

GHEORGHE-IOAN MIHALAŞ
ANCA TUDOR
SORIN PARALESCU

BIOINFORMATICA



Colecția

ȘTIINȚELE EXACTE ÎN CERCETAREA MEDICALĂ

GHEORGHE-IOAN MIHALAȘ
ANCA TUDOR
SORIN PARALESCU

BIOINFORMATICA



2011

Editura VICTOR BABEȘ

Piața Eftimie Murgu 2, cam. 316, 300041 Timișoara

Tel./ Fax 0256 495 210

e-mail: *evb@umft.ro, calaevb@umft.ro*

www.evb.umft.ro

Director general: Prof. univ. dr. Ștefan-Iosif Drăgulescu

Consilier editorial: Claudia Christian

Coperta: Caius Fericean

Colecția: ȘTIINȚELE EXACTE ÎN CERCETAREA MEDICALĂ

Coordonator colecție: Prof. univ. dr. Gheorghe-Ioan Mihalaș

Referent științific: Prof. univ. dr. Diana Lungeanu

© 2011 Toate drepturile asupra acestei ediții sunt rezervate.

Reproducerea parțială sau integrală a textului sau imaginilor, pe orice suport, fără acordul scris al autorilor, este interzisă și se va sancționa conform legilor în vigoare.

Descrierea CIP a Bibliotecii Naționale a României

MIHALAȘ, GHEORGHE-IOAN

Bioinformatică / Gheorghe-Ioan Mihalaș, Anca Tudor,

Sorin Paralescu. - Timișoara: Editura Victor Babeș, 2011

ISBN 978-606-8054-33-9

I. Tudor, Anca

II. Paralescu, Sorin

57:004

Tipărit la Tipografia Eurostampa

B-dul Revoluția din 1989 nr. 26, Timișoara

Tel. 0256- 204 816, edituraeurostampa@gmail.com

Prefață

Bioinformatica a apărut natural, ca disciplină înrudită cu informatica medicală (ulterior integrată în ea), dedicată stocării și prelucrării datelor din biologie. În ultimele decenii, aceste date au început să conțină tot mai multe detalii structurale, în special secvențe moleculare. Informațiile noi, care soseau cu o cadență accelerată impuneau cu stringență nu numai noi abordări în organizarea și sistematizarea datelor și algoritmi performanți de prelucrare, ci și necesitatea pregătirii unor specialiști în bioinformatică. Aplicațiile din ce în ce mai numeroase în domeniul medical și farmaceutic au ridicat și mai mult presiunea, astfel încât în ultimul deceniu al secolului trecut au apărut primele calificări, la nivel masteral, în bioinformatică. Trebuie însă remarcat caracterul interdisciplinar deosebit al bioinformaticii, acoperind o paletă mult mai largă decât alte discipline de graniță. Tendința de îngustare a domeniului de specializare, manifestată pregnant în învățământul superior în ultima parte a secolului XX a redus sensibil aria de recrutare a specialiștilor în toate ramurile multidisciplinare, iar pentru bioinformatică această reducere a devenit drastică. Comunitatea specialiștilor bioinformaticieni, deși într-o rapidă creștere, este încă destul de redusă – probabil așa se explică de ce această specialitate este cea mai bine remunerată dintre toate specialitățile informaticii medicale.

Ce este deosebit la această interdisciplinaritate? Așa cum am afirmat mai sus - paleta largă de intersecție a unor discipline, uneori destul de îndepărtate între ele. Într-adevăr, centrul atenției se îndreaptă spre principalele molecule de interes informațional: acizi nucleici și proteine. Deci este necesară o bază solidă de biologie moleculară, care este ea însăși un domeniu interdisciplinar, necesitând o serie de cunoștințe de biochimie și de biologie celulară. Pentru aplicațiile medicale și interpretarea hărților cromozomiale sunt necesare o serie de cunoștințe de genetică. Să ne îndreptăm privirea și în cealaltă parte, privind dezvoltarea metodelor de prelucrare. Sunt necesare cunoștințele solide de teoria probabilităților și statistică, topologie sau analiza seriilor. Determinarea structurilor tridimensionale și a aplicațiilor farmaceutice de “drug design” sunt aplicații ale chimiei cuantice. O serie de noțiuni de fizică și biofizică sunt de asemenea necesare, atât pentru înțelegerea proprietăților structurilor analizate cât și a metodelor utilizate pentru obținerea și interpretarea datelor experimentale. Să nu uităm nici capitolul de biologia sistemelor, în care sunt urmărite procesele de semnalizare celulară, studiate de fiziologia celulară și prelucrate prin simulare atât prin ecuații diferențiale cât și simulare stochastică. Pentru viziunea integrativă sunt necesare noțiuni din teoria sistemelor, cibernetică și teoria evoluționismului. Desigur, să adăugăm în final și câteva discipline fundamentale din știința calculatoarelor – baze de date și limbaje de programare.

Nu este deci de mirare că, la înscrierea la cursurile de bioinformatică, în universitățile ce au introdus aceste cursuri, se presupune că studenții au parcurs anterior cursuri de biofizică, biochimie, biologie moleculară și celulară sau genetică și au baza matematică adecvată.

Având în vedere complexitatea pregătirii în atâtea domenii adiacente, prezentul manual a fost alcătuit astfel încât și absolvenții unui profil tehnic să-și găsească aici și noțiunile fundamentale din aceste discipline. În această viziune am inclus, după primul capitol privitor la noțiuni generale de informatică medicală, un set de patru capitole dedicate prezentării sumare a noțiunilor utile pentru abordările din partea a doua.

În ultima parte sunt prezentate lucrările practice de laborator, esențiale pentru a simți adevărata dimensiune a imenselor baze de date moleculare și complexitatea metodelor de prelucrare.

Deși această carte – prima de acest fel din țară – a fost elaborată mai ales pentru studenții la nivel masteral în profilul “Sisteme informatice în îngrijirea sănătății” (SIIS) de la Universitatea Politehnica, respectiv profilul “Biologie moleculară și biotehnologii” (BMB), de la Universitatea de Medicină și Farmacie, ea se adresează de fapt tuturor celor interesați de acest domeniu fascinant, care va oferi încă numeroase surprize. Lumea vie are încă multe secrete nedescifrate!

Închei prefața citând un student (într-un sondaj despre curs): “Este mai interesant ca la ‘Discovery!’”

*Prof. dr. G. I. Mihalaș
membru al Academiei de Științe Medicale*

CUPRINS

Partea I (G.I.Mihalaş)

1. INTRODUCERE ÎN INFORMATICA MEDICALĂ	1
1.1. Obiectul Informaticii Medicale	1
1.1.1. Scurt istoric	1
1.1.2. Obiectul informaticii medicale	1
1.2. Teoria informaţiei	1
1.2.1. Noţiunea de informaţie	1
1.2.2. Proprietăţile informaţiei	1
1.2.3. Triada abordărilor complete	2
1.2.4. Cantitatea de informaţie	2
1.2.5. Relaţia între entropia informaţională şi entropia termodinamică	3
1.2.6. Redondanţă	3
1.3. Transmiterea informaţiei	4
1.3.1. Sisteme de comunicaţie	4
1.3.2. Transmiterea informaţiei în materia vie	5
1.4. Informatica medicală	6
1.4.1. Date şi cunoştinţe	6
1.4.2. Ciclul elementar al informaţiei medicale	7
1.4.3. Tipuri de date	7
1.4.4. Tipuri de cunoştinţe	7
1.4.5. Operaţii cu informaţii	7
1.5. Capitolele informaticii medicale şi structura cursului	8
1.5.1. Clasificarea informaţiei medicale pe nivele structurale	8
1.5.2. Bioinformatica – relaţii cu alte discipline	8
2. INTRODUCERE ÎN BIOFIZICĂ	9
2.1. Obiectul şi capitolele biofizicii	9
2.1.1. Ştiinţe interdisciplinare	9
2.1.2. Capitolele biofizicii	9
2.2. Structura atomului	10
2.2.1. Proprietăţi	10
2.2.2. Nucleul atomic	10
2.2.3. Norul electronic	11
2.2.4. Structura norului electronic – principii	12
2.2.5. Structura norului electronic al elementelor	12
2.3. Atomul de carbon	13
2.3.1. Structura norului electronic	13
2.3.2. Hibridizarea	14
2.4. Structura moleculei. Legături chimice	15
2.4.1. Stabilitatea atomilor	15
2.4.2. Legătura covalentă	16
2.4.3. Energia de legătură	17
2.4.4. Legătura ionică	17
2.5. Forţe intermoleculare	18
2.5.1. Legătura de hidrogen	19
2.5.2. Forţe Van der Waals	19

2.5.3. Forțe de dispersie.....	19
2.6. Molecula de apă.....	19
2.6.1. Structura moleculară.....	19
2.6.2. Legăturile de hidrogen ale moleculei de apă.....	20
2.6.3. Proprietățile apei.....	21
2.7. Soluții.....	21
2.7.1. Sisteme disperse.....	21
2.7.2. Concentrații.....	22
2.8. pH-ul soluțiilor.....	23
2.8.1. Disocierea electroliților.....	23
2.8.2. Produsul ionic al apei.....	23
2.8.3. Scara pH.....	24
2.9. Termodinamica biologică.....	25
2.9.1. Parametrii de stare ai unui sistem termodinamic.....	25
2.9.2. Procese termodinamice.....	26
2.9.3. Funcții de stare.....	26
2.9.4. Principiul al doilea al termodinamicii.....	27
2.10. Procese cuplate.....	27
2.10.1. Natura proceselor cuplate.....	27
2.10.2. Laturile metabolismului.....	27
2.10.3. Stocarea energiei pentru procesele biologice.....	28
2.11. Forțe termodinamice.....	28
2.12. Transport transmembrantar.....	28
2.12.1. Clasificare, proprietăți.....	28
2.12.2. Transport pasiv.....	28
2.12.3. Transport activ.....	29
3. NOȚIUNI DE BIOCHIMIE.....	31
3.1. Aminoacizii.....	31
3.1.1. Structură generală.....	31
3.1.2. Formele ionice ale aminoacizilor.....	32
3.1.3. Proprietățile aminoacizilor.....	33
3.1.4. Hidrofobicitatea.....	34
3.2. Proteine.....	36
3.2.1. Legătura peptidică.....	36
3.2.2. Structura primară a proteinelor.....	37
3.2.3. Structura secundară a proteinelor.....	38
3.2.4. Structura terțiară a proteinelor.....	40
3.2.5. Structura cuaternară.....	40
3.3. Componentele acizilor nucleici.....	41
3.3.1. Componentele unui nucleotid.....	41
3.3.2. Pentozele din acizii nucleici.....	41
3.3.3. Bazele azotate.....	41
3.3.4. Nucleozide.....	42
3.3.5. Nucleotide.....	43
3.4. Structura acizilor nucleici.....	43
3.4.1. Formarea lanțului polinucleotidic.....	43
3.4.2. Structura primară a acizilor nucleici.....	44
3.4.3. Structura secundară a acizilor nucleici.....	45

4. ELEMENTE DE BIOLOGIE CELULARĂ ȘI MOLECULARĂ	47
4.1. Structura generală a unei celule umane	47
4.1.1. Enumerarea principalelor componente ale celulei umane	47
4.1.2. Membrana celulară	48
4.1.3. Citoplasma	49
4.1.4. Nucleul celular	50
4.1.5. Organitele celulare	50
4.1.6. Mitocondria	51
4.1.7. Ribozomii	52
4.2. Diviziunea celulară	53
4.2.1. Mitoza	53
4.2.2. Meioza	53
4.3. Replicarea ADN	54
4.3.1. Noțiunea de replicare	54
4.3.2. Fazele replicării ADN	55
4.4. Sinteza proteinelor	56
4.4.1. Paradigma centrală a bioinformaticii	56
4.4.2. Transcripția	56
4.4.3. Codul genetic	57
4.4.4. ARN de transport	58
4.4.5. Structura ribozomilor	59
4.4.6. Mecanismul translației	59
4.5. Controlul sintezei proteinelor	60
5. ELEMENTE DE GENETICĂ.....	61
5.1. Scurt istoric. Capitolele geneticii	61
5.1.1. Epoca mendeliană	61
5.1.2. Teoria cromozomială a eredității	61
5.1.3. Genetica moleculară	61
5.1.4. Bioinformatica	61
5.2. Genetica mendeliană	62
5.2.1. Terminologie	62
5.2.2. Legea purității gameților	62
5.2.3. Dihibridismul și legea segregării independente a caracterelor	64
5.2.4. Alte tipuri de segregare	65
5.2.5. Principiul Hardy-Weinberg	66
5.3. Teoria cromozomială	66
5.3.1. Localizarea genelor	66
5.3.2. Morfologia cromozomilor	66
5.3.3. Transmiterea înlănțuită a genelor (linkage)	67
5.3.4. Schimbul reciproc de gene (crossing-over)	68
5.3.5. Hărți cromozomiale	70
5.3.6. Markerii pentru hărți genetice	71
5.4. Genetica moleculară	71
5.4.1. Corespondențe moleculare	71
5.4.2. Structura unui cromozom	71
5.4.3. Hărți fizice	71
5.4.4. Exemplu de hartă cromozomială	73

6. ANALIZA SECVENȚIALĂ INDIVIDUALĂ.....	75
6.1. Introducere	75
6.1.1. Obiectul bioinformaticii	75
6.1.2. Capitolele bioinformaticii.....	75
6.1.3. Obiectul analizei secvențiale	75
6.2. Analiza secvențială grafică	76
6.2.1. Baze teoretice	76
6.2.2. Aspecte structurale	77
6.3. Structura tridimensională	77
6.3.1. Roți elicoidale și elice amfipatice.....	77
6.3.2. Alte structuri.....	78
7. COMPARAREA A DOUĂ SECVENȚE	79
7.1. Introducere	79
7.1.1. Fundamentele analizei secvențiale	79
7.1.2. Evenimente	79
7.1.3. Termeni	79
7.2. Graficele de puncte „dot plots”	80
7.2.1. Principiul metodei	80
7.2.2. Filtre	80
7.3. „Distanțe” între secvențe.....	81
7.3.1. Noțiunea de aliniere.....	81
7.3.2. Distanțe.....	82
7.3.3. „Problemele” analizei secvențiale	82
7.4. Programare dinamică – Algoritmul Needleman – Wunsch.....	84
7.4.1. Alinierea globală	84
7.4.2. Principiul algoritmului.....	84
7.4.3. Marcarea traseului	86
7.4.4. Exemplu	86
7.5. Alinierea locală. Algoritmul Smith – Waterman.....	89
7.5.1. Deosebiri față de algoritmul NW.....	89
7.5.2. Descrierea algoritmului SW	90
7.6. Modele complexe.....	92
7.6.1. Potriviri repetate (Repeated matches).....	92
7.6.2. Potriviri suprapuse.....	93
7.6.3. Potriviri hibride (Hybrid match conditions)	94
7.7. Gap-uri afine	95
7.7.1. Baze teoretice	95
7.7.2. Tipuri de gap-uri.....	95
8. MATRICI DE SUBSTITUȚIE.....	97
8.1. Introducere	97
8.2. Matrici de substituție pentru proteine.....	97
8.2.1. Matrici PAM.....	97
8.2.2. Matrici BLOSUM.....	100
8.3. Matrici de substituție pentru acizi nucleici.....	101
8.3.1. Matricea Jukes – Cantor	101
8.3.2. Matricea Kimura.....	102
8.4. Testarea semnificației alinierii	103
8.4.1. Scoruri	103
8.4.2. Semnificația scorurilor	104

9. ALINIEREA MULTIPLĂ.....	105
9.1. Semnificația alinierii multiple.....	105
9.1.1. Factorii luați în considerare pentru MSA	105
9.1.2. Obiective în alinierea multiplă	105
9.2. Scop și motivație pentru alinierea multiplă.....	106
9.2.1. Punerea problemei.....	106
9.2.2. Motivație	106
9.3. Scoruri pentru alinierea multiplă.....	106
9.3.1. Privire generală	106
9.3.2. Suma perechilor.....	107
9.3.3. Entropia minimă	107
9.3.4. Reprezentări grafice	107
9.4. Algoritmi pentru alinierea multiplă.....	108
9.4.1. Programare dinamică.....	108
9.4.2. Metode de aliniere progresivă	109
9.4.3. Algoritmul MSA	109
9.5. Modele de ordonare	111
9.5.1. Modelul “stea”.....	111
9.5.2. Modelul “arbore”.....	112
10. LANȚURI MARKOV	113
10.1. Lanțuri Markov simple	113
10.1.1. Descrierea unui lanț Markov	113
10.1.2. Probabilitatea unei secvențe	114
10.2. Estimarea parametrilor modelului	115
10.2.1. Punerea problemei.....	115
10.2.2. Estimarea Laplace	116
10.2.3. Estimarea Laplace generalizată	117
10.2.4. Estimarea parametrilor de ordin I.....	117
10.2.5. Lanțuri Markov de ordin superior	118
10.3. “Open Reading Frames” (ORF).....	119
10.3.1. Noțiunea de “cadru de citire”	119
10.3.2. Metode pentru găsirea genelor	120
11. MODELE MARKOV ASCUNSE	121
11.1. Suportul biologic al abordării	121
11.1.1. Insulele CpG.....	121
11.1.2. Distincția “stare” - “simbol”.....	121
11.2. Enunțul problemei în HMM	122
11.2.1. Denumirea HMM (Hidden Markov Model).....	122
11.2.2. Reprezentarea schematică a unui HMM.....	122
11.3. Algoritmi de calcul pentru HMM	123
11.3.1. Algoritmul Viterbi.....	123
11.3.2. Exemplu	124
11.4. Aplicarea HMM pentru discriminare	125
11.4.1. Notății.....	125
11.4.2. Scoruri “log-odd”	126
11.5. Modele Markov ascunse cu inserții și deleții	127
11.5.1. Stări silențioase	127
11.5.2. Schema HMM generală.....	127

12. ANALIZA FILOGENETICĂ.....	129
12.1. Inferență filogenetică	129
12.1.1. Obiectivul analizei filogenetice	129
12.2. Noțiuni generale despre arbori	129
12.2.1. Terminologie	129
12.2.2. Utilitate, motivație	129
12.3. Proprietățile arborilor	131
12.3.1. Definiție, structură	131
12.3.2. Tipuri de arbori	131
12.3.3. Numărul de arbori	132
12.4. Construcția arborilor filogenetici - generalități	133
12.4.1. Date pentru construcția arborilor	133
12.4.2. Metode de construcție a arborilor filogenetici	133
12.4.3. Comparăție între metodele bazate pe distanțe și cele bazate pe parsimonie	133
12.5. Metode bazate pe distanțe	135
12.5.1. Proprietățile distanțelor, formularea problemei	135
12.5.2. Algoritmul UPGMA	135
12.5.3. Distanțe ultrametrice	136
12.5.4. Metoda Neighbor Joining	137
12.6. Metode bazate pe “parsimonie”	138
12.6.1. Enunțul problemei	138
12.6.2. Algoritmul lui Fitch	138
12.6.3. Parsimonie ponderată	140
12.7. Testarea arborilor - metoda „bootstrap”	143

Partea a II-a (A. Tudor, S. Paralescu)

1. BAZE DE DATE – BD integrată NCBI (I). Utilizarea motorului de căutare ENTREZ și modulul PUBMED	147
1.1. Obiectivele lucrării de laborator	147
1.2. Baze de date și instrumente de căutare folosite în Bioinformatică	147
1.3. Utilizarea Entrez	148
1.3.1. Calificatorii câmpului Entrez	150
1.3.2. Exemple	151
1.4. Utilizarea PubMed	152
1.4.1. Exemple	153
1.5. Exerciții propuse	154
2. BAZE DE DATE – BD INTEGRATĂ NCBI (II). PROGRAMUL BLAST	155
2.1. Obiectivele lucrării de laborator	155
2.2. Introducere	155
2.3. Programul BLAST	156
2.4. Exemple	160
2.5. Utilizarea Bazelor de date cu secvențe proteice	163
2.6. Exerciții propuse	166

3. STRUCTURA SECUNDARĂ, TERȚIARĂ ȘI CUATERNARĂ A PROTEINELOR. PROGRAMELE PDB, Cn3D ȘI RASMOL	167
3.1. Obiectivele lucrării de laborator	167
3.2. Noțiuni introductive	167
3.2.1. Proprietăți fizico-chimice	167
3.3. Determinarea structurii proteinelor prin tehnici experimentale și metode de modelare comparativă	172
3.4. Grafica moleculară	176
3.5. Exemplu de explorare a structurilor proteice cu PDB (Protein Data Bank)	178
3.6. Utilizarea Cn3D	179
3.7. Exemplu de utilizare a programului RasMol	181
3.8. Exerciții propuse	183
4. APLICAȚIA VECTOR NTI – PROGRAMUL ALIGN X	185
4.1. Obiectivele lucrării de laborator	185
4.2. Exemplu 1	186
4.3. Exemplu 2	190
4.4. Exemplu 3	191
4.5. Exerciții propuse	193
5. CLUSTAL X	195
5.1. Obiectivele lucrării de laborator	195
5.2. Introducere	195
5.3. Utilizarea Clustal X	195
5.4. Crearea fișierului de intrare pentru alinierea multiplă	196
5.4.1. Exemplu	197
5.4.2. Alinierea multiplă	198
5.4.3. Introducerea datelor în programul ClustalX	199
5.5. Setarea parametrilor aliniamentului	199
5.5.1. Parametrii alinierii pereche	200
5.5.2. Parametrii alinierii multiple	202
5.5.3. Formatul de ieșire al alinierii	202
5.6. Crearea alinierii	203
5.7. Scrierea aliniamentului ca fișier Postscript	204
5.8. Arbori filogenetici folosind ClustalX	205
5.9. Exerciții propuse	205
Bibliografie	209

Partea I

1. Introducere în Informatica Medicală

1.1. *Obiectul Informaticii Medicale*

1.1.1. *Scurt istoric*

A. Informatica medicală este o disciplină tânără, termenul apărând în cursul anilor '60. În accepțiunea inițială informatica medicală cuprindea *programele de calculator* cu aplicabilitate în domeniul medical. Progresul tehnic rapid a arătat însă că, pentru aceleași aplicații, atât programele cât și suportul fizic se schimbau; ceea ce rămânea la fel era modul în care era prelucrată *informația*.

B. Astfel, în accepțiunea actuală, centrul definiției s-a mutat de la „calculator” la informație. Coiera [1997] chiar atrage atenția în acest sens: „Informatica medicală se ocupă de calculatoare tot atât de mult cât se ocupă cardiologia de stetoscoape”.

1.1.2. *Obiectul informaticii medicale*

Accepțiunea clasică: totalitatea programelor de *calculator* cu aplicații în domeniul biomedical și sănătate.

Definiția actuală: disciplina care studiază întregul flux al *informației medicale:* generare, achiziție, stocare, transmitere, prelucrare și utilizare.

1.2. *Teoria informației*

1.2.1. *Noțiunea de informație*

Pentru a ne ocupa de *informația medicală*, să încercăm mai întâi să privim conceptul de *informație* la modul general.

A. Termenul de informație este folosit în mod curent în viața de zi cu zi, fiind cel mai adesea asociat cu aducerea unui element de noutate. Fiind un concept cu grad înalt de generalitate (categorie filosofică), informația nu poate fi definită în manieră clasică, pornind de la genul proxim și precizând diferențele specifice, ci prin proprietatea sa esențială – cea de a înlătura o nedeterminare.

B. Noțiunea de informație

Informația este un concept cu grad înalt de generalitate caracterizat prin proprietatea de a înlătura o nedeterminare (incertitudine).

1.2.2. *Proprietățile informației*

Informația nu este materie; totuși ea nu poate exista înafara materiei. Norbert Wiener spunea „Creierul nu secretă informație precum ficatul fiere”.

Informația nu este energie; totuși ea nu se poate transmite fără un suport energetic.

Nu este o relație directă între cantitatea de energie ce însoțește transmiterea unei informații și cantitatea de informație transmisă. De exemplu, energia unui trăsnet în

timpul unei furtuni este imensă, însă informația transmisă este neglijabilă; în schimb un foșnet într-o pădure, purtat de o energie infimă poate reprezenta o informație vitală – punând pe fugă un animal! Să semnalăm totuși că nu există relație nici între cantitatea de informație și efectele sale. De ex. Legenda lui Tezeu din metodologia greacă: Tezeu promisese tatălui său Aegeus, că dacă va învinge în luptă minotaurul va înlocui pânza neagră a corabiei cu pânza albă, dar a uitat și tatăl său s-a aruncat de pe stânci. Informația primită a fost doar 1 bit.

1.2.3. *Triada abordărilor complete*

A. Introducerea aspectelor informaționale în studiul materiei vii completează imaginea noastră privind complexitatea sistemelor biologice, actualmente considerându-se că o abordare completă trebuie să acopere atât aspectele materiale și energetice cât și cele informaționale.

B. Triada abordărilor complete:

- aspectul material – structura
- aspectul energetic – suportul funcțional
- aspectul informațional – mecanismul funcțional.

C. Valoarea utilă a informației depinde de receptor.

Sensul noțiunii de informație, așa cum a fost prezentat mai sus etc. legat de altă noțiune – nedeterminarea (sau incertitudinea) – vag definită la rândul său. Același mesaj poate să aibă valori informaționale diferite pentru diferiți receptori: pentru o persoană care deja știa conținutul său cantitatea de informație primită este zero, însă pentru receptorii care nu-i știau conținutul va putea fi evaluată cantitatea de informație primită.

1.2.4. *Cantitatea de informație*

Parcurerea acestui subiect necesită cunoștințe fundamentale de teoria probabilităților

A. Pornind de la proprietatea fundamentală a informației, cea de a înlătura o nedeterminare, Shannon a considerat că informația primită este invers proporțională cu probabilitatea de apariție a evenimentului: dacă se va întâmpla un eveniment cu probabilitate mare, informația primită este mică; în schimb primim o informație „mai mare” dacă apare un eveniment mai rar. (Ziariștii exploatează intens această relație!)

B. Relația propusă de Shannon pentru calculul cantității de informație care este primită când se petrece un eveniment cu probabilitatea p_i cuprinde logaritmul în baza 2 din inversul probabilității p_i .

