

Βιοτληροφορική

Στοίχιση αλληλουχιών

Είναι μια από τις βασικότερες διαδικασίες της μοριακής βιολογίας και η 'εκ των ουκ άνευ' της βιοπληροφορικής. Η μέθοδος είναι αναπόσπαστα δεμένη με την εξέλιξη (κληρονομούμενες αλλαγές πληροφορίας) ασχέτως με το εάν η εξέλιξη αποτελεί τμήμα του προβλήματος [«πριν από πόσα Myr διαχωρίστηκαν τα είδη αυτά ;»]
ή όχι
[«κωδικοποιεί το γονίδιο αυτό για μια πρωτεάση ;»].

Εφαρμογές στοιχίσεων

- Έρευνα βάσεων δεδομένων για αναζήτηση ομολόγων αλληλουχιών και (δυνητικά) ταυτοποίηση λειτουργίας.
- Συστηματική/ταξινομική αλληλουχιών, γονιδιωμάτων, οργανισμών.
- Συστηματική/ταξινομική οικογενειών πρωτεΐνων και ταυτοποίηση συντηρημένων (και συνεπώς, λειτουργικά/δομικά σημαντικών) καταλοίπων.
- Αναγνώριση και ταυτοποίηση μοτίβων σε πρωτεΐνικές οικογένειες.
- Πρόβλεψη της δομής μίας πρωτεΐνης με βάση την ομολογία της με άλλη πρωτεΐνη γνωστής δομής (homology modeling).

Ομολογία, αναλογία, ομοιότητα

Δύο αλληλουχίες είναι ομόλογες εάν έχουν αποκλίνει (εξελικτικά) από μία κοινή προγονική αλληλουχία.

Η ομολογία είναι μια απόλυτη δήλωση σχέσης : δύο αλληλουχίες ή είναι ή δεν είναι ομόλογες (η ομολογία δεν επιδέχεται ποσοτικούς προσδιορισμούς).

Δύο αλληλουχίες είναι ανάλογες εάν είναι προϊόντα συγκλίνουσας εξέλιξης, δηλαδή εάν έχουν αντίστοιχες δομές ή και λειτουργίες αλλά χωρίς να έχουν προέλθει (εξελικτικά) από κάποια κοινή προγονική αλληλουχία.

Η αναλογία είναι επίσης μια απόλυτη δήλωση απουσίας εξελικτικής σχέσης.

Ομολογία, αναλογία, ομοιότητα

Η ομοιότητα δύο αλληλουχιών είναι ένα ποσοτικό μέτρο των μεταξύ τους κοινών χαρακτηριστικών.

Ένα από τα ευρέως χρησιμοποιούμενα μέτρα είναι το επί τοις εκατό ποσοστό ταυτότητας μεταξύ των στοιχισμένων αλληλουχιών (% sequence identity), αν και υπάρχουν και άλλα μέτρα τα οποία λαμβάνουν υπόψη το κατά πόσο οι παρατηρούμενες αλλαγές είναι συντηρητικές ή όχι [το οποίο, για να επιτείνει τη σύγχυση, αναφέρεται από πολλά προγράμματα ως επί τοις εκατό «ομοιότητα» (% similarity), η οποία είναι διακριτή από την επί τοις εκατό ταυτότητα].

Ορθόλογες και παράλογες αλληλουχίες

Ορθόλογες (orthologues) ονομάζονται δύο ομόλογες αλληλουχίες οι οποίες έχουν την ίδια λειτουργία αλλά προέρχονται από διαφορετικά είδη. Οι ορθόλογες αλληλουχίες είναι το κυρίως υλικό για τη δημιουργία φυλογενετικών δένδρων.

Παράλογες (paralogues) ονομάζονται δύο ομόλογες αλληλουχίες ενός και του αυτού είδους οι οποίες έχουν παρόμοιες (αλλά όχι πτανομοιότυπες) λειτουργίες. Αυτές δίνουν πληροφορία για την εξελικτική πορεία γονιδίων μετά από φαινόμενα γονιδιακού διπλασιασμού.

Το πρόβλημα

Το ζητούμενο της στοίχισης αλληλουχιών είναι να αποκαλυφθούν οι μεταξύ τους εξελικτικές σχέσεις [δηλ. μεταλλάξεις & εισαγωγές/διαγραφές]. Όντως, μία οποιαδήποτε στοίχιση μεταξύ αλληλουχιών υποδηλώνει ένα εξελικτικό σενάριο π.χ.

**AGGACTGCG-TAG-TAC
AGG---TCGATAGGCAC**

ή

**AGGACT-GCGTA-GTAC
AGG---TCGA-TAGGCAC**

ή ...

Το πρόβλημα

Δυστυχώς η εξελικτική ιστορία των αλληλουχιών μας είναι άγνωστη (εάν δεν ήταν δεν θα χρειαζόμασταν τη διαδικασία στοίχισης). Αυτό έχει δύο συνέπειες :

1. Δεν γνωρίζουμε εάν όντως οι αλληλουχίες έχουν εξελικτική σχέση (δηλαδή εάν είναι όντως στοιχίσιμες).
2. Εάν έχουν εξελικτική σχέση, το μόνο που συνήθως έχουμε σαν βάση για να συνάγουμε τη στοίχιση είναι η παροντική μορφή των αλληλουχιών (εκτός και εάν υπάρχει διαθέσιμη η αλληλουχία από κάποιο κοινό εξελικτικό πρόγονο).

Η μηδενική υπόθεση

Επειδή δεν γνωρίζουμε εάν όντως οι αλληλουχίες έχουν εξελικτική σχέση, πρέπει να τροποποιήσουμε το στόχο της στοίχισης αλληλουχιών ως εξής :
Στόχος μας είναι να συγκρίνουμε δύο διαφορετικές υποθέσεις. Η μηδενική υπόθεση είναι ότι οι αλληλουχίες δεν έχουν εξελικτική σχέση (ότι δηλ. δεν έχουν προκύψει από κάποιο κοινό εξελικτικό πρόγονο). Αυτή η υπόθεση υπονοεί ότι η οποιαδήποτε ομοιότητα μεταξύ των αλληλουχιών είναι τυχαία και δεν οφείλεται σε εξελικτικές διαδικασίες.
Για παράδειγμα :

Η μηδενική υπόθεση

Με τη βοήθεια μίας γεννήτριας τυχαίων αριθμών παρήχθησαν οι κάτωθι αλληλουχίες A & B :

A : **CAGTATTGCTAGCATTG**

B : **GTATGCGGAACGTATTG**

Παρότι δεν υπάρχει καμία μεταξύ τους εξελικτική σχέση, οι δυο αυτές αλληλουχίες έχουν ομοιότητες, όπως φαίνεται από την παρακάτω στοίχιση

A : **CAGTATTGC---TAGCATTG**

B : **GTAT-GCGGAACGTATTG**

Σε αυτό το τεχνητό παράδειγμα γνωρίζουμε εκ των προτέρων ότι αυτές οι ομοιότητες είναι τυχαίες.

Αλλά εάν το μόνο που είχαμε σαν δεδομένο ήταν οι καθ'αυτό αλληλουχίες θα συναγάγαμε το ίδιο ;

Η εξελικτική υπόθεση

Η δεύτερη υπόθεση (την οποία επιθυμούμε να συγκρίνουμε με το μοντέλο της τυχαίας ομοιότητας) είναι οι αλληλουχίες να έχουν εξελικτική σχέση, δηλ. να έχουν προκύψει από έναν κοινό εξελικτικό πρόγονο μέσω μεταλλάξεων και indels. Σε αυτή την περίπτωση θα θέλαμε η στοίχιση να αποκαλύπτει αυτά τα εξελικτικά γεγονότα.

Για παράδειγμα :

Εξελικτική υπόθεση : παράδειγμα

Ξεκινώντας από την αλληλουχία ACGTACGT, και χρησιμοποιώντας σταθερούς ρυθμούς μεταλλάξεων και γεγονότων εισαγωγών/διαγραφών, προέκυψε μετά από ~9000 γενιές η αλληλουχία ACACCGGTCCCTAATAATGGCC. Επανάληψη της ίδιας διαδικασίας (ξεκινώντας από την ACGTACGT και για ίδιο αριθμό γενιών), έδωσε την αλληλουχία CAGGAAGATCTTAGTTC. Επειδή σε αυτή την περίπτωση η ιστορία των αλληλουχιών μας είναι γνωστή, μπορούμε να στοιχίσουμε την αρχική αλληλουχία με κάθε μία από τις τελικές :

Εξελικτική υπόθεση : παράδειγμα

--**ACG-T-A**--**CG-T**--
ACACGGTCCTAATAATGGCC

και

--**AC-GTA-C**--**G-T**--
CAG-GAAGATCTTAGTTC

και με υπέρθεση να βρούμε την εξελικτικά ορθή στοίχιση μεταξύ των τελικών αλληλουχιών :

-**ACAC-GGTCCCTAAT**--**AATGGCC**
CAG-GAA-G-AT--**CTTAGTTC**--
* * * *

Εξελικτική υπόθεση : παράδειγμα

--ACG-T-A---CG-T----
ACACGGTCCTAATAATGGCC

και

---AC-GTA-C--G-T--
CAG-GAAGATCTTAGTTC

και με υπέρθεση να βρούμε την εξελικτικά ορθή στοίχιση μεταξύ των τελικών αλληλουχιών :

-ACAC-GGTCC**TAAT**--AATGGCC
CAG-GAA-G-AT--CTTAGTTC--
* * * *

Μόνο που ένας αλγόριθμος θα έδινε άλλη στοίχιση ...

ACACG--GTCCTAATAATGGCC****
-CAGGAAGATCT--TAGTT--C
*** * * ** *** * *

Το πρόβλημα

Άρα το πρόβλημα της στοίχισης αλληλουχιών είναι

1. Να υπολογίσουμε την πιθανότητα η ομοιότητα των αλληλουχιών να οφείλεται σε εξελικτική σχέση, και,
2. Να προσδιορίσουμε τη στοίχιση εκείνη που καλύτερα αναπαριστά την εξελικτική σχέση (ιστορία) των αλληλουχιών.

Δυστυχώς, είναι αδύνατο να υπολογίσουμε το [1o] χωρίς να έχουμε ήδη κάνει μια στοίχιση. Αυτό συμβαίνει γιατί ο υπολογισμός της στατιστικής σημασίας μιας στοίχισης απαιτεί να υπάρχει η στοίχιση (ώστε να μπορούμε στη συνέχεια να ρωτήσουμε «ποιά είναι η πιθανότητα να έχει προκύψει αυτή η στοίχιση τυχαία;»).

Η Λύση

- Υποθέτουμε (a priori) πως οι αλληλουχίες έχουν εξελικτική σχέση.
- Προσδιορίζουμε την στοίχιση η οποία μεγιστοποιεί την μεταξύ τους (εξελικτική) ομοιότητα.
- Ελέγχουμε την στατιστική σημαντικότητα της προκύπτουσας στοίχισης :
 - Εάν η ομοιότητα της στοίχισης δεν είναι στατιστικά σημαντική, καταρρίπτουμε την αρχική μας υπόθεση και συνάγουμε ότι οι αλληλουχίες δεν είναι ομόλογες.
 - Στην αντίθετη περίπτωση πιθανολογούμε πως οι αλληλουχίες είναι ομόλογες και αναζητούμε επιπρόσθετες (μη υπολογιστικές;) ενδείξεις.

Τα υπολογιστικά προβλήματα

Το πρόβλημα μας τώρα είναι το εξής :
Δεδομένων των αλληλουχιών, πως θα προσδιορίσουμε
την στοίχιση εκείνη που μεγιστοποιεί την μεταξύ τους
(εξελικτική) ομοιότητα.

Το οποίο, με τη σειρά του, μας φέρνει στο πρόβλημα
του πως μετράμε την «εξελικτική ομοιότητα»
(ώστε να μπορούμε να την μεγιστοποιήσουμε).

Πίνακες βαθμολόγησης

Πίνακες βαθμολόγησης

Οι πίνακες βαθμολόγησης είναι απότειρες να κωδικοποιηθούν αριθμητικά οι μέσες (στατιστικά) προτιμήσεις της φυσικής επιλογής σε ό,τι αφορά τις μεταλλάξεις (αμινοξικές αλλαγές).

Οι τιμές που περιέχουν είναι (μεταφορικά) ενδεικτικές για το τι άποψη έχει η φυσική επιλογή για τα πετραγμένα της εξέλιξης. Δεν κωδικοποιούν τις συχνότητες συγκεκριμένων μεταλλάξεων, αλλά τις συχνότητες με τις οποίες οι διάφορες μεταλλάξεις γίνονται αποδεκτές.

Πίνακες βαθμολόγησης

Επειδή η φυσική επιλογή δρα στο επίπεδο των φαινοτύπων (και όχι των γονοτύπων), θα επικεντρωθούμε από εδώ και πέρα στις στοιχίσεις πρωτεϊνικών αλληλουχιών. Οι διαδικασίες και αλγόριθμοι που θα περιγραφούν μεταφέρονται σχεδόν αυτούσιες και στην περίπτωση αλληλουχιών νουκλεϊκών οξέων.

Πίνακες βαθμολόγησης

Ο πλέον απλοϊκός πίνακας βαθμολόγησης είναι :

Για κάθε θέση που οι αλληλουχίες ταυτίζονται,
αύξησε την βαθμολογία κατά ένα σταθερό θετικό
αριθμό, αλλιώς προχώρησε στην επόμενη θέση.
Αυτό σε μορφή πίνακα θα μπορούσε να γραφτεί :

A	1
C	0 1
M	0 0 1
P	0 0 0 1
F	0 0 0 0 1
.....	
A	C M P F ...

'Μοναδιαίοι πίνακες βαθμολόγησης'.

Πίνακες βαθμολόγησης

Ένας τέτοιου τύπου πίνακας έχει ελάχιστη ευαισθησία (διαγνωστική ισχύ) : αγνοεί πλήρως τις σχέσεις συγγένειας μεταξύ των διάφορων ομάδων αμινοξέων (π.χ. τπολικά, υδρόφοβα, φορτισμένα, ...) όπως αυτές προκύπτουν από τις φυσικοχημικές τους ιδιότητες (αλλά και από τον γενετικό κώδικα).

