSRCS SSPSAESQYLSSVDSFGSPPTAAASQECAGLGEMPGSFVPTVTA  60
*****

FOSB_MOUSE      ITTSQDLQWLQPTLISSMAQSQGQPLASQPPAVDPYDMPGTSYSTPGLSAYSTGGASGS  120
FOSB_HUMAN      ITTSQDLQWLQPTLISSMAQSQGQPLASQPPVDPYDMPGTSYSTPGMSGYSSGGASGS  120
*****

FOSB_MOUSE      GGPSTSTTTSGPV SARPARARPRRPREETLTPEEEKRRVRRERNKLA AAKCRNRRRELT  180
FOSB_HUMAN      GGPSTSGTTS GPGPARPARARPRRPREETLTPEEEKRRVRRERNKLA AAKCRNRRRELT  180
*****

FOSB_MOUSE      DRLQ AETDQLEEEKAELESEIAELQKEKERLEFVLVAHKPGCKI PYE EGP GPGPLAEVRD  240
FOSB_HUMAN      DRLQ AETDQLEEEKAELESEIAELQKEKERLEFVLVAHKPGCKI PYE EGP GPGPLAEVRD  240
*****

FOSB_MOUSE      LPGSTSAKEDGFGWLLPPPPPPPLPFQSSRDAPPNLTASLFT HSEVQV LGDPFPV VVSPSY  300
FOSB_HUMAN      LPGSAPAKEDGFSWLLPPPPPPPLPFQTSQDAPPNLTASLFT HSEVQV LGDPFPV VVNSPY  300
****

FOSB_MOUSE      TSSFVLTCP EVSAFAGAQR TSGSEQSDPLNSP SLLAL  338
FOSB_HUMAN      TSSFVLTCP EVSAFAGAQR TSGSDQSDPLNSP SLLAL  338
*****
```

Εικόνα 2: Αποτελέσματα ευθυγράμμισης του ClustalW2.





# Πολλαπλή ευθυγράμμιση (22/29)

## Clustal Ω: νέο Clustal

<http://www.clustal.org/>

<http://www.ebi.ac.uk/Tools/msa/clustalo/help/index.html>

<http://www.ebi.ac.uk/Tools/msa/clustalo/help/faq.html#11>



# Πολλαπλή ευθυγράμμιση (23/29)

## Προβλήματα ClustalW2

- όταν οι ακολουθίες μοιάζουν σε ορισμένα σημεία αλλά όχι σ' όλο τους το μήκος
- όταν μία από τις ακολουθίες έχει ένα μεγάλο τμήμα ενσωματωμένο σε σχέση με τις άλλες
