

Πρωτεϊνική Μηχανική

Ιωάννης Παυλίδης
Τμήμα Χημείας | Πανεπιστήμιο Κρήτης

E-mail: ipavlidis@uoc.gr

Τηλ.: +30 2810 54-5130 | Γραφείο Γ211

ΕΠΙΣΗΜΗ ΣΕΛΙΔΑ: [HTTP://WWW.CHEMISTRY.UOC.GR/PAVLIDIS/](http://www.chemistry.uoc.gr/pavlidis/)

Εισαγωγή στη βιοπληροφορική

Βιοπληροφορική

Η Βιοπληροφορική είναι η χρήση Η/Υ για την απόκτηση, διαχείριση και ανάλυση βιολογικών δεδομένων

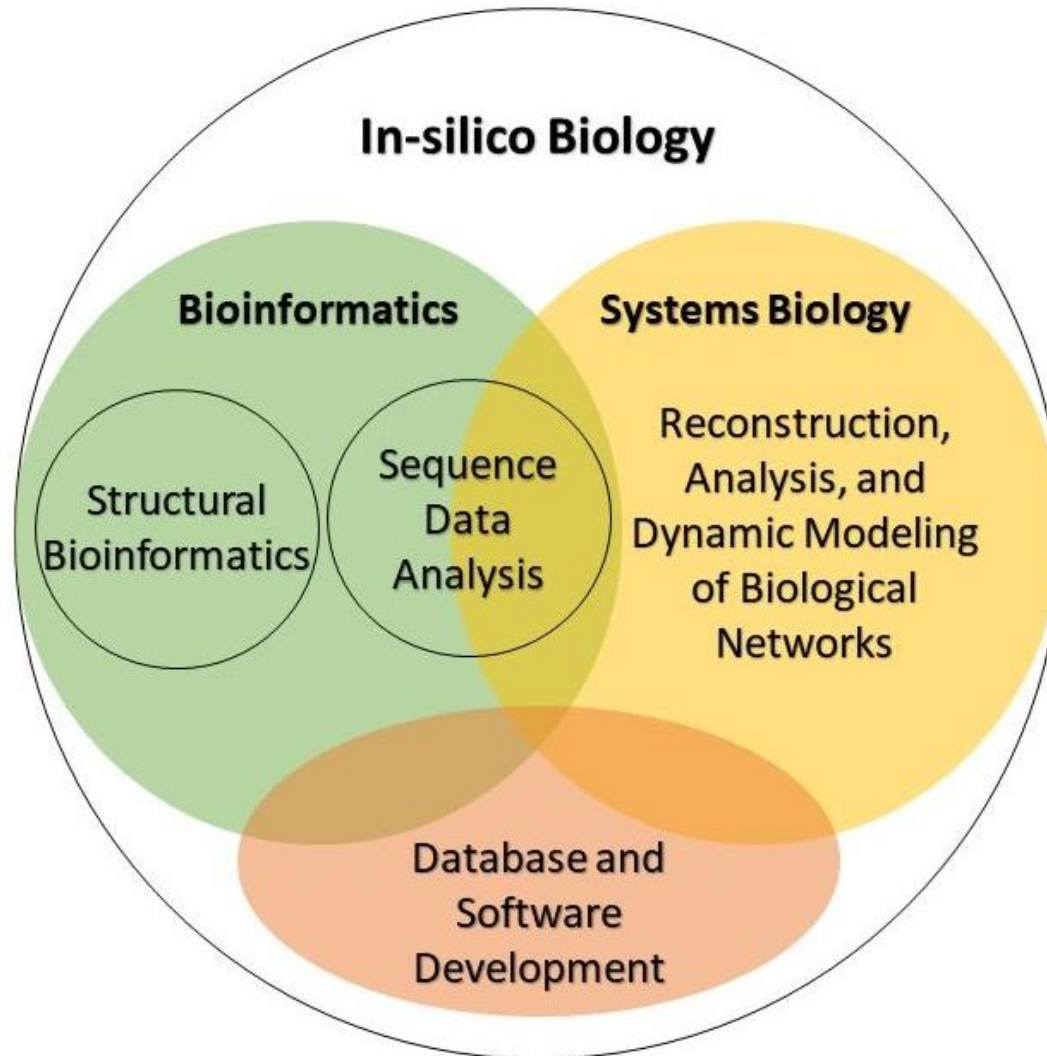
Περιλαμβάνει πολλά πεδία της επιστήμης υπολογιστών:

- ❖ Σχεδιασμός βάσεων δεδομένων
- ❖ Προγραμματισμός
- ❖ Διαχείριση δικτύων και πληροφοριακών συστημάτων