Πίνακες βαθμολόγησης

	T		C		A		G		
T	TTC	Phe	TCT	Ser	TAT	Try	TGT	Cys	T
	TTA	Leu	TCC		TAC		TGC		C
C	TTG		TCA	Pro	TAA	Stop	TGA	Stop	A
	CTT	Leu	CCT		TAG		TGG	Trp	G
	CTC		CCC		CAT	His	CGT	Arg	T
	CTA		CCA		CAC		CGC		C
A	CTG		CCG		CAA	Gln	CGA		A
	ATT	Ile	ACT	Thr	CAG		CGG		G
	ATC		ACC		AAT	Asn	AGT	Ser	T
	ATA		ACA		AAC		AGC		C
G	ATG	Met	ACG	Ala	AAA	Lys	AGA	Arg	A
	GTT	Val	GCT		AAG		AGG		G
G	GTC		GCC		GAT	Asp	GGT	Gly	T
	GTA		GCA		GAC		GGC		C
	GTG		GCG		GAA	Glu	GGA		A
G					GAG		GGG		G

ΟΙ ΤΤÍΝΑΚΕΣ PAM

PAM 250, Deyhoff, et al. (1978)

C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	
12																				
S	0	2																		
T	-2	1	3																	
P	-3	1	0	6																
A	-2	1	1	1	2															
G	-3	1	0	-1	1	5														
N	-4	1	0	-1	0	0	2													
D	-5	0	0	-1	0	1	2	4												
E	-5	0	0	-1	0	0	1	3	4											
Q	-5	-1	-1	0	0	-1	1	2	2	4										
H	-3	-1	-1	0	-1	-2	2	1	1	3	6									
R	-4	0	-1	0	-2	-3	0	-1	-1	1	2	6								
K	-5	0	0	-1	-1	-2	1	0	0	1	0	3	5							
M	-5	-2	-1	-2	-1	-3	-2	-3	-2	-1	-2	0	0	6						
I	-2	-1	0	-2	-1	-3	-2	-2	-2	-2	-2	-2	-2	2	5					
L	-6	-3	-2	-3	-2	-4	-3	-4	-3	-2	-2	-3	-3	4	2	6				
V	-2	-1	0	-1	0	-1	-2	-2	-2	-2	-2	-2	-2	2	4	2	4			
F	-4	-3	-3	-5	-4	-5	-4	-6	-5	-5	-2	-4	-5	0	1	2	-1	9		
Y	0	-3	-3	-5	-3	-5	-2	-4	-4	-4	0	-4	-4	-2	-1	-1	-2	7	10	
W	-8	-2	-5	-6	-6	-7	-4	-7	-7	-5	-3	2	-3	-4	-5	-2	-6	0	0	17

ΟΙ ΤΠÍΝΑΚΕΣ PAM

PAM είναι τα αρχικά του Point Accepted Mutation. Το PAM είναι ένα μέτρο της εξελικτικής απόκλισης μεταξύ δύο αλληλουχιών τέτοιο ώστε το 1 PAM να είναι η εξελικτική απόσταση που απαιτείται ώστε να μεταλλαχθεί το 1% των αμινοξέων μιας πρωτεΐνης. Μία απόσταση 100 PAM δεν υπονοεί ότι όλα τα αμινοξέα είναι διαφορετικά. Μερικά θα μεταλλαχθούν περισσότερες από μία φορές, άλλα καμία (ανάλογα με την πίεση της φυσικής επιλογής). Κατά μέσο όρο, μία απόσταση 250 PAM αντιστοιχεί σε ~20% ταυτότητα.

ΟΙ ΤΠÍΝΑΚΕΣ PAM

Εάν δεν υπήρχε η πίεση της φυσικής επιλογής (και εάν συνεπώς όλες μεταλλάξεις οι γινόντουσαν αποδεκτές), τότε η συχνότητα εμφάνισης των διάφορων αμινοξικών αλλαγών θα ήταν ανάλογη της συχνότητας εμφάνισης των συγκεκριμένων αμινοξέων στις πρωτεΐνες.

Παράδειγμα : εάν δεν υπήρχε η φυσική επιλογή, και εάν η συχνότητα εμφάνισης της τρυπποφάνης (W) είναι π.χ. μόνο 1%, τότε οι μεταλλάξεις προς την τρυπποφάνη θα είναι σχετικά σπάνιες. Οι αναμενόμενες συχνότητες εμφάνισης των μεταλλάξεων με βάση τις συχνότητες εμφάνισης των αμινοξέων (στις πρωτεϊνικές αλληλουχίες) ονομάζονται "συχνότητες υποβάθρου".

ΟΙ ΤΠÍΝΑΚΕΣ ΡΑΜ

Παρουσία της φυσικής επιλογής, οι συχνότητες εμφάνισης των διαφόρων μεταλλάξεων αλλάζουν. Η αλλαγή οφείλεται στο ότι επιλέγονται κατά προτίμηση οι μεταλλάξεις εκείνες οι οποίες είναι συμβατές με την δομή και λειτουργία της πρωτεΐνης (δηλ. εκείνες που δεν μειώνουν την λειτουργικότητα της). Με άλλα λόγια, επιλέγονται σημειακές μεταλλάξεις οι οποίες έγιναν αποδεκτές από τη φυσική επιλογή (εξ ου και 'point accepted mutations'). Οι συχνότητες εμφάνισης μεταλλάξεων όπως τις παρατηρούμε σε ομόλογες πρωτεΐνες (παρουσία της φυσικής επιλογής) ονομάζονται "παρατηρούμενες συχνότητες".

ΟΙ ΠÍΝΑΚΕΣ PAM

Η κάθε καταχώρηση στον πίνακα της Deyhoff, είναι ο λογάριθμος του πηλίκου της παρατηρούμενης συχνότητας δια την συχνότητα υποβάθρου της μετάλλαξης (log-odds ratio). Για παράδειγμα, η μετάλλαξη W \Leftrightarrow K (τρυπτοφάνη-λυσίνη) έχει τιμή -3. Άρα,

$$10 \log [P_{obs} / P_{back}] = -3$$

$$\Rightarrow P_{obs} / P_{back} = 10^{(-3/10)} = 0.5011$$

$$\Rightarrow P_{obs} = 0.5011 \cdot P_{back}$$

$$\Rightarrow P_{obs} = P_{back} / 2$$

το οποίο σημαίνει ότι παρουσία της φυσικής επιλογής, μια μετάλλαξη από τρυπτοφάνη προς λυσίνη συμβαίνει 2 φορές πιο σπάνια απ'ότι θα συνέβαινε κατά τύχη.

ΟΙ ΠΤÍΝΑΚΕΣ ΡΑΜ

Με βάση τις ιδιότητες των λογαρίθμων λοιπόν, όταν κάποια τιμή του πίνακα είναι θετική, τότε η αντίστοιχη μετάλλαξη παρατηρείται πιο συχνά απ'ότι θα περιμέναμε εάν συνέβαινε τυχαία. Άρα, οι θετικές τιμές του πίνακα αντιστοιχούν σε προτιμώμενες (από την φυσική επιλογή) μεταλλάξεις.

Κατ'ανalogία, οι αρνητικές τιμές του πίνακα αντιστοιχούν σε μεταλλάξεις που σπανίως γίνονται αποδεκτές, ενώ μία τιμή ίση με το μηδέν αντιστοιχεί σε ουδέτερες μεταλλάξεις.

ΟΙ ΤΤÍΝΑΚΕΣ PAM

PAM 250

C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	
12	0	2																		
S	0	2																		
T	-2	1	3																	
P	-3	1	0	6																
A	-2	1	1	1	2															
G	-3	1	0	-1	1	5														
N	-4	1	0	-1	0	0	2													
D	-5	0	0	-1	0	1	2	4												
E	-5	0	0	-1	0	0	1	3	4											
Q	-5	-1	-1	0	0	-1	1	2	2	4										
H	-3	-1	-1	0	-1	-2	2	1	1	3	6									
R	-4	0	-1	0	-2	-3	0	-1	-1	1	2	6								
K	-5	0	0	-1	-1	-2	1	0	0	1	0	3	5							
M	-5	-2	-1	-2	-1	-3	-2	-3	-2	-1	-2	0	0	6						
I	-2	-1	0	-2	-1	-3	-2	-2	-2	-2	-2	-2	-2	2	5					
L	-6	-3	-2	-3	-2	-4	-3	-4	-3	-2	-2	-3	-3	4	2	6				
V	-2	-1	0	-1	0	-1	-2	-2	-2	-2	-2	-2	-2	2	4	2	4			
F	-4	-3	-3	-5	-4	-5	-4	-6	-5	-5	-2	-4	-5	0	1	2	-1	9		
Y	0	-3	-3	-5	-3	-5	-2	-4	-4	-4	0	-4	-4	-2	-1	-1	-2	7	10	
W	-8	-2	-5	-6	-6	-7	-4	-7	-7	-5	-3	2	-3	-4	-5	-2	-6	0	0	17

ΟΙ ΤΠÍΝΑΚΕΣ ΡΑΜ

Ένας φαύλος κύκλος ;

Προκειμένου να υπολογίσουμε τις απαιτούμενες συχνότητες μεταλλάξεων, πρέπει να έχουμε ήδη βρει οικογένειες ομόλογων πρωτεϊνών και να τις έχουμε στοιχίσει. Άρα, για να δημιουργήσουμε τον πίνακα χρειαζόμαστε το προϊόν εφαρμογής του. Η λύση αυτού του φαύλου κύκλου είναι ότι για ομόλογες αλληλουχίες υψηλής ομοιότητας (αλληλουχίες που διαχωρίστηκαν πρόσφατα) η στοίχιση είναι προφανής και συνήθως περιορίζεται στην απλή παράθεση των αλληλουχιών. Από αυτές τις στενά συσχετιζόμενες αλληλουχίες μπορούμε εύκολα να υπολογίσουμε πίνακες όπως ο ΡΑΜ1.

ΟΙ ΠÍΝΑΚΕΣ PAM

Ένας φαύλος κύκλος :

Οι πίνακες για πιο μακρινές (εξελικτικά) σχέσεις (π.χ. PAM200 ή PAM250) προκύπτουν με απλό πολλαπλασιασμό πινάκων ξεκινώντας από τους χαμηλότερης 'τάξης' πίνακες. Ένα παράδειγμα για κάποιες υποθετικές αλληλουχίες οι οποίες αποτελούνται από μόνο δύο τύπων μονομερή (α, β) είναι :

PAM (x)

PAM (2x)

	α	β	\Rightarrow		α	β
α	2	-1		α	5	-5
β	-1	3		β	-5	10

ΟΙ ΤΠÍΝΑΚΕΣ BLOSUM

Ένας φαύλος κύκλος ;

Το πρόβλημα με αυτήν την προσέγγιση, είναι ότι για την στοίχιση απόμακρων εξελικτικά αλληλουχιών (που είναι και αυτό που περισσότερο μας ενδιαφέρει), δεν χρησιμοποιούμε "παρατηρούμενες" συχνότητες μεταλλάξεων, αλλά συχνότητες που τις συνάγαμε υπολογιστικά από αλληλουχίες με μεγάλη ομοιότητα.

ΟΙ ΠÍΝΑΚΕΣ BLOSUM

Henikoff & Henikoff (1992)

Οι πίνακες BLOSUM (για BLOcks SUbstitution Matrix) αποφεύγουν το προαναφερθέν πρόβλημα με το να υπολογίζουν τις παρατηρούμενες συχνότητες μεταλλάξεων από εξελικτικά απομακρυσμένες ομόλογες αλληλουχίες.

Το πρόβλημα αρχικής στοίχισης αυτών των αλληλουχιών αποφεύγεται με το να μη χρησιμοποιούνται ολόκληρες οι αλληλουχίες, παρά μόνο πρωτεϊνικά μοτίβα (χωρίς κενά) από τη βάση BLOCKS.

OI TTÍVAKΕΣ BLOSUM

OPN3_HUMAN/293-309

OPN3_MOUSE/291-307

OPN4_HUMAN/334-350

OPN4_MOUSE/331-347

OPS1_CALVI/311-327

OPS1_DROME/313-329

OPS1_DROPS/314-330

OPS1_HEMSA/319-335

OPS1_LIMPO/312-328

OPS1_PATYE/316-332

OPS1_SCHGR/317-333

OPS2_DROME/320-336

OPS2_DROPS/320-336

OPS2_HEMSA/319-335

OPS2_LIMPO/312-328

OPS2_PATYE/276-292

VSylfAKSNTvyNPviY

VSylfAKSSTvyNPviY

VPaviAKASAihNPiiY

VPaviAKASAihNPiiY

WGacfAKSAAcyNPivY

WGacfAKSAAcyNPivY

WGacfAKSAAcyNPivY

LPallAKSCScyNPfvY

WGsvfAKANScyNPivY

LPmmlAKSSSmhNPvvY

WGslfAKANAvfNPivY

WGatfAKTSAvyNPivY

WGatfAKTSAvyNPivY

LPallAKSCScyNPfvY

WGsvfAKANScyNPivY

LPtlfAKASCayNPfiY

BLOSUM62

A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
4																			
-1	5																		
-2	0	6																	
-2	-2	1	6																
0	-3	-3	-3	9															
-1	1	0	0	-3	5														
-1	0	0	2	-4	2	5													
0	-2	0	-1	-3	-2	-2	6												
-2	0	1	-1	-3	0	0	-2	8											
-1	-3	-3	-3	-1	-3	-3	-4	-3	4										
-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4									
-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5								
-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5							
-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6						
-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7					
1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4				
0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5			
-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11		
-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	
0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4

PAM250

C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
12	0	2																	
S	0	2																	
T	-2	1	3																
P	-3	1	0	6															
A	-2	1	1	1	2														
G	-3	1	0	-1	1	5													
N	-4	1	0	-1	0	0	2												
D	-5	0	0	-1	0	1	2	4											
E	-5	0	0	-1	0	0	1	3	4										
Q	-5	-1	-1	0	0	-1	1	2	2	4									
H	-3	-1	-1	0	-1	-2	2	1	1	3	6								
R	-4	0	-1	0	-2	-3	0	-1	-1	1	2	6							
K	-5	0	0	-1	-1	-2	1	0	0	1	0	3	5						
M	-5	-2	-1	-2	-1	-3	-2	-3	-2	-1	-2	0	0	6					
I	-2	-1	0	-2	-1	-3	-2	-2	-2	-2	-2	-2	-2	2	5				
L	-6	-3	-2	-3	-2	-4	-3	-4	-3	-2	-2	-3	-3	4	2	6			
V	-2	-1	0	-1	0	-1	-2	-2	-2	-2	-2	-2	-2	2	4	2	4		
F	-4	-3	-3	-5	-4	-5	-4	-6	-5	-5	-2	-4	-5	0	1	2	-1	9	
Y	0	-3	-3	-5	-3	-5	-2	-4	-4	-4	0	-4	-4	-2	-1	-1	-2	7 10	
W	-8	-2	-5	-6	-6	-7	-4	-7	-7	-5	-3	2	-3	-4	-5	-2	-6	0 0 17	
C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W

Βαθμολόγηση στοιχίσεων

Ο πίνακας βαθμολόγησης είναι ένα σημαντικό βήμα για να μπορούμε να βαθμολογήσουμε μια στοίχιση, άλλα δεν αρκεί. Αυτό που λείπει είναι ένας τρόπος να λάβουμε υπόψη τα γεγονότα εισαγωγών/διαγραφών. Για παράδειγμα, οι αλληλουχίες

α **GATGGCGTAGCTAGGATTAAACA**

β **AGTGCATTTAGATCAGGATCTTCTATAGC**

μπορούν να στοιχιθούν έτσι,

GACGGCGT--AG--CTAGCAT----TA-A-CA

AGTG-CATTTAG-ATCAGGATCTTCTATCGC-

* * * * ** ** * *

ή και έτσι,

GACGGCGTAGCTAGCATTAAACA-----

---AGTGCATTTAG-ATCAGGATCTTCTATCGC

* * * *** ** * *

Βαθμολόγηση στοιχίσεων

Η πρώτη στοίχιση έχει δύο παραπάνω ταυτότητες βάσεων, αλλά έχει εισαγάγει ένα μάλλον απίθανο αριθμό από εισαγωγές/διαγραφές για να επιτύχει τις δύο επιπλέον ταυτότητες. Η δεύτερη στοίχιση φαίνεται να έχει μεγαλύτερο βιολογικό περιεχόμενο, κυρίως γιατί τα γεγονότα εισαγωγών/διαγραφών που υπονοεί είναι συνεχή και όχι διάσπαρτα και μεμονωμένα.