- Cantitatea de informație eliberată de un eveniment a cărui probabilitate este p_i

$$I_i = -\log_2 p_i \quad (1.2.4.a)$$

C. Pe baza acestei relații stabilește unitatea de măsură pentru cantitatea de informație, numită *bit* (de la Binary digiT).

Definiție: un bit este cantitatea de informație primită când se înlătură o nedeterminare de $\frac{1}{2}$.

D. În mod uzual informația se transmite printr-o succesiune de evenimente, numită adesea *mesaj*, iar un eveniment într-un mesaj se mai numește *simbol*.

În cazul unui mesaj format din N evenimente, fiecare eveniment i apare de n_i ori, aducând de fiecare dată informația I_i , deci mesajul aduce informația.

$$I = n_1 I_1 + n_2 I_2 + \dots + n_k I_k = \sum_{i=1}^k n_i I_i \quad (1.2.4.b)$$

E. Valoarea medie a informației corespunzătoare unui eveniment într-un șir de N evenimente se mai numește „entropie informațională”, H, și se calculează astfel:

$$H = \frac{n_1 I_1 + n_2 I_2 + \dots + n_k I_k}{N} = \frac{n_1}{N} I_1 + \frac{n_2}{N} I_2 + \dots + \frac{n_k}{N} I_k$$

la limită (i.e. atunci când $N \rightarrow \infty$) relația devine:

$$H = p_1 I_1 + \dots + p_k I_k = \sum_{i=1}^k p_i I_i$$

Înlocuind I_i conform relației (1.2.4.a), obținem formula (1.2.4.c), formulă fundamentală în teoria informației, numită și formula lui Shannon pentru entropia informațională.

- Entropia informațională este cantitatea medie de informație per eveniment (simbol) într-un mesaj:

$$H = - \sum_{i=1}^k p_i \log_2 p_i \quad (1.2.4.c)$$

1.2.5. Relația între entropia informațională și entropia termodinamică

A. Termenul de „entropie” a fost introdus în termodinamică pentru enunțarea principiului al II-lea al termodinamicii: „În procesele termodinamice entropia nu poate să scadă: ea rămâne constantă în cadrul proceselor reversibile și crește în cazul proceselor ireversibile”.

B. Relația între entropia termodinamică și cea informațională poate fi înțeleasă pornind de la experimentul „ideal” propus de Maxwell pentru explicarea variației entropiei în cazul proceselor ireversibile, prezentat în cadrul cursului de biofizică.

C. Se vede deci că sistemul poate evolua în sens contrar celui dictat de al II-lea principiu al termodinamicii în cazul în care primește o informație. Acesta este mecanismul prin care sistemele vii evoluează spre stări tot mai organizate și deosebite de mediul înconjurător.

1.2.6. Redondanță

A. Entropia informațională are valoare maximă când evenimentele din mesaj sunt echiprobabile: $p_i = 1/k$;

atunci $H_{\max} = k \left(\frac{1}{k} \log \frac{1}{k} \right)$ de unde se obține relația (1.2.6.a)

$$H_{\max} = \log_2 k \quad (1.2.6.a)$$

B. Un exemplu ar fi cazul unui mesaj criptat, în care probabilitatea apariției unui simbol este (cel puțin aparent) independentă de simbolurile anterioare. În mesajele

reale probabilitatea unui simbol depinde de simbolurile anterioare; putem, în funcție de context, să „ghicim” ce urmează, putem folosi prescurtări, putem observa greșeli cum ar fi omisiunea unei litere etc. Deci informația nu este distribuită uniform în mesaj sau chiar în interiorul cuvintelor, cantitatea de informație transportată în realitate fiind inferioară celei maxime ce ar putea fi transmise folosind aceeași lungime a textului. Această diferență, între cantitatea maximă ce poate fi conținută în mesaj și cea reală se numește redondanță și reprezintă o parte din mesaj care... nu conține informație!

C. Relația de definiție a redondanței absolute este:

$$R = H_{max} - H_{real} \quad (1.2.6.b)$$

Raportând redondanța absolută la H_{max} se definește redondanța relativă.

$$R_r = R / H_{max} \quad (1.2.6.c)$$

D. *Rolul redondanței*

Aparent redondanța ar reprezenta o încărcătură inutilă în mesaj. Totuși, prezența ei diminuează rolul negativ al perturbațiilor ce apar în cursul transmiterii informației, putând deseori reconstitui mesajul inițial chiar dacă unele simboluri au fost perturbate.

1.3. Transmiterea informației

1.3.1. Sisteme de comunicație

A. Am precizat anterior că valoarea utilă a informației depinde de receptor, deci noțiunea de informație are sens doar dacă se transmite; altfel, rămâne în faza de „informație potențială”.

B. Transmiterea informației presupune o *sursă* a informației (emițător E) și un destinatar (receptor R). Spațiul dintre S și R reprezintă *canalul de comunicație* (C). Pe canalul de comunicație pot să apară diverse *zgomote* care perturbă sistemul de comunicație afectând calitatea transmisiei (figura 1.3.1.a).

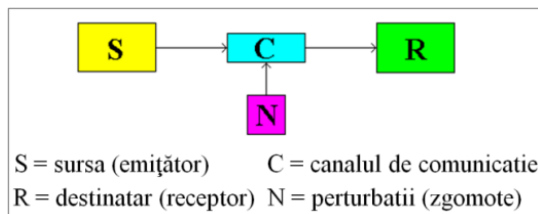


Fig. 1.3.1.a. Sistem de comunicație simplu

C. Să introducem doi termeni importanți în cadrul sistemelor de comunicație:

- *mesaj* – un termen pe care îl folosim când ne referim la conținutul informațional al transmisiei
- *semnal* – suportul fizic care transportă mesajul (sunet, curent electric etc.).

D. Pentru diminuarea efectelor perturbațiilor sau pentru a asigura transmiterea mesajului la distanțe foarte mari, se introduc pe canalul de transmisie niște dispozitive numite *traductori*. Un traductor schimbă suportul fizic al unui semnal (figura 1.3.1.b). De exemplu, în cazul unei convorbiri telefonice, microfonul este traductorul localizat lângă emițător, transformând sunetele (variații ale presiunii aerului) în variații ale unui curent electric. Canalul de comunicație este reprezentat de firele telefonice. La

destinatar un alt traductor, casca telefonică, transformă variațiile curentului electric în vibrații ale unei membrane elastice generând astfel sunete.

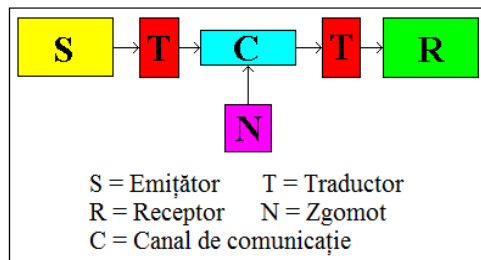


Fig. 1.3.1.b. Schema generală a unui sistem de comunicație

E. Există și alte dispozitive ce pot fi utilizate în sisteme de comunicație, de ex. *modem-ul*. Denumirea modem provine de la modulator/demolator. Modemul este un dispozitiv care asigură modularea semnalului, adică suprapunerea semnalului real peste un semnal purtător (undă purtătoare) care are caracteristici încât se diminuează efectul perturbațiilor (de ex. perturbațiile uzuale, de joasă frecvență, sunt eliminate dacă unda purtătoare are frecvență înaltă).

F. O altă transformare pe care o putem aplica semnalului pentru transmisie este *codificarea*. Mesajul este compus uzual dintr-o succesiune de simboluri. Totalitatea simbolurilor utilizate pentru a compune un mesaj se numește *alfabet*. Simbolurile alfabetului se mai numesc „*litere*”, iar cu literele putem construi *cuvinte*. Totalitatea cuvintelor cu sens reprezintă un *dicționar*, iar precizarea sensului cuvintelor se numește *semantică*. Cu ajutorul cuvintelor se pot construi propoziții; regulile de construcție a propozițiilor se numește *sintaxă*. Un dicționar împreună cu semantica și o sinteză reprezintă un *limbaj*. Noi folosim uzual pentru comunicație *limbaje naturale*, dar există posibilitatea utilizării unor *limbaje formale* sau *artificiale*. Diferitele componente ale sistemului de comunicație pot folosi diferite alfabet sau dicționare. Transpunerea unui mesaj dintr-o formă ce utilizează un alfabet într-o formă în alt alfabet, cu anumite reguli de corespondență se numește *codificare*. Operațiunea inversă se numește *decodificare*. Transpunerea unui mesaj dintr-un limbaj în altul se numește *traducere*.

G. Să mai menționăm legat de sistemele de comunicație că există o capacitate limitată de transmisie a informației pe canalul de comunicație, numită *viteză de transmisie*, măsurată în bit/secundă.

1.3.2. Transmiterea informației în materia vie

A. *Codul genetic*. Informația privind structura proteinelor ce pot fi sintetizate este stocată în molecula de ADN din nucleu. Acizii nucleici conțin 4 baze azotate: adenina A, timina T, citozina C și guanina G (în cazul ARN în loc de timină apare uracilul U). Proteinele sunt formate din 20 de aminoacizi esențiali. O succesiune de 3 baze azotate din ADN se numește *codon* și poartă informația pentru codificarea unui aminoacid într-o secvență proteică. Totalitatea corespondenților între codoni și aminoacizii corespunzători poartă denumirea de *cod genetic*. Porțiunea dintr-un lanț ADN care poartă informația pentru sinteza unei proteine se numește *genă*, iar ansamblul tuturor genelor unei specii se numește genom. Genomul uman conține circa 30.000 gene. Prezentul curs este dedicat acestui tip de informație.

B. *Codificarea informației în sistemul nervos.* Pe axoni informația este transmisă printr-o succesiune de impulsuri nervoase; fiecare impuls nervos este un potențial de acțiune care are întotdeauna aceeași amplitudine. Unui stimul mai intens îi corespunde o rată mai ridicată de formare a potențialelor de acțiune; spunem că pe axon informația privind intensitatea stimulului este *codificată în frecvență*.

La nivelul sinapselor are loc o descărcare a veziculelor cu mediator chimic în spațiul sinaptic, cantitatea de mediator descărcată fiind proporțională cu frecvența impulsurilor nervoase pe axon; spunem că în spațiul sinaptic informația privind intensitatea stimulului este *codificată în amplitudine*, aceasta fiind reprezentată de cantitatea de mediator descărcată.

La nivelul membranei postsinaptice, mediatorul se cuplează pe receptorii postsinaptici, se deschid canalele de sodium, membrana se depolarizează și apare un potențial care se propagă pe membrana corpului neuronal sau pe dendrite. Spunem că informația este *codificată în amplitudine*, aceasta fiind reprezentată de potențialul local.

1.4. Informatica medicală

După această incursiune în teoria informației putem reveni la noțiunea centrală din informatică medicală și anume *informația medicală*.

Ce este informația medicală și când apare ea?

1.4.1. Date și cunoștințe

Să încercăm să schițăm în cel mai simplificat mod actul medical primar și anume vizita pacientului la medic. Poziția centrală în activitate medicală este ocupată de *pacient*. Fără pacient nu există medicină! Actorul principal al activității medicale este *medicul*, dar în activitatea medicală sunt implicate numeroase alte persoane care aparțin așa-numitelor „*profesii aliate*”. Dialogul medic-pacient începe uzual cu expunerea de către pacient a motivelor pentru care s-a prezentat la medic. Această descriere reprezintă transmiterea unor informații de la pacient către medic. Informațiile care se transmit sau se utilizează într-un act medical (sau ca urmare a unui act medical) reprezintă *informația medicală*. Dialogul este succedat de către examenul obiectiv al pacientului, medicul colectând astfel și alte informații despre pacient. Să observăm că aceste informații au un caracter individual – sunt valabile pentru *acest pacient*. Aceste informații se numesc *date*. Uzual paleta datelor se completează cu informații provenind și din alte investigații (probe de laborator, explorări funcționale, radiografii etc.). Indiferent cât de complexe ar fi ca reprezentare, ele sunt „date”, fiind caracteristice unui anumit individ.

În același timp, medicul utilizează și alt fel de informații, numite *cunoștințe*. Acestea au un caracter general și sunt acumulate în cursul pregătirii profesionale precum și în experiența sa practică. Fără aceste cunoștințe informațiile sub formă de date nu pot fi interpretate (revenim la afirmația că valoarea utilă a informației depinde de receptor; practic, fără aceste cunoștințe receptorul datelor nu este „medic”). De aceea numeroși autori numesc *informație* doar *datele interpretate*. Pentru a evita confuzia între termenul *informație* folosit la modul general și *informație* pentru treapta de *date interpretate*, vom păstra termenul de *date interpretate* pentru acest nivel.

1.4.2. Ciclul elementar al informației medicale

Prin interpretarea datelor de către medic pe baza cunoștințelor sale, este generată o nouă informație numită diagnostic. Pe baza diagnosticului, folosind din nou cunoștințele sale, medicul stabilește un plan terapeutic pe care îl aplică pacientului cu scopul de a îmbunătăți starea pacientului. Urmărirea evoluției pacientului este însoțită de colectarea unor noi informații sub formă de date. Se observă că se încheie un ciclu al fluxului informațional în activitatea medicală, numit „ciclul elementar al informației medicale” (figura 1.4.2).

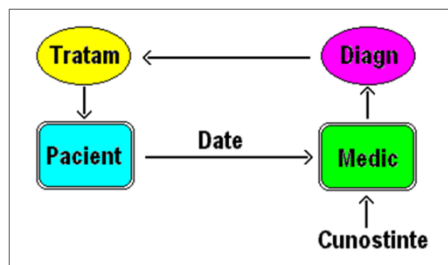


Fig. 1.4.2. Ciclul elementar al informației medicale

1.4.3. Tipuri de date

Informațiile culese despre starea pacientului, adică datele, pot îmbrăca diverse forme:

- *date calitative* – cu caracter descriptiv, așa cum apar în anamneză
- *date numerice* – forma uzuală de prezentare a rezultatelor de laborator
- *grafice* – modul de înregistrare a evoluției în timp a unor mărimi biologice (ex.: semnalul ECG, EEG etc.)
- *sunete* – de ex. fonocardiograma; modul de prelucrare este asemănător cu cel al altor semnale
- *imagini* – radiografia, tomografia, ecografia etc.
- *imagini dinamice* – filme.

Modul de achiziție, stocare și prelucrare este specific pentru fiecare tip de date și în cadrul cursului nostru le corespund capitole separate.

1.4.4. Tipuri de cunoștințe

Cunoștințele pot fi de mai multe feluri:

- *cunoștințe explicite* – care se pot formaliza, se pot exprima în propoziții, pot fi ușor transmise pe cale orală sau scrisă
- *abilități* sau *cunoștințe tacite* (limba engleza - *skill*) – cele câștigate prin experiență practică (de ex. îndemânarea unui chirurg sau a unui dentist); nu pot fi transmise ușor.

1.4.5. Operații cu informații

Urmărind ciclul de viață al informației, din momentul generării sale până în momentul utilizării, observăm că informația suferă o serie de operații:

- *achiziția (colectarea)* – presupune mijloace specifice tipului de informație
- *stocarea* – baze de date, respective baze de cunoștințe
- *transmitere* – căi, procedee

- *prelucrare* – cu o largă paletă de metode specifice, pentru a extrage elementele esențiale în vederea interpretării și utilizării
- *protecție* – măsurile ce se impun pentru asigurarea integrității informației stocate sau transmise, precum și a confidențialității acesteia
- *interpretare/utilizare* – pasul final, în care informația este integrată în acțiunile specifice nivelului.

1.5. Capitolele informaticii medicale și structura cursului

1.5.1. Clasificarea informației medicale pe nivele structurale

În ciclul elementar al informației medicale prezentat mai sus am luat în considerare informațiile care apar în activitatea medicală curentă, la nivelul individului, numit pacient. Totuși fenomenele care se petrec în materia vie (legate de starea de sănătate a pacientului) privesc deseori nivele infraindividuale, pornind de la nivelul molecular sau celular, urcând prin nivelul de țesut, organ sau sistem până la nivelul întregului organism sau nivelul individual.

Pe de altă parte, activitatea medicală este organizată în unități care prestează servicii pentru populație, deci putem urmări fluxul informațional și la nivel supraindividual, de comunitate. Corespunzător acestor nivele structurale avem diferite discipline biomedicale precum și diferite capitole corespunzătoare ale informaticii medicale.

1.5.2. Bioinformatica – relații cu alte discipline

Bioinformatica este un capitol de sine stătător al informaticii medicale. Însă, prin specificul obiectului său de a urmări modul în care este reprezentată informația în structurile vii, obiectul său este strâns legat de o serie de discipline privind structura și funcția structurilor din materia vie. Astfel, sunt necesare cunoștințele de biofizică, biochimie, biologie celulară și moleculară precum și genetică.

Din punct de vedere al abaterilor formale, legătura cu teoria probabilităților, biostatistica, teoria sistemelor și chimia cuantică este evidentă. De asemenea, se folosesc diverse limbaje de programare – R, PERL, C++, Java etc.

După cum vedem, se poate afirma – pe bună dreptate – că bioinformatica este una dintre cele mai complexe discipline, pregătirea specialiștilor în acest domeniu fiind cel mai adesea inclusă în categoria „studii avansate”.

2. Introducere în biofizică

2.1. *Obiectul și capitolele biofizicii*

2.1.1. *Științe interdisciplinare*

A. Evoluția cunoașterii umane a consacrat de-a lungul timpului conturarea unor domenii solide, bine definite mai ales în conținut decât în definiții formale. Astfel și-au stabilit o poziție solidă diferite științe: matematica, fizica, chimia, biologia etc. Ele au nu numai obiect diferit ci și mod de abordare specific, considerat adesea chiar „mod de gândire”. Însă complexitatea naturii și modul elaborat al gândirii umane au șters granițele dintre științele clasice, fertilizând un teritoriu numit astăzi științe interdisciplinare, ce au cunoscut o dezvoltare deosebită în secolul XX și care s-au dovedit a fi foarte prolifiche.

B. În acest context putem privi și intersecția între două domenii gigant – fizica, respectiv biologia. Mai multe discipline de graniță și-au găsit originea în această intersecție: biofizica, bionica, biocibernetica, biotehnologia, fizica medicală etc., fiecare cu obiect și metode specifice. Putem astfel preciza obiectul câtorva dintre aceste științe de graniță, cele mai bine conturate și dezvoltate.

Biofizica – este disciplina ce se ocupă cu studiul fenomenelor fizice care se petrec la nivelul materiei vii.

Bionica – se ocupă cu aplicarea în tehnică a unor „soluții” din natură (imitarea unor procese din materia vie). Ex.: forma elicopterului, profilul aripilor de avion, eclocația etc.

Biocibernetica – studiază mecanismele de reglare, control și comandă în materia vie precum și integrarea diferitelor nivele informaționale în sistemele biologice.

Biotehnologia – utilizarea organismelor vii sau bioprocесelor în inginerie, tehnologie sau medicină ex.: inginerie genetică, culturi de celule sau de țesuturi, etc.

Fizica medicală – aplicații ale fizicii în medicină, pentru diagnostic (toată aparatura medicală radiografia, ecografia etc.) sau tratament (cu ultrasunete, terapia prin radiații etc.).

2.1.2. *Capitolele biofizicii*

Există mai multe variante de împărțire a biofizicii pe capitole:

A. În funcție de fenomenul fizic urmărit de ex.: biomecanică, bioelectricitate, biomagnetism, fenomene optice în biologie etc.

B. În funcție de nivelul structural abordat astfel avem:

a) *biofizica moleculară* – care se ocupă cu studiul proprietăților moleculelor componente ale materiei vii,

b) *biofizica celulară* – în care se urmăresc fenomenele fizice la nivelul celular, inclusiv particularitățile de abordare a fenomenelor fizice din materia vie în general; în acest capitol sunt incluse: termodinamica biologică, fenomenele de transport – inclusiv

transportul transmembrantar, respectiv bioelectrogeneza – generarea fenomenelor electrice la nivelul membranelor celulare,

c) *biofizica sistemelor complexe* – aspecte fizice ale unor sisteme: aspecte mecanice ale sistemului osteomuscular, aspecte de dinamica fluidelor în sistemul circulator, aspecte specifice în organele de simț mecanismul văzului, auzului etc.,

d) *biofizica ambientală*, numită frecvent și interacțiunea factorilor fizici cu materia vie, capitol dedicat variatelor interacțiuni cu diverși factori fizici: presiune, curent electric, unde mecanice - vibrații, ultrasunete, unde electromagnetice toată gama, de la unde radio și microunde la radiații infraroșii, ultraviolete, X și gamma, radiații corpusculare neutroni etc.

C. Vom prezenta doar câteva noțiuni fundamentale din capitolele de biofizică moleculară și celulară utile pentru înțelegerea fenomenelor ce vor fi prezentate în capitolele următoare.

Din biofizica moleculară vom prezenta pe scurt:

- structura atomului
- atomul de carbon
- structura moleculară
- forțe intermoleculare
- molecula de apă
- soluții
- noțiuni generale despre pH.

Din biofizica celulară ne vom opri la:

- noțiuni de termodinamică biologică
- forțe termodinamice
- procese cuplate
- transport transmembrantar.

2.2. Structura atomului

2.2.1. Proprietăți

A. Atomul are dimensiuni de ordinul 10^{-10} m (10^{-10} m = 1Å Ångstrom). Masa se exprimă în unități atomice de masă sau daltoni (1u.a.m = 1u = 1Da = $1,66 \times 10^{-27}$ kg). Masa relativă a atomului este numărul care arată de câte ori este mai greu atomul respectiv față de u.a.m., de ex.: masa relativă a atomului de carbon ^{12}C este 12.

B. Atomul este neutru din punct de vedere electric.

C. Atomul este format din nucleu și nor electronic.

2.2.2. Nucleul atomic

A. Nucleul atomic are dimensiuni de ordinul 10^{-15} m și concentrează practic masa întregului atom.

B. Nucleul este încărcat electric cu sarcină pozitivă. Nucleul este format din protoni și neutroni. Protonii sunt particule elementare încărcate pozitiv cu sarcina +1e și masa aproximativ 1u; neutronii sunt particule elementare neutre electric, cu masa aproximativ 1u.

C. Nucleul este caracterizat prin două numere:

- *numărul de ordine* Z (sau *număr atomic*), care reprezintă numărul de protoni din nucleu și determină poziția atomului în tabelul periodic al elementelor. Sarcina nucleului este $+Ze$, unde „ e ” este sarcina electrică elementară ($1e = 1,6 \times 10^{-19}C$, $C = \text{coulomb}$)

- *numărul de masă* A , care reprezintă numărul total de protoni și neutroni din nucleu. Cum atât masa protonului cât și a neutronului sunt apropiate de $1u$, numărul de masă reprezintă aproximativ masa nucleului exprimată în u .

D. Simbolic un nucleu se notează ${}_Z^AX$, unde X este simbolul chimic al atomului iar Z și A , ca indici inferiori, respectiv superiori, sunt numărul de ordine, respectiv numărul de masă. (A și Z se pot scrie și ambii de aceeași parte a simbolului elementului).

2.2.3. Norul electronic

A. Nucleul este înconjurat de un nor electronic, alcătuit din Z electroni, fiecare cu sarcina negativă elementară.

a) Norul electronic este împărțit în nivele sau straturi. Un strat electronic este caracterizat prin numărul cuantic n , numit *număr cuantic principal*, care determină energia E_n și dimensiunea norului electronic r_n ;

- valori posibile $n = 1, 2, 3, \dots$ și corespund nivelelor energetice ale straturilor notate K, L, M, \dots ;
- energia nivelului $E_n \sim -1/n^2$, iar distanța medie (\sim raza orbitei) $r_n \sim n^2$.

Observație: energia electronului legat în atom este negativă, electronii din apropierea nucleului sunt puternic legați, cei periferici sunt slab legați; când $n \rightarrow \infty$, $E \rightarrow 0$ și electronul devine liber. Pentru smulgerea unui electron din atom este necesară o energie egală cu energia sa de legătură, iar atomul devine un ion pozitiv.

Ex.: pentru $Z = 1$, $E_1 = -13,6\text{eV}$ iar $r_1 = 0,53\text{Å}$ ($1\text{eV} = 1,6 \times 10^{-19}\text{J}$, $\text{J} = \text{Joule}$).

b) Fiecare strat (nivel) este împărțit în subnivele; un subnivel într-un nivel este caracterizat prin numărul cuantic l numit *număr cuantic orbital* sau azimutal, care determină forma norului electronic;

- valori posibile $l = 0, 1, \dots, n-1$, în total n valori, corespunzând subnivelelor notate în ordine s, p, d, f . Energia subnivelelor crește cu l , dar mai slab decât variația cu n ; un nivel are n subnivele;
- forma norului pentru:
 - $s = 0$ - formă sferică
 - $s = 1$ - formă bilobară
 - $s = 2$ - formă tetralobară etc.

c) Fiecare subnivel este împărțit în orbitali; un orbital într-un subnivel este caracterizat prin *numărul cuantic magnetic* m , care determină orientarea în spațiu a norului electronic;

- valori posibile $m = -l, \dots, -1, 0, +1, \dots, +l$, în total $2l+1$ valori. Toți orbitalii unui subnivel au aceeași energie și se numesc orbitali „degenerați”. Subnivelele s ($l = 0$) au un singur orbital (sferic); subnivelele p ($l = 1$) au trei orbitali (bilobari), orientați pe cele trei direcții din spațiu: p_x, p_y, p_z .

B. Notația simbolică a unui orbital cuprinde numărul nivelului (n) urmat de simbolul subnivelului (s, p, d, f), la care se adaugă un indice privind orientarea spațială determinată de numărul cuantic magnetic; (în cazul orbitalilor s , de formă sferică acest indice nu are sens).

Deci un orbital într-un atom este caracterizat prin 3 numere cuantice:

- nr cuantic principal n ($n = 1, 2, 3, \dots$, pentru nivelele K, L, M, \dots)
- nr cuantic orbital l ($l = 0, 1, \dots, n-1$, pentru subnivelele s, p, d, f)
- nr cuantic magnetic m ($m = -l, \dots, 0, \dots, +l$).