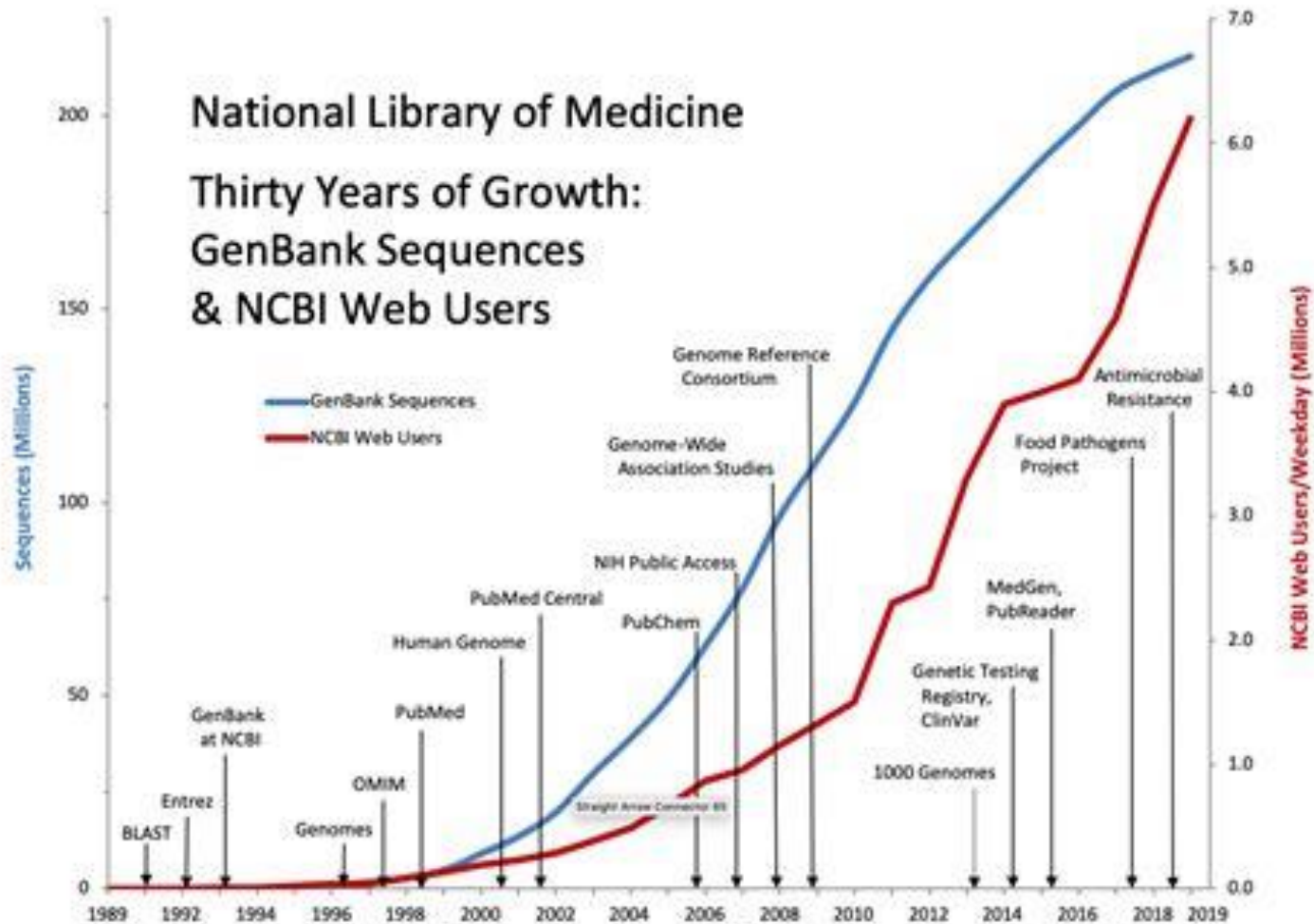
Μπορεί να εφαρμοστεί σε πολλά πεδία:

- ❖ Μοριακή βιολογία
- ❖ Γενετική
- ❖ Πρωτεϊνική βιοτεχνολογία

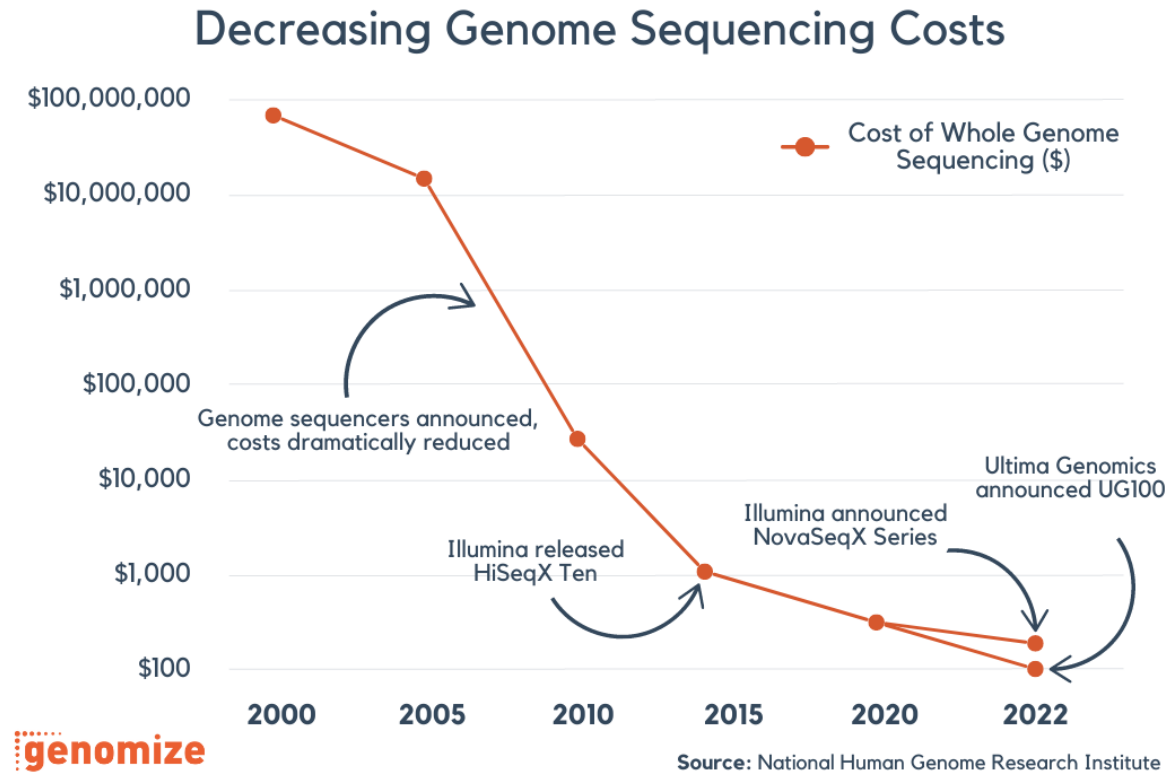
Βιοπληροφορική



Ραγδαία παραγωγή δεδομένων



Ραγδαία παραγωγή δεδομένων



Sequencing by Expansion (SBX) platform

«Αντιστροφή» της επιστημονικής μεθόδου

Η επιστημονική μέθοδος:

Σχηματισμός Υπόθεσης  Συλλογή δεδομένων για επιβεβαίωση

Η «νέα» επιστημονική μέθοδος με την βιοπληροφορική:

Ανάλυση (πολλών) δεδομένων  Δημιουργία υπόθεσης

Μία θεωρία γεννάται από τα δεδομένα

Data-mining: Νέα υπολογιστικά εργαλεία παρέχουν πρόσβαση και νέους τρόπους ανάλυσης της αχανούς πληροφορίας στις βάσεις δεδομένων

Ιστορία της βιοπληροφορικής

1960s: Πρώτα «βήματα» της βιοπληροφορικής

Η **Dayhoff** δημιούργησε την **πρώτη πρωτεϊνική βάση δεδομένων** και το πρώτο πρόγραμμα για την **συναρμολόγηση αλληλουχιών**

Δημιουργός και του κώδικα των αμινοξέων με ένα γράμμα



Margaret Oakley Dayhoff

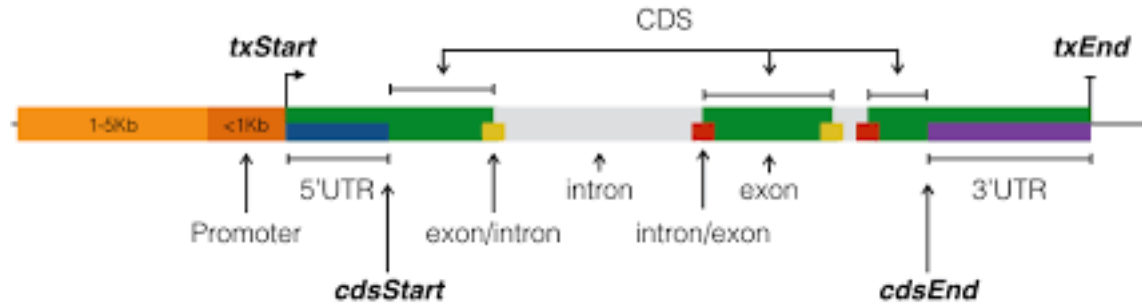


Αλγόριθμοι για ανάλυση δεδομένων

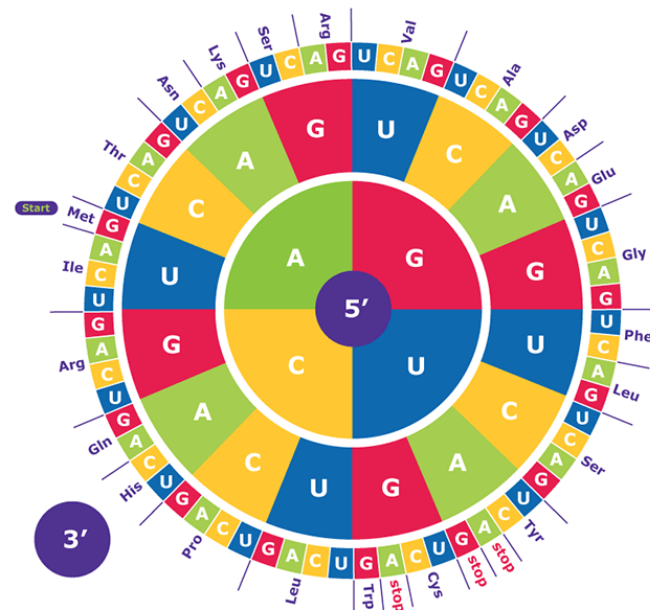


Ανάλυση αλληλουχιών DNA

Δομή του DNA και περιοχές στο γονιδίωμα



Μετάφραση της αλληλουχίας του DNA σε αμινοξική αλληλουχία



Ανάλυση αλληλουχιών DNA

```
CGCCTTTGACCGCCGCGGCTTTGGCCGCTCGGACCAACCCTGGACCGGCAACGACTACGACACCT
TCGCCGACGACATCGCCCAGTTGATCGAACACCTGGACCTCAAGGAGGTGACCCTGGTGGGCTTC
TCCATGGGCGGGCGGCGATGTGGCCCCGCTACATCGCCCCGCCACGGCAGCGCACGGGTGGCCGGCCT
GGTGTGCTGGGCGCCGTCACCCCGCTGTTTCGGCCAGAAGCCCGACTATCCGCAGGGTGTCCCGC
TCGATGTGTTTCGCAAGGTTCAAGACTGAGCTGCTGAAGGATCGCGCGCAGTTCATCAGCGATTTC
AACGCACCGTTCTATGGCATCAACAAGGGCCAGGTCTCTCCCAAGGCGTGCAGACCCAGACCCT
GCAAATCGCCCTGCTGGCCTCGCTCAAGGCCACGGTGGATTG
```

- ❖ Από που ξεκινάω να διαβάζω;
- ❖ Πόσους τρόπους έχω να διαβάσω την συγκεκριμένη αλληλουχία;
- ❖ Είναι προκαρυωτικό ή ευκαρυωτικό γονίδιο (αν είναι όντως κωδικοποιούσα περιοχή);

Hints:

Πάντα χρησιμοποιείτε την γραμματοσειρά **Courier** για αλληλουχίες!

Κρατά ίσο πλάτος για κάθε γράμμα

.fasta: Γενικός τύπος αρχείων για αλληλουχίες DNA/πρωτεϊνών.