Βαθμολόγηση στοιχίσεων

Από την άποψη των στοιχίσεων, είναι αδύνατο να διακριθούν τα γεγονότα εισαγωγής από αυτά των διαγραφών. Για το λόγο αυτό, αναφερόμαστε σε αυτά με τον όρο "κενά" (gaps), και στον τρόπο βαθμολόγησης τους ως "κόστος κενών" (gap penalty). Η πλέον συνηθισμένη μέθοδος βαθμολόγησης των κενών αποτελείται από δύο όρους :

- Ο πρώτος όρος αφαιρεί από την βαθμολογία της στοίχισης ένα σταθερό ποσό για κάθε κενό που δημιουργείται ασχέτως του μήκους του κενού.
- Ο δεύτερος όρος αφαιρεί ένα σταθερό ποσό για κάθε επέκταση του μήκους ενός κενού.

Βαθμολόγηση στοιχίσεων

Άρα, για να υπολογίσουμε την βαθμολογία μίας στοιχίσης :

- Για κάθε ζεύγος στοιχισμένων καταλοίπων, αυξάνουμε την βαθμολογία κατά όσο αναφέρεται στον πίνακα βαθμολόγησης που χρησιμοποιούμε (PAM, BLOSUM, ...).
- Για κάθε κενό μήκους N , αφαιρούμε $N \cdot \alpha$ όπου α είναι το gap penalty. Εάν χρησιμοποιούμε τον διπλό τρόπο βαθμολόγησης κενών (προηγούμενη διαφάνεια), τότε αφαιρούμε $[N \cdot \alpha + k]$ όπου k είναι το κόστος δημιουργίας του κενού. Για απλότητα, από εδώ και πέρα θα υποθέτουμε ότι $k=0$.

Άρα, το πρόβλημα είναι ...

Με δεδομένες δύο αλληλουχίες A & B, έναν πίνακα βαθμολόγησης K και ένα κόστος εισαγωγής κενού α,
βρείτε για ποιά στοίχιση ---από όλες τις δυνατές στοιχίσεις των αλληλουχιών A&B --- η βαθμολογία (της στοίχισης) μεγιστοποιείται.

Πόσες είναι "οι δυνατές στοιχίσεις" ;

Το μέγεθος του προβλήματος

WHAT
WIAT

WHAT **WH-AT** **WHA-T** **WH-A-T**
WIAT **W-IAT** **WI-AT** **W-I-AT**

W-HAT **WHA-T** **W--HAT**
WI-AT **W-IAT** **WIA--T** ...

Το μέγεθος του προβλήματος

Για δύο αλληλουχίες A & B, με αντίστοιχα μήκη m & n , η φαινομενική πολυπλοκότητα του προβλήματος είναι ανάλογη του $[m^n]$.

Για μία στοίχιση δύο πρωτεΐνων 200 αμινοξέων εκάστη, και εάν εξετάζαμε 1 δις στοιχίσεις ανά δευτερόλεπτο, μετά από 15 δις χρόνια (την ηλικία του σύμπαντος) θα είχαμε εξετάσει μόνο 10^{27} στοιχίσεις (δεν θα είχαμε καν ξεκινήσει).

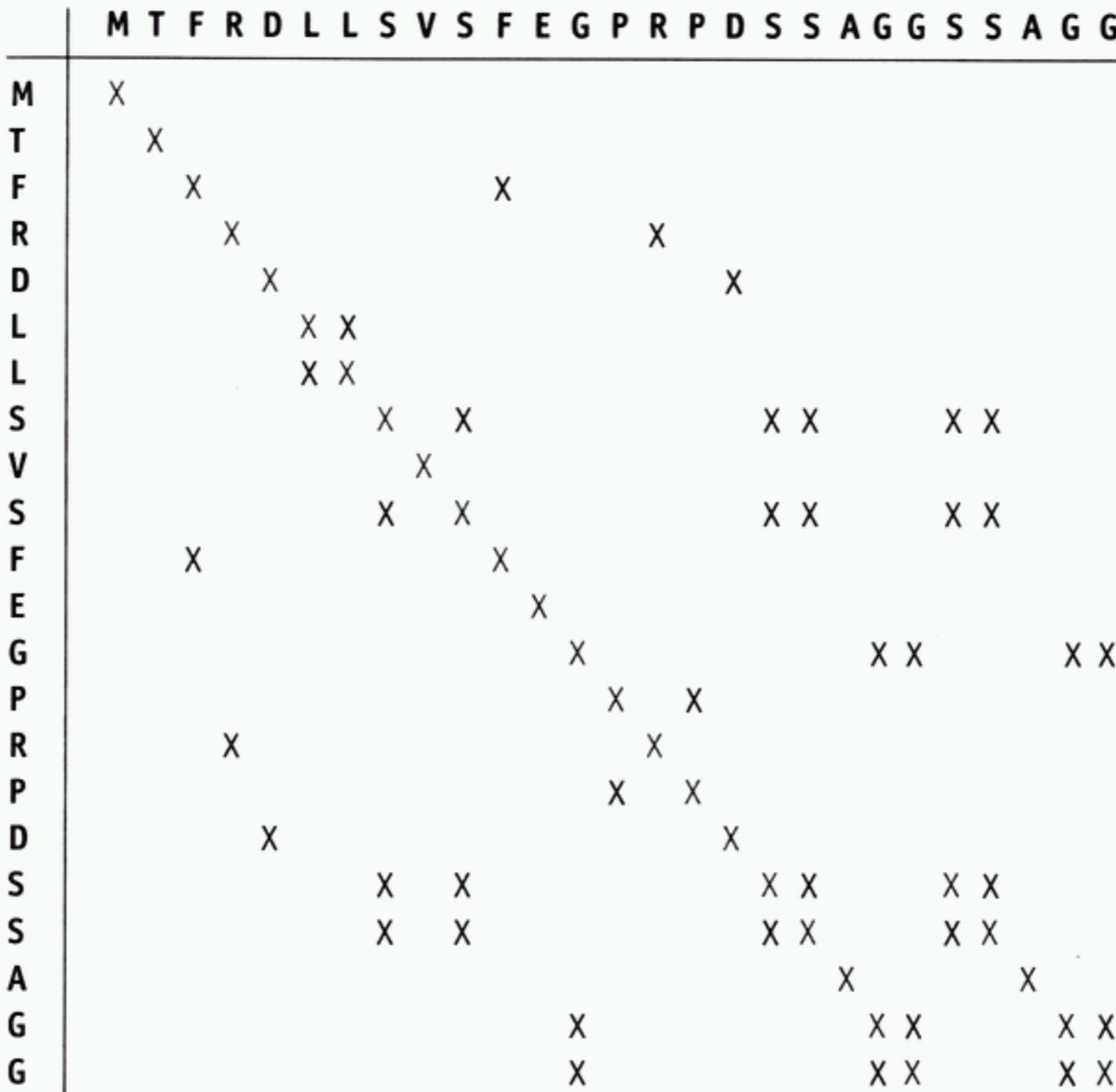
Στην πραγματικότητα, η πολυπλοκότητα του προβλήματος είναι μόνο $[m \cdot n]$. Για το παραπάνω πρόβλημα, θα είχαμε τελειώσει με τη στοίχιση σε 40 εκατομμυριοστά του δευτερολέπτου.

ΣΤΟΙΧΙΣΕΙΣ : dot plots

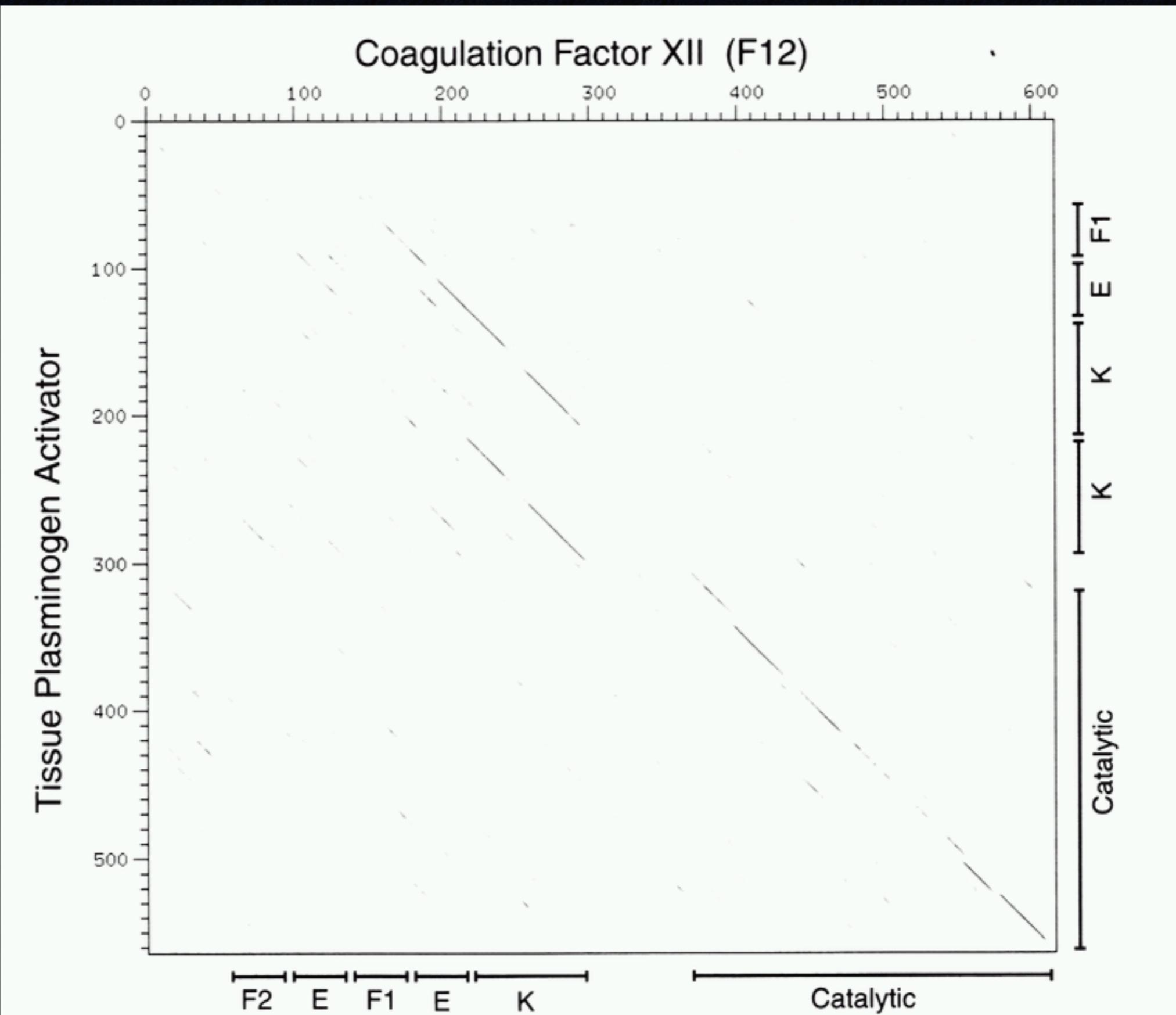
Μια μέθοδος που βοηθάει στην κατανόηση του γιατί η στοίχιση δυο αλληλουχιών έχει πολυπλοκότητα ανάλογη του $[m \cdot n]$ είναι τα λεγόμενα dot plots.

Τα dot plots (σημείο-διαγράμματα ;) δεν είναι ένας αλγόριθμος στοίχισης, αλλά ένα μέσο παρουσίασης των σχέσεων μεταξύ δυο αλληλουχιών. Πρόκειται για μια δισδιάστατη γραφική αναπαράσταση στην οποία κατά μήκος των δυο αξόνων αναπαρίστανται οι δυο αλληλουχίες και στον χώρο που ορίζεται από αυτές τοποθετούνται σημεία των οποίων η φωτεινότητα είναι ανάλογη της σχέσης των αμινοξέων (των αλληλουχιών) που αντιστοιχούν στις συντεταγμένες του σημείου.

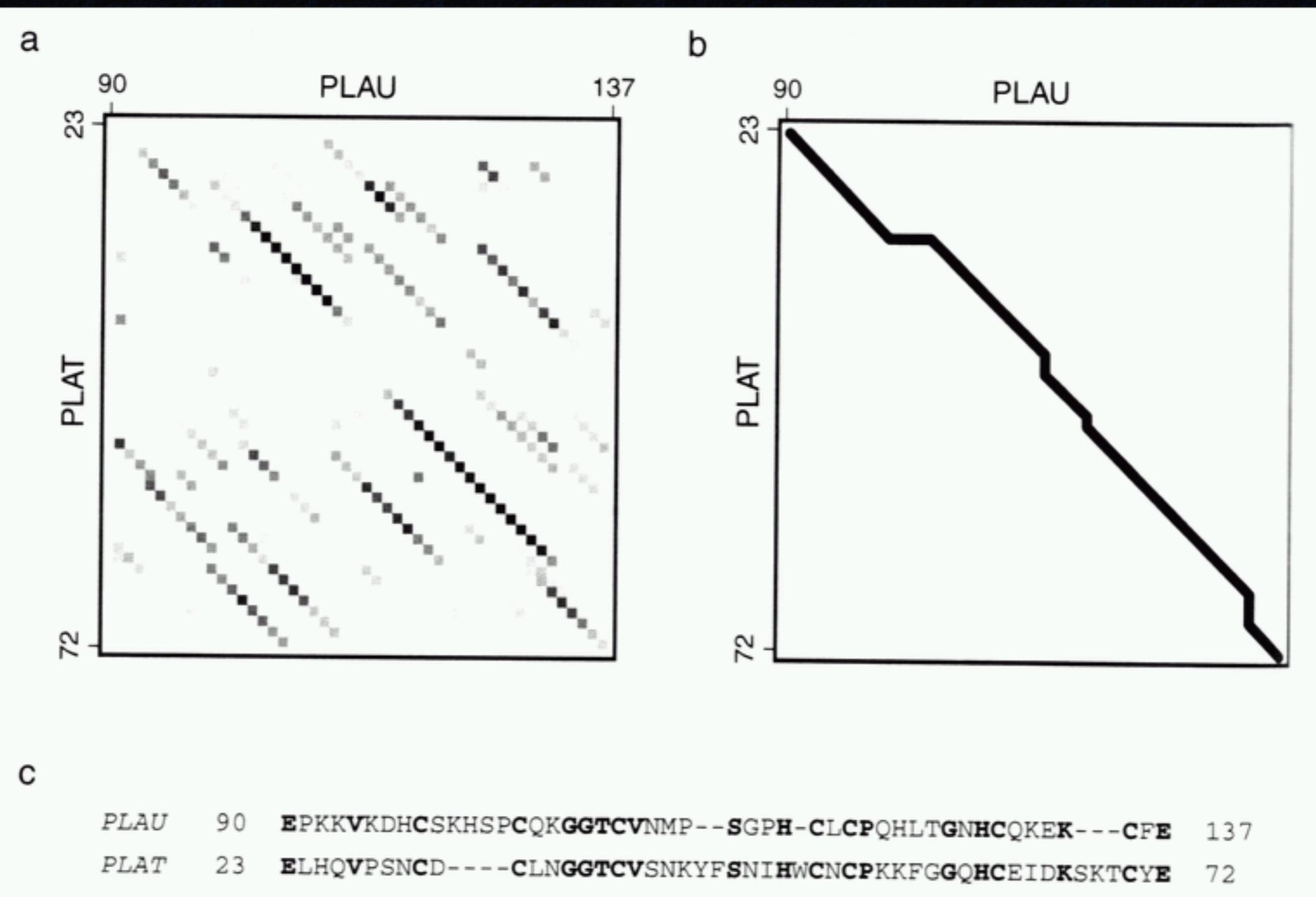
ΣΤΟΙΧΙΣΕΙΣ : dot plots



ΣΤΟΙΧΙΣΕΙΣ : dot plots



ΣΤΟΙΧΙΣΕΙΣ : dot plots



Αλγόριθμοι στοίχισης

Άρα, για δύο αλληλουχίες A & B, με αντίστοιχα μήκη m & n , το πρόβλημα είναι να βρεθεί το καλύτερο (υψηλότερης βαθμολογίας) μονοπάτι μέσω ενός πίνακα $[mn]$. Το μονοπάτι αυτό θα πρέπει να περνάει διαδοχικά από όλα τα αμινοξέα των δυο αλληλουχιών (ολική στοίχιση).