C. Un electron într-un atom este caracterizat prin 4 numere cuantice:

- 3 numere cuantice n, l, m ce determină orbitalul pe care se găsește
- numărul cuantic magnetic de *spin* S , care definește electronul pe orbital
- valori posibile: $s = +1/2, -1/2$, în total 2 valori.

D. Principiul lui Pauli

Enunț: Într-un atom nu pot exista mai mulți electroni cu aceleași valori pentru cele 4 numere cuantice.

Consecință: pe un orbital pot exista maxim 2 electroni, cu spin opus.

2.2.4. Structura norului electronic – principii

A. Atomii diferă între ei prin structura norului electronic, care determină proprietățile chimice ale atomului și principalele proprietăți fizice.

Structura norului electronic poate fi analizată ierarhic, pornind de la cel mai simplu atom – atomul de hidrogen, care are un singur electron și urmărind în continuare structura pentru atomii cu mai mulți electroni.

B. În completarea starturilor electronice pentru atomi vom lua în considerare următoarele principii și reguli:

- a) principiul energiei minime nivelele subnivelele orbitalii se completează începând cu orbitalii de cea mai joasă energie;
- b) principiul lui Pauli pe un orbital încap maxim doi electroni, cu spin diferit;
- c) regula lui Hund în cazul orbitalilor degenerați se completează întâi fiecare orbital cu câte un electron cu același spin, apoi cel de-al doilea electron pe fiecare orbital.

2.2.5. Structura norului electronic al elementelor

Cu principiile prezentate mai sus putem reprezenta schematic structura norului electronic al elementelor de la începutul tabelului periodic al elementelor.

În figura 2.2.5.a sunt prezentate simbolic structurile norului electronic pentru elementele cu Z până la 13.

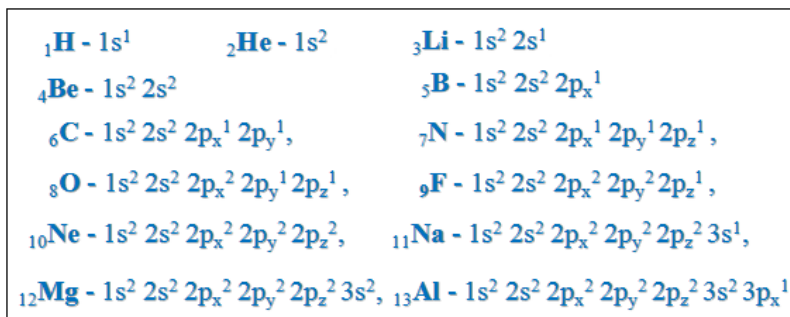


Fig. 2.2.5.a. Structura norului electronic al elementelor de la începutul tabelului periodic

În fig. 2.2.5.b. sunt prezentate pozițiile relative ale nivelelor și subnivelelor energetice.

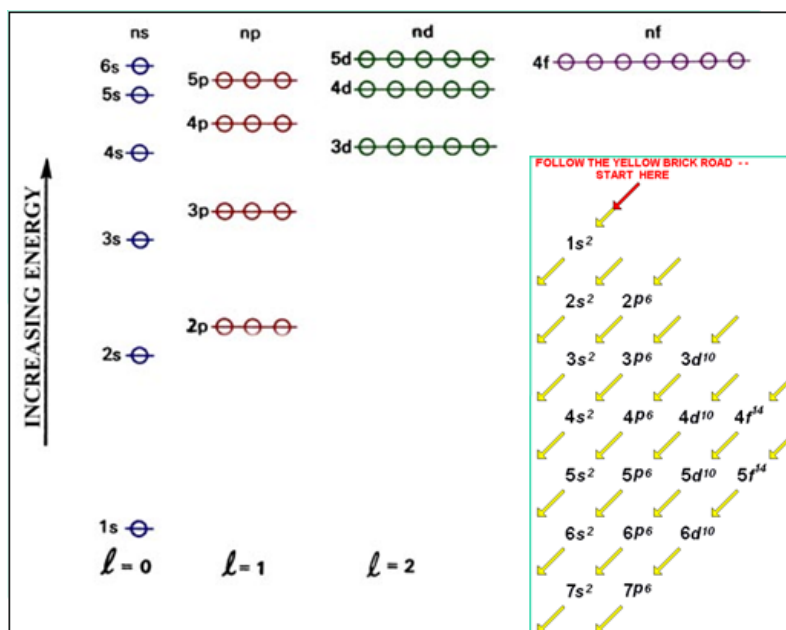


Fig. 2.2.5.b. Structura norului electronic și poziția relativă a subnivelelor energetice

2.3. Atomul de carbon

2.3.1. Structura norului electronic

A. Atomul de carbon ocupă poziția centrală în chimia organică, numeroase proprietăți ale moleculelor din materia vie fiind mai ușor înțelese dacă înțelegem configurația norului electronic al atomului de carbon.

Din cei 6 electroni, 2 sunt pe stratul K, aproape de nucleu, pe care nu îl părăsesc în nici un proces biochimic, deci îl neglijăm în comentariile ce urmează.

B. Cei 4 electroni de pe stratul L sunt situați în stare fundamentali în structura $2s^2, 2p_x^1, 2p_y^1$.

Configurația norului electronic al atomului de carbon:	1s	2s	2p _x	2p _y	2p _z
- sus: în stare fundamentală	↓↑	↓↑	↑	↑	
- jos: în stare excitată (hibridizare)	↓↑	↑	↑	↑	↑

Fig. 2.3.1. Structura norului electronic al atomului de carbon

Aceasta explică comportamentul bivalent în unele situații, de exemplu în molecula de CO, prin completarea orbitalilor 2p_x și 2p_y care erau incompleți.

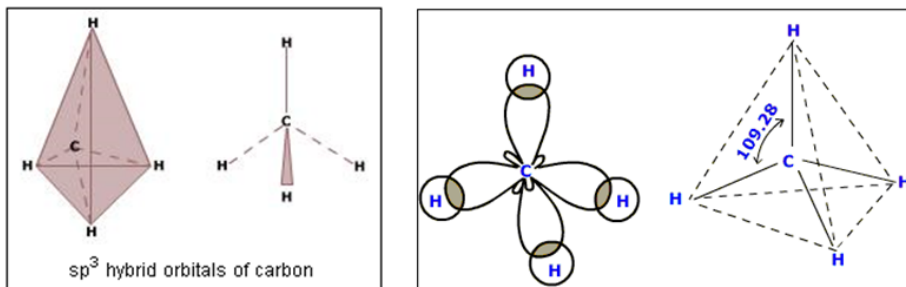
2.3.2. Hibridizarea

A. Diferența de energie între subnivelul 2s și 2p este foarte mică, astfel încât un electron de pe 2s sare pe orbitalul liber 2p_z, având acum 4 orbitali cu câte un electron fiecare, cu același spin (conform regulii lui Hund). Se explică astfel caracterul uzual tetravalent al atomului de carbon.

Saltul electronului de pe 2s pe 2p_z se numește hibridizare.

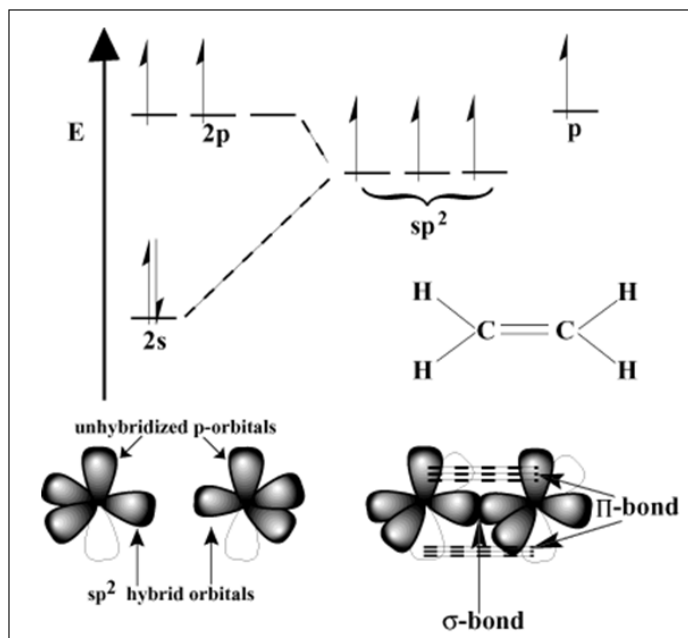
B. Hibridizarea este însoțită de o reorientare în spațiu a orbitalilor. Aceasta se poate realiza în trei feluri:

a) *Hibridizare sp³*: toți orbitalii se reorientează simetric în spațiu, cu densitate maximă orientată spre vârfurile unui tetraedru și toți au aceeași energie, caz în care atomul de carbon are cele 4 valențe echivalente și le poate angaja în diferite legături chimice.

Fig. 2.3.2.a. Hibridizarea sp³

Prin această orientare în spațiu a valențelor se pot explica și o serie de proprietăți, cum ar fi chiralitatea izomeria sterică, sau izomeria optică, în funcție de atomii cu care se realizează legăturile.

b) *Hibridizarea sp²*. În cazul în care un electron de pe un orbital, (ex.: 2p_x), era angajat într-o legătură chimică, (de exemplu cu un alt atom de carbon), reorientarea va cuprinde numai orbitalul s și ceilalți 2 orbitali p (fig. 2.3.2.b).

Fig. 2.3.2.b. Hibridizarea sp^2

În acest caz poate să apară și o a doua legătură între doi atomi, prin întrepătrunderea intersecția norului electronic ai orbitalilor care nu erau orientați de-a lungul axei dintre atomi ci într-un plan perpendicular. Evident, cele două legături nu sunt echivalente; prima, generată la intersecția orbitalilor de-a lungul axei este puternică și se numește legătură σ , iar a doua este mai slabă și se numește legătură π . (Comparativ, pentru legătura simplă C-C, energia legăturii este 347 kJ/mol, iar a legăturii π din legătura dublă este 263kJ/mol).

Important de remarcat că ceilalți orbitali se orientează simetric față de norul format.

c) *Hibridizarea sp* . În cazul în care ambii electroni de pe 2p sunt angajați în legături chimice, reorientarea se face între orbitalul s și orbitalul $2p_z$ rămas gol.

Într-o astfel de situație apare și legătura triplă, în care apare o legătură σ și două legături π , (mai slabe în medie decât când era doar una).

C. Orientarea norului electronic în moleculele organice (în special proteine) este un subiect important, prin aceasta explicându-se proprietățile lor. Există software dedicat pentru astfel de analize.

2.4. Structura moleculei. Legături chimice

2.4.1. Stabilitatea atomilor

A. Atomii au tendința de a avea un nor electronic complet. În funcție de structura norului electronic putem estima gradul de stabilitate și proprietățile diverselor clase de atomi. Astfel, atomii cu norul electronic complet sunt foarte stabili (grupa 18) - gazele nobile (inerte).

B. Atomii cu puțini electroni pe ultimul strat (electroni de valență) au tendința de a-i ceda mai ușor (caracter electropozitiv), în timp ce atomii cărora le lipsesc puțini electroni pentru completarea subnivelului periferic au tendința de a acapara electroni (caracter electronegativ). Cu cât aceste tendințe sunt mai puternice, atomii respectivi sunt mai puțini stabili (mai reactivi).

2.4.2. Legătura covalentă

A. Cea mai evidentă manifestare a tendințelor atomilor către configurația electronică stabilă apare prin formarea legăturilor chimice, iar dintre acestea, cea mai frecventă (și mai generală) formă este legătura covalentă.

Legătura covalentă se realizează prin punerea în comun a electronilor de pe orbitalii incompleți (cu câte 1 electron). Orbitalul nou format se numește orbital molecular și se redistribuie spațial între cei doi atomi, existând posibilitatea – prin chimie cuantică - de a calcula densitatea norului electronic nou format.

B. Legătura covalentă nepolară

În cazul în care legătura covalentă se realizează între atomi identici, norul electronic molecular va fi simetric distribuit peste ambii atomi, centrul sarcinilor pozitive (aflat la mijlocul distanței între cele două nuclee) va coincide cu centrul sarcinilor negative; molecula formată este nepolară (exemple: H_2 , Cl_2 , figura 2.4.2. sus). Schematic perechea de electroni puși în comun se notează cu o linie „-”, cu semnificația de legătură simplă (o singură pereche de electroni); în mod uzual electronii neparticipanți nici nu se reprezintă.

C. Legătura covalentă polară

În cazul în care atomii sunt diferiți, atomul mai electronegativ atrage norul electronic mai mult spre sine și va deveni un centru al sarcinilor negative, în timp ce atomul mai slab electronegativ va rămâne parțial denudat de nor electronic, devenind un centru al sarcinilor pozitive. Molecula obținută este o moleculă polară, este caracterizată printr-un moment de dipol (μ) și se va orienta de-a lungul liniilor de câmp în cazul plasării într-un câmp electric. De asemenea, fiind un dipol, se va implica într-o serie de fenomene electrice cu caracter molecular. Exemple: HCl , prezentat în figura 2.4.2. jos, H_2O etc.

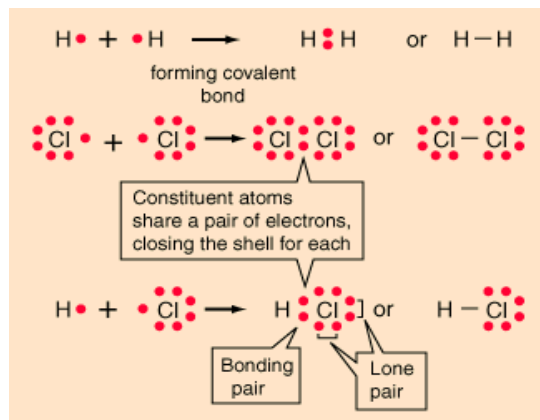


Fig. 2.4.2. Legătura covalentă polară și nepolară

2.4.3. Energia de legătură

A. Energia de legătură a unei legături chimice este definită ca energia necesară pentru ruperea sa (deci este o energie negativă) și este egală cu energia eliberată la formarea moleculei. Se exprimă uzual în kJ/mol (se mai folosește kcal/mol).

B. Un grafic al energiei potențiale a unui sistem de doi atomi este prezentat în figura 2.4.3. și poartă numele de groapă de potențial. Distanța între atomi pentru energia potențială minimă se numește distanța de echilibru.

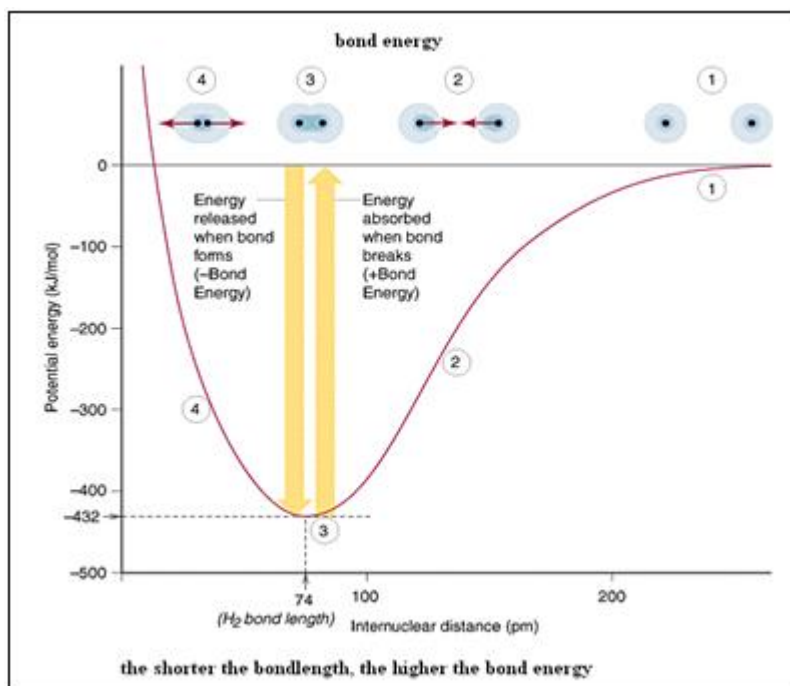


Fig. 2.4.3. Energia de legătură

C. Atomii nu sunt fișii la distanța de echilibru ci execută mișcări de oscilație vibrație în jurul poziției de echilibru. Amplitudinea vibrațiilor crește cu temperatura. Nivelele energetice de vibrație sunt cuantificate. Fiecare moleculă are un spectru de vibrație specific, cu benzi situate uzual în regiunea infraroșie a spectrului.

D. Stabilitatea unei molecule depinde de energia de legătură: cu cât energia de legătură este mai mare (la formare s-a eliberat mai multă energie), cu atât molecula este mai stabilă.

E. Există o relație între distanța de echilibru și energia de legătură pentru energii de legătură mai mari, forțele de atracție sunt mai puternice și distanța de echilibru lungimea legăturii este mai mică.

2.4.4. Legătura ionică

A. Un caz limită al legăturii covalente polare este cazul în care elementul mai electronegativ preia integral electronul pus în comun, devenind iar negativ, iar celălalt atom a devenit un ion pozitiv. Ex.: NaCl, MgO (figura 2.4.4.a).

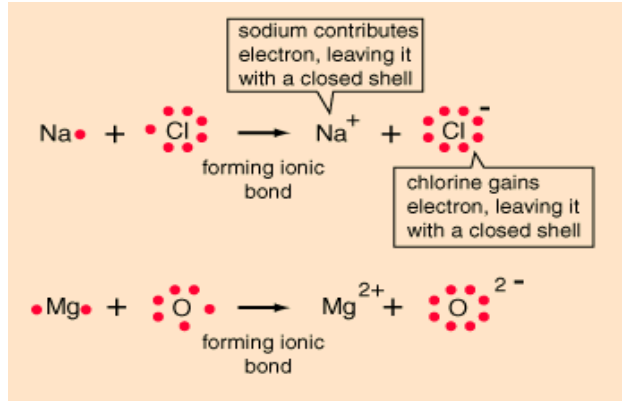


Fig. 2.4.4.a. Legătura ionică

B. Între cei doi ioni se exercită o forță electrostatică coulombiană, iar energia de legătură va fi chiar energia potențială corespunzătoare acestei interacțiuni (figura 2.4.4.b)

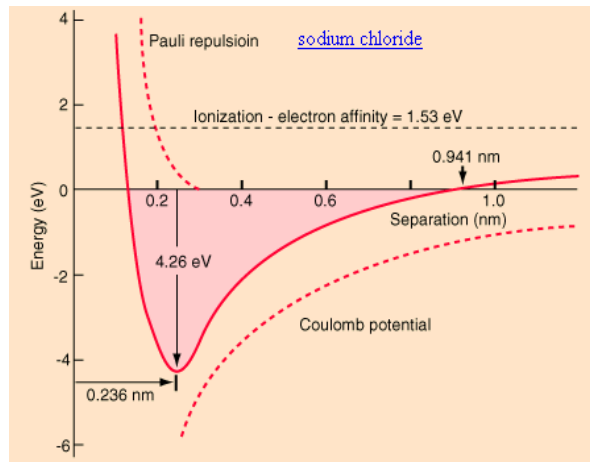


Fig. 2.4.4.b. Energia legăturii ionice

C. În cazul legăturii ionice, noțiunea de moleculă devine formală, în structuri avem doar ioni independenți, pozitivi și negativi. Aceasta este situația și în cristale, (de ex.: NaCl) și în soluție; să facem totuși observația că ionii în soluție sunt hidratați, (moleculele de apă fiind polare, sunt atrase cu capătul de sarcină opusă către ion).

2.5. Forțe intermoleculare

Moleculele din structuri interacționează între ele. Aceste interacțiuni sunt mai slabe decât legăturile chimice, însă existența forțelor intermoleculare determină o serie de proprietăți importante ale moleculelor, ex.: solubilitatea.

Ele pot fi:

- legătura de hidrogen
- forțe Van der Waals
- forțe de dispersie.

2.5.1. Legătura de hidrogen

A. Hidrogenul participant în legăturile covalente polare este adesea victima unei interacțiuni cu un atom mai electronegativ, rămânând parțial privat de nor electronic și un centru de sarcină pozitivă. În aceste condiții el poate fi atras de un nor electronic complet al unui orbital, creat de o pereche de electroni neparticipanți la vreo legătură chimică (deci nor dens în jurul unui atom electronegativ - centru de sarcină negativă).

B. Legătura de hidrogen este definită ca interacțiune de natură electrostatică între hidrogenul unei molecule cu electronii neparticipanți ai unui atom din altă moleculă, de obicei un atom de oxigen sau azot.

C. Energia de legătură în cazul punților de hidrogen este de 10 ori mai slabă decât a legăturilor covalente (4,5 kcal/mol, față de 110 kcal/mol în legătura O-H). Exemple: H₂O, NH₃, A-T/G-C etc. (figura 2.5.1).

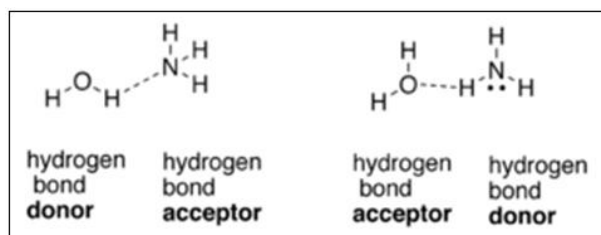


Fig. 2.5.1. Punți de hidrogen

2.5.2. Forțe Van der Waals

A. Forțele van der Waals sunt forțe de atracție slabe, ce se manifestă doar pe distanțe foarte mici, scăzând foarte repede cu distanța dintre molecule. Ele pot fi datorate interacțiunilor de origine electrostatică între dipolii moleculari.

B. Astfel putem avea:

- interacțiuni dipol - dipol, care apar între centrul sarcinilor pozitive ale unei molecule polare și centrul sarcinilor negative al altei molecule polare;
- interacțiuni dipol - dipol indus, care pot implica și molecule nepolare, dar care, sub acțiunea câmpului electric al unei molecule puternic polare, își deformează norul electric și apare un dipol indus care interacționează cu molecula polară.

2.5.3. Forțe de dispersie

Forțele de dispersie sunt și mai slabe ca forțele Van der Waals, acționează doar la distanțe foarte mici și cresc odată cu masa moleculară. Le luăm în considerare doar dacă celelalte forțe sunt toate mici.

2.6. Molecula de apă

2.6.1. Structura moleculară

A. Apa este un component esențial al materiei vii, având o serie de proprietăți care trebuie luate în considerare pentru a înțelege o suită de fenomene din lumea vie.

B. Structura moleculară a apei este bine cunoscută, H_2O . Norul electronic este deplasat către oxigen, care devine centru de sarcini negative, centrul sarcinilor pozitive fiind la semidistanța între atomii de hidrogen. Datorită respingerii electrostatice dintre nucleele de hidrogen, unghiul între valențe crește de la 90° la 105° , iar orbitalii cu electroni neparticipanți ai oxigenului se reorientează spre vârfurile unui tetraedru (neregulat, nu identic cu norul carbonului, dar asemănător). În figura 2.6.1. sunt prezentate schematic ideile redată aici, formarea legăturilor covalente O-H, unghiul între valențele O-H, tetraedrul cu poziția atomilor de hidrogen și densitățile maxime ale norului electronilor neparticipanți.

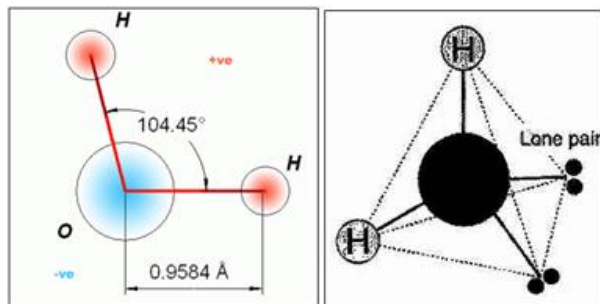


Fig. 2.6.1. Structura moleculei de apă

2.6.2. Legăturile de hidrogen ale moleculei de apă

A. O moleculă de apă poate prezenta 4 legături de hidrogen

– două legături, stabilite de cei doi atomi de hidrogen, care se îndreaptă, fiecare, spre o altă moleculă de apă, mai precis către electronii neparticipanți ai altei molecule de apă;

– alte două legături prin cei doi lobi corespunzători orbitalului $2p$ care hibridizează cu $2s$, formând un nor cu densități maxime către celelalte două vârfuri ale tetraedrului.

B. O reprezentare sugestivă, folosită frecvent în chimia organică și biochimie, este prezentată în figura 2.6.2, în care cu culoare închisă este redat atomul de oxigen, iar cu culoare deschisă atomul de hidrogen. Sunt evidențiate cele 4 legături posibile.