- αν μια ακολουθία έχει ένα επαναλαμβανόμενο τμήμα πολλές φορές



# Πολλαπλή ευθυγράμμιση (24/29)

## T-Coffee [Tree based Consistency Objective Function For alignment Evaluation]

- Πραγματοποιεί πολλαπλή ευθυγράμμιση ακολουθιών
- Μπορεί να ευθυγραμμίσει ακολουθίες από πρωτεΐνες, DNA και RNA



# Πολλαπλή ευθυγράμμιση (25/29)

πως δουλεύει το T-Coffee;

συνδυάζει αποτελέσματα από διαφορετικές μεθόδους ευθυγράμμισης: πχ. ClustalW2, δομική ευθυγράμμιση (structural alignment)

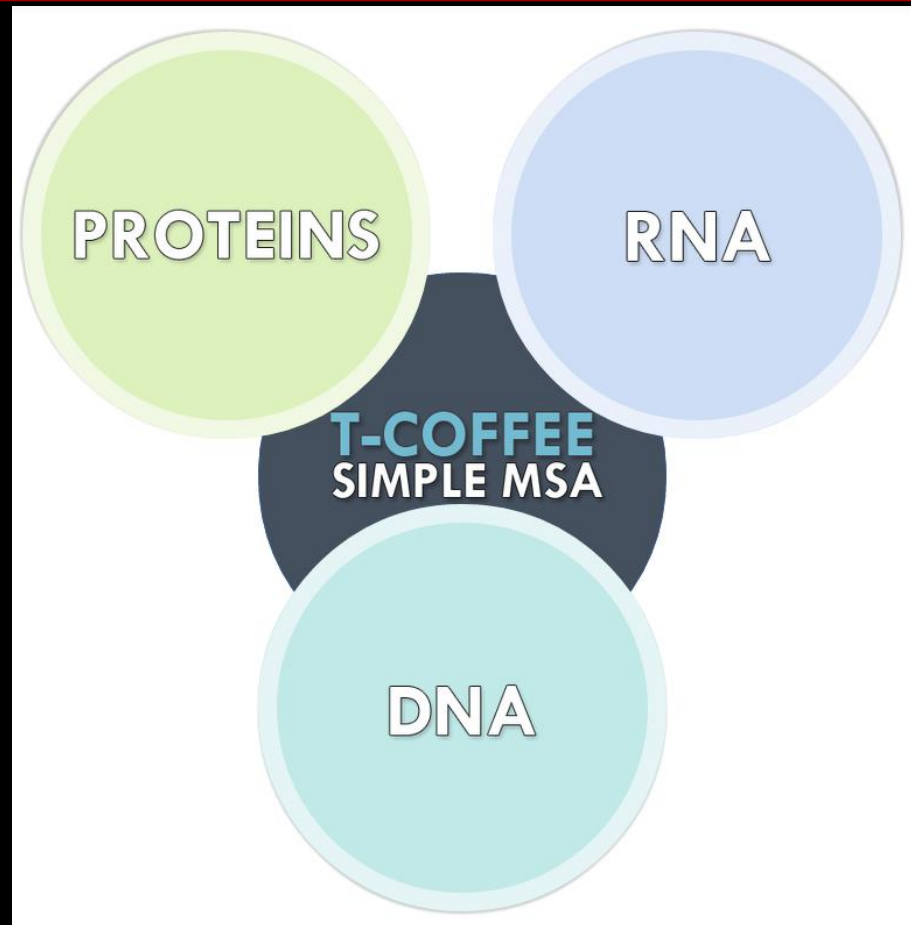
αν δώσουμε απλά κάποιες ακολουθίες τότε το T-Coffee θα τις συγκρίνει ανά δυο (global & local alignments )

θα συνδυάσει τις ευθυγραμμίσεις σε μία πολλαπλή ευθυγράμμιση

αντί της χρήσης πίνακα δημιουργείται μια βιβλιοθήκη που προκύπτει από την τοπική και ολική στοίχιση όλων των ζευγών ακολουθιών. Η βιβλιοθήκη αυτή, που τροποποιείται καθώς ο αλγόριθμος προχωρά, χρησιμοποιείται για την βαθμολόγηση



# Πολλαπλή ευθυγράμμιση (26/29)



Εικόνα 3: το user interface του T-Coffee

<http://tcoffee.crg.cat/apps/tcoffee/index.html>



# Πολλαπλή ευθυγράμμιση (27/29)

Εισάγουμε τις ακολουθίες που θέλουμε να ευθυγραμμίσουμε (<http://tcoffee.org.cat/apps/tcoffee/do:regular>) και το T-coffee κάνει την ευθυγράμμιση

```
MSA
The multiple sequence alignment result as produced by T-coffee.