Needleman & Wunsch

Ο αλγόριθμος των Needleman & Wunsch (N & W), ανήκει στην κατηγορία των λεγομένων μεθόδων 'δυναμικού προγραμματισμού' (dynamic programming). Η βασική παρατήρηση στην οποία στηρίζεται ο αλγόριθμος είναι ότι οποιοδήποτε υποσύνολο της βέλτιστης στοίχισης, θα πρέπει επίσης να είναι βέλτιστο (ειδάλλως η στοίχιση θα μπορούσε να βελτιστοποιηθεί και άλλο μέσω βελτίωσης της υποστοίχισης).

N & W : η μέθοδος

Υποθέστε ότι έχουμε ήδη επιτύχει μία βέλτιστη στοίχιση μεταξύ των (κ) πρώτων αμινοξέων από την αλληλουχία A με τα (μ) πρώτα αμινοξέα από την B. Εάν βρούμε ένα τρόπο να επεκτείνουμε τη στοίχιση κατά μία θέση προς τα δεξιά έτσι ώστε η προκύπτουσα στοίχιση να είναι επίσης βέλτιστη, έχουμε τελειώσει : εφαρμόζοντας επαναληπτικά την μέθοδο αυτή, θα στοιχίσουμε (βέλτιστα) ολόκληρες τις αλληλουχίες.

N & W : η μέθοδος

Για την επέκταση της στοίχισης κατά μία θέση υπάρχουν μόνο τρεις τρόποι :

- Στοιχίζουμε το επόμενο αμινοξύ της A, δηλαδή το $A(\kappa+1)$, με το επόμενο της B, το $B(\mu+1)$.
- Προσθέτουμε ένα αμινοξύ από την A, το $A(\kappa+1)$, και εισαγάγουμε ένα κενό στη B.
- Προσθέτουμε ένα αμινοξύ από την B, το $B(\mu+1)$, και εισαγάγουμε ένα κενό στη A.

Αυτό που θέλουμε είναι να διαλέξουμε εκείνο τον τρόπο επέκτασης που μεγιστοποιεί την βαθμολογία της τελικής στοίχισης.

N & W : η μέθοδος

Έστω, λοιπόν, ότι η βαθμολογία της βέλτιστης στοίχισης των (κ) αμινοξέων της A, με τα (μ) της B, είναι $S(\kappa, \mu)$.

Τότε :

Εάν στοιχίσουμε το επόμενο αμινοξύ της A, δηλαδή το $A(\kappa+1)$, με το επόμενο της B, το $B(\mu+1)$, η βαθμολογία της στοίχισης θα αυξηθεί κατά τόσο όσο η βαθμολογία υποκατάστασης των αμινοξέων $A(\kappa+1)$ και $B(\mu+1)$ όπως αυτή δίδεται από τον πίνακα βαθμολόγησης K. Άρα σε αυτή την περίπτωση :

$$S(\kappa+1, \mu+1) = S(\kappa, \mu) + K[A(\kappa+1), B(\mu+1)]$$

N & W : η μέθοδος

Εάν προσθέσουμε το επόμενο αμινοξύ της A, δηλαδή το $A(\kappa+1)$, και εισαγάγουμε ένα κενό στη B, η τελική βαθμολογία θα είναι ίση με τη βαθμολογία της βέλτιστης στοίχισης των $(\kappa+1)$ αμινοξέων της A με τα (μ) της B μείον το gap penalty (α). Δηλαδή :

$$S(\kappa+1, \mu+1) = S(\kappa+1, \mu) - \alpha$$

N & W : η μέθοδος

Εάν προσθέσουμε το επόμενο αμινοξύ της B, δηλαδή το $B(\mu+1)$, και εισαγάγουμε ένα κενό στη A, η τελική βαθμολογία θα είναι ίση με τη βαθμολογία της βέλτιστης στοίχισης των (κ) αμινοξέων της A με τα $(\mu+1)$ της B μείον το gap penalty (α). Δηλαδή :

$$S(\kappa+1, \mu+1) = S(\kappa, \mu+1) - \alpha$$

N & W : η μέθοδος

Η καινούργια βέλτιστη στοίχιση θα είναι αυτή για την οποία η τιμή του $S(\kappa+1, \mu+1)$ μεγιστοποιείται.

Η ολική βαθμολογία της καινούργιας βέλτιστης στοίχισης θα είναι προφανώς :

$$S(\kappa+1, \mu+1) = \max \{ S(\kappa, \mu) + K[A(\kappa+1), B(\mu+1)], \\ S(\kappa+1, \mu) - \alpha, \\ S(\kappa, \mu+1) - \alpha \}$$

N & W : η μέθοδος

Άρα, εάν έχουμε τιμές για τα $S(\kappa, \mu)$, $S(\kappa+1, \mu)$ και $S(\kappa, \mu+1)$, τότε μπορούμε να υπολογίζουμε το $S(\kappa+1, \mu+1)$ (και την αντίστοιχη στοίχιση) απλά με το να βρίσκουμε το μέγιστο των τριών αριθμητικών εκφράσεων που δίδονται στις αγκύλες της προηγούμενης εξίσωσης. Γράφοντας αυτά τα αποτελέσματα σε μορφή πίνακα δείχνει ότι εάν έχουμε τιμές βαθμολογίας για την πρώτη γραμμή και την πρώτη στήλη του, μπορούμε επαναληπτικά να τον συμπληρώσουμε ολόκληρο :

N & W : η μέθοδος

S (κ, μ)

S (κ, μ+1)

S (κ, μ+2)

S (κ, μ+3)

S (κ+1, μ)

S (κ+2, μ)

S (κ+3, μ)

N & W : η μέθοδος

$S(\kappa, \mu)$

$S(\kappa, \mu+1)$

$S(\kappa, \mu+2)$

$S(\kappa, \mu+3)$

$S(\kappa+1, \mu)$

$S(\kappa+1, \mu+1)$

$S(\kappa+2, \mu)$

$S(\kappa+3, \mu)$

$$S(\kappa+1, \mu+1) = \max \{ S(\kappa, \mu) + K[A(\kappa+1), B(\mu+1)], \\ S(\kappa+1, \mu) - \alpha, \\ S(\kappa, \mu+1) - \alpha \}$$

N & W : η μέθοδος

$S(\kappa, \mu)$

$S(\kappa, \mu+1)$

$S(\kappa, \mu+2)$

$S(\kappa, \mu+3)$

$S(\kappa+1, \mu)$

$S(\kappa+1, \mu+1)$

$S(\kappa+1, \mu+2)$

$S(\kappa+1, \mu+3)$

$S(\kappa+2, \mu)$

$S(\kappa+3, \mu)$

$$S(\kappa+1, \mu+1) = \max \{ S(\kappa, \mu) + K[A(\kappa+1), B(\mu+1)], \\ S(\kappa+1, \mu) - \alpha, \\ S(\kappa, \mu+1) - \alpha \}$$

N & W : η μέθοδος

$S(\kappa, \mu)$

$S(\kappa, \mu+1)$

$S(\kappa, \mu+2)$

$S(\kappa, \mu+3)$

$S(\kappa+1, \mu)$

$S(\kappa+1, \mu+1)$

$S(\kappa+1, \mu+2)$

$S(\kappa+1, \mu+3)$

$S(\kappa+2, \mu)$

$S(\kappa+2, \mu+1)$

$S(\kappa+2, \mu+2)$

$S(\kappa+2, \mu+3)$

$S(\kappa+3, \mu)$

$S(\kappa+3, \mu+1)$

$S(\kappa+3, \mu+2)$

$S(\kappa+3, \mu+3)$

$$S(\kappa+1, \mu+1) = \max \{ S(\kappa, \mu) + K[A(\kappa+1), B(\mu+1)], \\ S(\kappa+1, \mu) - \alpha, \\ S(\kappa, \mu+1) - \alpha \}$$

N & W : η μέθοδος

Άρα, για να βρούμε την βέλτιστη στοίχιση των δύο αλληλουχιών A & B, με μήκη m & n, χρειαζόμαστε τιμές για τα $S(0,0)$ $S(1,0)$ $S(2,0)$... $S(m,0)$ και για τα $S(0,1)$ $S(0,2)$... $S(0,n)$.

Για τον αλγόριθμο των N & W η αρχικοποίηση των τιμών είναι απλή (και προφανής) :

$$S(0,0) = 0$$

$$S(\kappa,0) = -\kappa \cdot \alpha$$

$$S(0,\mu) = -\mu \cdot \alpha$$

N & W : παράδειγμα 1ο

Θέλουμε να στοιχίσουμε τις αλληλουχίες ASPERA και APTERA. Ο (μοναδιαίος) πίνακας βαθμολόγησης είναι :

και το gap penalty έχει τιμή 1 (ανά κενό).

N & W : παράδειγμα 1ο

	A	S	P	E	R	A
0	-1	-2	-3	-4	-5	-6
A	-1					
P	-2					
T	-3					
E	-4					
R	-5					
A	-6					

Αρχικοποίηση του πίνακα.

N & W : παράδειγμα 1ο

	A	S	P	E	R	A
0	-1	-2	-3	-4	-5	-6
A	-1	3				
P	-2					
T	-3					
E	-4					
R	-5					
A	-6					

Συμπλήρωση.

N & W : παράδειγμα 1ο

	A	S	P	E	R	A
0	-1	-2	-3	-4	-5	-6
A	-1	3	2	1	0	-1
P	-2	2				
T	-3					
E	-4					
R	-5					
A	-6					

Συμπλήρωση.

N & W : παράδειγμα 1ο

	A	S	P	E	R	A
0	-1	-2	-3	-4	-5	-6
A	-1	3	2	1	0	-1
P	-2	2	3	5	4	3
T	-3	1	2	4	5	4
E	-4	0	1	3	7	6
R	-5	-1	0	2	6	10
A	-6	-2	-1	1	5	9
						13

Συμπλήρωση.

N & W : παράδειγμα 1ο

	A	S	P	E	R	A
0	-1	-2	-3	-4	-5	-6
A	-1	3	2	1	0	-1
P	-2	2	3	5	4	3
T	-3	1	2	4	5	4
E	-4	0	1	3	7	6
R	-5	-1	0	2	6	10*
A	-6	-2	-1	1	5	9
						13*

Μονοπάτι βέλτιστης στοίχισης.

N & W : παράδειγμα 1ο

	A	S	P	E	R	A
0	-1	-2	-3	-4	-5	-6
A	-1	3*	2*	1	0	-1
P	-2	2	3	5*	4	3
T	-3	1	2	4*	5	4
E	-4	0	1	3	7*	6
R	-5	-1	0	2	6	10*
A	-6	-2	-1	1	5	9
	A S P - E R A					
	A - P T E R A					

A S P - E R A
A - P T E R A

N & W : παράδειγμα 2ο

Θέλουμε να στοιχίσουμε τις αλληλουχίες ASPERA και APTERA. Ο (μοναδιαίος) πίνακας βαθμολόγησης είναι :

A	3				
C	0 3				
M	0 0 3				
P	0 0 0 3				
F	0 0 0 0 3				
.....					
A	C	M	P	F

και το gap penalty έχει τιμή 2 (ανά κενό).

N & W : παράδειγμα 2ο

	A	S	P	E	R	A
0	-2	-4	-6	-8	-10	-12
A	-2					
P	-4					
T	-6					
E	-8					
R	-10					
A	-12					

N & W : παράδειγμα 2ο

	A	S	P	E	R	A
0	-2	-4	-6	-8	-10	-12
A	-2	3	1			
P	-4					
T	-6					
E	-8					
R	-10					
A	-12					

N & W : παράδειγμα 2ο

	A	S	P	E	R	A
0	-2	-4	-6	-8	-10	-12
A	-2	3	1	-1	-3	-5
P	-4	1	3	4	2	0
T	-6	-1	1	3	4	2
E	-8	-3	-1	1	6	4
R	-10	-5	-3	-1	4	9
A	-12	-7	-5	-3	2	7
						12

N & W : παράδειγμα 2ο

	A	S	P	E	R	A	
0	-2	-4	-6	-8	-10	-12	
A	-2	3*	1	-1	-3	-5	-7
P	-4	1	3*	4	2	0	-2
T	-6	-1	1	3*	4	2	0
E	-8	-3	-1	1	6*	4	2
R	-10	-5	-3	-1	4	9*	7
A	-12	-7	-5	-3	2	7	12*

A S P E R A
A P T E R A

ΤΟΠΙΚΕΣ VS ΟΛΙΚΕΣ ΣΤΟΙΧΙΣΕΙΣ

Ο αλγόριθμος των Needleman & Wunsch οδηγεί στην εύρεση της βέλτιστης ολικής στοίχισης μεταξύ δύο αλληλουχιών, δηλ. όλα τα αμινοξέα της μίας στοιχίζονται με όλα τα αμινοξέα της άλλης.

Άρα, η χρήση του αλγορίθμου υπονοεί πως γνωρίζουμε (ή υποθέτουμε) ότι οι αλληλουχίες είναι ομόλογες καθ'όλο το μήκος τους.

Συνηθέστατα, αυτό δεν ισχύει και η ομολογία μεταξύ των αλληλουχιών δεν είναι είναι καθολική αλλά τοπική (π.χ. δύο multi-domain πρωτεΐνες που έχουν ομολογία μόνο σε ένα από τα domains τους).

ΤΟΤΙΚΕΣ VS ΟΛΙΚΕΣ ΣΤΟΙΧΙΣΕΙΣ

Σε τέτοιες περιπτώσεις αυτό που χρειαζόμαστε είναι μια μέθοδο προσδιορισμού της τοπικά βέλτιστης στοίχισης. Δηλ. ζητάμε να βρούμε για ποιές υπακολουθίες των αλληλουχιών η βαθμολογία της μεταξύ τους στοίχισης μεγιστοποιείται (και, βέβαια, ποιά είναι αυτή η στοίχιση).

Η αλγορίθμική υλοποίηση αυτής της ιδέας είναι πολύ απλούστερη από την εκφώνηση της : Η λύση είναι μια τροποποιημένη μορφή του αλγόριθμου των Needleman-Wunsch για ολικές στοιχίσεις. Αυτός ο τροποποιημένος αλγόριθμος είναι γνωστός ως ο αλγόριθμος των Smith & Waterman για τοπικές στοιχίσεις.

Smith & Waterman

Για τον αλγόριθμο των Smith & Waterman οι συνθήκες αρχικοποίησης είναι :

Αρχικοποίηση :

$$S(0, 0) = 0$$

$$S(\kappa, 0) = 0$$

$$S(0, \mu) = 0$$

Smith & Waterman

Για τον αλγόριθμο των Smith & Waterman η συνθήκη επανάληψης είναι :

$$S(\kappa+1, \mu+1) = \max \{ S(\kappa, \mu) + K[A(\kappa+1), B(\mu+1)], \\ S(\kappa+1, \mu) - \alpha, \\ S(\kappa, \mu+1) - \alpha, \\ 0 \} \\ }$$

Δηλαδή, η βαθμολογία δεν μπορεί να γίνει αρνητική (καθώς π.χ. ένα κενό αυξάνει σε μήκος).