Συνήθως χρησιμοποιούνται κεφαλαία για αμινοξέα, μικρά για DNA

.gb: Γενικός τύπος αρχείων για αλληλουχίες DNA με σημειώσεις

Ανάλυση αλληλουχιών DNA

Πως αναγνωρίζω ένα ανοιχτό πλαίσιο ανάγνωσης;

- 1. Kozak consensus sequence:** Αλληλουχία κοντά στο κωδικόνιο έναρξης στους ευκαρυώτες CC(A/G)CCATGG
- 2. RBS:** Αλληλουχία πρόσδεσης του ριβοσώματος στους προκαρυώτες
- 3. Χρήση κωδικονίων** διαφορετική μεταξύ κωδικοποιούσας και μη κωδικοποιούσας περιοχής στο DNA
- 4. Θέσης Wobble (3η)** προτίμηση σε G/C ή A/T ανά οργανισμό
- 5. Στοιχίση** με γνωστές αλληλουχίες

Στοίχιση αλληλουχιών

Πρακτικές αναζήτησης σε βάσεις δεδομένων:

- Αναζήτηση κειμένου:

“show all adrenaline receptors”

- Αναζήτηση με ομοιότητα αλληλουχιών:

Βρες όλες τις παρόμοιες αλληλουχίες με την αλληλουχία του ανθρώπινου υποδοχέα αδρεναλίνης

- ❖ Ο δεύτερος τρόπος είναι πιο ορθολογικός στην επιστήμη υπολογιστών
- ❖ Οι σημειώσεις μπορεί να είναι λάθος

Σύγκριση αλληλουχιών

Hamming distance: Σε δεδομένη αλληλουχία ίσου μήκους, είναι ο αριθμός των διαφορών (χαμηλότερο σκορ είναι καλύτερο)

Χωρίς στοίχιση, HD:3

AGGVLI IQV

AGGVLI QVG

Με στοίχιση, HD:1

AGGVLI IQV

AGGVLI -QVG

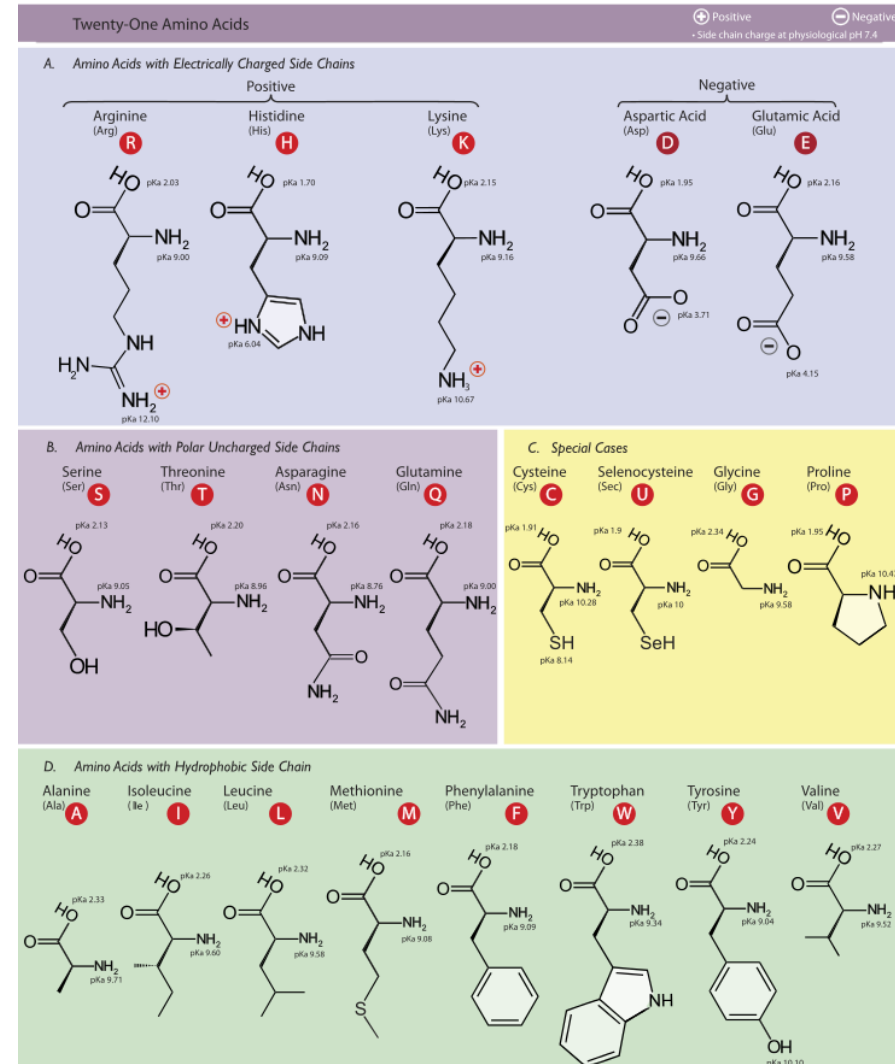
Σε πραγματικές αλληλουχίες, το ερώτημα είναι πιο περίπλοκο, καθώς οι αλληλουχίες είναι μεγαλύτερες και όχι ίδιου μεγέθους

Ταυτότητα και ομοιότητα

Η εύρεση ταυτότητας σε αλληλουχίες είναι εύκολη υπολογιστικά

Όταν μελετάμε και την ομοιότητα, το ερευνητικό ερώτημα γίνεται πιο περίπλοκο

- ❖ Πως στοιχίζουμε;
- ❖ Που επιτρέπουμε μία ασυμφωνία (mismatch) για να βελτιστοποιήσουμε το σκορ;



Σκορ στοίχισης

Το σκορ της στοίχισης έρχεται από συγκεκριμένους πίνακες:

- όταν ταιριάζουν οι θέσεις, έχει θετικό σκορ
- όταν δεν ταιριάζουν, έχει ένα αρνητικό σκορ

Ίση επίδραση

	A	C	T	G	
A	1	-1	-1	-1	-2
C	-1	1	-1	-1	-2
T	-1	-1	1	-1	-2
G	-1	-1	-1	1	-2
	-2	-2	-2	-2	

value of matching G with A

value of matching C with gap

value of matching G with G

Πιθανή ασυμμετρία

Letters from string 2

	A	C	T	G	
A	1	-2	-1	-1	-2
C	-2	2	-1	-1	-2
T	-1	-1	3	-4	-7
G	-1	-1	-3	1	-2
	-2	-2	-6	-2	

Different match costs

Different mismatch costs

Different gap cost (depends on which string)

Letters from string 1

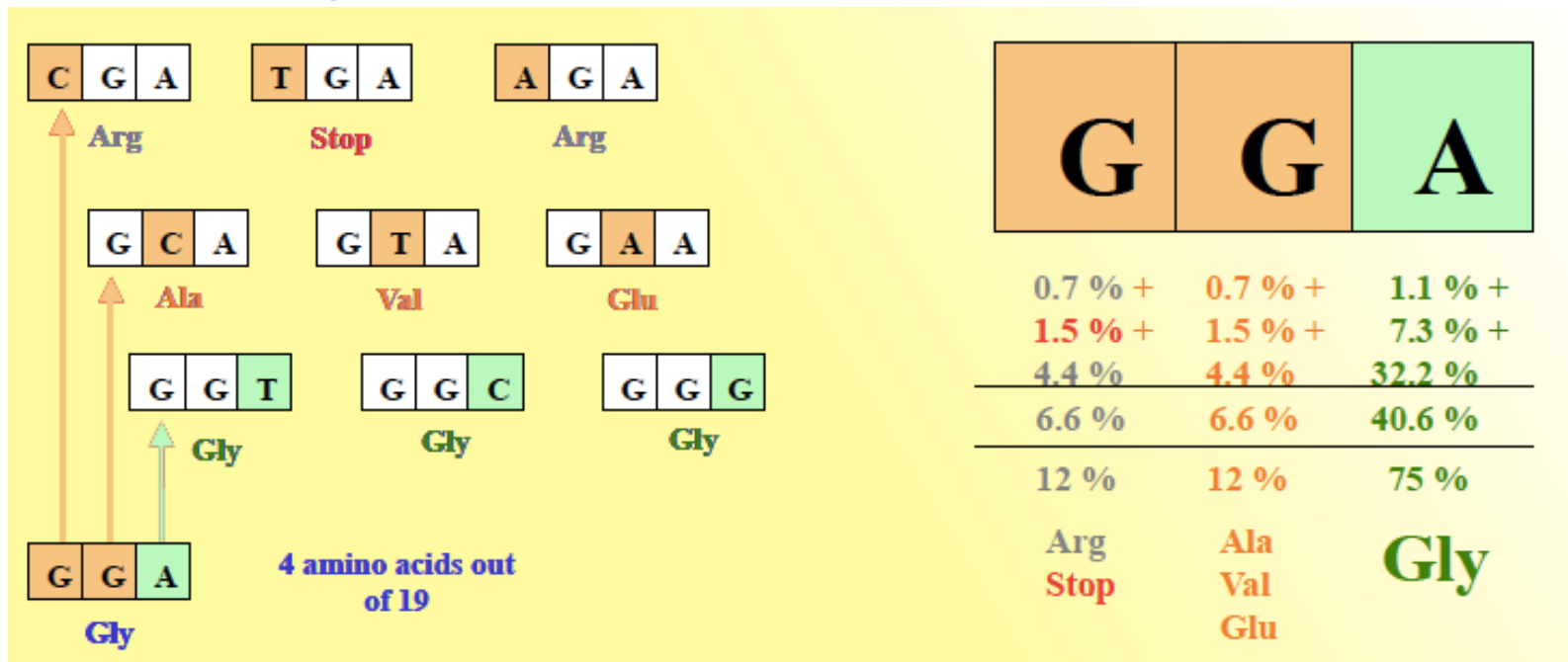
Αλγόριθμοι στοίχισης

Αλγόριθμος Dayhoff : βασισμένος στο PAM matrix

PAM: Point accepted mutation

Πιθανότητα της αλλαγής ενός αμινοξέος σε ένα άλλο σε 100 θέσεις

Πίνακες για μεγαλύτερες αλληλουχίες έχουν προετοιμαστεί με πολλαπλασιασμό πινάκων



Στοιχίση βασισμένη στη δομή

Στοιχίστε τις παρακάτω αλληλουχίες:

PI**F**ENHGT

PI**F**ENHGT

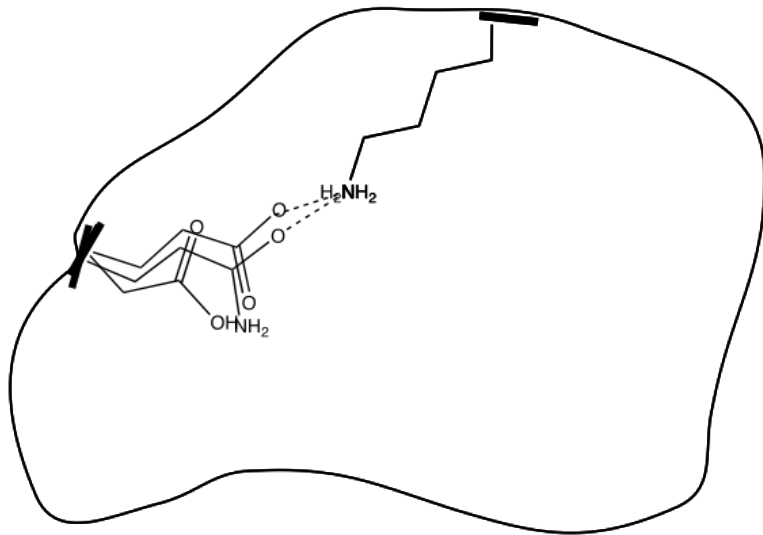
PI**F**ENHGT

PI-**DQ**HGT

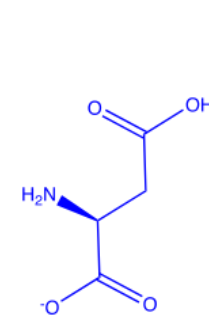
PI**D**-**Q**HGT

PI**DQ**-HGT

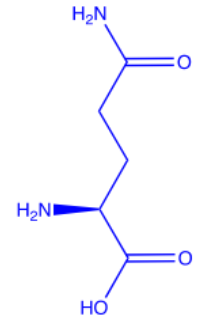
Αν το γλουταμικό οξύ παρέχει ένα δεσμό υδρογόνου σημαντικό για τη σωστή αναδίπλωση της πρωτεΐνης;