T-COFFEE, Version_11.00.8cbe486 (2014-08-12 22:05:29 - Revision 8cbe486 - Build 477)
Cedric Notredame
SCORE=993
*
BAD AVG GOOD
*
E.coli : 99
S.flexneri : 99
S.dysenteriae : 99
S.sonnei : 99
S.enterica : 99
S.glossinidius : 98
P.aeruginosa : 97
P.luminescens : 99
Y.pseudotubercu : 98
Y.pestis : 98
cons : 99

E.coli ATGATC-----CCCTTACAACATGGACTGATCCTCGCGGCAATCTTATTTCGTTCTTGGCTTA
S.flexneri ATGATC-----CCCTTACAACATGGACTGATCCTCGCGGCAATCTTATTTCGTTCTTGGCTTA
S.dysenteriae ATGATC-----CCCTTACAACATGGACTGATCCTCGCGGCAATCTTATTTCGTTCTTGGCTTA
S.sonnei ATGATC-----CCCTTACAACATGGACTGATCCTCGCGGCAATCTTATTTCGTTCTTGGCTTA
S.enterica ATGATC-----CCCTTACAACATGGACTGATCCTCGCGGCAATCTTATTTCGTTCTTGGCTTA
S.glossinidius ATG-----ATCCCGTATACACACGGTCTATCCTGGCCGCTATCCTGTTTGTGCTGGGCTG
P.aeruginosa ATGAACGCAATACCACTGGAAACACGGCCTGGCCCTCGCCAGCGTCTGTTTCGCCCTCGGACTG
P.luminescens ATGATA-----CCTCTTACAACATGGACTGATTTTGGCGGCAATCTTATTTCGTTTGGGCTG
Y.pseudotubercu ATG-----ATCCCTCTACAACATGGCCTGATTCCTGGCGGCCATTCTGTTTGTGCTGGGCTA
Y.pestis ATG-----ATCCCTCTACAACATGGCCTGATTCCTGGCGGCCATTCTGTTTGTGCTGGGCTA

cons *** ** * ** * ** * ** * ** * ** *
```

Εικόνα 4: Αποτελέσματα ευθυγράμμισης του T-coffee.



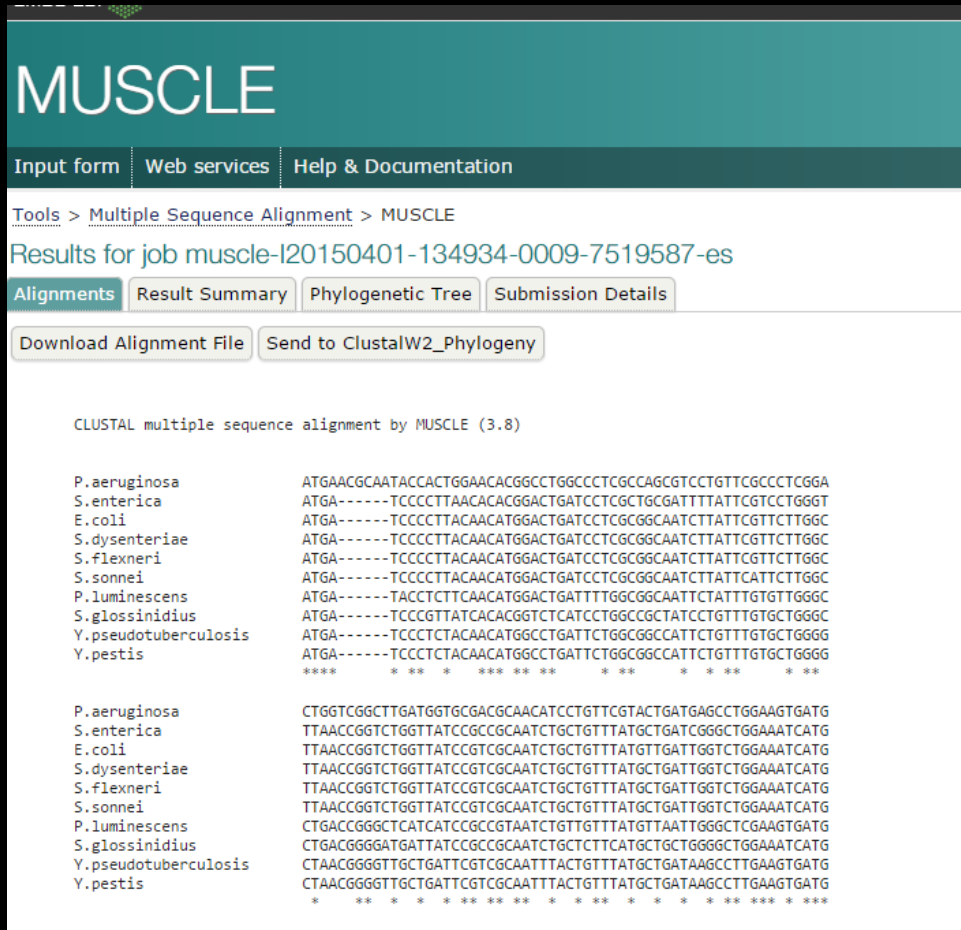
# Πολλαπλή ευθυγράμμιση (28/29)

## **MUSCLE** [Multiple Sequence Comparison by Log-Expectation]

- ακρίβεια που σε περιπτώσεις ξεπερνάει ClustalW & T-Coffee
- πολλαπλή ευθυγράμμιση πρωτεϊνών
- αρκετά γρήγορος αλγόριθμος
- 2 - 3 φορές πιο γρήγορος από CLUSTALW2 κατά μέσο όρο
- μπορεί να ευθυγραμμίσει 1000 πρωτεΐνες μέσου μήκους 282 αμινοξέων σε 21 sec. σε ένα μέσης υπολογιστικής ισχύος μηχάνημα



# Πολλαπλή ευθυγράμμιση (29/29)



MUSCLE

Input form | Web services | Help & Documentation

Tools > Multiple Sequence Alignment > MUSCLE

Results for job muscle-I20150401-134934-0009-7519587-es

Alignments | Result Summary | Phylogenetic Tree | Submission Details