Smith & Waterman

Άρα, για την εύρεση μίας τοπικά βέλτιστης στοίχισης μεταξύ δυο αλληλουχιών, ουσιαστικά εφαρμόζουμε αυτούσιο τον αλγόριθμο των N & W, με τη διαφορά ότι αρχικοποιούμε την πρώτη στήλη και γραμμή του πίνακα στο μηδέν, και ότι όποτε η βαθμολογία (με βάση τον αλγόριθμο των N & W) θα έπρεπε να γίνει αρνητική, την αντικαθιστούμε με το μηδέν.

S & W : πταράδειγμα

Θέλουμε να στοιχίσουμε τις αλληλουχίες APTERA και PERASPERAADASTRA. Ο (μοναδιαίος) πίνακας βαθμολόγησης είναι :

A	2				
C	0	2			
M	0	0	2		
P	0	0	0	2	
F	0	0	0	0	2
.....					
A	C	M	P	F	...

και το gap penalty έχει τιμή 1 (ανά κενό).

S & W : παράδειγμα

	P	E	R	A	S	P	E	R	A	A	D	A	S	T	R	A
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0															
P	0															
T	0															
E	0															
R	0															
A	0															

Αρχικοποίηση του πίνακα.

S & W : παράδειγμα

	P	E	R	A	S	P	E	R	A	A	D	A	S	T	R	A
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0													
P	0															
T	0															
E	0															
R	0															
A	0															

Συμπλήρωση.

S & W : παράδειγμα

	P	E	R	A	S	P	E	R	A	A	D	A	S	T	R	A	
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
A	0	0	0	0	2	1	0	0	0	2	2	1	2	1	0	0	2
P	0	2	1	0	1	2	3	2	1	1	2	2	1	2	1	0	1
T	0																
E	0																
R	0																
A	0																

Συμπλήρωση.

S & W : παράδειγμα

	P	E	R	A	S	P	E	R	A	A	D	A	S	T	R	A
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	0	2	1	0	0	0	2	2	1	2	1	0	0
P	0	2	1	0	1	2	3	2	1	1	2	2	1	2	1	0
T	0	1	2	1	0	1	2	3	2	1	1	2	2	1	4	3
E	0	0	3	2	1	0	1	4	3	2	1	1	2	2	3	4
R	0	0	2	5	4	3	2	3	6	5	4	3	2	2	2	5
A	0	0	1	4	7	6	5	4	5	8	7	6	5	4	3	4

Συμπλήρωση.

S & W : παράδειγμα

	P	E	R	A	S	P	E	R	A	A	D	A	S	T	R	A
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	0	2	1	0	0	0	2	2	1	2	1	0	0
P	0	2	1	0	1	2	3	2	1	1	2	2	1	2	1	0
T	0	1	2	1	0	1	2	3	2	1	1	2	2	1	4	3
E	0	0	3	2	1	0	1	4	3	2	1	1	2	2	3	4
R	0	0	2	5	4	3	2	3	6	5	4	3	2	2	2	5
A	0	0	1	4	7	6	5	4	5	8*	7	6	5	4	3	4

Μονοπάτι μίας τοπικά βέλτιστης στοίχισης.

S & W: πταράδειγμα

	P	E	R	A	S	P	E	R	A	A	D	A	S	T	R	A
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	0	2*	1	0	0	0	2	2	1	2	1	0	0
P	0	2	1	0	1	2*	3	2	1	1	2	2	1	2	1	0
T	0	1	2	1	0	1	2*	3	2	1	1	2	2	1	4	3
E	0	0	3	2	1	0	1	4*	3	2	1	1	2	2	3	4
R	0	0	2	5	4	3	2	3	6*	5	4	3	2	2	2	5
A	0	0	1	4	7	6	5	4	5	8*	7	6	5	4	3	4

Μονοπάτι μίας τοπικά βέλτιστης στοίχισης.

P E R A S P E R A A D A S T R A
A P T E R A

S & W : πταράδειγμα

	P	E	R	A	S	P	E	R	A	A	D	A	S	T	R	A	
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
A	0	0	0	0	2*	1*	0	0	0	2	2	1	2	1	0	0	2
P	0	2	1	0	1	2	3*	2	1	1	2	2	1	2	1	0	1
T	0	1	2	1	0	1	2*	3	2	1	1	2	2	1	4	3	2
E	0	0	3	2	1	0	1	4*	3	2	1	1	2	2	3	4	3
R	0	0	2	5	4	3	2	3	6*	5	4	3	2	2	2	5	4
A	0	0	1	4	7	6	5	4	5	8*	7	6	5	4	3	4	7

Μονοπάτια τοπικά βέλτιστων στοιχίσεων ?

P E R A S P - E R A A A D A S T R A
A - P T E R A

S & W: πταράδειγμα

	P	E	R	A	S	P	E	R	A	A	D	A	S	T	R	A
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	0	2	1	0	0	0	2	2	1	2	1	0	0
P	0	2*	1	0	1	2	3	2	1	1	2	2	1	2	1	0
T	0	1*	2	1	0	1	2	3	2	1	1	2	2	1	4	3
E	0	0	3*	2	1	0	1	4	3	2	1	1	2	2	3	4
R	0	0	2	5*	4	3	2	3	6	5	4	3	2	2	2	5
A	0	0	1	4	7*	6	5	4	5	8	7	6	5	4	3	4

Μονοπάτι μίας 'πλησίου της βέλτιστης' στοίχισης.

P - E R A S P E R A A D A S T R A
P T E R A

S & W : πταράδειγμα

	P	E	R	A	S	P	E	R	A	A	D	A	S	T	R	A
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	0	2	1	0	0	0	2	2	1	2*	1	0	0
P	0	2	1	0	1	2	3	2	1	1	2	2	1	2*	1	0
T	0	1	2	1	0	1	2	3	2	1	1	2	2	1	4*	3
E	0	0	3	2	1	0	1	4	3	2	1	1	2	2	3*	4
R	0	0	2	5	4	3	2	3	6	5	4	3	2	2	2	5*
A	0	0	1	4	7	6	5	4	5	8	7	6	5	4	3	4
																7*

Μονοπάτι μίας άλλης 'πλησίου της βέλτιστης' στοίχισης.

P E R A S P E R A A D A S T - R A
A P T E R A

Στοίχιση πολλών αλληλουχιών

Η ανάγκη στοίχισης πολλών αλληλουχιών προκύπτει σε μια πληθώρα περιπτώσεων :

- Ταυτοποίηση συγγενών αλληλουχιών (για τη δημιουργία οικογενειών).
- Εύρεση φυλογενετικών σχέσεων μεταξύ των μελών μιας οικογένειας αλληλουχιών (σε συνδυασμό με τους αλγόριθμους φυλογενετικών δένδρων).
- Εύρεση συντηρημένων μοτίβων και αυστηρά συντηρημένων καταλοίπων.
- Βελτίωση της ευαισθησίας της έρευνας των βάσεων δεδομένων μέσω της χρήσης ολόκληρων στοιχίσεων (αντί απλών αλληλουχιών).

Στοίχιση πολλών αλληλουχιών

Μία στοίχιση πολλών αλληλουχιών έχει τη μορφή δισδιάστατου πίνακα στον οποίο οι γραμμές αντιστοιχούν στις αλληλουχίες και οι στήλες στις αμινοξικές θέσεις. Για παράδειγμα :

A	S	P	-	E	R	A
A	-	P	T	E	R	A
A	S	-	T	-	R	A

Όπως και για τη στοίχιση δύο αλληλουχιών, στόχος της στοίχισης πολλών αλληλουχιών είναι να αποκαλύψει τις μεταξύ τους εξελικτικές σχέσεις.

Ομοιότητα με τη στοίχιση δύο αλληλουχιών

Το πρόβλημα της στοίχισης πολλών αλληλουχιών φαίνεται να είναι μια άμεση επέκταση του προβλήματος της στοίχισης δύο αλληλουχιών : αντί να βρίσκουμε το βέλτιστο μονοπάτι σε ένα δισδιάστατο πίνακα [όπως κάναμε με τους αλγόριθμους δυναμικού πργραμματισμού (N&W, S&W)], τώρα θα αναζητούμε το βέλτιστο μονοπάτι σε ένα πίνακα υψηλότερης διάστασης (σε ένα τρισδιάστατο πίνακα για τρεις αλληλουχίες, σε ένα τετραδιάστατο πίνακα για τέσσερις, κοκ).

Υπολογιστικό πρόβλημα

Το πρόβλημα με την προσέγγιση αυτή είναι το υπολογιστικό κόστος : για δύο αλληλουχίες με μήκη m το κόστος εύρεσης της βέλτιστης στοίχισης είναι ($m \cdot m$). Για τρεις θα είναι ($m \cdot m \cdot m$), και για K αλληλουχίες θα είναι (m^K). Έτσι επιστρέψαμε σε ένα υπολογιστικό κόστος εκθετικά ανάλογο του αριθμού των αλληλουχιών, μόνο που σε αυτή την περίπτωση δεν είναι 'φαινομενικό' (όπως είχαμε δείξει για την περίπτωση της στοίχισης δυο αλληλουχιών), αλλά πραγματικό. Έτσι, το υπολογιστικό κόστος της στοίχισης τεσσάρων αλληλουχιών μήκους 100 κατάλοιπων είναι το ίδιο με το κόστος της πραγματοποίησης 10000 στοιχίσεων μεταξύ δύο τέτοιων αλληλουχιών.

Υπολογιστικό πρόβλημα

Για μικρό αριθμό αλληλουχιών (π.χ. τρεις αλληλουχίες) υπάρχουν αλγόριθμοι (και προγράμματα) για την εύρεση της βέλτιστης μεταξύ τους στοίχισης. Αυτός ο περιορισμός σε αριθμό αλληλουχιών μειώνει κατά πολύ τη βιολογική χρησιμότητα αυτών των αλγορίθμων.

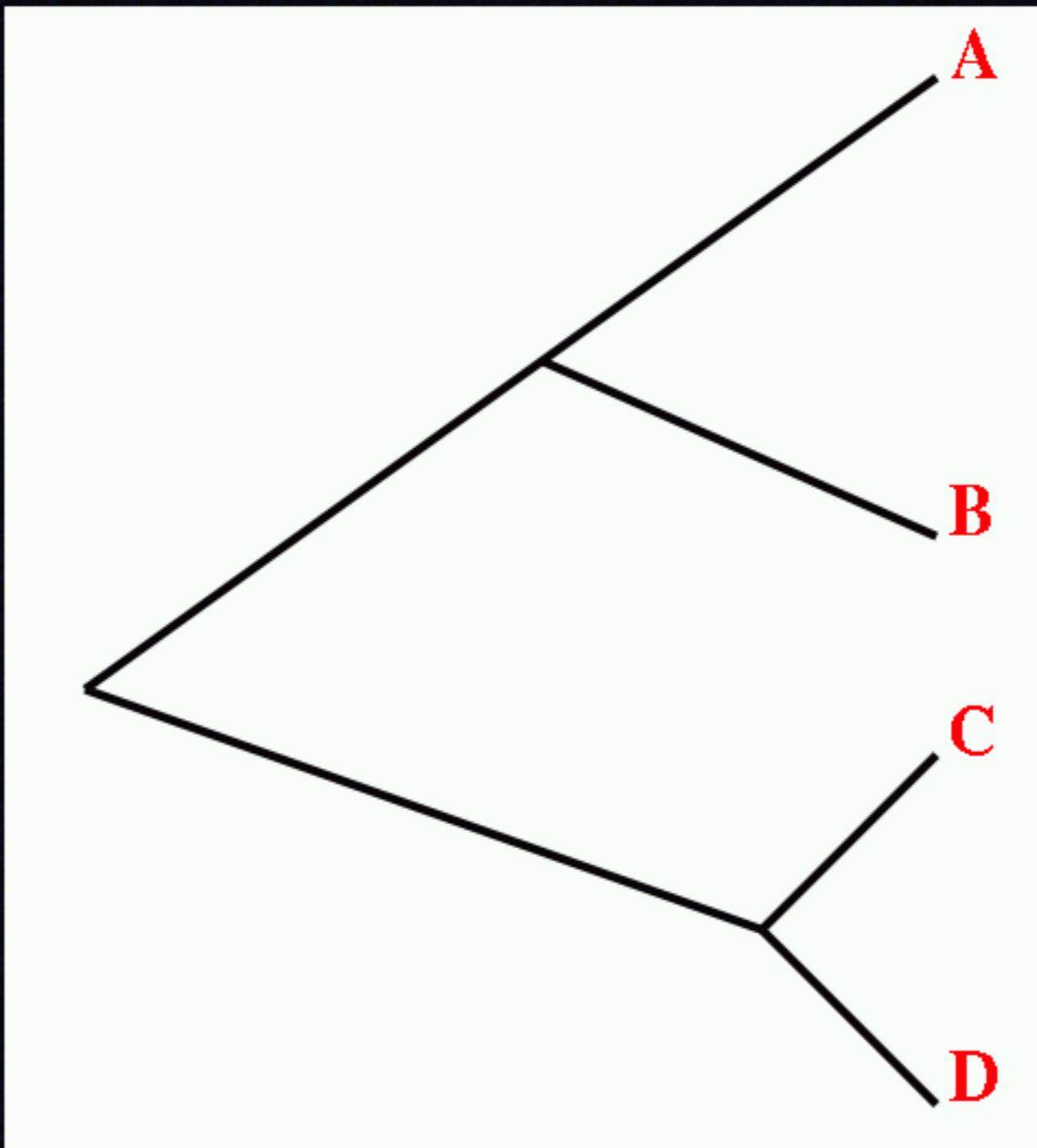
Για αυτό το λόγο, ο χώρος της στοίχισης πολλών αλληλουχιών κυριαρχείται από ευρεστικούς αλγόριθμους οι οποίοι χρησιμοποιούν επιπλέον υποθέσεις και περιορισμούς προκειμένου να μειώσουν το υπολογιστικό κόστος. Οι πλέον γνωστοί αλγόριθμοι για στοίχιση πολλών αλληλουχιών ανήκουν στις λεγόμενες προοδευτικές (progressive) μεθόδους (επίσης γνωστές και ως ιεραρχικές).

Ιεραρχικές μέθοδοι

Η βασική ιδέα πίσω από αυτές τις μεθόδους είναι ότι η ύπαρξη εξελικτικών σχέσεων μεταξύ των αλληλουχιών καθιστά περιττή την αναζήτηση μιας καθολικά βέλτιστης στοίχισης : η στοίχιση που αναζητούμε είναι αυτή που καλύτερα αναπαριστά τις μεταξύ τους εξελικτικές σχέσεις και όχι αυτή που αποδίδει το σύνολο των μεταξύ τους ομοιοτήτων.

Έτσι, εάν για παράδειγμα γνωρίζαμε ότι για τέσσερις αλληλουχίες A, B, C, D ισχύει το δένδρο :

Ιεραρχικές μέθοδοι



τότε, η άμεση αναζήτηση σχέσεων μεταξύ των αλληλουχιών A και D περιττεύει. Καλύτερο είναι το :

Ιεραρχικές μέθοδοι

- Βρες την εξελικτική σχέση ανάμεσα στις αλληλουχίες A & B. Η μεταξύ τους ομολογία αναπαριστά με πληρότητα ό,τι μπορούμε να συνάγουμε για την κοινή τους προγονική αλληλουχία (τον κόμβο του δένδρου από τον οποίο προήλθαν).
- Βρες την εξελικτική σχέση ανάμεσα στις αλληλουχίες C & D. Η μεταξύ τους ομολογία αναπαριστά με πληρότητα ό,τι μπορούμε να συνάγουμε για την κοινή τους προγονική αλληλουχία (τον κόμβο του δένδρου από τον οποίο προήλθαν).
- Χρησιμοποίησε ό,τι γνωρίζεις για τους δύο ενδιάμεσους κόμβους για να συνάγεις την μεταξύ τους εξελικτική σχέση (και με τη ρίζα του δένδρου).