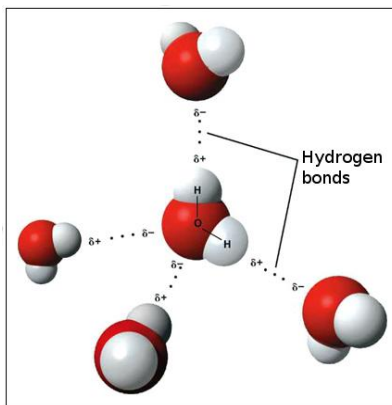


Fig. 2.6.2. Legăturile de hidrogen ale moleculei de apă

2.6.3. Proprietățile apei

A. Prezența legăturilor de hidrogen face ca apa să ocupe un loc excepțional comparativ cu alte molecule, favorizând poziția sa privilegiată ca și component esențial al materiei vii. Fără a intra în detalii explicative, prezentăm în continuare o listă a proprietăților apei, în care, în bună parte, explicațiile invocă prezența punților de hidrogen:

- gheața are structură cristalină hexagonală, cu spații libere, având densitatea mai mică decât apa lichidă,
- la 0°C numai 15% din legăturile de hidrogen se rup; apa lichidă are structură cvasicristalină, formată în special din trimeri și dimeri,
- densitatea apei este maximă la 4°C,
- tensiunea superficială a apei este foarte mare (cca. 72 N/m, față de cca. 22 N/m pentru alcool),
- căldura specifică (implicit capacitatea calorică) mare, acumulând o parte din energie pentru ruperea unor legături de hidrogen, deci ca energie potențială și nu ca energie cinetică de vibrație a moleculelor; inerția termică este considerată ca unul din factorii cheie în evoluția materiei vii ($c = 4,18 \text{ kJ/kg}\cdot\text{grd} = 1 \text{ kcal/kg}\cdot\text{grd}$),
- temperatura de topire $t_t = 0^\circ\text{C}$, temperatura de fierbere $t_f = 100^\circ\text{C}$ (față de alcool: $t_t = -38^\circ\text{C}$, $t_f = 78^\circ\text{C}$),
- conductibilitatea termică ridicată ($\sim 0,0013 \text{ cal/cm}^2\cdot\text{grd}$),
- căldurile latente foarte mari: căldura latentă de topire $\lambda_t \sim 80 \text{ kcal/kg} = 335 \text{ kJ/kg}$, căldura latentă de fierbere $\lambda_f \sim 580 \text{ kcal/kg} = 2258 \text{ kJ/kg}$,
- conductibilitatea electrică foarte mică (rezistivitate mare, $\sigma \sim 11 \cdot 10^{-6} \Omega^{-1}\cdot\text{m}^{-1}$),
- constanta dielectrică relativă foarte mare ($\epsilon = 80$), datorită efectului dipolilor moleculari,
- proprietăți optice: transparentă, indice de refracție $n = 4/3$, nu absoarbe în domeniul vizibil, absoarbe moderat radiațiile UV și este opacă la radiațiile IR.

B. Vâscozitatea – este o proprietate ce reflectă gradul de frecare internă între molecule în cazul curgerii lichidului. Vâscozitatea scade cu temperatura. În cazul apei se constată o scădere aproximativ liniară de la 4°C la 35°C, apoi o scădere bruscă în intervalul 35 - 40°C și din nou o scădere liniară până la fierbere. Aceasta datorită faptului că în intervalul 35 - 40°C are loc ruperea masivă a unor punți de hidrogen, nu mai rămân în soluție trimeri (grupe de trei molecule de apă) și mulți dimeri se desfac. Temperatura de 37°C se mai numește al doilea punct de topire al apei și nu întâmplător animalele homeoterme au ca temperatură de echilibru tocmai această valoare.

Proprietățile enumerate mai sus, chiar dacă nu au fost explicate în detaliu, joacă un rol esențial în înțelegerea unui mare număr de procese biologice și vom face referiri specifice unde este cazul.

2.7. Soluții

2.7.1. Sisteme disperse

A. În natură nu ne întâlnim cu sisteme pure, formate dintr-un singur fel de molecule, ci cu sisteme formate din mai multe tipuri de molecule. Aceste sisteme poartă numele generic de sisteme disperse.

Repartiția spațială a tipurilor de molecule poate fi uniformă (sisteme omogene) sau neuniformă (sisteme heterogene). Pentru moment vom lua în considerare doar sistemele omogene.

B. Compoziția sistemelor disperse

În general un sistem dispers are două componente principale, numite *faze*:

- o fază continuă, numită *solvent*, alcătuind componenta majoritară a sistemului,
- o fază discontinuă discretă, numită *solut*, reprezentată de substanța dizolvată.

C. Cel mai adesea considerăm cele două faze ca fiind reprezentate de starea de agregare lichidă. Totuși, la modul general, putem defini sisteme disperse cu toate combinațiile posibile pentru solvent și solut.

D. Clasificarea sistemelor disperse după diametrul particulelor solutului.

a) *soluții moleculare*, cu diametrul sub 10 Å; cum de obicei masa moleculară este proporțională cu volumul, deci cu puterea a treia a diametrului, uzual se consideră în categoria soluții moleculare, soluțiile care au pentru substanța dizolvată masa moleculară $M < 1000$,

b) *soluții coloidale*, cu diametrul între 10 și 1000 Å și masa moleculară peste 10^3 , dar sub 10^6 ,

c) *dispersii* – medii cu particule cu diametrul 1000 Å sau $M > 10^6$.

E. O altă clasificare posibilă ia în considerare procesul de disociere; astfel soluțiile pot fi:

a) electroliți – în cazul în care moleculele disociază să specificăm aici că electroliții pot fi la rândul lor:

- i. electroliți tari, când disocierea este totală este cazul tuturor sărurilor, dar și acizii și bazele tari;
- ii. electroliți slabi, când disocierea este parțială, acizi sau baze slabe în aceste situații trebuie definită și o constantă de disociere, ca raportul între numărul moleculelor disociate și numărul total de molecule dizolvate.

b) ne-electroliți – când substanțele dizolvate sunt formate din molecule care nu disociază, ex.: glucoză, uree etc.

2.7.2. Concentrații

A. O caracteristică esențială a soluțiilor este concentrația acestora, exprimată ca măsură a proporției moleculelor substanței dizolvate în soluție.

Există mai multe modalități de a exprima concentrația unei soluții:

a) Concentrația procentuală – exprimă cantitatea de substanță dizolvată, exprimată în grame, la 100 ml soluție;

b) Concentrația molară – exprimă numărul de moli de substanță dizolvată în 1 l soluție ($1M = 1 \text{ mol} =$ cantitatea de substanță, exprimată în grame, numeric egală cu masa moleculară exprimată în u.a.m.; 1 mol are N_A molecule, N_A este numărul lui Avogadro $N_A = 6,023 \cdot 10^{23}$ molecule/mol);

c) Concentrația normală – exprimă numărul de echivalenți gram de substanță dizolvați în 1 l de soluție (1 echivalent-gram = 1 mol/z, unde z este numărul de sarcini electrice de un semn, ce se obține la disocierea unei molecule).

B. Relații

$$C(M/l) = m(g)/V(l) \cdot M \quad (2.7.2)$$

unde C este concentrația, în mol/litru, V = volumul soluției în l, M este masa moleculară

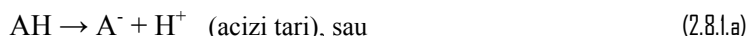
2.8. pH-ul soluțiilor

2.8.1. Disocierea electroliților

A. Electroliții disociază în soluție. Disocierea poate fi totală sau parțială. În cazul disocierii parțiale se definește constanta de disociere k .

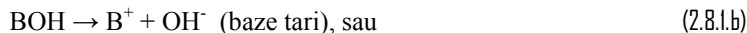
B. Acizi

Substanțele care prin disociere pun în libertate protoni ioni H^+ se numesc acizi. Acizii tari disociază total, acizii slabi disociază parțial. Ex.: HCl , H_2CO_3



C. Baze

Substanțele care prin disociere pun în libertate ioni OH^- se numesc baze (substanțe alcaline). Bazele tari disociază total, bazele slabe disociază parțial. Ex.: $NaOH$, NH_4OH .



D. Disocierea apei

Molecula de apă disociază parțial în soluție eliberând un proton de H^+ și un ion oxidril (hidroxil) OH^- .



E. Constanta de echilibru

În cazul disocierii parțiale, putem explica legea acțiunii maselor constanta de echilibru se numește constantă de disociere. Iată, pentru acizii slabi și bazele slabe.



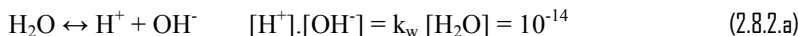
F. Bazele ca acceptori de protoni

O substanță poate avea caracter bazic fără a elibera direct ioni OH^- în soluție apoasă. Este suficient a accepta protoni (ca opus al eliberării de protoni de către acizi) protonii acceptați provin dintr-o moleculă de apă, astfel încât în soluție, în mod indirect apare un ion OH^- . Ex.: amoniacul NH_3 are caracter alcalin:



2.8.2. Produsul ionic al apei

A. În cazul apei, în condiții normale ($p = 1 \text{ atm}$, $t = 20^\circ C$) produsul $[H^+] \cdot [OH^-]$ se numește produs ionic al apei și are valoarea 10^{-14} .



B. Cum fiecare moleculă de apă eliberează în mod egal un ion de H^+ și unul de OH^- rezultă că:

$$[H^+] = [OH^-] = 10^{-7} \quad (2.8.2.b)$$

Cu alte cuvinte 1 moleculă din 10 milioane este disociaată!

C. Ionul de hidroniu

De menționat că ionul H^+ nu se găsește niciodată liber în soluție, ci este întotdeauna cuplat pe o moleculă de apă formând ionul de hidroniu H_3O^+ . Pentru simplitate în abordarea formală păstrăm însă notația de H^+ .

2.8.3. Scara pH

A. Definiția pH-ului

Exprimarea concentrației ionilor H^+ într-o soluție se face prin definirea scării pH.

Definiție: pH-ul unei soluții este logaritmul cu semn schimbat al concentrației ionilor de hidrogen.

$$pH = -\log [H^+] \quad (2.8.2.c)$$

Se aplică logaritmi zecimali. În cazul apei cunoscând $[H^+] = 10^{-7}$ rezultă $pH = 7$. Pentru un acid tare cu concentrația 1 mol/l, $pH = 0$.

B. Definiția pOH-ului

Similar, logaritmul cu semn schimbat al concentrației ionilor OH^- se numește pOH.

$$pOH = -\log [OH^-] \quad (2.8.2.d)$$

C. Relația pH – pOH

Prin logaritizarea formulei 2.8.2.a se obține relația:

$$pH + pOH = 14 \quad (2.8.2.e)$$

Din acest motiv, nu este necesar a lucra cu ambele mărimi, fiind suficientă una. În practică se folosește doar scara pH.

D. Descrierea scării pH

Caracterul acid, neutru sau alcalin al unei soluții se poate recunoaște după valoarea pH-ului. Scara pH, stabilită de Sorensen are valori între 0 și 14 și este ilustrată în figura 2.8.3, în care sunt redată și câteva exemple de valori de pH ale unor soluții mai des întâlnite.

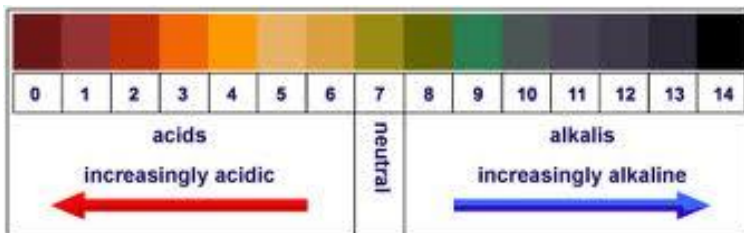


Fig. 2.8.3. Scara pH

E. Importanța pH-ului în materia vie

Există substanțe care au caracter amfoter, adică disociază atât ca acid cât și ca bază. În această categorie intră și amino - acizii. Gradul de disociere al grupărilor carboxil sau

amino va depinde de pH-ul soluției. Însă gradul de disociere determină și sarcina electrică a moleculei, deci o serie de proprietăți. De aceea, organismele vii au mecanisme de reglare pentru menținerea pH-ului la valori relativ constante.

F. Soluții tampon

Soluțiile tampon sunt soluții care mențin constantă valoarea pH-ului. Ele sunt de regulă compuse dintr-un acid slab și o sare a sa, sau dintr-o bază slabă și o sare a sa. Să luăm primul caz:



Când în soluție se adaugă un acid deci apare o abundență de ioni H^+ , echilibrul primei reacții se mută la stânga, în soluție existând suficienți ioni A^- proveniți din sare. În consecință se formează noi molecule neutre AH și nu crește concentrația ionilor H^+ .

În cazul în care se adaugă o bază, procesul este invers, ionii H^+ neutralizați de ionii OH^- din bază, generând molecule de apă, sunt refăcuți prin disocierea altor molecule de AH, echilibrul mutându-se la dreapta.

pH-ul unei soluții tampon se calculează cu relația Henderson-Hasselbach:

$$\text{pH} = \text{pk}_a + \log \left(\frac{[\text{f.deprot.}]}{[\text{f.protonata}]} \right) \quad (2.8.2.g)$$

unde [f. deprot.] reprezintă concentrația formei deprotonate, iar [f. protonata] reprezintă concentrația formei protonate.

Dacă luăm $\text{pk}_a = -\log k_a$, observăm că atunci când sistemul este în echilibru, $\text{pH} = \text{pk}_a$ iar concentrația formei protonate este egală cu cea a formei deprotonate.

În exemplul de mai sus, forma protonată este AH iar forma deprotonată este A^- sau AMe (sărurile disociază total).

2.9. Termodinamica biologică

Unul dintre capitolele importante din biofizica celulară este cel referitor la termodinamica biologică. Vom urmări în continuare câteva deosebiri esențiale între materia nevie și materia vie din punct de vedere termodinamic. Să revedem însă câteva noțiuni introductive de termodinamică.

2.9.1. Parametrii de stare ai unui sistem termodinamic

A. Obiectul de studiu în termodinamică se numește sistem termodinamic. El este definit ca o porțiune finită din univers, delimitată fizic sau imaginar de restul universului, numit mediu exterior.

B. Un sistem termodinamic este caracterizat la un moment dat printr-un set de parametri numiți parametri de stare. De ex.: un gaz într-un cilindru cu piston este un sistem termodinamic, caracterizat prin parametrii de stare presiune p , temperatură T și volum V .

C. Parametrii de stare pot fi:

- *extensivi* – care depind de dimensiunea sistemului, de ex.: volumul, masa, numărul de moli, energia internă, etc.,

- *intensivi* – care nu depind de dimensiunea (masa) sistemului de ex.: temperatura, presiunea, densitatea, etc.

Parametrii extensivi au proprietatea de aditivitate, cei intensivi nu. Dacă asociem două sisteme 1 și 2 pentru a crea un nou sistem, 3, vom avea masa finală $m_3 = m_1 + m_2$ energia internă $U_3 = U_1 + U_2$ ș.a.m.d.

D. Sistem omogen, gradienti, sistem izolat

Când parametrii intensivi au aceeași valoare în orice punct din sistem spunem că sistemul este omogen. Când sistemul nu este omogen, putem defini variația unui parametru intensiv Y de-a lungul unei axe x prin gradientul parametrului respectiv dY/dx .

E. Un sistem termodinamic care nu poate schimba cu exteriorul nici substanță, nici energie se numește sistem izolat. În cazul în care sistemul nu schimbă substanță, dar poate schimba energie sub formă de lucru mecanic sau căldură sistemul se numește închis. În contrast cu sistemul închis este sistemul deschis, care permite și schimb de substanță. Un caz aparte de sistem închis este sistemul izolat adiabatic, care nu schimbă substanță, iar energia poate fi schimbată numai sub formă de lucru mecanic, nu și sub formă de căldură.

2.9.2. *Procese termodinamice*

A. Stare de echilibru. Starea unui sistem în care parametrii de stare rămân constanți în timp și prin sistem nu circulă fluxuri se numește stare de echilibru.

B. Stare staționară - este starea unui sistem care are parametrii de stare constanți în timp, dar sistemul este traversat de fluxuri constante.

C. Trecerea unui sistem dintr-o stare în altă stare se numește transformare sau proces termodinamic.

D. Un proces este caracterizat prin mărimi de proces, numite variații. Dacă (1) și (2) sunt două stări distincte ale sistemului în două momente, t_1 și t_2 , atunci variația unui parametru de stare X va fi $\Delta X = X_2 - X_1$, care pentru intervale scurte de timp devine dX/dt . Variațiile pot fi definite atât pentru parametrii intensivi cât și pentru cei extensivi.

E. Există procese în care unii parametri rămân constanți; iată denumirile câtorva procese particulare:

- transformarea izotermă – când temperatura rămâne constantă,
- transformarea izobară – când presiunea rămâne constantă,
- transformarea izocoră – când volumul rămâne constant.

F. Procese reversibile și ireversibile

Un proces termodinamic între două stări 1 și 2 se numește reversibil dacă sistemul poate reveni din starea 2 în starea 1 prin aceleași stări intermediare. În caz contrar este proces ireversibil. Procesele reversibile sunt ideale, procesele reale sunt ireversibile, dar putem avea procese care se apropie destul de mult de procesele reversibile.

2.9.3. *Funcții de stare*

Procesele termodinamice sunt mai ușor descrise cu ajutorul unor anumite funcții de stare, pe care le enumerăm fără a le defini sau analiza proprietățile:

- entropia termodinamică S (exprimată în J/K)
- energia internă U (exprimată în J)
- energia liberă $F = U - T.S$ (J)
- entalpia $H = U + p.V$ (J)
- entalpia liberă $G = U + p.V - T.S$ (J).

2.9.4. Principiul al doilea al termodinamicii

Cel mai controversat aspect al termodinamicii biologice este legat de al doilea principiu al termodinamicii. Vom face câteva scurte comentarii legate de acest principiu.

A. Enunțul său din fizica clasică are mai multe formulări echivalente.

a) Căldura nu poate trece de la sine de la un corp cu temperatură mai scăzută la unul cu temperatură mai ridicată

b) În procesele termodinamice entropia nu poate scădea; variația entropiei este zero în procese reversibile și pozitivă în procesele ireversibile.

B. Entropia termodinamică poate fi corelată intuitiv (demonul lui Maxwell) cu gradul de ordine la scară moleculară. În procesele ireversibile crește dezordinea moleculară (numărul de stări în care putem aranja moleculele componente) asociată cu creșterea entropiei. Acest lucru este valabil pentru sisteme nevii.

C. Însă procesele din materia vie sfidează (aparent) al doilea principiu al termodinamicii: structurile evoluează către stări tot mai ordonate!

Explicația poate fi dată prin legătura care se face între entropia termodinamică și cea informațională. Un proces termodinamic însoțit de scăderea entropiei termodinamice ar fi posibil dacă sistemul își crește entropia informațională.

2.10. Procese cuplate

2.10.1. Natura proceselor cuplate

A. Procesele din materia vie nu sunt procese izolate ci procese cuplate. Este deci posibil a avea cuplate două procese – unul în care entropia scade (reacții de sinteză) cuplat cu unul în care entropia crește (un proces catabolic, spontan), respectând și principiul al II – lea prin creșterea globală a entropiei în ansamblul proceselor cuplate.

B. Din punct de vedere energetic un proces care necesită energie (endoenergetic) trebuie cuplat cu un proces ce eliberează energie (exoenergetic).

2.10.2. Laturile metabolismului

Procesele cuplate le întâlnim ca laturi principale ale metabolismului.

A. *Catabolismul* – cuprinde procese exoenergetice, de degradare moleculară, în care entropia crește sunt procese spontane.

B. *Anabolismul* – cuprinde latura specifică a materiei vii, procese endoenergetice în care se produc structuri mai ordonate și mai bogate în energie în aceste procese entropia scade (reacții de sinteză, endoenergetice).

2.10.3. Stocarea energiei pentru procesele biologice

Pentru a asigura energia necesară reacțiilor anabolice, în materia vie energia este stocată în molecule de ATP (acid adenzin-trifosforic). Aceste molecule, numite și molecule macroergice, sunt sintetizate în organellele celulare numite mitocondrii, printr-un proces numit fosforilare oxidativă. Descifrarea mecanismelor sintezei ATP prin pompa de protoni de către Mitchell a fost răsplătită cu un premiu Nobel, în 1974.

De fapt importanța proceselor termodinamice, înțelegerea lor în contextul materiei vii, a generat o bună conturare a termodinamicii biologice, numită și termodinamica proceselor ireversibile, pentru care s-a acordat premiul Nobel lui Prigogine, în 1971.

2.11. Forțe termodinamice

A. Forțe și fluxuri termodinamice

În termodinamica biologică se generalizează termenul de forță, numită forță termodinamică, definită ca gradient al unui parametru intensiv. Efectul produs de o forță termodinamică se numește flux conjugat și se definește ca transferul unei mărimi printr-o secțiune unitară, în unitatea de timp.

B. Lista forțelor termodinamice și fluxurile lor conjugate

Tabelul 2.11. Forțe și fluxuri termodinamice

Forța	Param. intensiv	Flux conjugat	Flux
Δp	gradient de presiune	flux de volum (curgere)	J_V
ΔT	gradient de temperatură	flux de căldură	J_Q
ΔE	gradient de potențial electric	curent electric	I
ΔC	gradient de concentrație (fiecare subst. i)	difuziune	J_i
$\Delta \pi$	gradient de presiune osmotică	osmoză	J_w
$\Delta \mu$	gradient de potențial chimic	rata reacției chimice	v

2.12. Transport transmembrantar

2.12.1. Clasificare, proprietăți

O serie de molecule importante în studiile de bioinformatică le întâlnim în structurile care asigură transportul substanțelor prin membranele biologice.

Transportul prin membranele biologice poate fi:

- transport pasiv în sensul gradientului electrochimic, fără consum energetic
- transport activ contra gradientului electrochimic, cu consum energetic.

2.12.2. Transport pasiv

Transportul pasiv se poate realiza prin:

- a) difuziune – în cazul gazelor sau substanțelor solubile în lipide
- b) prin canale ionice – există structuri moleculare de natură proteică, ce traversează membrana celulară, având forma unui canal, adesea având situri de cuplare temporară a ionilor. În figura 2.12.2 este prezentat un canal de Na^+ .

c) transport facilitat – cu ajutorul unor molecule ”carrier” pentru molecule pentru care nu există nici canale, nici nu pot difuza, de ex. transportul glucozei prin membrana eritrocitară.

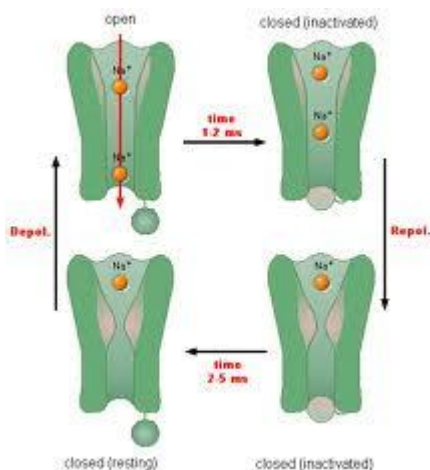


Fig. 2.12.2. Canal de Na⁺

2.12.3. Transport activ

A. Transportul activ are un rol important în menținerea parametrilor fizico-chimici ai celulelor. El este asigurat de structuri specializate numite pompe, fiind frecvente pompele pentru ioni.

Structura moleculelor ce formează pompele este mai complexă decât a canalelor, ele având și situsuri de cuplare a moleculelor ce asigură suportul energetic. De obicei procesul începe cu o reacție de fosforilare, din partea unei molecule de ATP.

B. În figura 2.12.3 este prezentată pompa Na⁺/K⁺, întâlnită cel mai frecvent în celule pentru a asigura concentrațiile normale intra și extra celulare de Na⁺ și K⁺.

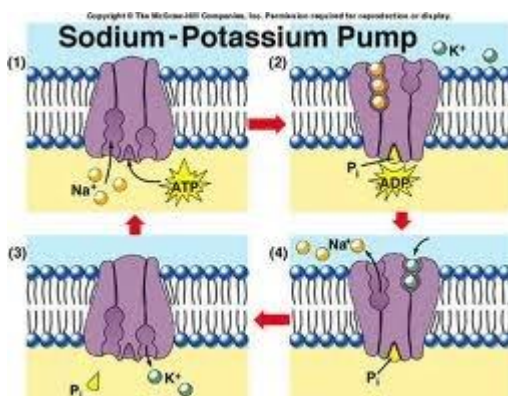


Fig. 2.12.3. Pompa Na⁺/K⁺

3. Noțiuni de biochimie

În acest capitol vom trece în revistă structura și proprietățile a două mari clase de molecule ce constituie de fapt obiect de studiu pentru bioinformatică:

- proteinele, compuse din aminoacizi,
- acizi nucleici, compuși din nucleotide.

3.1. Aminoacizii

3.1.1. Structură generală

A. Aminoacizii (AA) sunt molecule organice care conțin o grupare amino ($-\text{NH}_2$) și o grupare acidă carboxil ($-\text{COOH}$).

B. Formula generală este $\text{H}_2\text{N}-\text{CH}-\text{R}-\text{COOH}$. Atomul de carbon primar (numit și C_α) are cele 4 valențe ocupate astfel: una leagă gruparea amino, a doua leagă gruparea carboxil, a treia leagă un proton ($-\text{H}$) iar a patra leagă o grupare numită radical sau reziduu ($-\text{R}$). Aminoacizii diferă între ei prin acest radical (figura 3.1.1.a).

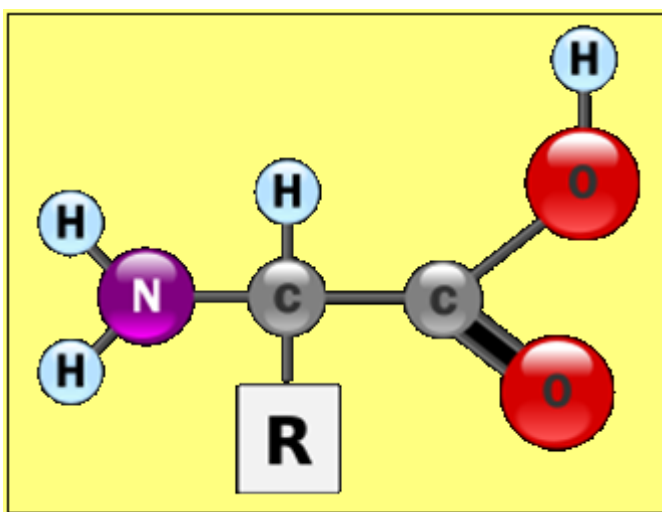


Fig. 3.1.1.a. Formula generală a aminoacizilor

C. În materia vie întâlnim doar 20 AA, aceiași atât în regnul vegetal cât și animal, la toată scara evolutivă, demonstrând elocvent unitatea materiei vii pe pământ. Însă numărul practic infinit al combinațiilor posibile, acoperă cu prisosință orice paletă imaginativă, explicând toate diversitățile, la fel cum cu numărul restrâns al literelor alfabetului se pot crea orice cuvinte.