Download Alignment File | Send to ClustalW2\_Phylogeny

CLUSTAL multiple sequence alignment by MUSCLE (3.8)

```
P.aeruginosa      ATGAACGCAATACCACTGGAACACGGCCTGGCCCTCGCCAGCGTCTGTTCCGCCCTCGGA
S.enterica        ATGA-----TCCCCTTAACACACGGACTGATCCTCGCTGCGATTTTATTCGCTCGGGT
E.coli            ATGA-----TCCCCTTACAACATGGACTGATCCTCGCGGCAATCTTATTCGTTCTTGGC
S.dysenteriae    ATGA-----TCCCCTTACAACATGGACTGATCCTCGCGGCAATCTTATTCGTTCTTGGC
S.flexneri       ATGA-----TCCCCTTACAACATGGACTGATCCTCGCGGCAATCTTATTCGTTCTTGGC
S.sonnei         ATGA-----TCCCCTTACAACATGGACTGATCCTCGCGGCAATCTTATTCATTCTTGGC
P.luminescens    ATGA-----TACCTCTTCAACATGGACTGATTTGGCGGCAATCTATTGTGTGGG
S.glossinidius   ATGA-----TCCCGTTATCACACGGTCTCATCTGGCCGCTATCCTGTTTGTGCTGGG
Y.pseudotuberculosis ATGA-----TCCCTCTACAACATGGCCTGATCTGGCGGCAATCTGTTTGTGCTGGGG
Y.pestis         ATGA-----TCCCTCTACAACATGGCCTGATCTGGCGGCAATCTGTTTGTGCTGGGG
**** * * * * * * * * * * * * * * * * * * * * * * * * * * * * * *

P.aeruginosa      CTGGTCGGCTTGATGGTGCAGCAGCAACATCCTGTTCTGTAATGATGAGCCTGGAAGTGATG
S.enterica        TTAACCGGTCTGGTTATCCGTCGCAATCTGCTGTTTATGCTGATCGGGCTGGAAATCATG
E.coli            TTAACCGGTCTGGTTATCCGTCGCAATCTGCTGTTTATGTTGATTGGTCTGGAAATCATG
S.dysenteriae    TTAACCGGTCTGGTTATCCGTCGCAATCTGCTGTTTATGCTGATTGGTCTGGAAATCATG
S.flexneri       TTAACCGGTCTGGTTATCCGTCGCAATCTGCTGTTTATGCTGATTGGTCTGGAAATCATG
S.sonnei         TTAACCGGTCTGGTTATCCGTCGCAATCTGCTGTTTATGCTGATTGGTCTGGAAATCATG
P.luminescens    CTGACGGGATGATTATCCGTCGCAATCTGCTTCTCATGCTGCTGGGGCTGGAAATCATG
S.glossinidius   CTGACGGGATGATTATCCGTCGCAATCTGCTTCTCATGCTGCTGGGGCTGGAAATCATG
Y.pseudotuberculosis CTAACGGGTTGCTGATTGCTCGCAATTTACTGTTTATGCTGATAAGCCTTGAAGTGATG
Y.pestis         CTAACGGGTTGCTGATTGCTCGCAATTTACTGTTTATGCTGATAAGCCTTGAAGTGATG
* * * * * * * * * * * * * * * * * * * * * * * * * * * * * *
```

Εικόνα 5: Αποτελέσματα ευθυγράμμισης του MUSCLE.