Ιεραρχικές μέθοδοι

Άρα, το "στοίχισε τις αλληλουχίες A,B,C,D"
μετασχηματίστηκε στο :

- Στοίχισε τις A & B => AB
- Στοίχισε τις C & D => CD
- Στοίχισε τις AB & CD => ABCD

Η διαφορά από άποψη υπολογιστικού κόστους είναι
τεράστια : εάν οι τέσσερις αλληλουχίες είχαν μήκος
200 καταλοίπων, η καινούργια μέθοδος θα ήταν
~13000 φορές πιο γρήγορη από μια σχολαστική
μέθοδο στοίχισης.

Ακόμη ένας φαύλος κύκλος ;

Η μέθοδος που παρουσιάστηκε έχει ένα ουσιώδες πρόβλημα : για να βρούμε τις φυλογενετικές σχέσεις μεταξύ των αλληλουχιών χρειαζόμαστε την μεταξύ τους στοίχιση (όλων των αλληλουχιών). Άρα, για να βρούμε τη στοίχιση χρειαζόμαστε το προϊόν της. Το αδιέξοδο αυτό αίρεται μέσω της παραδοχής ότι ένα δενδρόγραμμα-οδηγός για τις εξελικτικές σχέσεις μεταξύ των αλληλουχιών μπορεί να δημιουργηθεί λαμβάνοντας υπόψη μόνο τις ομοιότητες μεταξύ ζευγών αλληλουχιών (χωρίς τη δημιουργία μιας στοίχισης όλων των αλληλουχιών).

Ακόμη ένας φαύλος κύκλος ;

Για παράδειγμα, υποθέστε ότι για τρεις αλληλουχίες A,B,Γ πραγματοποιήσαμε όλες τις δυνατές ανά ζεύγη στοιχίσεις, και ελέγξαμε κάθε μία από αυτές με βάση το Z-τεστ που αναφέρθηκε στην προηγούμενη διάλεξη. Τα αποτελέσματα (με τη μορφή πίνακα) ήταν :

	A	B	Γ
A	-	9	4
B	9	-	5
Γ	4	5	-

Από αυτόν το πίνακα μπορούμε χωρίς καμία επιπλέον ανάλυση να συνάγουμε ότι οι αλληλουχίες A και B είναι πιο στενά συνδεδεμένες μεταξύ τους απ' ότι κάθε μία από αυτές με την αλληλουχία Γ.

Ακόμη ένας φαύλος κύκλος ;

Συνεπώς, αυτό που θα κάναμε σε αυτή την περίπτωση είναι να στοιχίσουμε τις A & B (χρησιμοποιώντας για παράδειγμα τον αλγόριθμο των N & W), και στη συνέχεια, να στοιχίσουμε (πάλι με τον N & W) την αλληλουχία Γ με την προϋπάρχουσα στοίχιση των A & B.

Το οποίο μας φέρνει στην ερώτηση πως στοιχίζουμε βέλτιστα (κατά N & W) όχι δυο αλληλουχίες, αλλά μια αλληλουχία και μια ολόκληρη στοίχιση, ή ακόμα και δύο στοιχίσεις μεταξύ τους.

ΣΤΟΙΧΙΣΗ ΣΤΟΙΧΙΣΕΩΝ

Η βασική απλοποίηση του προβλήματος προκύπτει από την απαίτηση ότι οι ήδη στοιχισμένες αλληλουχίες θα χρησιμοποιηθούν όχι ως ανεξάρτητες αλληλουχίες, άλλα ως μια (μικτή) αλληλουχία : το ποιο αμινοξύ της μιας είναι στοιχισμένο με ποιο αμινοξύ της άλλης πρόκειται να μείνει αμετάβλητο. Εάν απαιτηθεί η προσθήκη κενών, αυτά (τα κενά) μπαίνουν ταυτόχρονα σε όλες τις ήδη στοιχισμένες αλληλουχίες. Π.χ.

ASFKLMEMNERA

+

ASPERA

APTERA

==>

ASFKLMEMTERA

AS-----PERA

AP-----TERA

Στοίχιση στοιχίσεων

Η ουσιαστική διαφορά από τον απλό αλγόριθμο του N & W, έγκειται στον τρόπο βαθμολόγησης τόσο για τις στοιχίσεις μεταξύ αμινοξέων όσο και για την εισαγωγή κενών. Για πληρότητα θα αναφέρουμε έναν από τους πλέον απλοϊκούς αθροιστικούς τρόπους βαθμολόγησης :

Η βαθμολογία της στοίχισης K αμινοξέων από μία θέση μίας στοίχισης A, με Λ αμινοξέα από μία στοίχιση B, είναι ίση το άθροισμα των βαθμολογιών υποκατάστασης κάθε αμινοξέος (από τα K) της στοίχισης A με κάθε ένα (από τα Λ) της στοίχισης B. Για παράδειγμα :

ΣΤΟΙΧΙΣΗ ΣΤΟΙΧΙΣΕΩΝ

Έστω δύο στοιχίσεις τεσσάρων αλληλουχιών,

ASPERA καὶ **AMEMPTA**
APTERA **ASE-PTA**

Η βαθμολογία της στοίχισης των P-T (τρίτη θέση πρώτης στοίχισης) με τα M-S (δεύτερη θέση δεύτερης στοίχισης) θα είναι

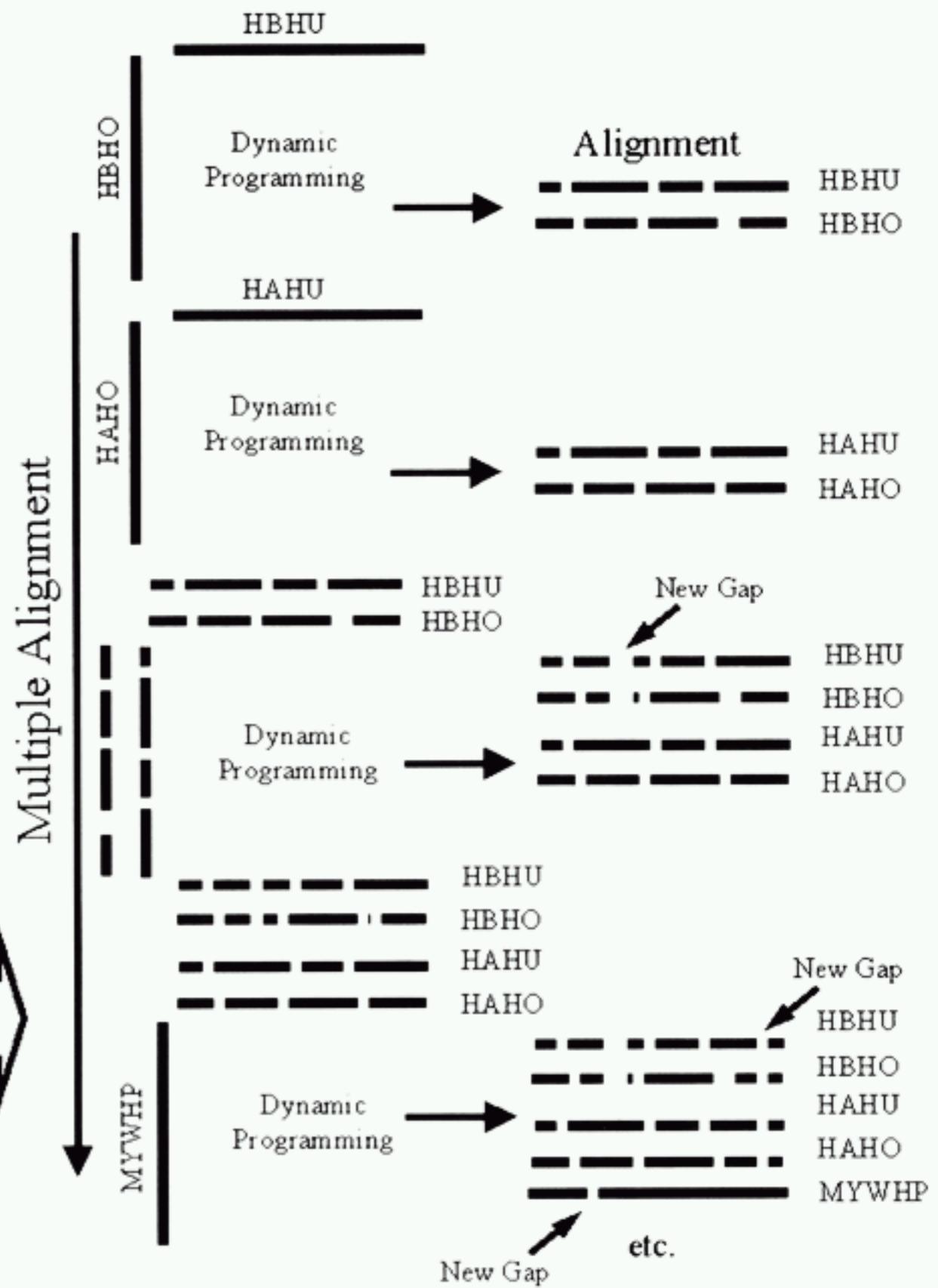
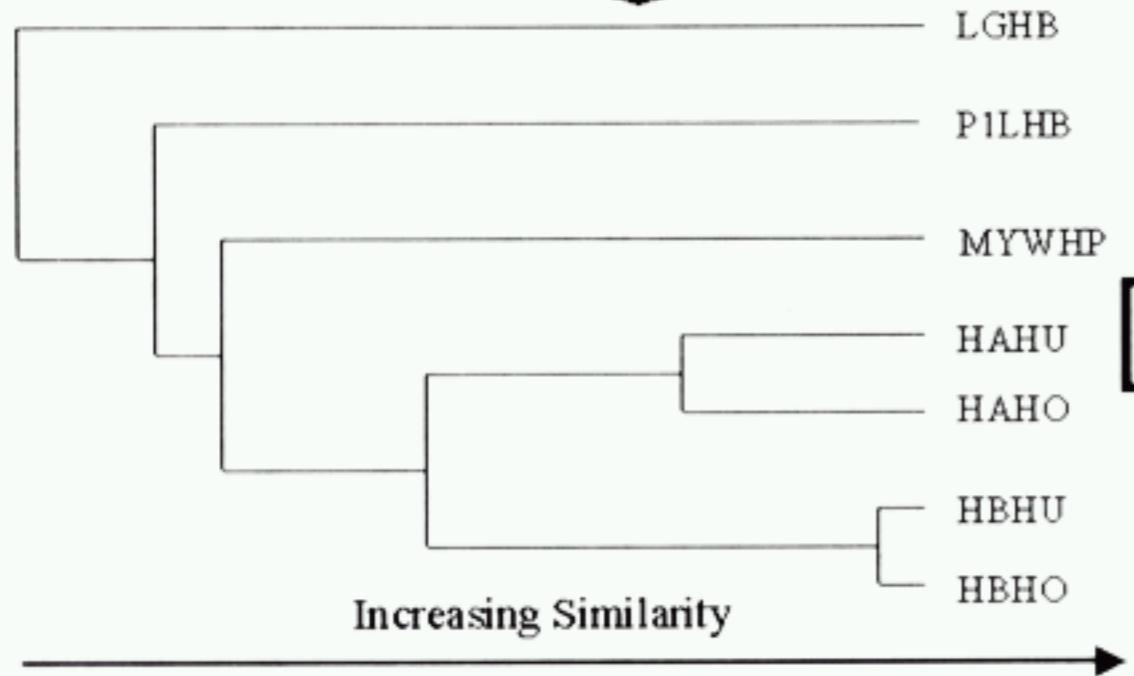
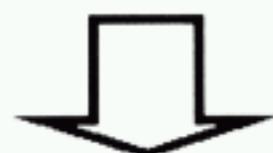
$$\Sigma(P,M) + \Sigma(P,S) + \Sigma(T,M) + \Sigma(T,S)$$

όπου Σ είναι ο πίνακας βαθμολόγησης (π.χ. PAM250).

Ιεραρχικές μέθοδοι

	HAHU	HBHU	HAHO	HBHO	MYWHP	P1LHB	LGHB
HAHU							
HBHU	21.1						
HAHO	32.9	19.7					
HBHO	20.7	39.0	20.4				
MYWHP	11.0	9.8	10.3	9.7			
P1LHB	9.3	8.6	9.6	8.4	7.0		
LGHB	7.1	7.3	7.5	7.4	7.3	4.3	

Cluster Analysis



Παράδειγμα

Θέλουμε να στοιχίσουμε τις αλληλουχίες ASPERA
ASTRA και APTERA. Ο (μοναδιαίος) πίνακας
βαθμολόγησης είναι :

A	3						
C	0	3					
M	0	0	3				
P	0	0	0	3			
F	0	0	0	0	3		
.....							
A	C	M	P	F	...		

και το gap penalty έχει τιμή 2 (ανά κενό).

Παράδειγμα

Η μέθοδος επίλυσης είναι :

- Στοιχίζουμε τις αλληλουχίες ASPERA και APTERA
=> Βαθμολογία Σ1
- Στοιχίζουμε τις αλληλουχίες ASPERA και ASTRA
=> Βαθμολογία Σ2
- Στοιχίζουμε τις αλληλουχίες APTERA και ASTRA
=> Βαθμολογία Σ3
- Για το ζεύγος αλληλουχιών με την υψηλότερη βαθμολογία βρίσκουμε τη βέλτιστη στοίχιση.
- Τέλος, στοιχίζουμε το ήδη στοιχισμένο ζεύγος με την τρίτη αλληλουχία.