D. În figura 3.1.1.b sunt prezentate structurile moleculare ale aminoacizilor în forma lor uzuală.

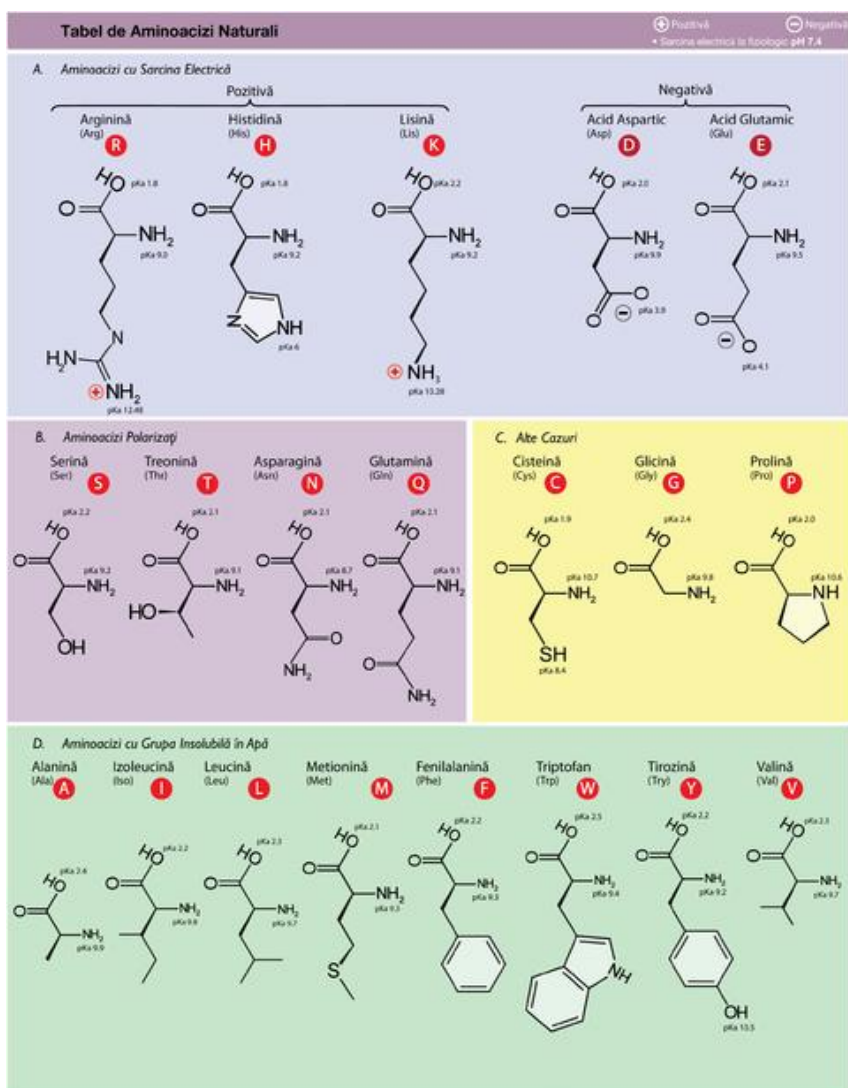


Fig. 3.1.1.b. Structurile moleculare ale aminoacizilor

3.1.2. Formele ionice ale aminoacizilor

A. Gruparea carboxil are caracter acid, putând elibera un proton și trecând în forma $-\text{COO}^-$, cu sarcină negativă.

B. Gruparea amino are ca caracter bazic, putând accepta un proton și trecând în forma $-\text{NH}_3^+$, cu sarcină pozitivă.

C. Având caracter dublu, atât acid cât și alcalin, spunem că are un caracter amfoter.

D. Formele ionice care se pot forma, fiind atât pozitive cât și negative, cu posibilitatea ca o moleculă să prezinte simultan ambele sarcini, poartă denumirea de “zwitter-ioni”, și sunt prezentate în figura 3.1.2.

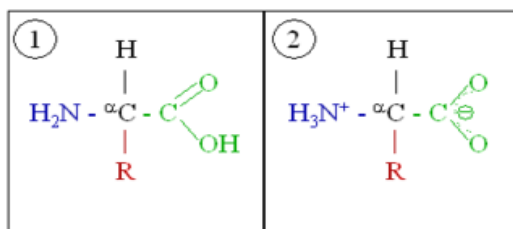


Fig. 3.1.2. Forme ionice ale aminoacizilor

E. Sarcina pe care o are un aminoacid la un moment dat depinde de pH-ul mediului. În mediu acid, unde avem o abundență de protoni H⁺, gruparea carboxil nu disociază, în schimb gruparea amino se va găsi în formă ionică, iar aminoacidul va avea sarcină pozitivă. În schimb în mediu bazic, gruparea amino va fi neutră, iar cea carboxil ionizată, având o sarcină globală negativă.

F. Pentru fiecare aminoacid există o anumită valoare a pH-ului, numită *punct izoelectric* pentru care gradul de disociere ca acid este egal cu cel de disociere ca bază, numărul sarcinilor pozitive este egal cu cel al sarcinilor negative, iar molecula în ansamblu este neutră.

G. Dependența sarcinii de pH stă la baza unei metode de separare a proteinelor numită *electroforeză*. O bandă de hârtie de filtru îmbibată în soluție salină are marginile în contact cu doi electrozi între care se aplică o diferență de potențial de ordinul sutelor de volți. Pe bandă se depune o picătură din soluția unei proteine. În prezența câmpului electric moleculele se deplasează (migreză) în funcție de sarcina lor. Viteza de deplasare este invers proporțională cu masa moleculară. Această deplasare datorată câmpului electric se numește *electroforeză*.

3.1.3. Proprietățile aminoacizilor

A. Abundență

Cei 20 AA nu au o răspândire uniformă în proteinele din natură. Abundența lor (în procente) este prezentată în tabelul 3.1.3.

B. Izomerie structurală

Dacă se trasează un plan imaginar creat de carbonul α, carbonul din carboxil și azotul din amino, atunci pentru fiecare AA avem două posibilități de poziționare a radicalului R - fie de partea stângă, fie de partea dreaptă a planului, cu alte cuvinte avem 2 izomeri sterici, care poartă numele S și R, iar proprietatea de a prezenta izomeri sterici se numește *chiralitate*.

Aminoacizii din proteinele din natură sunt cu toții izomeri S, acceptând cisteina (Cys) care este R și glicina (Gly) care este nonchirală.

C. Izomerie optică

Compușii care prezintă izomeri structurali, prezintă de obicei și izomerie optică. O serie de substanțe (numite „optic active”) au proprietatea de a roti planul luminii polarizate când această lumină le traversează. În funcție de sensul în care este rotit planul luminii polarizate substanțele pot fi D (dextrogire) - care rotesc planul la dreapta și L (levogire) care rotesc planul la stânga. În aminoacizii de sinteză cele două forme se produc în proporții egale. În aminoacizii naturali întâlnim numai forma L. Fenomenul încă nu are o explicație unanim acceptată.

Tabel 3.1.3. Răspândirea aminoacizilor

Denumirea (Residue)	cod 3-litere	cod 1 literă code	Abundență />(%) E.C.
Alanină	ALA	A	13.0
Arginină	ARG	R	5.3
Asparagină	ASN	N	9.9
Aspartat	ASP	D	9.9
Cisteină	CYS	C	1.8
Acid glutamic	GLU	E	10.8
Glutamină	GLN	Q	10.8
Glicină	GLY	G	7.8
Histidină	HIS	H	0.7
Isoleucină	ILE	I	4.4
Leucină	LEU	L	7.8
Lizină	LYS	K	7.0
Metionină	MET	M	3.8
Fenilalanină	PHE	F	3.3
Prolină	PRO	P	4.6
Serină	SER	S	6.0
Treonină	THR	T	4.6
Triptofan	TRP	W	1.0
Tirosină	TYR	Y	2.2
Valină	VAL	V	6.0

D. Aminoacizii esențiali

Majoritatea aminoacizilor pot fi sintetizați în organismul uman. Totuși există unii AA care nu pot fi sintetizați – aceștia poartă numele de AA esențiali. Ei sunt în număr de 9: Val, Leu, Ile, Lys, Thr, Met, His, Trp, Phe. AA sunt sintetizați de plante și sunt necesari în aportul alimentar, ei se găsesc și în produse animale.

E. Notății

Se folosesc frecvent niște notații prescurtate, formate din 3 litere convenabile pentru recunoașterea denumirii lor. Totuși în analiza secvențială este preferată notația simbolică cu câte 1 literă.

3.1.4. Hidrofobicitatea

A. Grupările din radicali R sunt foarte variate, cu proprietăți diferite, ceea ce conferă aminoacizilor – și prin ei, proteinelor – o paletă largă de proprietăți. În funcție de aceste grupări putem face o clasificare a AA (figura 3.1.4.a)

- AA cu grupări nepolare: Gly, Ala, Val, Leu, Ile, Phe,
- AA cu grupări polare neîncărcate electric: Ser, Thr, Cys, Met, Asn, Gln, Tyr, Trp, Pro,
- AA cu grupări polare cu caracter acid: Asp, Glu,
- AA cu grupări polare cu caracter bazic: Arg, Lys.

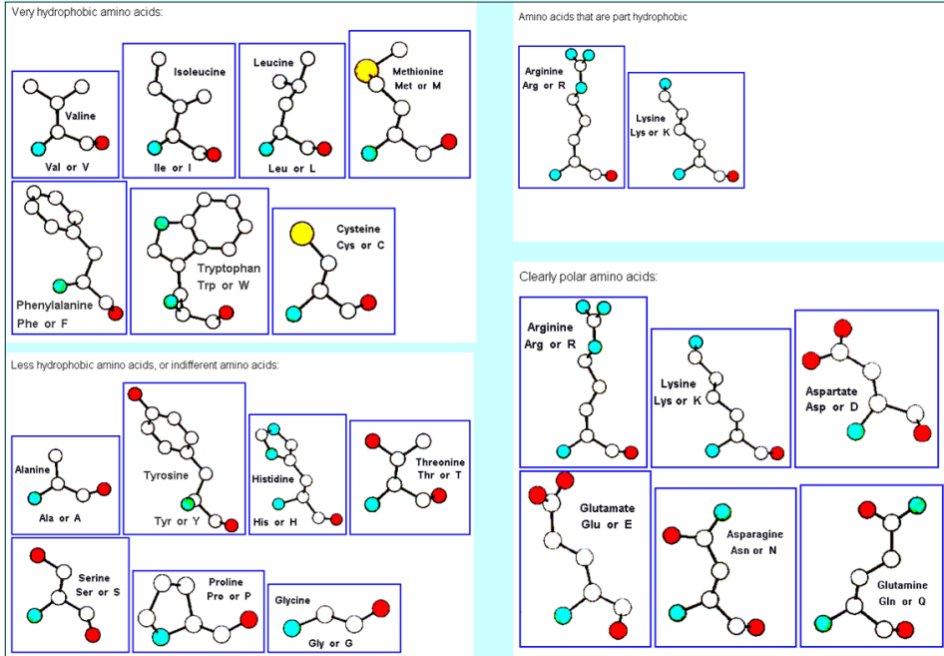


Fig. 3.1.4.a. Clasificarea aminoacizilor

B. Polaritatea grupărilor determină și solubilitatea lor în apă. Astfel AA cu grupări polare sunt hidrofobi, cei cu sarcină sunt hidrofilii, iar cei cu grupări polare neîncărcate electric sunt relativ solubili prin punțile de hidrogen care se pot forma.

C. Scări de hidrofobicitate

Evaluarea gradului de hidrofobicitate al unei molecule este dificilă, fiind cel mai adesea bazată pe partiția moleculei între diferiți solvenți. În funcție de solvenții folosiți s-au stabilit diverse scări care ierarhizează aminoacizii figura 3.1.4.b, dreapta.

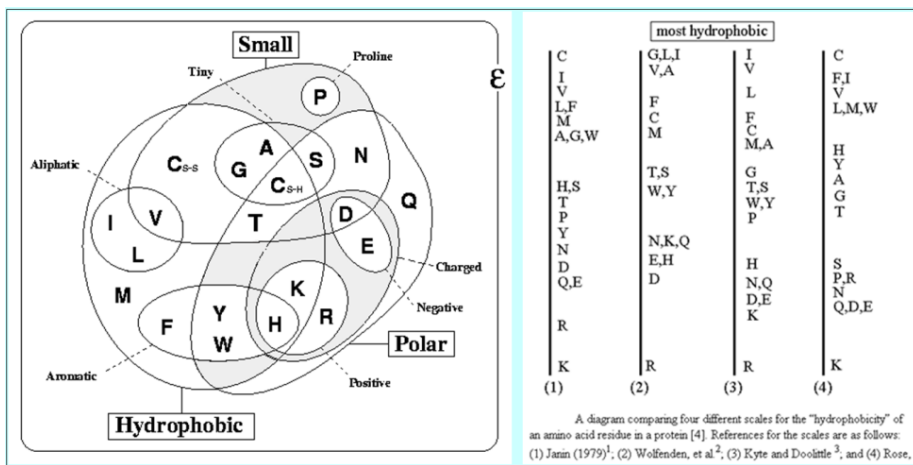


Fig. 3.1.4.b. Hidrofobicitatea aminoacizilor

D. Dacă ținem cont de dimensiunea moleculei, pe lângă hidrofobicitate și caracterul polar/nepolar sau sarcina moleculei la pH fiziologic, putem prezenta schematic o diagramă care sintetizează toate aceste proprietăți figura 3.1.4.b stânga.

E. Pentru a avea o imagine mai detaliată asupra proprietăților aminoacizilor, să mai adăugăm determinările unor proprietăți fizice: densitate, volum, suprafață, constante de echilibru de disociere a grupărilor din radical (unde este cazul). Aceste proprietăți sunt sintetizate în tabelul 3.1.4.

Tabelul 3.1.4. Proprietățile aminoacizilor solubilitate, densitate, volum, suprafață și pK

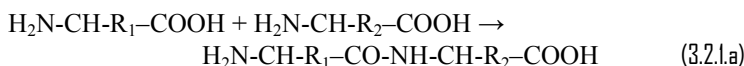
Name	Solubility (g/100g, 25 °C)	Crystal Density (g/ml)	pI at 25 °C	Residue Volume	Surface Area	Side Chain pKa
Alanine	16.65	1.401	6.107	88.6	115	-
Arginine	15	1.1	10.76	173.4	225	~12
Aspartic Acid	0.778	1.66	2.98	111.1	150	4.5
Asparagine	3.53	1.54	-	114.1	160	-
Cysteine	very	-	5.02	108.5	135	9.1-9.5
Glutamic Acid	0.864	1.460	3.08	138.4	190	4.6
Glutamine	2.5	-	-	143.8	180	-
Glycine	24.99	1.607	6.064	60.1	75	-
Histidine	4.19	-	7.64	153.2	195	6.2
Isoleucine	4.117	-	6.038	166.7	175	-
Leucine	2.426	1.191	6.036	166.7	170	-
Lysine	very	-	9.47	168.6	200	10.4
Methionine	3.381	1.340	5.74	162.9	185	-
Phenylalanine	2.965	-	5.91	189.9	210	-
Proline	162.3	-	6.3	112.7	145	-
Serine	5.023	1.537	5.68	89.0	115	-
Threonine	very	-	-	116.1	140	-
Tryptophan	1.136	-	5.88	227.8	255	-
Tyrosine	0.0453	1.456	5.63	193.6	230	9.7
Valine	8.85	1.230	6.002	140.0	155	-

a-amino pKa = 6.8 - 7.9, a-carboxyl pKa = 3.5 - 4.3

3.2. Proteine

3.2.1. Legătura peptidică

A. Doi aminoacizi se pot lega chimic prin reacția între gruparea carboxil a unui AA cu gruparea amino a celui alt, cu eliminarea unei molecule de apă:



B. Legătura formată se numește legătură peptidică. Reacția este bine ilustrată în figura 3.2.1.a, iar proprietățile legăturii peptidice sunt prezentate în figura 3.2.1.b.

C. Ne putem ușor imagina că dipeptidul format, având un capăt amino și un capăt carboxil, poate la rândul său forma o nouă legătură peptidică și așa mai departe. O astfel

de moleculă se numește polipeptid sau proteină. Pentru un număr mic de AA în structură se pot folosi prefixele corespunzătoare: dipeptid, tri~, tetra~, penta~ ș.a.m.d.

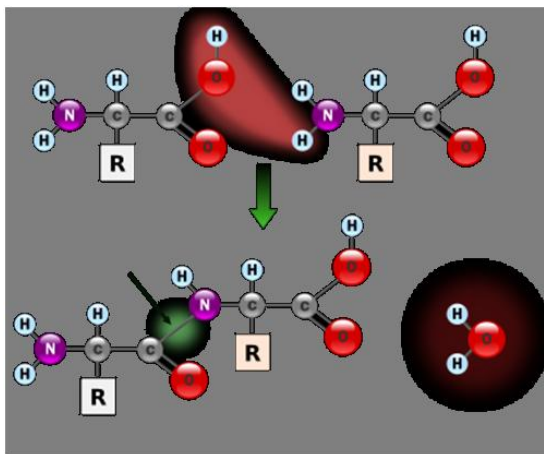


Fig. 3.2.1.a. Formarea legăturii peptidice

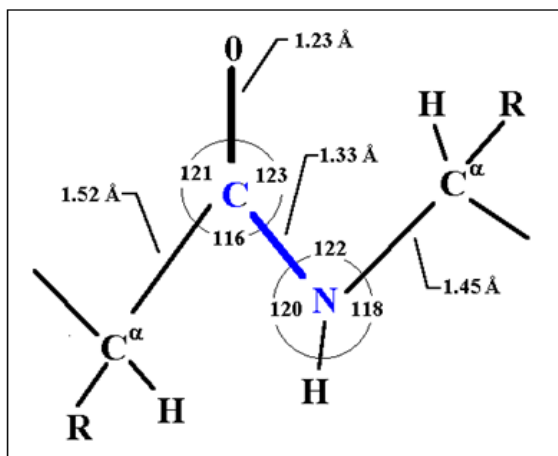


Fig. 3.2.1.b. Proprietățile legăturii peptidice

3.2.2. Structura primară a proteinelor

A. O primă caracteristică, de importanță majoră în bioinformatică este reprezentată de succesiunea de AA care formează molecula proteică.

B. Această succesiune poartă numele de „structură primară” a proteinelor. Pozițiile aminoacizilor se numerotează începând de la capătul cu gruparea amino liberă.

C. Determinarea structurii primare a proteinelor

Există mai multe metode de determinare a structurii primare a proteinelor, pe care nu le vom trata în detaliu ci doar le vom enumera aici:

- a) metode biochimice de determinare a secvenței, cuprinzând:
 - degradarea Edman,
 - fracționarea prin digestie cu enzime proteolitice,

- spectrometrie de masă,
- hidroliză în acizi (doar determinarea cantităților de aminoacizi).
 - b) deducerea din secvența genelor.

D. Exemple de structuri primare

Una dintre primele structuri primare determinate a fost cea a insulinei, prezentată în figura 3.2.2, de către Sanger în 1951, pentru care i s-a acordat premiul Nobel.

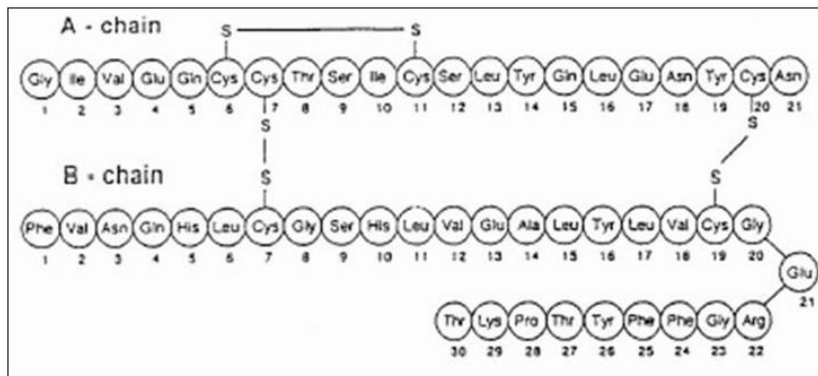


Fig. 3.2.2. Structura primară a proteinelor: insulina

3.2.3. Structura secundară a proteinelor

A. O moleculă proteică poate avea de la câteva zeci de aminoacizi până la câteva sute. Prin proprietățile pe care le au radicalii aminoacizilor din lanț, se formează adesea niște structuri spațiale specifice ușor de recunoscut. Corey și Pauling (1943) au fost primii care au semnalat aceste structuri denumite generic „structura secundară” a proteinelor.

B. Structura – helix

Cea mai bine cunoscută formă specifică este cea de helix pe care o întâlnim frecvent la proteine. Unii aminoacizi, mai ales Met, Ala, Leu, Glu și Lys se dispun în forma unor spirale ce s-ar înfășura pe un cilindru virtual, numit - helix. Alți AA însă nu se dispun în helix, apariția lor (de exemplu Gly și Pro) generând discontinuitatea spiralei, de aceea se și numesc ”helix breakers”. Procentul de AA angajați în structurile secundare variază. În tabelul 3.2.3. sunt trecute proporțiile de AA cuprinși în structuri de tip - helix.

Tabelul 3.2.3. Procent de AA în elice

Proteina	% elice α
Mioglobina	70
Insulina	38
Ovalbumina	31
Serumalbumina	46
Pepsina	31
Ribonucleaza	16
Chimotripsina	15

C. Să trecem în revistă câteva proprietăți ale unui α – helix (figura 3.2.3.a)

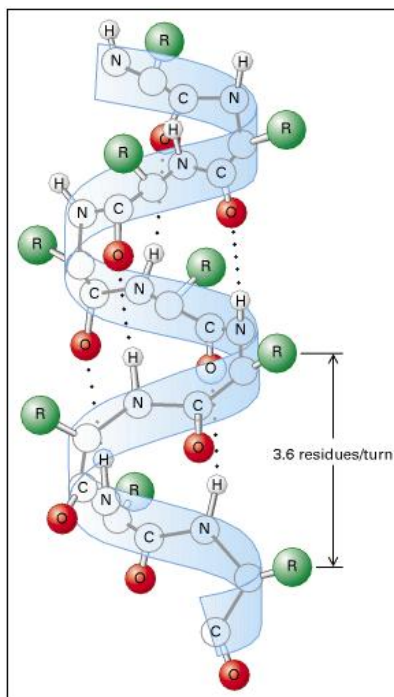


Fig. 3.2.3.a. Structura α helix a proteinelor

- o spiră cuprinde 3,6 AA,
- pasul spirei 5,21 Å,
- forma spiralei ca un șurub cu pasul spre dreapta,
- catenele (radicalii R) sunt orientate spre exteriorul cilindrului.

D. Determinarea structurii secundare

Cea mai potrivită metodă este difracția cu raze X, (utilizată de Perutz și Kendrew pentru structura mioglobinei), radiațiile X având lungimea de undă de ordinul de mărime al distanțelor interatomice (\sim Å). O limitare a metodei este că poate fi folosită numai pe cristale, deci moleculele care nu pot fi aduse în formă cristalină nu pot fi studiate prin această metodă.

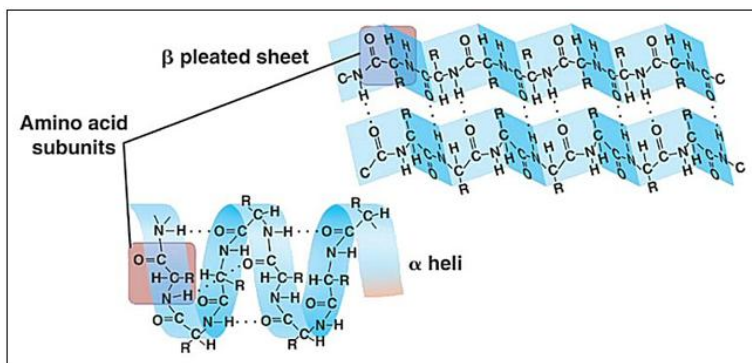


Fig. 3.2.3.b. Structura pliuri β a proteinelor

E. Fâșiile β

O altă formă întâlnită ca structură secundară este reprezentată de „fâșiile β ” sau „pliurile β ” (figura 3.2.3.b). Alternanța unor AA care generează alternanța orientării valențelor duce la formarea unor regiuni cu formă specifică de panglică pliată. Aceste fâșii apar în special în regiunile cu AA aromatici (Trp, Tyr și Phe). De asemenea, forma β -C este generată în zonele cu Ile, Val, Thr.

3.2.4. Structura terțiară a proteinelor

A. Pe lângă structura secundară, între diferitele părți ale unei molecule proteice pot apărea interacțiuni:

- legături între părți prin punți de hidrogen, legături disulfidice sau legături Van der Waals
- plieri, încovoieri.

Toate acestea generează o structură tridimensională complexă.

3.2.5. Structura cuaternară

A. În cazul moleculelor foarte mari se pot realiza chiar legături între mai multe unități, formând așa numita structură cuaternară. Procesul acesta este catalizat de enzime specifice numite „holoenzime”.

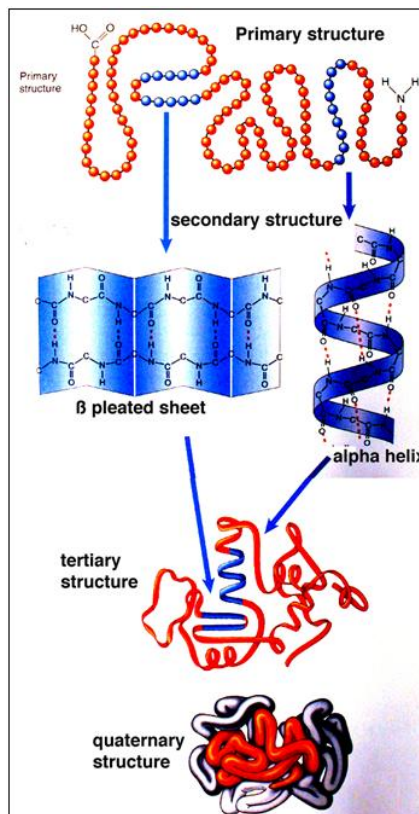


Fig. 3.2.5. Structura proteinelor – viziune de ansamblu

B. Exemple

Structură cuaternară întâlnim în special la proteinele globulare. Un exemplu: hemoglobina, care are 4 subunități, „molecula” fiind de fapt un tetramer.