# Οπτικοποιώντας μια ευθυγράμμιση (1/4)

## JalView

- γραφική αναπαράσταση της ευθυγράμμισης
- δυνατότητα επεξεργασίας της ευθυγράμμισης

<http://www.jalview.org/>



# Οπτικοποιώντας μια ευθυγράμμιση (2/4)

JalView

πολλές λειτουργίες

τα λογισμικά τείνουν να κάνουν  
το ένα τη «δουλειά» του άλλου



# Οπτικοποιώντας μια ευθυγράμμιση (3/4)

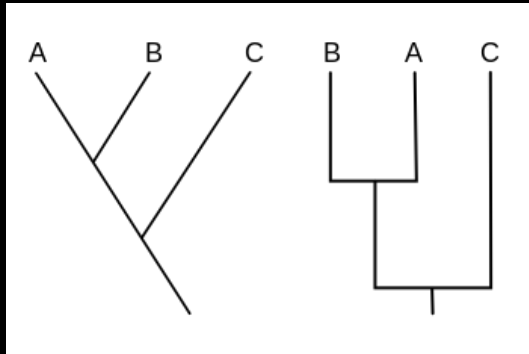
## JaView

- επεξεργασία αλληλουχιών
- πλοήγηση στις αλληλουχίες
- “conservation colouring”
- απεικόνιση τρισδιάστατων δομών πρωτεϊνών



# Οπτικοποιώντας μια ευθυγράμμιση (4/4)

## λίγο πριν απ' τα δέντρα

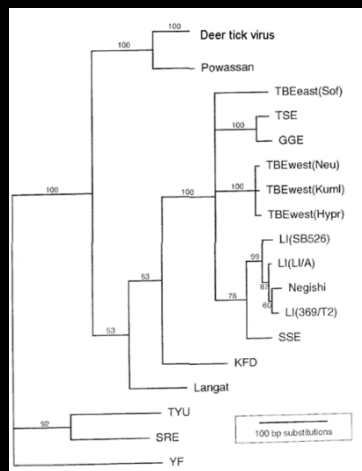


**Εικόνα 6:** κλαδόγραμμα : όλοι οι κλάδοι έχουν το ίδιο μήκος,

[http://commons.wikimedia.org/wiki/File:Identical\\_cladograms.svg](http://commons.wikimedia.org/wiki/File:Identical_cladograms.svg)

by Alexei Kouprianov, CC-BY-SA-3.0,

<https://creativecommons.org/licenses/by-sa/3.0/deed.en> .



**Εικόνα 7:** φυλόγραμμα : οι κλάδοι έχουν διαφορετικό μήκος (δείχνει και χρονική μεταβολή),

[http://commons.wikimedia.org/wiki/File:Tick-borne\\_encephalitis\\_phylogram\\_simplified.gif](http://commons.wikimedia.org/wiki/File:Tick-borne_encephalitis_phylogram_simplified.gif)

by Telford et al.