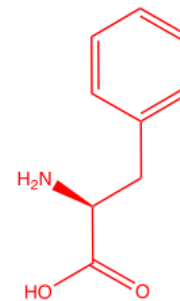
D



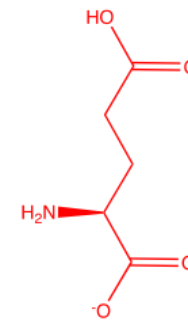
Q



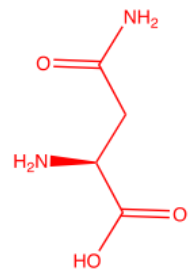
F



E



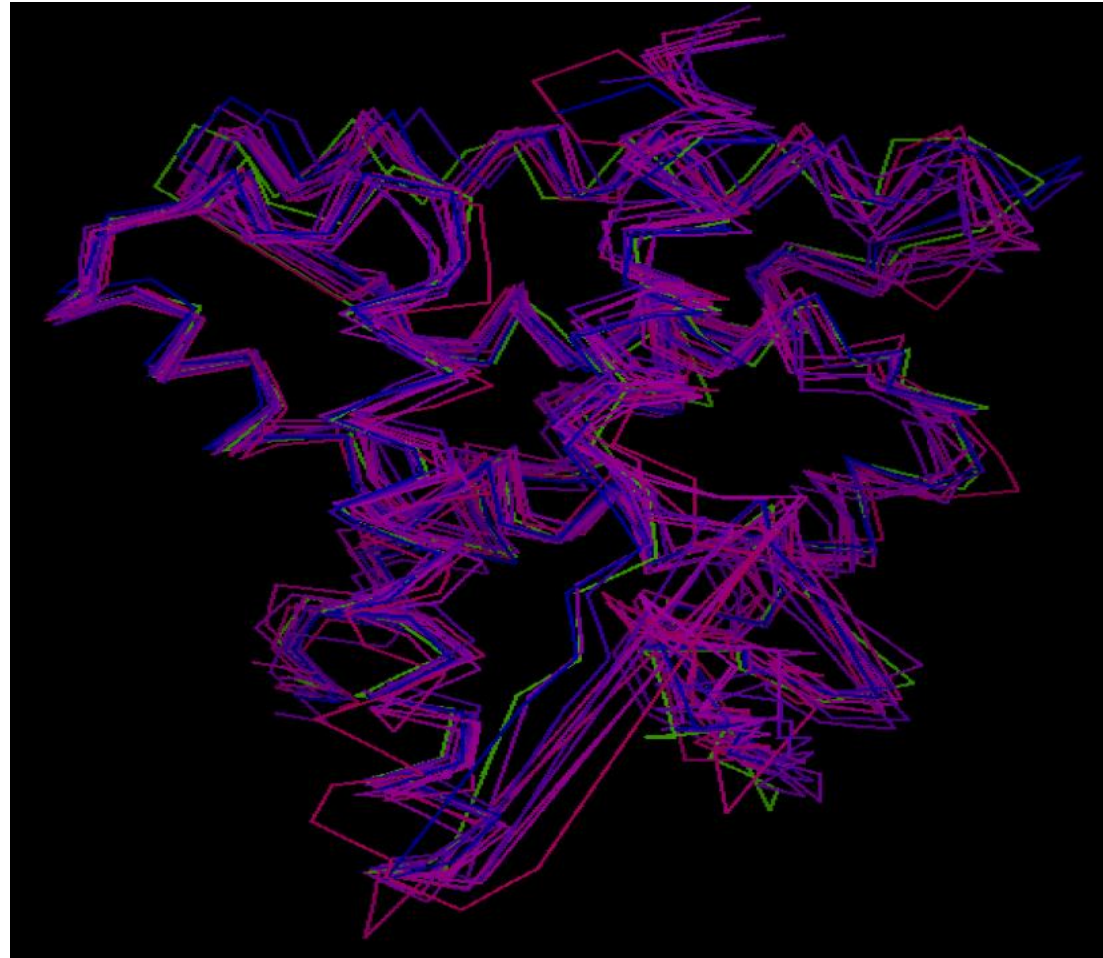
N



Στοιχίση βασισμένη στη δομή

Στις στοιχίσεις πάντα θα πρέπει να λαμβάνουμε υπόψη την δομή

Η στοιχίση 2 πρωτεϊνών συγκρίνει ταυτόχρονα και την 3D δομή τους



Στοίχιση βασισμένη στη δομή

Συντηρημένα αμινοξέα – λειτουργικότητα πρωτεϊνών

Υπάρχουν συντηρημένες και μη συντηρημένες περιοχές

Χρησιμοποιώντας αυτή τη γνώση, μπορούμε να κάνουμε στοίχιση αλληλουχιών με πολύ χαμηλή ταυτότητα

```
QLIPPLINLLMSIEPDVIYA LLTSLNQLGERQLLSVVKWSKSLPGFRNLHIDDQITLIQYSWMSLMVFLGW
DEEWELIKTVTEAHVATNAQ AFSHFTKIITPAITRVVDFAKKLPMFCELPCEDDQIILLKGCCMEIMSLRAAV
PEVGELIEKVRKAHQETFPA LWDKFSELSTKCIIKTVEFAKQLPGFTTLTIADQITLLKAACLDILILRICT
PQLEELITKVSKAHQETFPS LWDKFSSELATKCIIKIVEFAKRLPGFTGLSIADQITLLKAACLDILMLRICT
ANEDMPVERILEAELAVEPK PVTNICQAADKQLFTLVEWAKRIPHFSELPLDDQVILIRAGWNELLIASFH
EEQRMMIRELMDARMKTFDT LLPHMADMSTYMFKGIISYAKVISYFRDLPIEDQISLRKGAAFELCQLRFNT
ADLKSLAKRIYEAYLKNFNM IFHCCQCTSVETVTELTELAKAIPGFANLDLNDQVTLLKYGVYEAIFAMLSS
KPYNKIVSHLLVAEPEKIYA ALTTLCDLADRELVVIIGWAKHIPGFSTLSLADQMSLIQSAWMEILILGVVY
QLTPTLVSLLEVIEPEVLYA IMTTLNMLGGRQVIAAVKWAKAIPGFRNLHLDDQMTLIQYSWMSLMFAFALGW
LSPEQLVLTLLAEPPHVL I MMMSLTKLADKELVHMISWAKKIPGFVELSLFDQVRLIESCWMEVLMMLMW
SEIDRIAQNIIKSHLETCQY LWQQCAIQITHAIQYVVEFAKRITGFMELCQNDQIILLKSGCLEVVLVRMCR
PEEWDLIHIATEAHRSTNAK AFSEFTKIITPAITRVVDFAKKLPMFSELPCEDQIILLKGCCMEIMSLRAAV
```

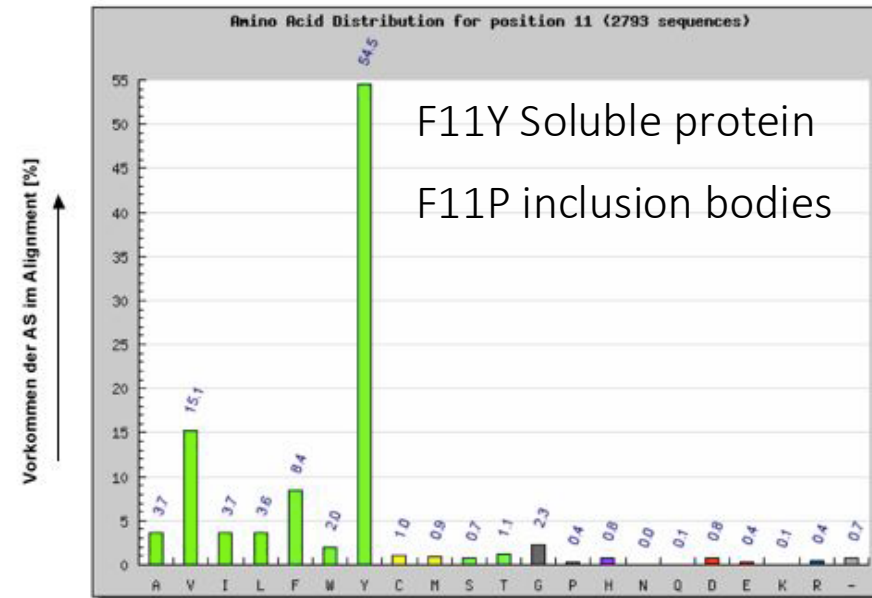
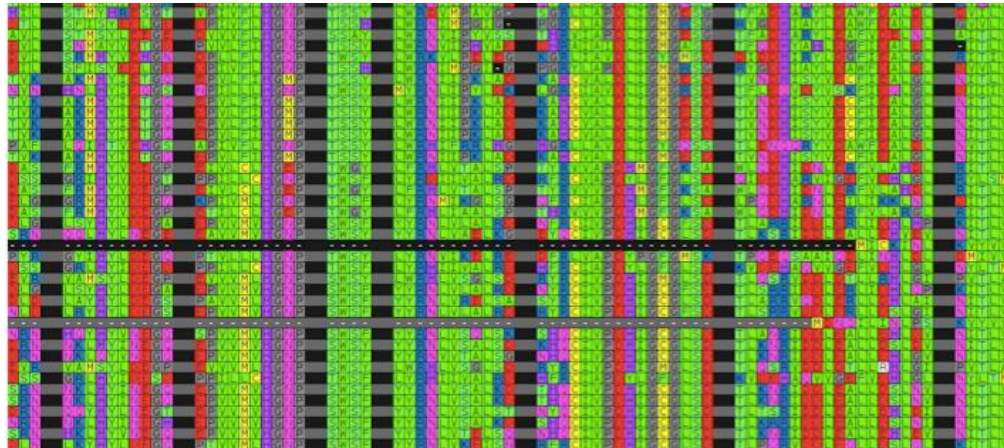
Στοίχιση βασισμένη στη δομή