În aceeași clasă mai putem include ADN polimeraza precum și proteinele componente ale canalelor ionice.

În figura 3.2.5 este prezentată o sinteză privind structura proteinelor.

3.3. Componentele acizilor nucleici

3.3.1. Componentele unui nucleotid

Acizii nucleici au și ei o structură secvențială, asemănătoare întrucâtva cu cea a proteinelor. Unitatea de secvență, echivalentă aminoacizilor, se numește „nucleotid”. Un nucleotid are 3 componente principale:

- o pentoză,
- un acid fosforic,
- o bază azotată.

3.3.2. Pentozele din acizii nucleici

A. În acizii nucleici întâlnim două tipuri de pentoze:

- riboza - în acidul ribonucleic (ARN),
- de(z)oxiriboza - în acidul de(z)oxiribonucleic ADN.

B. Structura chimică și numerotarea atomilor sunt prezentate în figura (3.2.2).

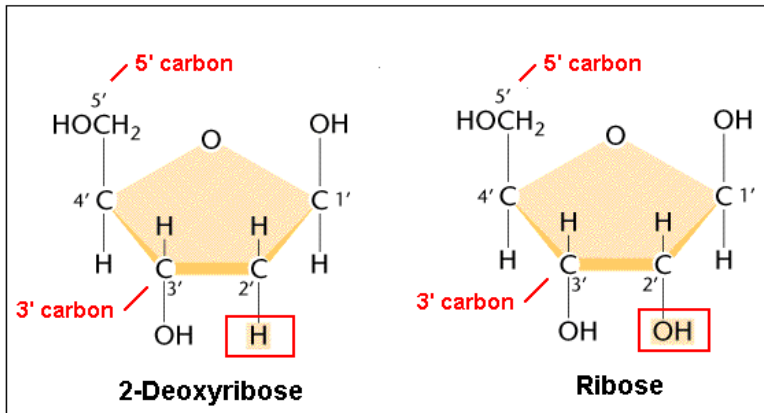


Fig. 3.2.2. Structura ribozei și dezoxiribozei

Observăm că atomii C_{1'} - C_{4'} și oxigenul formează un plan, în timp ce carbonul C_{5'} este înafara planului.

De asemenea, gruparea -OH de la carbonul 1' este de aceeași parte cu carbonul 5', în timp ce grupările -OH din pozițiile 2' și 3' sunt de cealaltă parte a planului.

C. În moleculele de dezoxiriboză lipsește gruparea -OH din poziția 2'.

3.3.3. Bazele azotate

A. Fiecare tip de acid nucleic conține 4 tipuri de baze azotate: două baze purinice și două baze pirimidinice.

Bazele purinice sunt adenina (A) și guanina (G), iar cele pirimidinice sunt citozina (C) și Timina (T) sau uracilul (U). În timp ce bazele A, G și C se găsesc și în ADN și în ARN, cea de a patra este diferită: în ADN întâlnim timina, în timp ce în ARN întâlnim uracilul.

B. Numerotarea atomilor din bazele azotate se face numai pentru atomii din cicluri și este redată împreună cu structura moleculelor de purină și pirimidină, precum și bazele azotate derivate: A și G, respectiv C, T și U în figura 3.3.3.

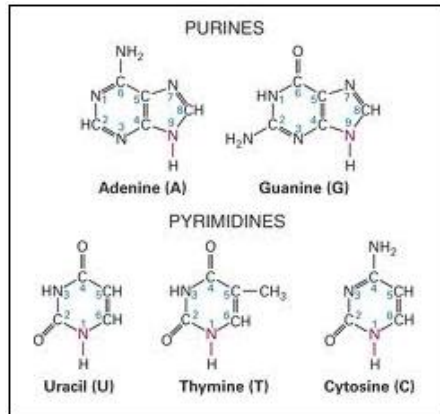


Fig. 3.3.3. Bazele azotate din acizii nucleici

3.3.4. Nucleozide

A. Bazele azotate se leagă de pentoze printr-o legătură glicozidică, formând o structură numită „nucleozid”. Denumirea nucleozidelor pornește de la numele bazei azotate; de ex.: de la adenină se formează adenozina (prin cuplarea de riboză) sau dezoxiadenozina (prin legarea de dezoxiriboză); similar avem guanină/guanozină, dezoxiguanozină, citozină/citidină/dezoxicitidină, timină/timidină (cu dezoxiriboză), respectiv uracil/uridină (cu riboză).

Legarea bazelor azotate se face întotdeauna la carbonul 1 al pentozei. Bazele purinice se leagă prin azotul din poziția 9, iar cele pirimidinice prin azotul din poziția 1.

B. Structura nucleozidelor este prezentată în figura 3.3.4.

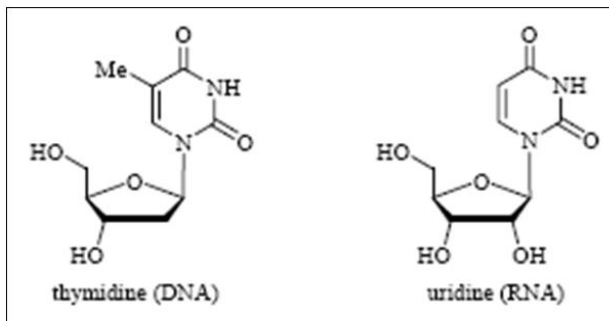


Fig. 3.3.4. Structura nucleozidelor

3.3.5. Nucleotide

Pentru formarea lanțurilor acizilor nucleici, nucleozidele sunt fosforilate. Cuplarea moleculei de acid fosforic se face la carbonul 5' al pentozei, structura nou formată fiind numită „nucleotid”. Deși nucleotidele, ca entități separate au denumiri specifice (acid adenilic, guanilic, citidilic sau timidilic, uneori sub formă de radical adenilat, guanilat, citidilat, timidilat) în cazul precizării unui nucleotid într-o secvență de acid nucleic, se folosește doar numele bazei azotate din componența sa. În figura 3.3.5 este redată structura nucleotidelor.

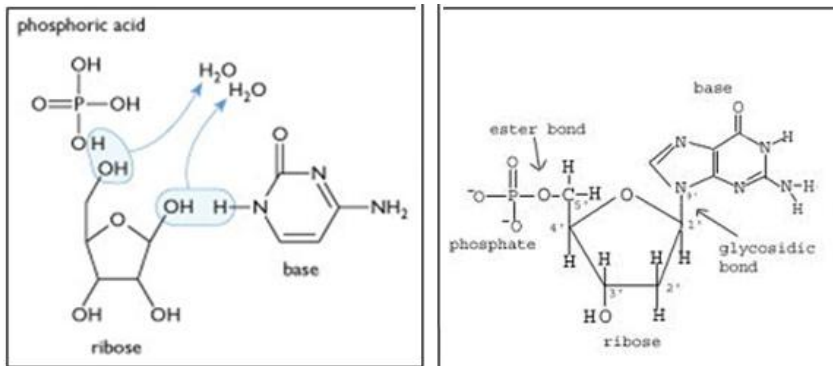


Fig. 3.3.5. Structura nucleotidelor

3.4. Structura acizilor nucleici

3.4.1. Formarea lanțului polinucleotidic

A. Nucleotidele se pot lega unele de altele prin eliminarea unei molecule de apă între o grupare - OH a fosforului legat de carbonul 5' și gruparea OH a carbonului 3' din pentoză (figura 3.4.1.a).

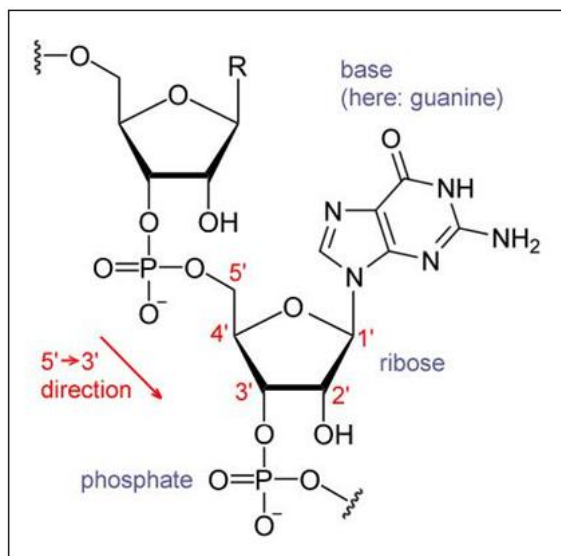


Fig. 3.4.1.a. Formarea unui dinucleotid

B. Prin repetarea procesului se poate obține un lanț care să cuprindă de la câteva zeci până la câteva mii de nucleotide (figura 3.4.1.b).

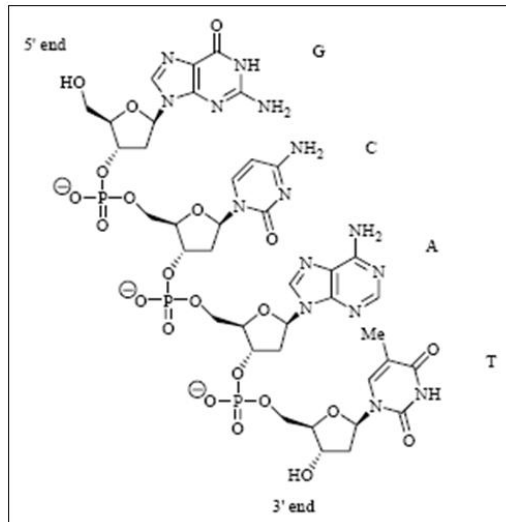


Fig. 3.4.1.b. Formarea lanțului polinucleotidic GCAT

3.4.2. Structura primară a acizilor nucleici

A. Toate lanțurile de acizi nucleici au două capete libere, unul la atomul 5', altul la atomul 3'. O secvență de lanț al unui acid nucleic se citește în ordinea 5' → 3'

B. Secvența bazelor azotate într-un lanț de acid nucleic poartă denumirea de structură primară a acizilor nucleici (figura 3.4.2).

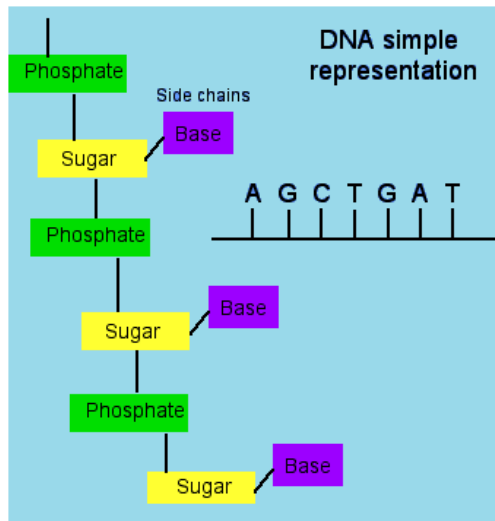


Fig. 3.4.2. Structura primară a unei secvențe ADN

3.4.3. Structura secundară a acizilor nucleici

A. Descifrarea structurii acizilor nucleici a reprezentat un pas important în înțelegerea unui întreg șir de fenomene din biologie în general și din genetică în special. Acumularea unor date experimentale privind constanța rapoartelor A/T și G/C, sau considerente teoretice privind structura tridimensională a moleculelor, au condus în cele din urmă la stabilirea modelului „dublu-helix” al moleculei de ADN de către Watson și Crick în 1953, laureați ai premiului Nobel.

B. Moleculele de ADN formează o dublă elice înfășurată pe un cilindru virtual, având orientate spre interior bazele azotate. Între două baze azotate din cele două lanțuri se stabilesc niște punți de hidrogen. Prin configurația spațială a norului electronic sunt posibile doar perechi între o adenină de pe un lanț cu o timină pe lanțul opus (prin 2 punți de hidrogen) sau între o guanină și o citozină (prin 3 punți de hidrogen), așa cum se poate vedea în figura 3.4.3.a.

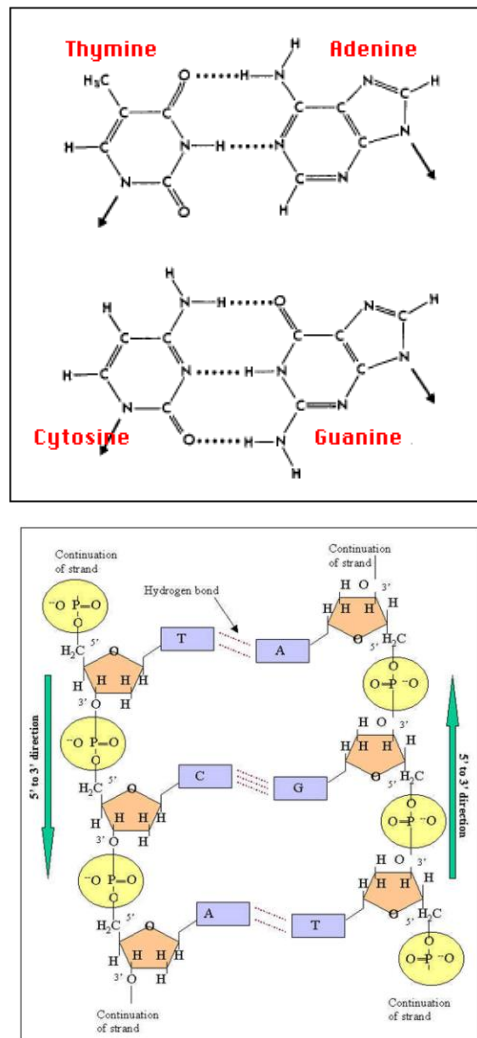


Fig. 3.4.3.a. Perechile de baze azotate în ADN

4. Elemente de biologie celulară și moleculară

O caracteristică esențială a materiei vii este reprezentată de structura sa celulară, alături de încă două proprietăți fundamentale: metabolismul și reproducerea.

Deși există mari diferențe între celulele diverselor specii, sau chiar între celulele aceluiași organism, există și o serie de caracteristici comune pe care le vom trece în revistă.

Disciplina care se ocupă cu studiul celulei și a fenomenelor la nivelul celular se numește „biologie celulară”. Din biologia celulară vom extrage sintetic câteva noțiuni privind:

- structura celulei umane (membrană, citoplasmă, nucleu, organite celulare – mitocondria și ribozomii),
- diviziunea celulară mitoza și meioza.

Explicarea fenomenelor biologice pe baze moleculare face obiectul unei discipline înrudite cu biologia celulară, numită „biologie moleculară”. Din această disciplină ne vom opri la două mecanisme:

- replicarea ADN,
- sinteza proteinelor.

4.1. Structura generală a unei celule umane

4.1.1. Enumerarea principalelor componente ale celulei umane

- a) membrana,
- b) citoplasma,
- c) nucleu,
- d) organite celulare.

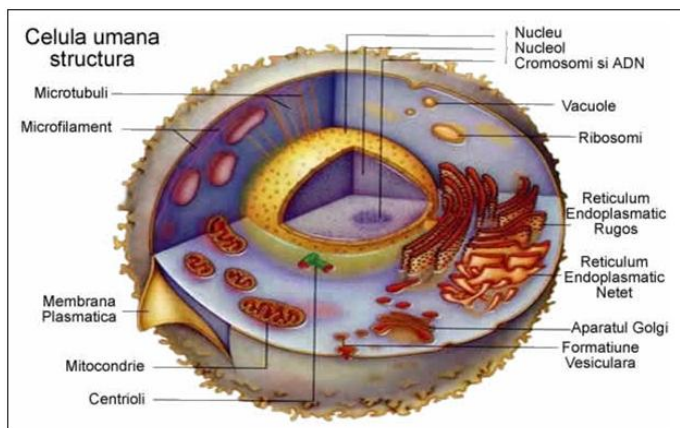


Fig. 4.1.1. Structura celulei umane

4.1.2. Membrana celulară

A. Membrana celulară delimitează celula la exterior, are rol de susținere și separare, fiind sediul unor procese de schimb de substanțe între celulă și mediul său extern.

B. Membrana celulară este alcătuită dintr-un strat dublu fosfolipidic, conținând în forme și proporții diferite o serie de structuri proteice (modelul „mozaic lichid” Singer-Nicholson, 1972), (figura 4.1.2.a).

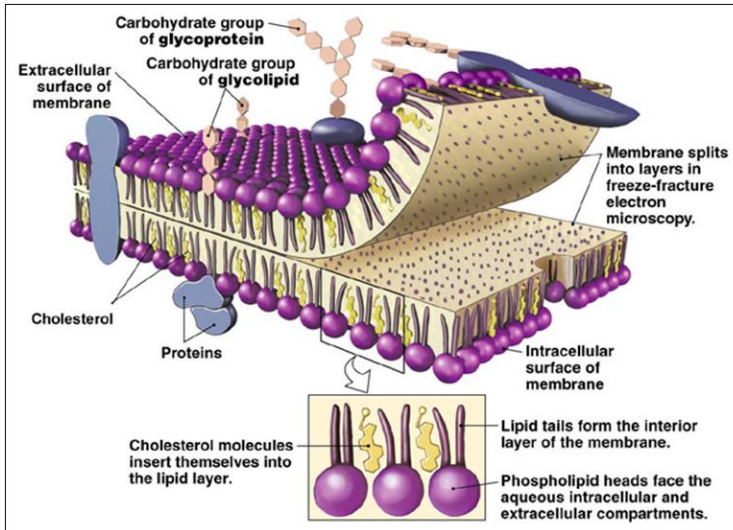


Fig. 4.1.2.a. Membrana celulară-modelul „mozaic lichid”

C. Proprietăți fizice

Membrana celulară are o grosime de 3-10 nm, cu slabă conductibilitate electrică, dar cu mare capacitate electrică, greu permeabilă pentru apă și ioni.

D. Matricea membranei celulare este alcătuită din fosfolipide ce formează un strat dublu.

Fosfolipidele au caracter amfipatic, având un capăt hidrofil și un capăt hidrofob (figura 4.1.2.b).

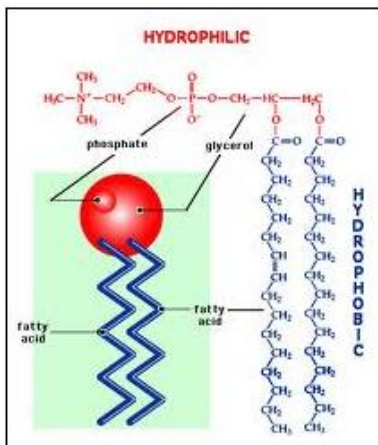


Fig. 4.1.2.b. Structura generală a unui fosfolipid

Acizii grași au 16-18 atomi de carbon în lanț și, având caracter hidrofob, se orientează spre interiorul membranei, în timp ce restul fosforic, având grupări OH este hidrofil, fiind orientat spre fețele exterioare ale membranei, fie spre interiorul celulei, fie spre exterior.

E. Proteinele membranare formează „insule” în matricea fosfolipidică. Ele pot fi:

- a) proteine intrinseci, care pot să fie la rândul lor:
- proteine ce traversează membrana (canale ionice pentru transport pasiv sau pompe ionice pentru transport activ), figura 4.1.2.c,
 - proteine parțial înglobate, ce formează structuri numite „receptori membranari”, cu rol foarte important în funcționarea membranei (de receptori se pot lega diferite molecule-hormoni, medicamente etc.).

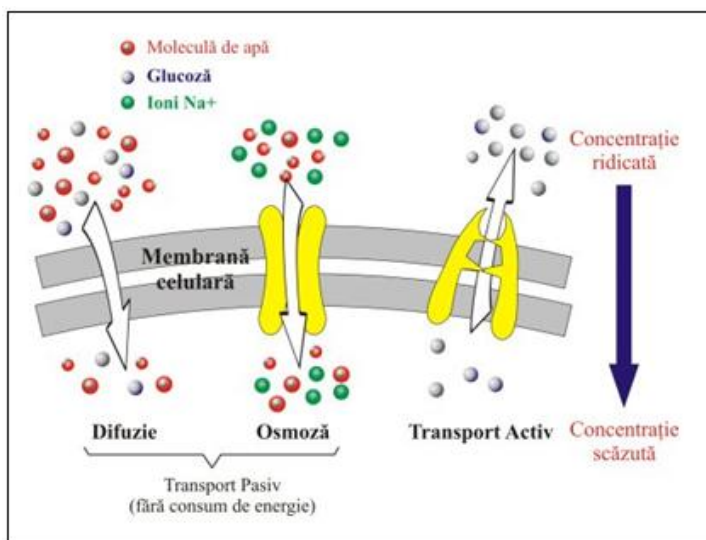


Figura 4.1.2.c. Canale și pompe ionice

b) Proteine extrinseci, legate slab la suprafața membranei.

F. Membrana celulară este totodată sediul fenomenelor electrice, de producere a potențialului de membrană.

G. Membrana are un rol important în semnalizarea inter și intracelulară privind transferul de informații în aceste procese de semnalizare sunt implicați și receptori membranari.

H. Să încheiem cu rolul membranei în apoptoză - fenomenul de moarte naturală programată.

4.1.3. Citoplasma

A. Citoplasma este o soluție apoasă ce reprezintă mediul de dispersie pentru substanțele dizolvate:

- ioni: [K⁺] 150 mM, [Na⁺] 10 mM, [Cl⁻] 100mM,
- molecule mici neutre [glucoză],
- macromolecule (proteine, polizaharide etc).

B. Citoplasma este străbătută de fire ce se sprijină pe membrană, nucleu și organelle celulare mari, ce formează un citoschelet.

C. Proprietățile citoplasmei sunt sintetizate în tabelul 4.1.3.

Tabel 4.1.3. Proprietățile fizice ale citoplasmei

Ion	Concentration in cytosol (millimolar)	Concentration in blood (millimolar)
Potassium	139	4
Sodium	12	145
Chloride	4	116
Bicarbonate	12	29
Amino acids in proteins	138	9
Magnesium	0.8	1.5
Calcium	<0.0002	1.8

4.1.4. Nucleul celular

A. Nucleul celular este o formațiune sferică sau ovoidală, cu diametrul de 4-6 μm , alcătuit în cea mai mare parte din ADN.

B. Structura nucleului (figura 4.1.4.a)

În interior se găsește o formațiune mai densă numită „nucleol”, iar la exterior este mărginit de o membrană nucleară, cu pori prin care pot trece molecule mici. Transportul prin membrana nucleară (pentru ARN și proteine) este asigurat de molecule numite „cargo-GTP-aze.”

C. Nucleul este sediul de stocare a informației genetice și joacă un rol important în procesul de diviziune celulară și în procesul de sinteză a proteinelor.

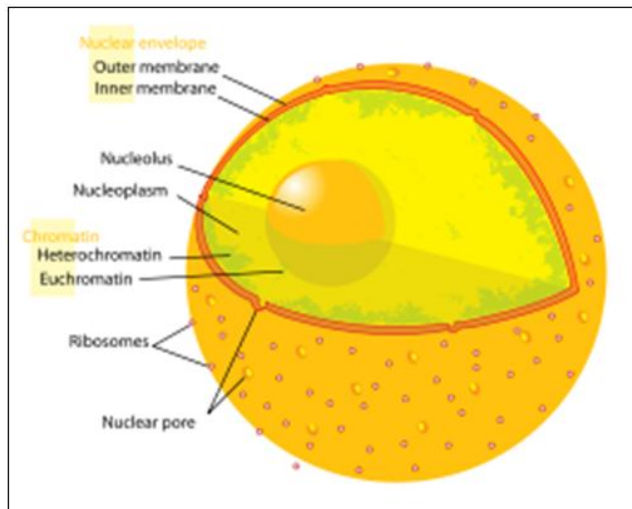


Fig. 4.1.4. Structura nucleului

4.1.5. Organitele celulare

În citoplasmă se găsesc numeroase organite celulare, fiecare cu un rol bine definit în funcționarea celulei:

- reticulul endoplasmic,
- ribozomi,
- mitocondrii,
- aparatul Golgi,

- lizozomi,
- centrioli,
- microtubuli, microfilamente,
- vacuole, vezicule.

Ne vom opri pe scurt doar la ribozomi și mitocondrii, având în vedere importanța lor în studiile de bioinformatică.

4.1.6. Mitocondria

A. Mitocondria este un organit celular care este sediul sintezei moleculelor macroergice de ATP.

B. Are formă alungită având în interior „creste” formate prin pliarea membranei interne. Membrana internă conține molecule ce asigură un transfer de electroni (NADH, FAD, ubiquinonă și citocrom c), însoțit de eliminare de protoni. Procesul se numește „fosforilare oxidativă”, având ca rezultat sinteza unei molecule de ATP. Funcționarea este asigurată de o pompă de protoni descrisă sub numele de ”teoria chemiosmotică” de către Mitchell, laureat al premiului Nobel în 1978.

Schematic, transferul electronilor în membrana internă a mitocondriei este prezentată în figura 4.1.6.a, iar principiul pompei de protoni este sintetizat în figura 4.1.6.b.