# Ανακεφαλαίωση (1/2)

## να θυμάστε:

- ✓ κάθε πρόγραμμα = υλοποίηση συγκεκριμένου αλγορίθμου
- ✓ οι αλγόριθμοι έχουν συγκεκριμένα χαρακτηριστικά, π.χ. πολυπλοκότητα, ολοκλήρωση, κτλ.
- ✓ πολλαπλή ευθυγράμμιση = ευθυγράμμιση πολλών ακολουθιών ( $>2$ ) ταυτόχρονα
- ✓ κυριότερα προγράμματα για πολλαπλή ευθυγράμμιση: ClustalW2, T-Coffee, MUSCLE

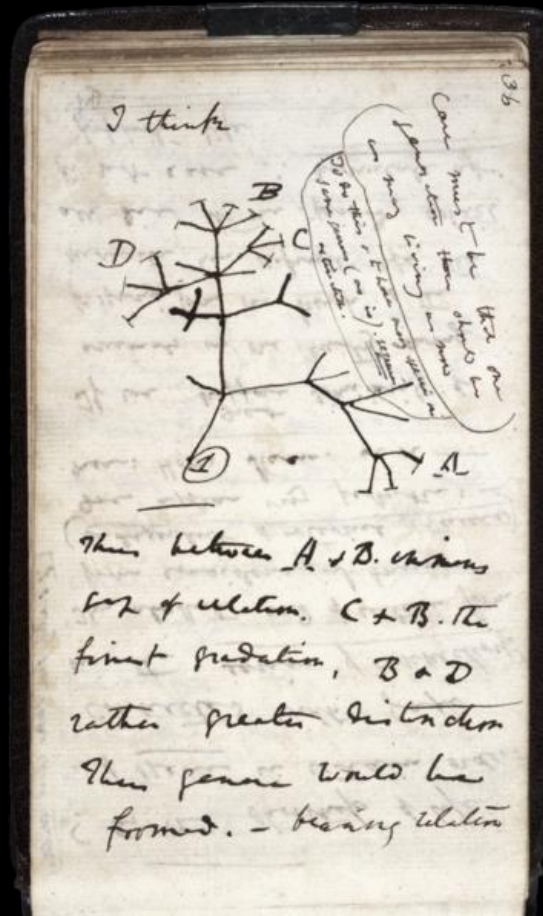


# Ανακεφαλαίωση (2/2)

## συνήθη βήματα στην π.ε.:

- ✓ ανάκτηση αλληλουχιών
- ✓ ClustalW2 ή T-coffee ή MUSCLE
- ✓ JalView ή αντίστοιχο εργαλείο
- ✓ εξαγωγή συμπερασμάτων
- ✓ φυλογενετικό δέντρο





είστε έτοιμοι για το πρώτο φυλογενετικό σας δέντρο!