Παράδειγμα

ASPERA-APTERA

	A	S	P	E	R	A
0	-2	-4	-6	-8	-10	-12
A	-2					
P	-4					
T	-6					
E	-8					
R	-10					
A	-12					

Παράδειγμα

ASPERA-APTERA

	A	S	P	E	R	A	
A	0	-2	-4	-6	-8	-10	-12
P	-2	3	1	-1	-3	-5	-7
T	-4	1	3	4	2	0	-2
E	-6	-1	1	3	4	2	0
R	-8	-3	-1	1	6	4	2
R	-10	-5	-3	-1	4	9	7
A	-12	-7	-5	-3	2	7	12

Παράδειγμα

ASPERA-ASTRA

	A	S	P	E	R	A	
A	0	-2	-4	-6	-8	-10	-12
S		-2					
T			-4				
R				-6			
A					-8		
						-10	

Παράδειγμα

ASPERA-ASTRA

	A	S	P	E	R	A
0	-2	-4	-6	-8	-10	-12
A	-2	3	1	-1	-3	-5
S	-4	1	6	4	2	0
T	-6	-1	4	6	4	2
R	-8	-3	2	4	6	7
A	-10	-5	0	2	4	6
						10

Παράδειγμα

APTERA-ASTRA

	A	P	T	E	R	A	
A	0	-2	-4	-6	-8	-10	-12
A	-2						
S	-4						
T	-6						
R	-8						
A	-10						

Παράδειγμα

APTERA-ASTRA

	A	P	T	E	R	A
0	-2	-4	-6	-8	-10	-12
A	-2	3	1	-1	-3	-5
S	-4	1	3	1	-1	-3
T	-6	-1	1	6	4	2
R	-8	-3	-1	4	6	7
A	-10	-5	-3	2	4	6
						10

Παράδειγμα

Βέλτιστη στοίχιση ASPERA-APTERA

	A	S	P	E	R	A	
0	-2	-4	-6	-8	-10	-12	
A	-2	3*	1	-1	-3	-5	-7
P	-4	1	3*	4	2	0	-2
T	-6	-1	1	3*	4	2	0
E	-8	-3	-1	1	6*	4	2
R	-10	-5	-3	-1	4	9*	7
A	-12	-7	-5	-3	2	7	12*

A S P E R A
A P T E R A

Παράδειγμα

ASPERA/APTERA vs ASTRA

	A/A	S/P	P/T	E/E	R/R	A/A
0	-2	-4	-6	-8	-10	-12
A	-2					
S	-4					
T	-6					
R	-8					
A	-10					

Παράδειγμα

ASPERA/APTERA vs ASTRA

	A/A	S/P	P/T	E/E	R/R	A/A	
0	-2	-4	-6	-8	-10	-12	
A	-2	6	4	2	0	-2	-4
S	-4	4	9	7	5	3	1
T	-6	2	7	12	10	8	6
R	-8	0	5	10	12	16	14
A	-10	-2	3	8	10	14	22

Παράδειγμα

ASPERA/APTERA vs ASTRA

	A/A	S/P	P/T	E/E	R/R	A/A	
0	-2	-4	-6	-8	-10	-12	
A	-2	6*	4	2	0	-2	-4
S	-4	4	9*	7	5	3	1
T	-6	2	7	12*	10*	8	6
R	-8	0	5	10	12	16*	14
A	-10	-2	3	8	10	14	22*

A S P E R A
A P T E R A
A S T - R A

Παράδειγμα 2

Εάν η τιμή του gap penalty ήταν ίση με 1 η βέλτιστη στοίχιση των ASPERA και APTERA θα ήταν :

A S P - E R A
A - P T E R A

οπότε ο πίνακας στοίχισης του ASTRA με αυτή τη στοίχιση θα ήταν :

Παράδειγμα 2

ASPERA/APTERA vs ASTRA

	A/A	S/-	P/P	-/T	E/E	R/R	A/A
0	-1	-2	-3	-4	-5	-6	-7
A	-1						
S	-2						
T	-3						
R	-4						
A	-5						

Παράδειγμα 2

ASPERA/APTERA vs ASTRA

	A/A	S/-	P/P	-/T	E/E	R/R	A/A
0	-1	-2	-3	-4	-5	-6	-7
A	-1	6	5	4	3	2	1
S	-2	5	9	8	7	6	5
T	-3	4	8	9	11	10	9
R	-4	3	7	8	10	11	16
A	-5	2	6	7	9	10	15
							22

Παράδειγμα 2

ASPERA/APTERA vs ASTRA

	A/A	S/-	P/P	-/T	E/E	R/R	A/A
0	-1	-2	-3	-4	-5	-6	-7
A	-1	6*	5	4	3	2	1
S	-2	5	9*	8*	7	6	5
T	-3	4	8	9	11*	10*	9
R	-4	3	7	8	10	11	16*
A	-5	2	6	7	9	10	15
							22*

A S P - E R A
A - P T E R A
A S - T - R A

Ιεραρχικές μέθοδοι

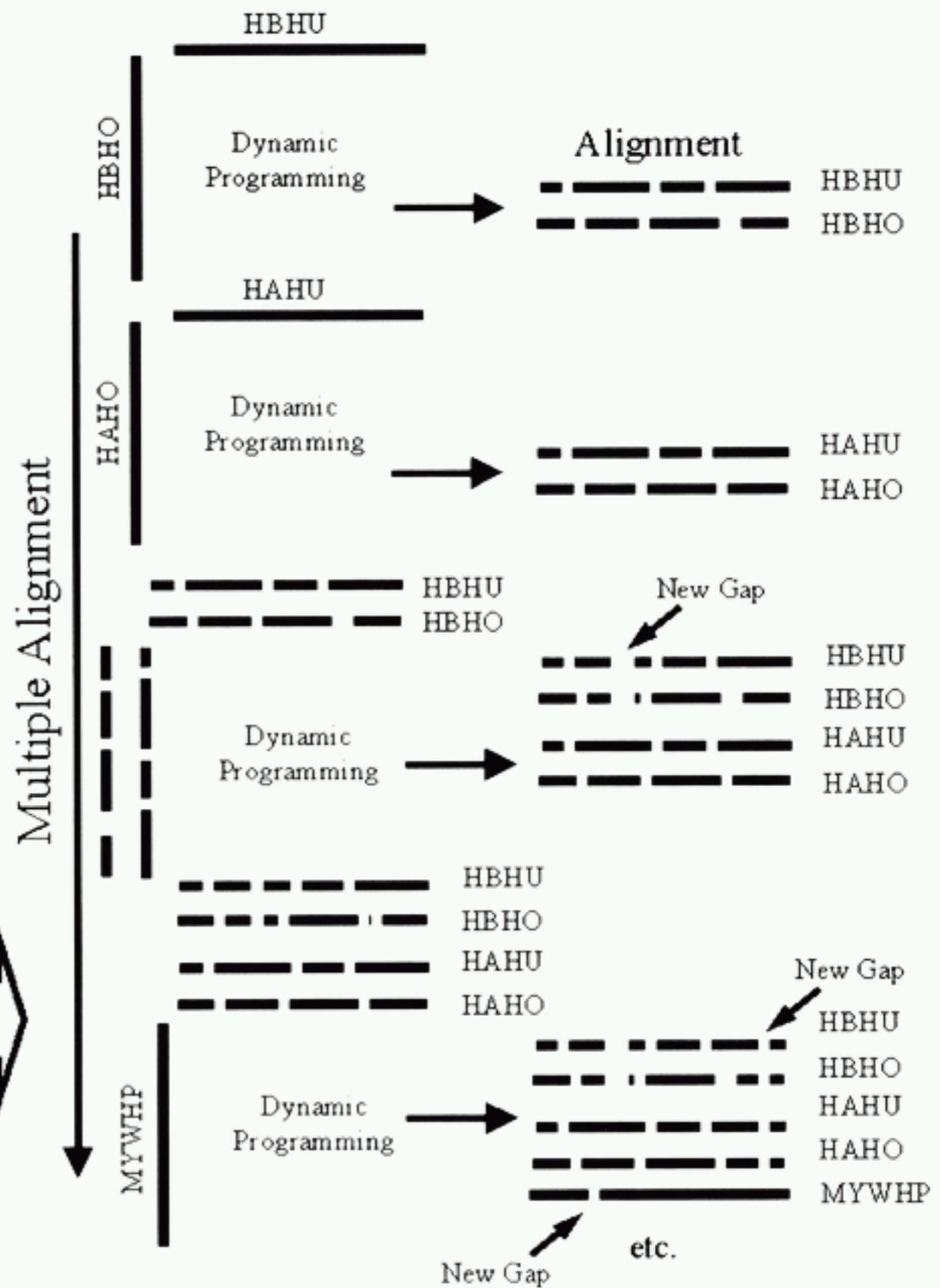
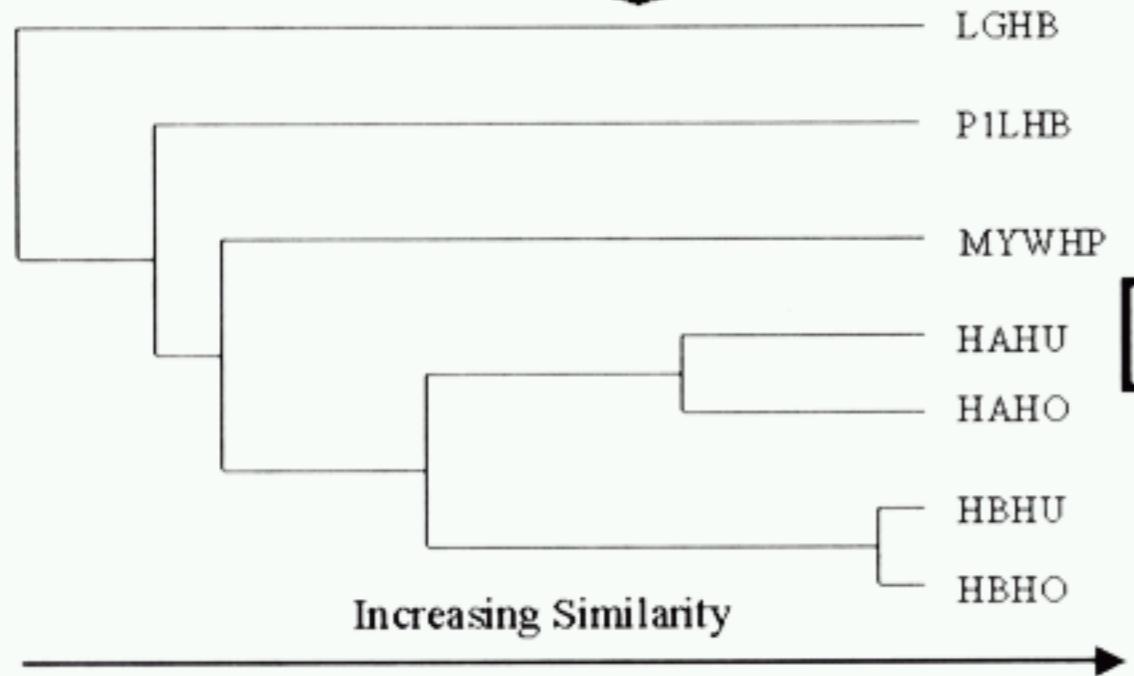
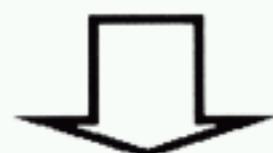
Προφανώς, το απλοϊκό αυτό παράδειγμα δεν μπορεί να καλύψει το εύρος των επιπλέον τεχνικών που χρησιμοποιούνται από τα διάφορα προγράμματα για περαιτέρω βελτίωση των προκυπτόντων στοιχίσεων. Σημειώστε επίσης ότι δεν προσπαθήσαμε καν να συνάγουμε (όπως θα έπρεπε) ένα δενδρόγραμμα από τις ανά ζεύγη στοιχίσεις. Ο λόγος είναι ότι η φυλογενετική ανάλυση αλληλουχιών θα αναλυθεί εκτενέστερα σε μελλοντική διάλεξη.

Παρ' όλα αυτά, ο βασικός αλγόριθμος για την ιεραρχική στοίχιση πολλών αλληλουχιών είναι αυτός που ήδη αναφέρθηκε :

Ιεραρχικές μέθοδοι

	HAHU	HBHU	HAHO	HBHO	MYWHP	P1LHB	LGHB
HAHU							
HBHU	21.1						
HAHO	32.9	19.7					
HBHO	20.7	39.0	20.4				
MYWHP	11.0	9.8	10.3	9.7			
P1LHB	9.3	8.6	9.6	8.4	7.0		
LGHB	7.1	7.3	7.5	7.4	7.3	4.3	

Cluster Analysis



Βελτιώσεις και επτεκτάσεις

Ο βασικός αλγόριθμος που περιγράψαμε μπορεί να βελτιωθεί. Μερικές από τις πλέον δημοφιλείς βελτιώσεις είναι :

- Επειδή τυχόν λάθη που θα γίνουν στα αρχικά στάδια της στοίχισης, θα συντηρηθούν και στα επόμενα, αρκετά προγράμματα κάνουν μετά το πέρας της αρχικής στοίχισης, επαναστοίχιση των αλληλουχιών. Υποθέστε για παράδειγμα ότι η βέλτιστη στοίχιση του πρώτου ζεύγους αλληλουχιών ήταν :

... **DEFLMPEF** ...
... **DEEKSTEF** ...

άλλα μετά το τέλος της στοίχισης όλες οι υπόλοιπες αλληλουχίες είχαν το μοτίβο :

Βελτιώσεις και επτεκτάσεις

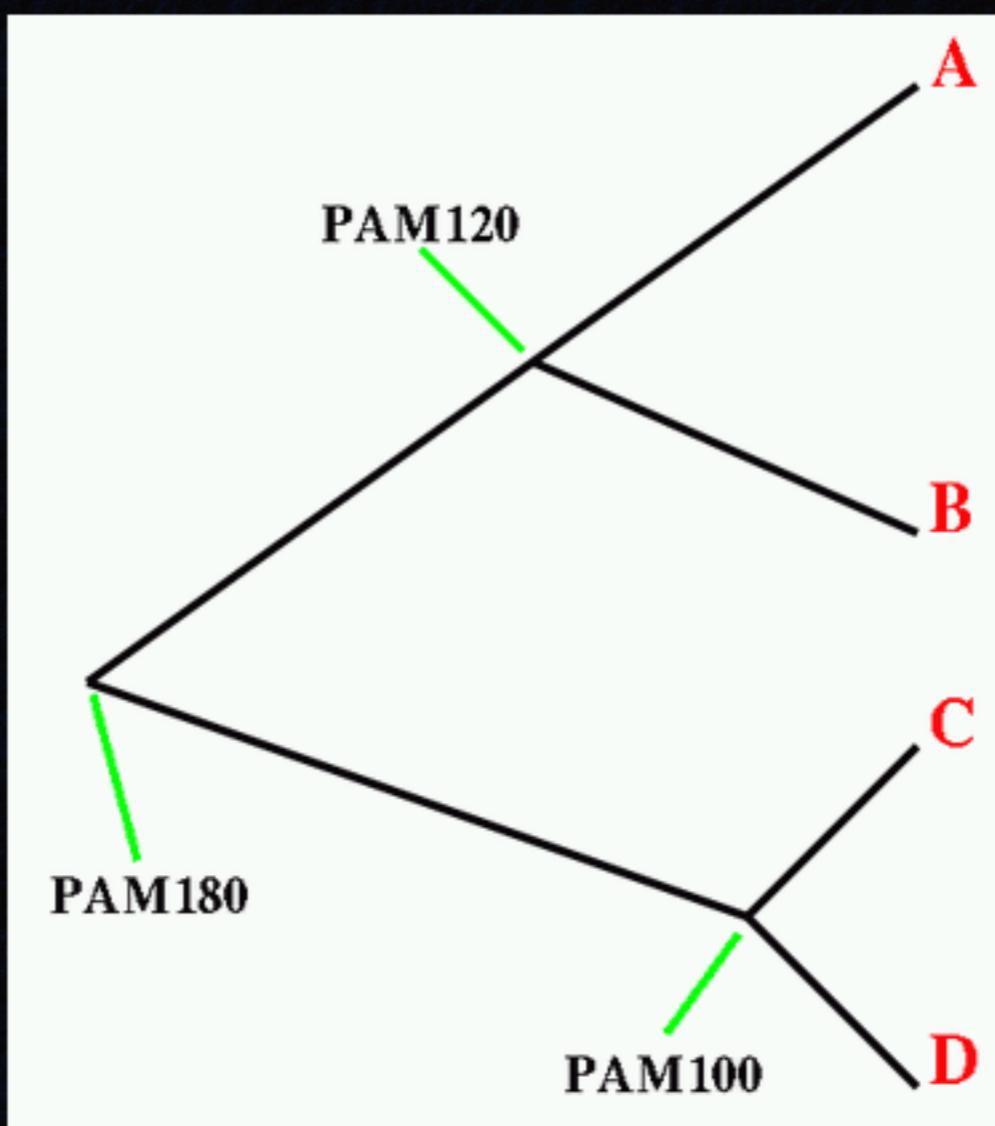
.... DEEKSTFLMPEF
.... DEEKSTFLMPEF
.....