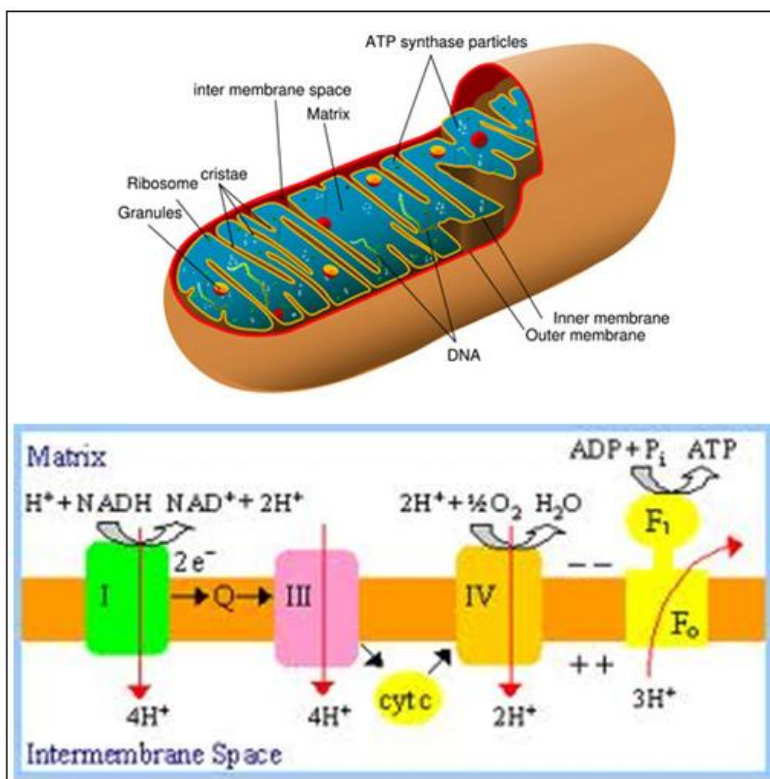


Fig. 4.1.6.a. Mecanismul fosforilării oxidative

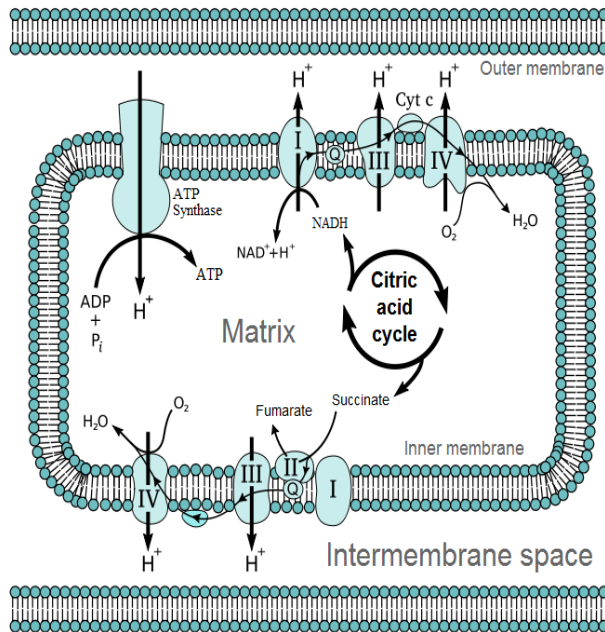


Fig. 4.1.6.b. Teoria chemiosmotică a pompei de protoni

4.1.7. Ribozomii

- A. Ribozomii sunt cele mai mici organite celulare, descoperite de către George Palade, de origine română, premiul Nobel în 1974.
- B. Sunt alcătuiți din ARN ribozomal și sunt sediul sintezei proteinelor.
- C. Un ribozom are trei situs-uri, notate E, P și A, în care se poate lega o moleculă de ARNt (ARN “de transport” sau „de transfer”, figura 4.1.7.)

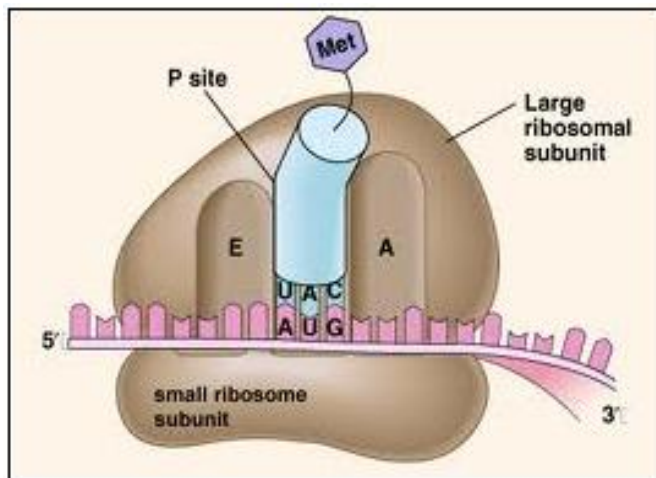


Fig. 4.1.7. Structura ribozomilor

4.2. Diviziunea celulară

Majoritatea celulelor umane (sunt doar unele excepții) se divid în cursul vieții lor. Procesul de diviziune se petrece diferit la celulele somatice, față de cele sexuale și vom trata separat cele două cazuri.

4.2.1. Mitoza

A. Mitoza este procesul de diviziune al celulelor somatice.

B. Celulele somatice sunt diploide – au fiecare cromozom în două exemplare. În cazul celulelor umane sunt 23 de perechi de cromozomi (unele detalii vor fi prezentate în cursul de genetică).

C. Ciclul de viață al unei celule somatice, adică intervalul de timp dintre două diviziuni, cuprinde mai multe faze, în care se sintetizează la început ARNm și proteine, apoi se sintetizează și ADN – practic se dublează cantitatea lui prin procesul de replicare în nucleu.

D. Diviziunea celulară ocupă cca 10% din ciclul celular și are mai multe faze:

- profaza – cromozomii devin vizibili, fiecare fiind dublat longitudinal, apar 2 centrioli ce migrează spre polii celulei, dispare membrana nucleară și se formează fusul de diviziune,
- metafaza – cromozomii se aranjează în regiunea ecuatorului și se leagă de fibrele fusului în zona centromerului,
- anafaza – deplasarea cromatidelor fiecărui cromozom spre polii celulei,
- telofaza – formarea la fiecare pol a câte unui nucleu, cu formarea de membrană nucleară.

În final se divizează și citoplasma formându-se două celule fiice cu nuclee identice.

Fenomenele descrise mai sus sunt prezentate sintetic în figura 4.2.1

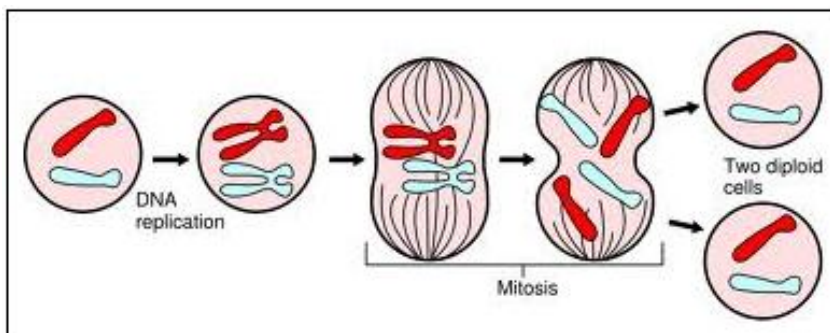


Fig. 4.2.1. Mitoza

4.2.2. Meioza

A. Diviziunea prin meioză este tipică pentru celulele germinale, ovocitul și spermatoцитul, în zona de maturizare ce formează celule sexuale, numite și „gameți”.

B. Meioza cuprinde două diviziuni succesive:

- a. diviziunea reduțională – când dintr-o celulă germinală diploidă se formează două celule haploide (cu fiecare cromozom într-un singur exemplar),

b. diviziunea eucariotă – asemănătoare mitozei, cu deosebirea că fiecare celulă care se divide acum este o celulă haploidă.

C. În final rezultă 4 celule haploide. În figura 4.2.2 este reprezentată în partea din dreapta meioza, comparativ cu mitoză reprezentată schematic în partea din stânga, fiind vizibile deosebirile din toate fazele diviziunii.

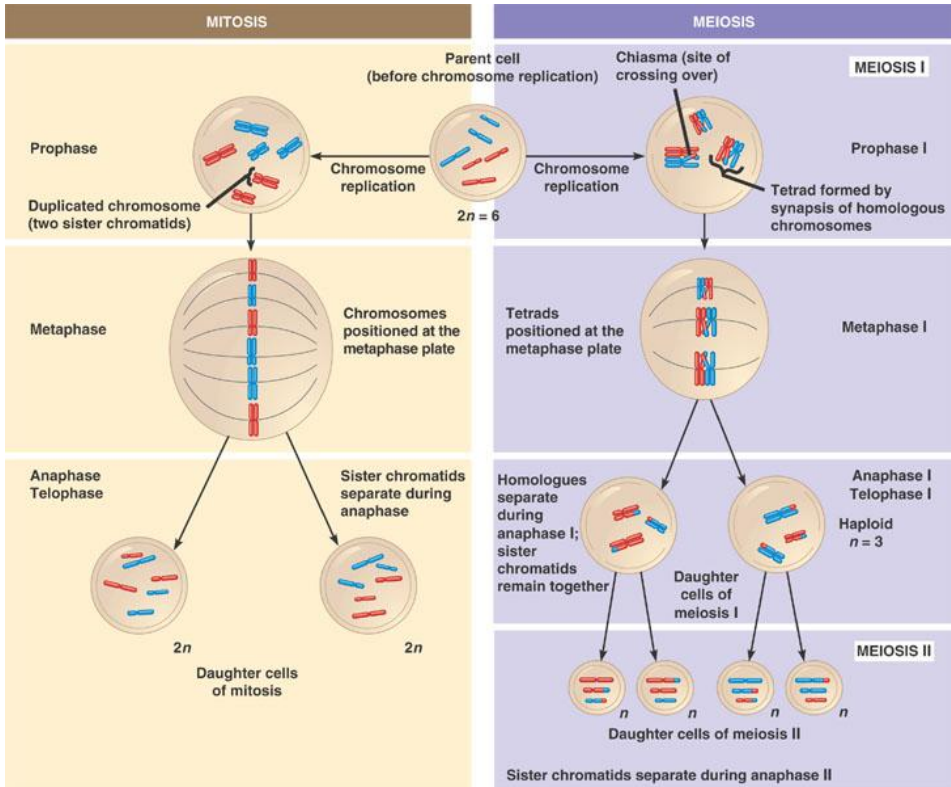


Fig. 4.2.2. Meioza

Să urmărim în continuare două procese esențiale din biologia moleculară: replicarea ADN și sinteza proteinelor.

4.3. Replicarea ADN

4.3.1. Noțiunea de replicare

A. Descifrarea mecanismului replicării ADN a avut o importanță deosebită prin aportul său la înțelegerea mecanismelor genetice de transmitere ereditară a informației.

B. Replicarea reprezintă procesul de producere a unei „replici”, adică o copie identică cu originalul. În cazul ADN este vorba de producere a două molecule de ADN identice pornind de la una singură.

4.3.2. Fazele replicării ADN

A. Replicarea este inițiată într-un punct al dublului helix, numit „origine” - o regiune cu secvență recunoscută de o proteină inițiatoare a replicării aceste regiuni sunt de obicei bogate în perechi AT, mai ușor de desfăcut.

B. Pornind de la origine, dublul helix este despiralat, cu o enzimă - topoisomeraza, apoi este desfăcut asemănător cu un fermoar, fiind rupte toate punțile de hidrogen pe o anumită porțiune din helix de către o enzimă - helicaza. Sunt astfel expuse liber spre exterior bazele azotate, în succesiunea lor, ale ambelor lanțuri ale helixului. Lanțul care este desfăcut de la 3'→5' încât completarea să apară în succesiunea firească 5'→3' se numește „lanț conducător” („leading strand”), iar celălalt este „lanț întârziat” („lagging strand”) (figura 4.3.2.)

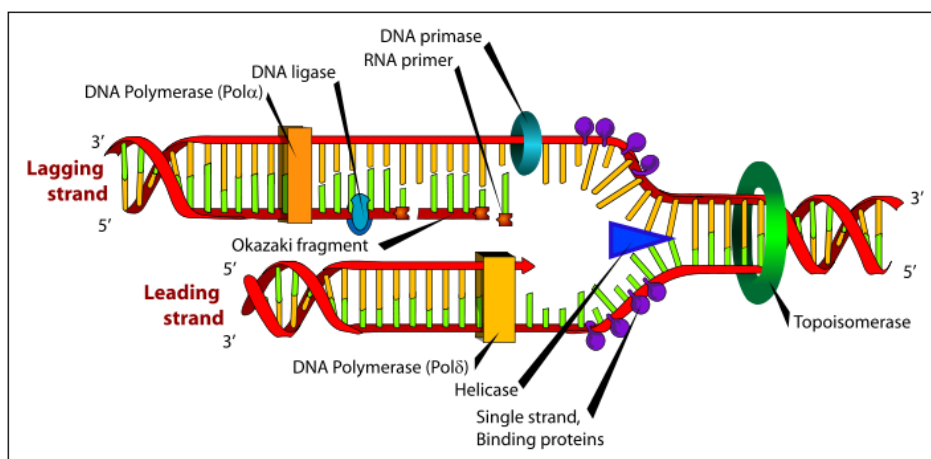


Fig. 4.3.2. Mecanismul replicării ADN

Forma de „furcă” a dublului helix desfăcut de helicază a generat și numele structurii în această fază: ”furcă de replicare” („replication fork”).

C. Pe lanțul conducător pozițiile expuse sunt ocupate succesiv de către nucleotidele corespunzătoare (conform perechilor posibile A-T și G-C), proces condus de o enzimă numită ADN-polimerază. Însă sinteza pe celălalt lanț nu poate fi realizată natural, succesiunea de sinteză fiind inversă, de la 3'→5'. De aceea, informația pentru un set de baze azotate libere este întâi transpusă în ordinea posibilă 5'→3' pe niște fragmente denumite Okazaki, fragmente ce sunt apoi cuplate pe lanț cu ajutorul enzimei ADN-ligaza.

D. Fiecare din lanțurile inițiale este acum completat cu un lanț pereche, care este spiralat și devine un dublu helix identic cu originalul.

E. În eucariote procesul de replicare poate începe în mai multe poziții din dublul helixurilor de ADN.

F. Replicarea poate să fie incompletă și să nu ajungă chiar până la capătul fiecărui cromozom, numit „telomer”, rezultând o scurtare generație după generație, fenomen considerat a avea un rol în procesul de îmbătrânire, în apoptoză sau în apariția unor boli (inclusiv unele forme de cancer).

4.4. Sinteza proteinelor

4.4.1. Paradigma centrală a bioinformaticii

A. Încercările de a explica marea variabilitate biologică fenotipică în contextul unității moleculare a structurilor biologice a condus la plasarea acesteia într-o poziție cheie pentru bioinformatică. Întrebarea la care dorim să răspundem este: Cum este transmisă informația stocată în structura moleculelor de ADN din nucleele celulelor, către celule, pentru a coordona practic toate procesele celulare, dar în special sinteza proteinelor, care sunt atât elemente structurale cât și elemente de control (enzime) a proceselor celulare. Mecanismul sintezei proteinelor este azi în bună măsură cunoscut, deși există încă elemente nu pe deplin elucidate. La descifrarea mecanismelor și-au adus contribuția cercetători din diverse domenii, un rol important avându-l cei din domeniul bioinformaticii.

B. Schema din figura 4.4.1. reprezintă sintetic principalele procese ce au loc pentru sinteza proteinelor.

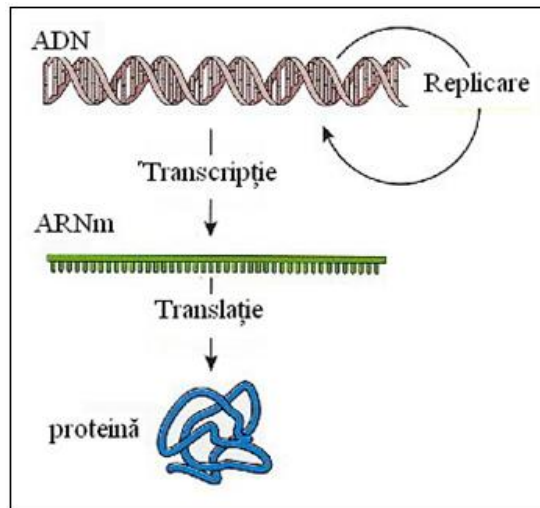


Fig. 4.4.1. Sinteza proteinelor – dogma centrală a bioinformaticii

Vom descrie în continuare aceste procese.

4.4.2. Transcripția

A. Regiunea din ADN care conține informația privind sinteza unei proteine se numește genă de sinteză. În amonte față de gena de sinteză se găsește o genă promotor, cu rol în declanșarea procesului de copiere a informației. Zona promotor mai este numită 5' UTR – 5' untranslated region, iar după gena de sinteză mai este o porțiune ce nu codifică secvență proteică numită 3' UTR.

B. Mecanismul transcripției

La activarea zonei promotor, dublul helix ADN din regiunea genei de sinteză este desfăcut în sensul 3'→5' de către helicază, apoi intervine ARN-polimeraza, care facilitează formarea unui lanț ARNm în sensul normal 5'→3'. Un singur lanț din cele două ale ADN este utilizat pentru citire, numit „lanț matriță” („template strand”), celălalt fiind numit lanț de codificare („coding strand”), deoarece conține secvența exact

în forma în care apare în molecula de ARN sintetizată, cu singura deosebire că, în ARN, în loc de timină apare uracilul. Transcripția este prezentată schematic în figura 4.4.2.

Procesul de transcripție are trei faze: prima este „inițierea”, prin activarea promotorului, a doua se numește „elongare”, reprezentată prin adăugarea, nucleotid cu nucleotid, a componentelor lanțului de ARN, iar ultima fază este „terminarea” (formarea unei bucle bogată în G-C, urmată de o succesiune de U).

Să menționăm că există și transcripție reversă, adică din ARN în ADN, întâlnită în unele cazuri patologice, cum ar fi infecția cu HIV.

C. Migrarea ARNm

Lanțul ARN format în nucleu se numește ARN mesager și se notează ARNm. În nucleu se mai găsesc niște organite numite „spliceosomi” care asigură un proces numit „splicing”. Prin acest proces sunt eliminate din molecula de ARNm formată, regiunile fără rol în codificarea secvenței proteice, numite „introni”, rămânând numai regiunile codificatoare, numite „exoni”. Molecula nou formată va traversa membrana nucleară și va ajunge la ribozomi, unde se va desfășura procesul de translație.

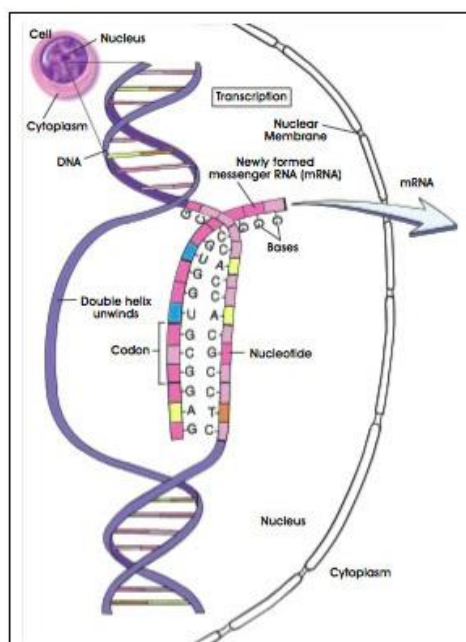


Fig. 4.4.2. Transcripția

4.4.3. Codul genetic

A. Pentru a putea folosi informația purtată într-un alfabet cu 4 litere, cum este cazul acizilor nucleici, la codificarea a 20 aminoacizi, vom avea nevoie de o succesiune de cel puțin 3 litere (George Gamow), care asigură $4^3 = 64$ combinații posibile (cu 2 litere am fi putut obține numai $4^2 = 16$ combinații. Această ipoteză, că o succesiune de 3 baze azotate într-o secvență ADN codifică 1 aminoacid într-o secvență proteică a fost confirmată experimental (Nirenberg, 1961). A fost astfel introdus termenul de „codon”, definit ca un triplet de baze azotate într-un lanț de acid nucleic pentru codificarea unui aminoacid într-o secvență polipeptidică.

		Second Letter				
		U	C	A	G	
1st letter	U	UUU Phe UUC UUA Leu UUG	UCU Ser UCC UCA UCG	UAU Tyr UAC UAA Stop UAG Stop	UGU Cys UGC UGA Stop UGG Trp	U C A G
	C	CUU Leu CUC CUA CUG	CCU Pro CCC CCA CCG	CAU His CAC CAA Gln CAG	CGU Arg CGC CGA CGG	U C A G
	A	AUU AUC Ile AUA AUG Met	ACU ACC ACA ACG	AAU Asn AAC AAA Lys AAG	AGU Ser AGC AGA AGG Arg	U C A G
	G	GUU Val GUC GUA GUG	GCU Ala GCC GCA GCG	GAU Asp GAC GAA Glu GAG	GGU Gly GGC GGA GGG	U C A G

Fig. 4.4.3. Codul genetic

B. Cercetări ulterioare (Khorana, Holley, premiul Nobel 1968) au permis stabilirea codului genetic prezentat în figura 4.4.3.

C. Observăm că pentru majoritatea aminoacizilor există mai multe codificări posibile. De asemenea există un codon de star, precum și trei codoni de stop – cărora nu le corespunde nici un aminoacid.

4.4.4. ARN de transport

A. În mecanismul sintezei proteinelor este necesară implicarea unor molecule care să asigure corespondența stabilită în codul genetic. Acestea sunt moleculele numite ARN de transport sau de transfer, notate ARNt (studiate în detaliu de Ochoa - premiul Nobel 1959).

B. Structura ARNt

Molecula de ARNt este destul de mică în comparație cu alte tipuri de ARN. Ea are 4 brațe dintre care două formează legături în timpul translației (figura 4.4.4)

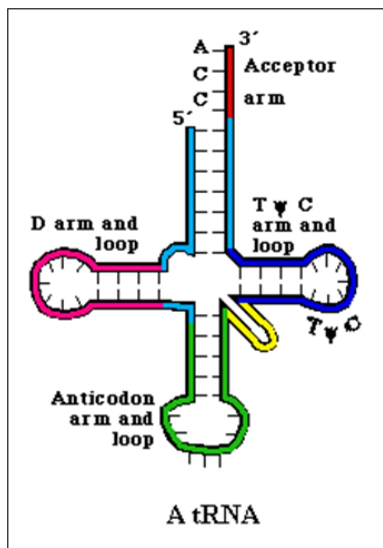


Fig. 4.4.4. Structura moleculei de ARNt

C. Observăm că brațul superior are o porțiune cu capătul 5' liber, care este de obicei fosforilat, iar la capătul 3' se poate atașa aminoacidul – atașarea este specifică.

D. Brațul opus are o buclă ce conține o secvență de trei baze azotate complementare codificării aminoacidului, secvență care se numește „anticodon”. Această porțiune este cea cu care molecula de ARNt se leagă de secvența potrivită din matrița ARNm.

E. Celelalte două ramuri, buclele D și T facilitează mecanismul de cuplare pe ribozom.

4.4.5. Structura ribozomilor

Revenim la descrierea structurii ribozomilor, începută în 4.1.7.

Ribozomii au două regiuni, o regiune care poate cupla molecule de ARNm venită din nucleu și o regiune care conține 3 situs-uri, notate E, P și A, regiune în care se leagă ARNt.

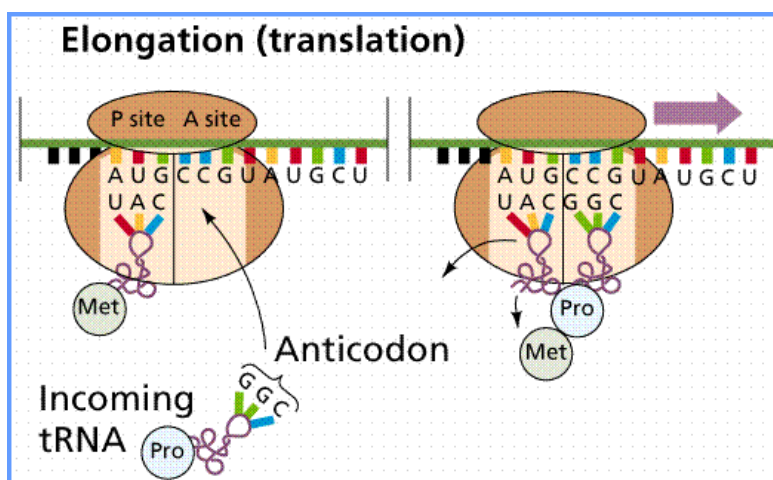


Fig. 4.4.5. Translația

4.4.6. Mecanismul translației

A. Când ARNm expune spre situl P porțiunea cu codonul de start (AUG), situl poate fi ocupat de un ARNt ce codifică metionina „Met” (conform codului genetic).

B. În situl A, încă liber, se poate lega de un ARNt care să aibă anticodonul corespunzător următorului codon din ARNm, codon expus către situl A.

C. Aminoacidul legat de ARNt de pe situl A se cuplează printr-o legătură peptidică de aminoacidul legat de ARNt de pe situl P.

D. La formarea acestei legături peptidice, AA din P se desprinde de ARNt de care era legat și întreaga structură formată se mută cu un codon spre în față: ARNt din P trece în E, ARNt din A trece în P, cu tot cu aminoacizii legați de el și în paralel se deplasează și ARNm în unitatea inferioară. Energia necesară procesului este furnizată de o moleculă de GTP, procesul fiind susținut de o altă proteină numită „factor de elongație”; de fapt, această fază se și numește „elongație”.

E. Situl A devenit vacant poate fi ocupat de un alt ARNt cu anticodon corespunzător următorului codon din matrița ARNm. Totodată ARNt de pe situl E se desprinde de lanțul ARNm. Aminoacidul legat de ARNt care a ocupat acum situl A

poate forma o legătură peptidică cu AA din poziția P, urmată de desprinderea acestuia de ARNT care l-a purtat.

Procesul se reia pas cu pas, generându-se un polipeptid, care iese din ribozom.

F. Procesul de translație se încheie când pe matrița ARNM apare un codon stop, care nu are echivalent în ARNT, deci situl A nu se mai ocupă. În acest moment, lanțul polipeptidic format se desprinde de ribozom.

Fazele sintezei proteinelor sunt prezentate schematic în figura 4.4.6.