αν και το πρώτο φυλογενετικό δέντρο φτιάχτηκε πριν εφευρεθούν οι υπολογιστές και πριν ανακαλυφθούν τα γονίδια

Εικόνα 8: Το δέντρο του Δαρβίνου,

[http://commons.wikimedia.org/wiki/File:Darwin\\_Tree\\_1837.png](http://commons.wikimedia.org/wiki/File:Darwin_Tree_1837.png)



# Βιβλιογραφία - Πηγές

1. Αρχεία Βοήθειας (Help files) των προγραμμάτων που παρουσιάστηκαν
2. Lecture notes by Per Kraulis, KTH Bioinformatics Oct-Dec 2001, <http://www.avatar.se/molbioinfo2001/index.html>
3. An Introduction to Bioinformatics Algorithms (Computational Molecular Biology) by Neil C. Jones, Pavel A. Pevzner, ISBN: 978-0262101066
4. Robert C Edgar, «MUSCLE: a multiple sequence alignment method with reduced time and space complexity», BMC Bioinformatics 2004 vol. 5, pp. 113
5. Notredame C., Higgins D., Heringa J., «T-Coffee: A novel method for multiple sequence alignments». *Journal of Molecular Biology* 2000, 302: 205-217
6. L. Bromham. 2008. Reading the story in DNA. Oxford University Press
7. Cédric Notredame, T-Coffee: What's New in The Grinder  
<http://www.google.gr/url?sa=t&rct=j&q=&esrc=s&source=web&cd=1&ved=0CDqQFjAA&url=http%3A%2F%2Fwww.tcoffee.org%2FPresentations%2Fmcoffee.ppt&ei=fAhzT67zBtDZ4QSyloGHDw&usq=AFQjCNFWjw-WddSAqN4TF81VLzQquHzlKq>
8. ClustalW and ClustalX version 2 Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H, Valentin F, Wallace IM, Wilm A, Lopez R, Thompson JD, Gibson TJ and Higgins DG *Bioinformatics* 2007 23(21): 2947-2948





# Σημείωμα Αναφοράς

Copyright Αριστοτέλειο Πανεπιστήμιο Θεσσαλονίκης,  
Σπύρος Γκέλης, «Βιοπληροφορική, Πολλαπλή  
ευθυγράμμιση». Έκδοση: 1.0. Θεσσαλονίκη 2014.  
Διαθέσιμο από τη δικτυακή διεύθυνση:  
[http://opencourses.auth.gr/eclass\\_courses](http://opencourses.auth.gr/eclass_courses).



# Σημείωμα Αδειοδότησης

Το παρόν υλικό διατίθεται με τους όρους της άδειας χρήσης Creative Commons Αναφορά - Παρόμοια Διανομή [1] ή μεταγενέστερη, Διεθνής Έκδοση. Εξαιρούνται τα αυτοτελή έργα τρίτων π.χ. φωτογραφίες, διαγράμματα κ.λ.π., τα οποία εμπεριέχονται σε αυτό και τα οποία αναφέρονται μαζί με τους όρους χρήσης τους στο «Σημείωμα Χρήσης Έργων Τρίτων».



Ο δικαιούχος μπορεί να παρέχει στον αδειοδόχο ξεχωριστή άδεια να χρησιμοποιεί το έργο για εμπορική χρήση, εφόσον αυτό του ζητηθεί.

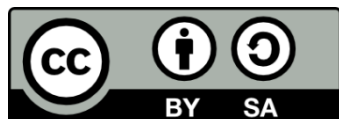
[1] <http://creativecommons.org/licenses/by-sa/4.0/>





# Σας ευχαριστώ!

Σπύρος Γκέλης  
Τμήμα Βιολογίας Α.Π.Θ.



Ευρωπαϊκή Ένωση  
Ευρωπαϊκό Κοινωνικό Ταμείο



ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ ΚΑΙ ΘΡΗΣΚΕΥΜΑΤΩΝ  
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης



ΕΥΡΩΠΑΪΚΟ ΚΟΙΝΩΝΙΚΟ ΤΑΜΕΙΟ



# Τέλος ενότητας

Επεξεργασία: Στυλιανή Μηνούδη  
Θεσσαλονίκη, Εαρινό εξάμηνο 2014



Ευρωπαϊκή Ένωση  
Ευρωπαϊκό Κοινωνικό Ταμείο



ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ ΚΑΙ ΘΡΗΣΚΕΥΜΑΤΩΝ  
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης



ΕΥΡΩΠΑΪΚΟ ΚΟΙΝΩΝΙΚΟ ΤΑΜΕΙΟ

# Διατήρηση Σημειωμάτων

Οποιαδήποτε αναπαραγωγή ή διασκευή του υλικού θα πρέπει να συμπεριλαμβάνει:

- το Σημείωμα Αναφοράς
- το Σημείωμα Αδειοδότησης
- τη δήλωση Διατήρησης Σημειωμάτων

μαζί με τους συνοδευόμενους υπερσυνδέσμους.