Προφανώς αυτό που απαιτείται είναι μια διόρθωση της αρχικής στοίχισης ώστε να γίνει :

.... DE----FLMPEF
.... DEEKST---EF
.... DEEKSTFLMPEF
.... DEEKSTFLMPEF
.....

Βελτιώσεις και επτεκτάσεις

Μια άλλη συνηθισμένη βελτίωση είναι η χρήση διαφορετικών πινάκων βαθμολόγησης σε διαφορετικά στάδια ανάπτυξης της στοίχισης. Αυτή η βελτίωση έχει προφανές βιολογικό περιεχόμενο :



Βελτιώσεις και επτεκτάσεις

Εάν η δευτεραγής δομή μίας εκ των αλληλουχιών είναι γνωστή, τότε μπορεί να υλοποιηθεί μία από τις βασικότερες βελτιώσεις (η οποία μάλιστα μπορεί να μεταφερθεί αυτούσια και στους αλγόριθμους στοίχισης δύο αλληλουχιών). Η βελτίωση αφορά την χρήση διαφορετικών τιμών για το gap penalty ανάλογα με το εάν το κενό πρόκειται να εισαχθεί σε μια περιοχή της αλληλουχίας που αντιστοιχεί σε ένα στοιχείο δευτεραγούς δομής (α, β) ή όχι (στροφές και random coil). Η διόρθωση αυτή έχει προφανές βιολογικό περιεχόμενο (indels στο μέσο μιας στροφής απορροφούνται πιο εύκολα απ' ότι indels στο μέσο ενός στοιχείου δευτεραγούς δομής).

Βελτιώσεις και επτεκτάσεις

Μια άλλη τεχνική που αφορά διαφορική εφαρμογή των gap penalties, στηρίζεται στις παρακάτω υποθέσεις :

- Οι περιοχές των πλέον όμοιων αλληλουχιών οι οποίες συσσωρεύουν κενά, είναι και οι περιοχές στις οποίες το gap penalty για τις υπόλοιπες (λιγότερο συγγενείς αλληλουχίες) θα πρέπει να έχει μικρότερες τιμές.
- Οι περιοχές των πλέον όμοιων αλληλουχιών οι οποίες δείχνουν τη μεγαλύτερη ποικιλότητα, είναι επίσης περιοχές για τις οποίες το κόστος εισαγωγής κενών (για τις υπόλοιπες, λιγότερο συγγενείς αλληλουχίες) θα πρέπει να μειωθεί.

Παράδειγμα προγράμματος: ClustalW

```
*****  
***** CLUSTAL W (1.82) Multiple Sequence Alignments *****  
*****
```

- 1. Sequence Input From Disc
- 2. Multiple Alignments
- 3. Profile / Structure Alignments
- 4. Phylogenetic trees

- S. Execute a system command
- H. HELP
- X. EXIT (leave program)

Your choice: ■

Βάσεις δεδομένων

Υπάρχει μια πληθώρα βάσεων δεδομένων που περιέχουν στοιχίσεις πολλών αλληλουχιών. Ο στόχος αυτών των βάσεων είναι να συγκεντρώσουν και να ομαδοποιήσουν τις υπάρχουσες πρωτοταγείς αλληλουχίες σε οικογένειες. Παρουσίαση της δομής των καταχωρήσεων αυτών των βάσεων μία-προς-μία είναι άσκοπη. Αυτό που αξίζει να αναφερθεί είναι η διάκριση ανάμεσα σε βάσεις που βασίζονται σε αυτόματες στοιχίσεις και σε βάσεις στις οποίες οι στοιχίσεις ελέγχονται από τους φροντιστές της βάσης.

Βάσεις δεδομένων : Pfam

A1AA_HUMAN	SLKYP....AI..MTER.KAA..AILALL.WVV.AL.VVSVGP.LLG...WKEPV..PPDE....RF		
A1AA_RAT	SLKYP....AI..MTER.KAA..AILALL.WAV.AL.VVSVGP.LLG...WKEPV..PPDE....RF		
A1AB_CANFA	SLQYP....TL..VTRR.KAI..LALLGV.WVL.ST.VISIGP.LLG...WKEPA..PNDD....KE		
A1AB_HUMAN	SLQYP....TL..VTRR.KAI..LALLSV.WVL.ST.VISIGP.LLG...WKEPA..PNDD....KE		
A1AB_MESAU	SLQYP....TL..VTRR.KAI..LALLSV.WVL.ST.VISIGP.LLG...WKEPA..PNDD....KE		
A1AB_RAT	SLQYP....TL..VTRR.KAI..LALLSV.WVL.ST.VISIGP.LLG...WKEPA..PNDD....KE		
A1AC_BOVIN	PLRYP....TI..VTQK.RGL..MALLCV.WAL.SL.VISIGP.LFG...WRQPA..PEDE....T		
A1AC_HUMAN	PLRYP....TI..VTQR.RGL..MALLCV.WAL.SL.VISIGP.LFG...WRQPA..PEDE....T		
A1AC_RAT	PLRYP....TI..VTQR.RGV..RALLCV.WVL.SL.VISIGP.LFG...WRQPA..PEDE....T		
OAR_DROME	PINYA....QK..RTVG.RVL..LLISGV.WLL.SL.LISSPP.LIG...W.NDW..PDEFT..SAT		
D1DR_CARAU	PFRYE....RK..MTPR.VAF..VMISGA.WTL.SV.LISFIPVQLK...WHKAQ..PIGFL..EVN		
D1DR_FUGRU	PFRYE....RK..MTPK.VAC..LMISVA.WTL.SV.LISFIPVQLN...WHKAQ..TASYVELNGT		
DADR_DIDMA	PFRYE....RK..MTPK.AAF..ILISVA.WTL.SV.LISFIPVQLN...WHKARPLSSPDG..NVS		
...	...		
A1AA_HUMANCGI..TE.....EAG...YA....VF.....SS		
A1AA_RATCGI..TE.....EVG...YA....IF.....SS		
A1AB_CANFACGV..TE.....EPF...YA....LF.....SS		
A1AB_HUMANCGV..TE.....EPF...YA....LF.....SS		
A1AB_MESAUCGV..TE.....EPF...YA....LF.....SS		
A1AB_RATCGV..TE.....EPF...YA....LF.....SS		
A1AC_BOVINICQI..NE.....EPG...YV....LF.....SA		
A1AC_HUMANICQI..NE.....EPG...YV....LF.....SA		
A1AC_RATICQI..NE.....EPG...YV....LF.....SA		
OAR_DROMEPCEL..TS.....QRG...YV....IY.....SS		
D1DR_CARAUASRR..DLPTDNC.....DSSL...NRT...YA....IS.....SS		
D1DR_FUGRUYAGD..LPPDNCD.....SSL...NRT...YA....IS.....SS		
DADR_DIDMASQDE...TMDNCD.....SSL...SRT...YA....IS.....SS		
...	...		
A1AA_HUMAN	V.CSFY...LPMAVIV.VMY.CRV.YVV....ARS..TTRSLEA.GVKR		
A1AA_RAT	V.CSFY...LPMAVIV.VMY.CRV.YVV....ARS..TTRSL....EA		
A1AB_CANFA	L.GSFY...IPLAVIL.VMY.CRV.YIV....AKR..TTKNL....EA		
A1AB_HUMAN	L.GSFY...IPLAVIL.VMY.CRV.YIV....AKR..TTKNL....EA		
A1AB_MESAU	L.GSFY...IPLAVIL.VMY.CRV.YIV....AKR..TTKNL....EA		
A1AB_RAT	L.GSFY...IPLAVIL.VMY.CRV.YIV....AKR..TTKNL....EA		
A1AC_BOVIN	L.GSFY...VPLTIIL.VMY.CRV.YVV....AKR..ESRGLKS.GLKT		
A1AC_HUMAN	L.GSFY...LPLAIIL.VMY.CRV.YVV....AKR..ESRGLKS.GLKT		
A1AC_RAT	L.GSFY...VPLAIIL.VMY.CRV.YVV....AKR..ESRGLKS.GLKT		
OAR_DROME	L.GSFF...IPLAIMT.IVY.IEI.FVA....TRR..RLRERA..RANK		
D1DR_CARAU	L.ISFY...IPVAIMI.VTY.TQI.YRI....AQK..QIRRIS..ALER		
D1DR_FUGRU	L.ISFY...IPVAIMI.VTY.TRI.YRI....AQK..QIRRIS..ALER		
DADR_DIDMA	L.ISFY...IPVAIMI.VTY.TRI.YRI....AQK..QIRRIS..ALER		
...	...		

Βάσεις δεδομένων : PRINTS

OAR_DROME	PINYAQ...KRTVGRVLLLISGVWLLSLLISSLSP.PLIG.WND.....WPDEFTSATP...
D2DR_RAT	PMLYNTR..YSSKRRVTVMIAIVWVLSFTISCP.LLFG.LNN...T.DQNE.....
D3DR_RAT	PVHYQHGTGQSSCRRVALMITAVWVLAFAVSCP.LLFG.FNT...TGDPSI.....
DADR_RAT	PFQYER...KMTPKAAFILISVAWTLSVLISFI.PVQLSWHKAK.PTWPLDGNFTSLEDT
DBDR_RAT	PFRYER...KMTQRVALVMVGLAWTLSILISFI.PVQLNWHRDKAGSQGQEGLLSNGTPW
A1AB_RAT	SLQYPT...LVTRRKAILALLSVWVLSTVISIG.PLLG.WKE.....PAPNDDKE....
B1AR_RAT	PFRYQS...LLTRARARALVCTVWAISALVSFL.PILMHWW.....RAESD.EARRCYND
B2AR_HUMAN	PKFYQS...LLTKNKARVIILMVWIVSGLTSFL.PIQMHWY.....RATHQ.EAINCYAN
B2AR_RAT	PKFYQS...LLTKNKARVVILMVWIVSGLTSFL.PIQMHWY.....RATHK.QAIDCYAK
B3AR_RAT	PLRYGT...LVTKRRARAADVLLWIVSATVSFA.PIMSQWW.....RVGADAEAQECHSN
5HTA_RAT	PIDYVN...KRTPRRAAALISLTWLGFLISIP.PMLG.WRTPEDRSDPDA.....
5HTD_RAT	ALEYSK...RRTAGHAAAMIAAVWAISICISIP.PLF..WRQ..ATAHEEMSD.....
5HT2_RAT	PIHHSR...FNSRTKAFLKIIAVWTISVGISMPIPVFG.LQDDSKVFKEGS.....
...
OAR_DROMECELTSQRG.....YVIYSSLGSFFIPLAIMTIVYIEIFVATRRRLRERARANK
D2DR_RATCIIANPA.....FVVYSSIVSFYVPFIVTLLVYIKIYIVLRKRRKR.....
D3DR_RATCSISNPD.....FVIYSSVVSFYVPFGVTVLVYARIYIVLRQRQRK.....
DADR_RAT	ED.....DNCDTRLSRT.....YAISSSLISFYIPVAIMIVTYTSIYRIAQKQIRR.....
DBDR_RAT	EEGWELEGRTENCDSSLNRT.....YAISSSLISFYIPVAIMIVTYTRIYRIAQVQIRR.....
A1AB_RATCGVTEEPF.....CALFCSLGSFYIPLAVILVMYCRVYIVAKRTTKN.....
B1AR_RAT	PK.....CCDFVTNRA.....YAIASSVVSFYVPLCIMAFVYLRFREAQKQVKK.....
B2AR_HUMAN	ET.....CCDFFTNQA.....YAIASSIVSFYVPLVIMVFVYSRWFQEAKRQLQK.....
B2AR_RAT	ET.....CCDFFTNQA.....YAIASSIVSFYVPLVVMVFVYSRWFQVAKRQLQK.....
B3AR_RAT	PR.....CCSFASNMP.....YALLSSSVSFYLPLLVMFLFVYARVFVVAKRQRRF.....
5HTA_RATCTISKDHG.....YTIYSTFGAFYIPLLLMLVLYGRIFRAARFRIRK.....
5HTD_RATCLVNTSQIS.....YTIYSTCGAFYIPSILLIILYGRIVVAARSRLN.....
5HT2_RATCLLADDN.....FVLIGSFVAFFIPLTIMVITYFLTIKSLQKEATL.....
...

Βάσεις δεδομένων

Η σύγκριση αυτή δείχνει ότι καθώς η ομοιότητα των αλληλουχιών μειώνεται, η βιολογική σημασία των στοιχίσεων που προκύπτουν από τις αυτόματες μεθόδους φαίνεται επίσης να μειώνεται (συνήθως λόγω υπερβολών στη χρήση των κενών).

Έτσι, καθώς η ομοιότητα των αλληλουχιών μειώνεται, αυξάνει η ανάγκη ανθρώπινης παρέμβασης (με χρήση κάποιου multiple sequence alignment editor).

Για αλληλουχίες σχετικά υψηλής ομοιότητας (π.χ. $Z>6$), αλγόριθμοι και άνθρωποι δίνουν συγκρίσιμα αποτελέσματα (παρόμοιες στοιχίσεις).

Εφαρμογές στην έρευνα βάσεων δεδομένων

Μια εμφανής χρήση της ύπαρξης μίας στοίχισης πολλών αλληλουχιών είναι η δυνατότητα να ερευνηθούν οι βάσεις δεδομένων για συγγενείς αλληλουχίες χρησιμοποιώντας όχι μια αλληλουχία, αλλά μία ολόκληρη στοίχιση. Η χρήση της στοίχισης αναμένεται να βελτιώσει το λόγο σήματος προς θόρυβο για την έρευνα λόγω του υψηλότερου πληροφοριακού περιεχομένου της στοίχισης (π.χ. πληροφορία για το ποια αμινοξέα είναι αποδεκτά σε κάποια θέση της αλληλουχίας). Μια τέτοιου τύπου έρευνα είναι ανάλογη με την έρευνα των βάσεων με ένα μοτίβο.

PSI-BLAST

Ο αλγόριθμος που χρησιμοποιείται από το PSI-BLAST για την έρευνα των βάσεων δεδομένων είναι άξιος ξεχωριστής μνείας. To Position-Specific Iterated BLAST αντιπροσωπεύει ένα υβρίδιο ανάμεσα στις ανά ζεύγη μεθόδους στοίχισης και τις στοιχίσεις πολλών αλληλουχιών. Η κεντρική ιδέα είναι η εξής : μετά από μία αρχική έρευνα των βάσεων δεδομένων (με μία αλληλουχία-στόχο), όσες νέες αλληλουχίες δείχνουν αρκετή ομοιότητα προς την αλληλουχία-στόχο χρησιμοποιούνται για τη δημιουργία ενός μοτίβου (μέσω της στοίχισης με την αρχική αλληλουχία). Το μοτίβο αυτό χρησιμοποιείται εκ νέου για την επανεξ-ταση των βάσεων δεδομένων μέχρις συγκλίσεως.

PSI-BLAST

Το αποτέλεσμα (από τη μεριά του τελικού χρήστη) είναι ότι μπορεί να πραγματοποιηθεί μια έρευνα των βάσεων δεδομένων με την ευαισθησία που περιμένουμε από τη χρήση μοτίβων, αλλά χωρίς τη χρονοβόρα και απαιτητική διαδικασία της δημιουργίας τους. Το πρόβλημα είναι ότι το όριο για την εισαγωγή μίας νέας αλληλουχίας στο 'μοτίβο' θα πρέπει να είναι αρκετά αυστηρό ώστε να εξασφαλίζει ότι δεν θα υπεισέλθουν ψευδή θετικά στην διαδικασία αναζήτησης. Αυτό είναι ιδιαίτερα σημαντικό για την περίπτωση που η αλληλουχία-στόχος περιέχει LCR περιοχές.