Μας βοηθάει στον ήμι-ορθολογικό σχεδιασμό πρωτεϊνών

Μαθαίνουμε από την φυσική εξέλιξη των πρωτεϊνών

Μπορούμε να μεταφέρουμε μεταλλάξεις μεταξύ ομόλογων πρωτεϊνών

Μπορούμε να παρατηρήσουμε δίκτυα συσχετισμένων θέσεων



Βιοπληροφορικά εργαλεία για ανάλυση πρωτεϊνών

Βιοπληροφορικά εργαλεία για πρωτεΐνες

Μελέτες σχέσης δομής και λειτουργίας

Ανάπτυξη μοντέλων νέων πρωτεϊνών ή μεταλλαγμάτων

- πρόβλεψη / κατανόηση καταλυτικής συμπεριφοράς
- ορθολογικός σχεδιασμός πρωτεϊνών

Μοντέλα ομολογίας

Ανάλυση docking

- docking μικρών μορίων στις πρωτεΐνες

Ειδικά εργαλεία για την καθοδήγηση της εξέλιξης πρωτεϊνών

- βάσεις δεδομένων με δεδομένα στοίχισης με βάση τη δομή

NCBI

National Center for Biotechnology Information

<https://www.ncbi.nlm.nih.gov>

- ❖ Κύρια βάση δεδομένων για αλληλουχίες (γονίδια, γονιδιώματα, πρωτεΐνες)
- ❖ Κατάθεση δημοσιευμένων αλληλουχιών
- ❖ Διασύνδεση με τη βιβλιογραφία (**PubChem**) και εργαλεία (**BLAST**)
- ❖ Μοναδικός κωδικός για κάθε αλληλουχία
- ❖ Πηγή εύρεσης νέων γονιδίων/πρωτεϊνών από γνωστά μοτίβα

Protein Data Bank

<http://www.rcsb.org/pdb>

- ❖ Κύρια βάση δεδομένων για 3D δομές πρωτεϊνών
- ❖ Περίπου 170.000 δομές (κατάθεση ~10.000 / έτος)
- ❖ Πηγή pdb αρχείων - βασική μορφή αρχείου για τέτοια δεδομένα
- ❖ Σύνδεση με βιβλιογραφία
- ❖ Δεδομένα ποιότητας της δομής
- ❖ 3D viewer
- ❖ Δομικές πληροφορίες και πληροφορίες για την αλληλουχία
- ❖ Εύρεση παρόμοιων πρωτεϊνών
- ❖ Κωδικός πάντα με 4 ψηφία!

ExPASy

Expert Protein Analysis System

<https://www.expasy.org>

Μία τεράστια βιβλιοθήκη βιοπληροφορικών εργαλείων για πρωτεΐνες και DNA

Στα πρωτεωμικά εργαλεία:

- ❖ **UniProt**: Λειτουργικές πληροφορίες ενζύμων
- ❖ **ENZYME**: Κατηγοριοποίηση ενζύμων
- ❖ **Biochemical pathways**
- ❖ **Peptide cutter**: προβλέπει θέσεις κοπής
- ❖ **Phyre2**: Μοντέλα ομολογίας

Other databases

UniProt: <https://www.uniprot.org>

Πληροφορίες σχετικά με την λειτουργία των πρωτεϊνών

Διαφορετικοί κωδικοί από την NCBI!

Δίνει σημαντικά σημεία στην αλληλουχία

KEGG: <https://www.genome.jp/kegg/>

Κυρίως δίνει δεδομένα σε σχέση με συγκεκριμένα μεταβολικά μονοπάτια

BRENDA: <https://www.brenda-enzymes.org>

Δομικά και λειτουργικά δεδομένα

Συσχέτιση και με την SCOP (Structural Classification of Proteins)

Εύρεση πρωτεϊνών με παρόμοια δομή

MolProbity

<http://molprobity.biochem.duke.edu>

- ❖ Έλεγχος ποιότητας της δομής
- ❖ Μπορεί να σχεδιάσει διάγραμμα Ramachandran
- ❖ Δίνει λίστα με δομικά προβλήματα όπως στερικές παρεμποδίσεις, περίεργες γωνίες δεσμών κτλ
- ❖ Μπορεί να στρίψει πλευρικές ομάδες με ασυμμετρία (His, Asn, Gln)
- ❖ Πολύ σημαντικό εργαλείο, ειδικά για την μελέτη ποιότητας μοντέλων ομολογίας

Protein viewer

[PyMol 0.99 \(newer versions are not freeware\)](#)

Το λογισμικό που χρησιμοποιούν οι περισσότεροι
Πλήρης ρύθμιση – πολλά plug-in για να δουλέψεις

[Chimera](#)

Freeware – χρησιμοποιείται αρκετά τα τελευταία χρόνια

[YASARA](#)

Freeware ως viewer

Πέρα της επισκόπησης της δομής, μπορεί να χρησιμοποιηθεί και για molecular docking, δημιουργία μοντέλων ομολογίας κτλ

Bio-product και 3DM

<https://www.bio-product.com>

Η 3DM είναι μία πλατφόρμα που συλλέγει, αποθηκεύει, συνδέει και οπτικοποιεί δεδομένα από μία υπεροικογένεια και χρησιμοποιεί μοναδικό αριθμό θέσης για τις αλληλουχίες.

1. *Αποδελτίωση δεδομένων από βιβλιογραφία*
2. *Κατηγοριοποίηση σε μία υπο-οικογένεια, γύρω από μία ομόλογη PDB δομή*
3. *Δημιουργία στοίχισης βασισμένης στη δομή για όλες τις αλληλουχίες*
4. *Δημιουργία δεδομένων από την μετα-ανάλυση, όπως συσχετισμός θέσεων, συντήριση δεδομένων, μοντέλα ομολογίας κτλ*
5. *Οπτικοποίηση στο YASARA*

Βιβλιογραφία

- ❖ Attwood and Parry-Smith (1999) Introduction to bioinformatics. Adison Wesley Longman
- ❖ Misener and Krawetz (2000) Methods and protocols. Humana press
- ❖ Higgins and Taylor (2000) Bioinformatics: Sequence, structure and databanks. Oxford University Press