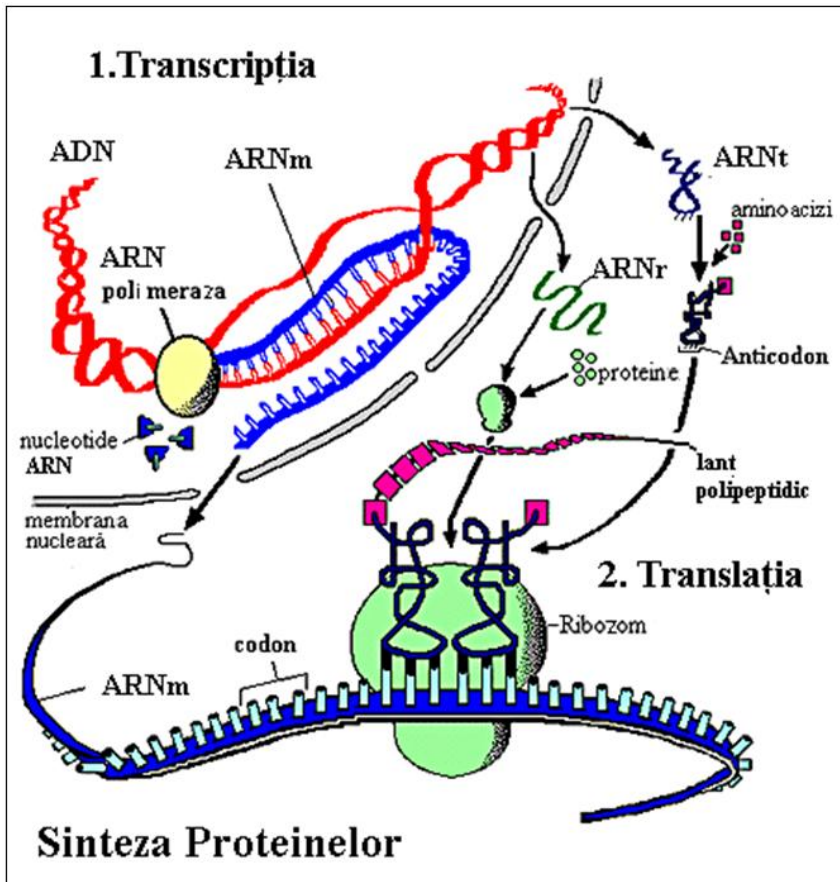


Fig. 4.4.6. Mecanismul translației

4.5. Controlul sintezei proteinelor

Mecanismele moleculare au un sistem de control care dirijează declanșarea, desfășurarea respectiv terminarea tuturor proceselor.

Controlul sintezei proteinelor la procariote este un mecanism de feedback studiat de Jacob, Monod și Lwoff (premiul Nobel 1965).

În cazul eucariotelor, mecanismul este mai complicat, descris prin modele cu feedback pozitiv. Nu vom intra aici în detalii privind aceste mecanisme de reglare.

5. Elemente de genetică

A. Genetica este o ramură a biologiei care a cunoscut una dintre cele mai spectaculoase evoluții în ultimele decenii, iar o serie de noțiuni din genetică sunt indispensabile pentru înțelegerea bioinformaticii.

B. Ca obiect de studiu, genetica studiază ereditatea și variabilitatea organismelor. Diversitatea lumii vii, reprezentată prin diferite caractere ale organismelor, este determinată genetic prin existența unor „factori ereditari” sau „gene” care se modifică sau rearanjează în cursul evoluției.

5.1. Scurt istoric. Capitolele geneticii

5.1.1. Epoca mendeliană

A. Gregor Mendel este recunoscut ca părinte al geneticii; în 1865, în urma studiilor pe diverse soiuri de mazăre, formulează legile eredității și introduce termenul de „factor ereditar”.

B. În 1900 legile lui Mendel sunt redescoperite și reformulate de Hugo de Vries, Tschermak și Correns.

C. În 1909, Johanson introduce termenul de „genă” pentru „factorul ereditar”.

5.1.2. Teoria cromozomială a eredității

A. La începutul sec. XX, Sutton și Boveri arată că genele sunt situate în cromozomi.

B. Părintele teoriei cromozomiale a eredității este Morgan, autor a numeroase studii pe *Drosophila melanogaster* (musculița de oțet), care au condus la localizarea genelor pe cromozomi (devenite ulterior hărți cromozomiale).

C. Sunt identificate tipurile de diviziune, fenomenele de linkage și cross over și recombinarea, explicând variabilitatea din natură.

5.1.3. Genetica moleculară

A. Începutul geneticii moleculare a fost reprezentat de descoperirea rolului genetic al ADN (Avery, 1944). În 1953, Watson și Crick descifrează structura dublu helix a ADN, arătând că genele sunt fragmente de ADN.

B. O suită de descoperiri importante au continuat să se acumuleze, aducând geneticii un palmares important de premii Nobel. Enumerăm aici cele mai semnificative descoperiri: codul genetic, secvențialitatea ADN, reglajul genetic al activității celulare, mecanismul sintezei proteinelor, descifrarea genomului uman.

5.1.4. Bioinformatica

A. De la bun început s-a văzut că funcția principală studiată în genetică este de fapt cea de stocare și transmitere de informație. S-a conturat astfel o nouă disciplină,

care urmărește aceleași fenomene ca și genetica în particular, sau biologia în general, însă din punct de vedere informațional – bioinformatica. Termenul a fost introdus în 1979 de către Hogeweg.

B. Dezvoltarea bioinformaticii este indisolubil legată de crearea bazelor de date moleculare și vital asociată de facilitățile oferite de internet.

C. Pentru analiza datelor acumulate, bioinformatica și-a dezvoltat propriile instrumente de lucru – elaborarea unor sofisticate algoritmi de analiză și comparare a secvențelor sau a arborilor filogenetici.

D. Biologia sistemelor, ca important capitol al bioinformaticii, tinde să se contureze ca disciplină separată, dedicată proceselor integrative – urmărirea interacțiunilor între componentele sistemelor biologice, pe baza informației purtate de fiecare componentă. În biologia sistemelor se face intensiv uz de metodele de modelare matematică și simulare pe calculator a proceselor biologice.

5.2. *Genetica mendeliană*

5.2.1. *Terminologie*

A. Genetica și-a dezvoltat o terminologie specifică, ajungându-se cu timpul să se găsească echivalenții moleculari pentru diverși termeni. Vom trece aici peste considerentele istorice, adăugând pentru principalii termeni și explicațiile moleculare adecvate. În plus, adesea termenii nu vor fi trecuți ca într-un dicționar, ci în mod firesc, cum apar în descrierea fenomenelor și interpretarea rezultatelor. O serie de termeni noi vor fi introduși ulterior (în secțiunea adecvată).

B. Mendel a lucrat pe plante, în special mazăre (*Pisum sativum*), urmărind diferite caractere distincte la urmași după diferite încrucișări. Câțiva termeni specifici aici:

a) *hibridare* - încrucișarea a doi indivizi ce au una sau mai multe caracteristici diferite. Caracteristicile au adesea manifestări bivalente (în perechi), uneori și plurivalente. De exemplu: talia plantei: înalte/pitice, culoarea bobului: galben/verde etc. Ulterior s-au găsit explicațiile moleculare ale multor caracteristici. De exemplu: pentru suprafața bobului netedă/zbârcită: prezența amidonului generează suprafața netedă, conținut ridicat de dextrină determină suprafața zbârcită. În funcție de numărul caracteristicilor diferite la încrucișare avem:

- *monohibridare* – încrucișare între părinți cu o singură pereche de caractere diferite,
- *dihibridare* – cazul a două perechi de caractere diferite,
- *polihibridare* – mai multe.

Să mai introducem aici și termenii:

- *hibrid*: un descendent rezultat în urma unei încrucișări (popular „hibrid”= soi (la plante) sau „rasă” la animale,
- *generație*: ansamblul descendenților obținuți din perechi de părinți similari.

b) *fenotip* – (termen introdus mai târziu): ansamblul caracteristicilor morfologice, fiziologice, biochimice și comportamentale ale unui organism.

5.2.2. *Legea purității gameților*

A. Observațiile lui Mendel

a) Prin încrucișarea plantelor (generația părinților P) cu bob neted, cu plante cu bob zbârcit, în prima generație F₁ s-au obținut numai plante cu bob neted. Acest aspect

(bob neted) al unui caracter (suprafața bobului), poartă denumirea de *aspect dominant* (numit și caracter dominant), notat uzual cu litere mari (A). Aspectul „bob zbârcit”, care nu apare în prima generație a fost numit „aspect recesiv” (caracter recesiv), notat uzual cu litere mici (a).

b) Prin încrucișarea plantelor din generația F₁, în a doua generație F₂ apar atât plante cu bob neted cât și cu bob zbârcit. Raportul între numărul descendenților cu aspect dominant (D) și recesiv (R) a fost [D:R= 3:1].

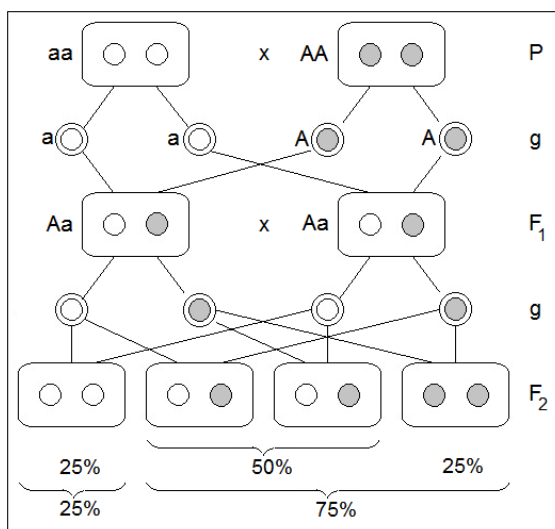


Fig. 5.2.2. Monohibridismul

B. Explicația lui Mendel

a) Mendel a intuit că în celulele somatice factorii ereditari se găsesc în pereche, iar în celulele sexuale sub formă simplă, (ulterior dovedit prin numărul de cromozomi: dublu în celulele somatice – celule diploide, simplu în cele sexuale – celule haploide).

b) Celulele sexuale, numite *gameți*, poartă caracterele părinților (în porțiuni numite ulterior „gene” – echivalent cu „factor ereditar” și aproximativ „1 genă = 1 proteină.”

c) Celula-ou, din care se dezvoltă viitorul organism apare prin unirea celor 2 gameți, unul de la tată și unul de la mamă.

d) Legea purității gameților afirmă că: Gameții sunt puri din punct de vedere genetic, (conțin doar un aspect al unui caracter = factor ereditar). Prin combinarea la întâmplare (probabilistă) a acestor gameți apare în generația a doua fenomenul segregării în proporția: „3D: 1R”.

C. Homozigoți și heterozigoți

a) organisme care au un singur tip de factori ereditari sunt pure din punct de vedere genetic, numite *homozigoți*. În organismele homozigote, gameții provenind de la cei doi părinți poartă același aspect al unui caracter. Homozigoții pot fi purtători fie ai aspectului dominant (notați AA), fie ai aspectului recesiv (notat aa).

- În observațiile menționate anterior s-a presupus că părinții (P) erau „homozigoți”, aspect justificat prin omogenitatea generațiilor la încrucișări exclusiv între organisme de același fel.

b) Organismele impure din punct de vedere genetic, adică cele ce posedau ambele aspecte ale unui caracter – factor ereditar, au primit numele de *heterozigoți*. Însă din punct de vedere al manifestării externe, este vizibil caracterul dominant (figura 5.2.2.a).

c) *Genotip* – totalitatea factorilor ereditari conținuți într-un organism. Termenul devine important în contrast cu „fenotipul”, în care surprindem doar aspectul extern, generat de factorii determinanți, deși același aspect poate fi generat de genotip-uri diferite.

5.2.3. Dihibridismul și legea segregării independente a caracterelor

A. Un pas important, efectuat tot de Mendel, a fost analiza aspectelor/caracterelor la descendenți provenind din părinți diferiți prin două perechi de caractere. Pentru abordarea formală să notăm caracterele cu A sau a (dominant = A, recesiv = a, ex.: bob neted A/zbârcit a), respectiv B sau b (bob galben B/bob verde b).

B. În prima generație toate plantele aveau fenotipic boabe netede și galbene, deși genotipic o parte erau hibride.

C. În a doua generație au apărut patru feluri de descendenți:

- două caractere dominante – proporție 9/16 (AB)
- un caracter dominant, altul recesiv – proporție 3/ 16 (Ab)
- celălalt caracter dominant, celălalt recesiv – 3/16 (aB)
- ambele recesive – 1/16 (ab).

D. În figura 5.2.3. este prezentată formarea generației F₂ în cazul dihibridismului.

Dihybrid cross
A and a represent one trait, and B and b represent a different trait that is linked to inheritance of A or a.

	AB	Ab	aB	ab
AB	AABB	AABb	AaBB	AaBb
Ab	AABb	AAbb	AaBb	Aabb
aB	AaBB	AaBb	aaBB	aaBb
ab	AaBb	Aabb	aaBb	aabb

Dominant for A and B = 9/16 Dominant for A, recessive for b = 3/16
 Recessive for a, dominant for B = 3/16 Recessive for a, recessive for b = 1/16

Fig. 5.2.3. Dihibridismul

E. Pe baza constatărilor experimentale, Mendel a enunțat și legea segregării independente a caracterelor: „Factorii ereditari pereche segregă (se separă) independent de alte perechi de factori ereditari. Raportul de segregare este 3D:1R pentru fiecare pereche, iar în cazul a două perechi este 9:3:3:1”.

F. Legea segregării independente a caracterelor a fost susținută de numeroase evidențe experimentale atât în regnul vegetal cât și animal, venind în sprijinul ipotezei unității lumii vii.

5.2.4. Alte tipuri de segregare

A. Noțiuni de alelă

Conceptului de „factor ereditar”, care determină un caracter, i s-a asociat și termenul de „alelă”, căruia i s-a atașat ulterior și baza moleculară – gena – reprezentată printr-o porțiune dintr-o moleculă de ADN. Astfel, spunem că aspectul de bob neted sau zbârcit este determinat de o pereche de gene alele, una dominantă, cealaltă recesivă. Organismele homozigote au doar o alelă, sub formă de pereche, fie dominantă (AA), fie recesivă (aa), iar cele heterozigote au ambele gene alele (Aa), fenotipic având aspectul dat de alela dominantă.

B. Relațiile interalelice s-au dovedit mai complexe decât forma simplă dominantă/recesivitate prezentată mai sus. Amintim aici, fără a intra în detalii, și alte tipuri de segregare.

a) Dominanța incompletă – destul de frecventă, în care organismele heterozigote nu manifestă fenotipic caracterul dominant ci prezintă o formă intermediară; de exemplu, la planta barba-împăratului (*Mirabilis jalapa*), la încrucișarea varietății cu flori roșii (AA), cu varietatea cu flori albe (aa), în F_1 apar numai flori roz (Aa), iar în F_2 , apar flori roșii/roz/albe în raportul 1:2:1, deci organismele heterozigote prezintă caracterul intermediar.

b) Supradominanța - organismele heterozigote au mai dezvoltate însușirile biologice (fertilitate, talie, etc), decât părinții homozigoți.

c) Gene letale – gene care în stare homozigotă determină moartea individului înainte de maturitatea sexuală. De exemplu: la încrucișarea a doi șoareci galbeni rezultă șoareci galbeni și șoareci de altă culoare în raport 2:1, șoarecii galbeni rezultați fiind toți heterozigoți, deci forma homozigotă este letală.

d) Polialelia – același caracter apare în mai mult de două forme, deci sunt mai mult decât 2 gene alele (A și a) pentru același „locus”. De ex.: musculițe (*Drosophila m*) cu ochi albi, roz, purpurii, corai, etc. tipul normal (nemutant) fiind ochi cărămizii.

e) Codominanța - Grupa sanguină la om este determinată de 3 gene polialele: L^A , L^B și I, unde L^A și L^B sunt dominante față de I, însă împreună sunt codominante. În tabelul 5.2.4. sunt prezentate schematic genotipurile corespunzătoare.

Tabelul 5.2.4. Fenotipurile și genotipurile grupelor sanguine la om

Grupa sanguină (fenotip)	Genotipuri posibile
AB	$L^A L^B$
A	$L^A L^A$ sau $L^A I$
B	$L^B L^B$ sau $L^B I$
0	II

Aceste relații sunt utile atât pentru transfuzie cât și în stabilirea paternității. Cunoscând grupa sanguină a copilului (de ex.: A) și a mamei (de ex.: 0), se pot stabili grupele sanguine ale tatălui (numai AB sau A).

f) Poligenia – adesea un caracter fenotipic este rezultatul interacțiunii mai multor gene nealele.

5.2.5. Principiul Hardy-Weinberg

A. Genofond

Între indivizii unei populații au loc schimburi de gene. Genotipurile parentale se desfac și se amestecă la nivelul descendenților, formând un tot unitar la nivelul populației, numit „fond genetic” sau „genofond” - totalitatea genelor unei populații.

B. În 1908 a fost enunțată legea stabilității frecvențelor genice numită și legea/principiul Hardy-Weinberg.

a) Enunț: Frecvențele genice rămân constante de la o generație la alta, frecvența genotipurilor fiind și ea constantă, determinată de frecvența genelor populației parentale.

Legea este valabilă pentru populații panmictice mari (cu încrucișare aleatoare), în care nu acționează forțe modificatoare ale frecvențelor genetice.

b) Din punct de vedere formal, pentru o pereche de gene alele A_1 și A_2 (nu trebuie precizată dominanța), care au frecvențele p și q (frecvența \approx probabilitatea), atunci organismele unei populații pot avea trei genotipuri, A_1A_1 , A_1A_2 și A_2A_2 cu frecvențele p^2 , $2pq$ și q^2 .

Sintetic este prezentat acest lucru în tabelul 5.2.5.

Tabelul 5.2.5. Ilustrarea legii Hardy-Weinberg

$\begin{matrix} \text{♂} \\ \diagdown \\ \text{♀} \end{matrix}$	$A_1 (p)$	$A_2 (q)$
$A_1 (p)$	$A_1A_1 (p^2)$	$A_1A_2 (pq)$
$A_2 (q)$	$A_1A_2 (pq)$	$A_2A_2 (q^2)$

5.3. Teoria cromozomială

5.3.1. Localizarea genelor

Lucrările lui Mendel au fost uitate, principiile sale fiind redescoperite la sfârșitul sec. XIX. De menționat formularea mai clară a lui Hugo de Vries: celulele au câte două copii ale aceluiași gene, câte una de la fiecare părinte. La începutul sec. XX s-au acumulat o serie de date privind diviziunea celulară, stabilindu-se legături evidente între datele studiilor de genetică (privind ereditatea) și localizarea cromozomială a informației genetice.

5.3.2. Morfologia cromozomilor

A. Fiecare specie de viețuitoare are caracteristic numărul și forma cromozomilor care pot fi vizualizați și identificați în timpul diviziunii celulare.

B. Un cromozom este alcătuit din două *cromatide*, unite într-un loc numit centromer. Ramurile mai lungi se notează cu „p”, iar cele mai scurte cu „q”.

C. După poziția centromerului cromozomii pot fi:

- metacentric – centromerul aproximativ median,
- submetacentric – centromerul depărtat de mijloc,
- subtelocentric – destul de aproape de un capăt,
- telocentric – centromerul situat terminal.

D. Cariotip, idiogramă

Ansamblul cromozomilor se numește *cariotip* și este caracteristic fiecărei specii. Dacă se adaugă măsurători pentru fiecare cromozom și se plasează în ordinea mărimii, se obține „idiograma” speciei respective.

E. Heterozomi, autozomi

Cromozomii somatici poartă denumirea de *autozomi* și perechile sunt uzual numerotate începând cu cei mai mari.

În afară de cromozomii perechi, celulele conțin și doi cromozomi ai sexului, denumiți *heterozomi*. Celulele organismelor de sex feminin au o pereche de cromozomi omologi, notați XX, în timp ce cele de sex masculin au un cromozom X și unul mai mic, notat Z.

F. În timpul meiozei, la formarea gameților, organismele feminine produc doi gameți X, iar cele masculine produc un gamet X și unul Z.

5.3.3. Transmiterea înlănțuită a genelor (*linkage*)

A. *Drosophila* ca model în genetică

Școala lui Morgan a utilizat muscuțita de oțet (*Drosophila melanogaster*) ca model de studiu, având ciclul de viață potrivit (2 săptămâni/generație), cariotip simplu, (4 perechi de cromozomi), cu cromozomi uriași în glandele salivare (ușor de văzut la microscop), fiind cunoscute cca 500 de mutații (modificări ale materialului genetic), însoțite de manifestări fenotipice.

B. Notății

Genele de tip normal (numit și „sălbatic” se notează cu 1 sau 2 litere urmate de semnul „+”, în timp ce alelele corespunzătoare mutante sunt notate fără „+”.

Vom folosi în continuare câteva exemple în care am folosit notațiile din tabelul de mai jos 5.3.3.

Tabelul 5.3.3. Exemple de gene alele la *Drosophila*

Caracterul	Gene alele Forma sălbatică	Forma mutantă
Culoarea ochilor	bw^{+} - ochi roșii [B]	bw ochi maro [b]
Culoarea corpului	b^{+} - gri [C]	b - negru [c]
Forma aripilor	vg^{+} normale [A]	vg – vestigiale (mici) [a]

C. Plasarea liniară a genelor

Este evident că numărul genelor (caracterelor) este mult mai mare decât numărul cromozomilor, deci mai multe gene trebuie să fie plasate pe un cromozom.

O ipoteză centrală în teoria cromozomială a eredității este plasarea liniară a genelor în cromozom.

D. Linkage

Culegând date de la sute de mii de indivizi, de-a lungul mai multor generații, urmărind diferite caractere, s-au constatat numeroase abateri de la legea segregării independente a caracterelor. Au apărut adesea descendenți cu transmiterea unor caractere „în bloc”, fenomen numit „linkage”. Un exemplu este redat în figura 5.3.3.a, din care reiese că descendenții din generația a 2-a (F_2) sunt doar de două tipuri, în proporție 1:1.

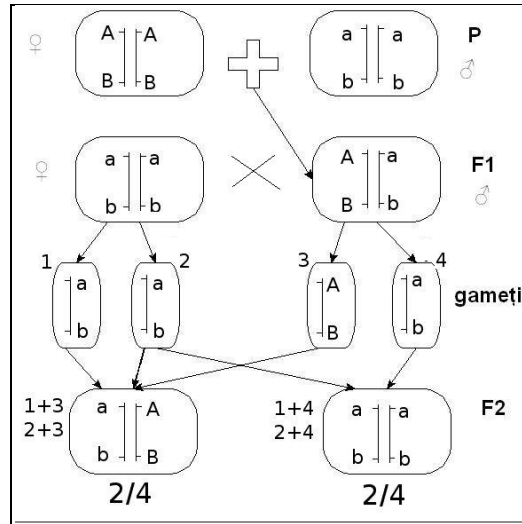


Fig. 5.3.3.a. Evidențierea linkage-ului

Conform legii segregării independente a caracterelor trebuia să avem 4 tipuri de descendenți în proporția 1:3:3:9. (figura 5.3.3.b).

Explicația: genele responsabile de cele două caractere sunt „legate” (se comportă ca una singură), fiind plasate pe același cromozom (cromozomul 2), apropiate una de alta. În cazul genelor plasate pe cromozomi diferiți legea segregării independente s-a respectat.

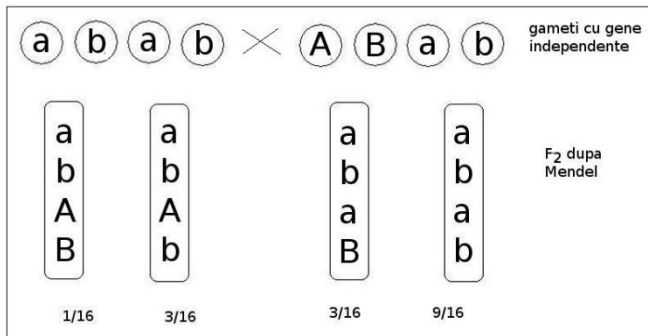


Fig. 5.3.3.b. Varianta mendeliană a exemplului 5.3.3.a

5.3.4. Schimbul reciproc de gene (crossing-over)

A. O altă constatare experimentală importantă a evidențiat că în general, și fenomenul de linkage are numeroase excepții, deci apar generații cu caractere segregate, dar nu în proporție mendeliană ci mult mai apropiată de cea explicată prin linkage.

B. Un exemplu tipic este redat în figura 5.3.4.a. Din părinți homozigoți mutanți, cu corp mutant negru (cc) și aripi normale (AA), încrucișați cu musculițe cu corp normal gri (CC) și aripi vestigiale (aa), apar în prima generație heterozigoți cu fenotip, normal (corp gri, aripi normale). Prin încrucișarea lor cu dubla mutantă homozigotă apar în generația a doua în procent ridicat (41,5%) variantele explicabile prin linkage,

însă apar și variante în care cele două caractere s-au segregat, în proporție scăzută, trecute în figură în dreptunghiuri cu linie punctată.

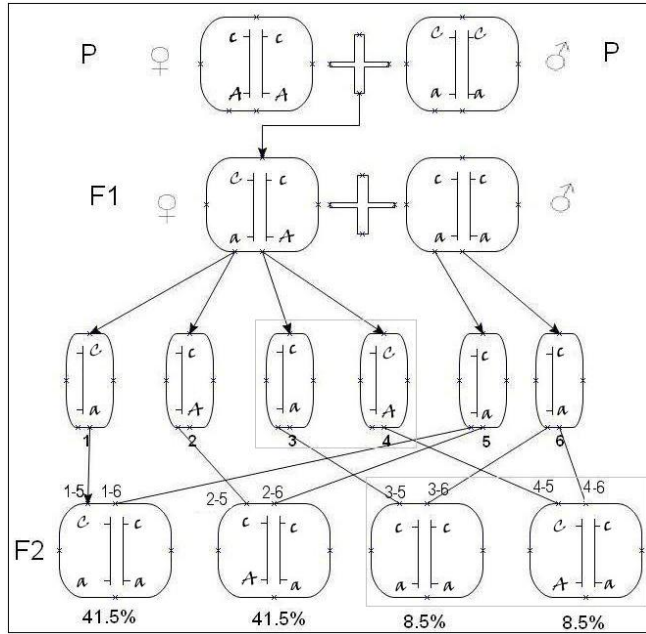


Fig. 5.3.4.a. Exemplu de crossingover

C. Crossing-over

Explicația lui Morgan, dovedită ulterior experimental este următoarea: în timpul meiozei cromozomii omologi se apropie mult, se ating, în unul sau mai multe puncte, iar la separare își pot schimba între ei anumite segmente cromatidice, ce duce la schimb de gene.

Ilustrarea acestui fenomen este dată în figura 5.3.4.b.

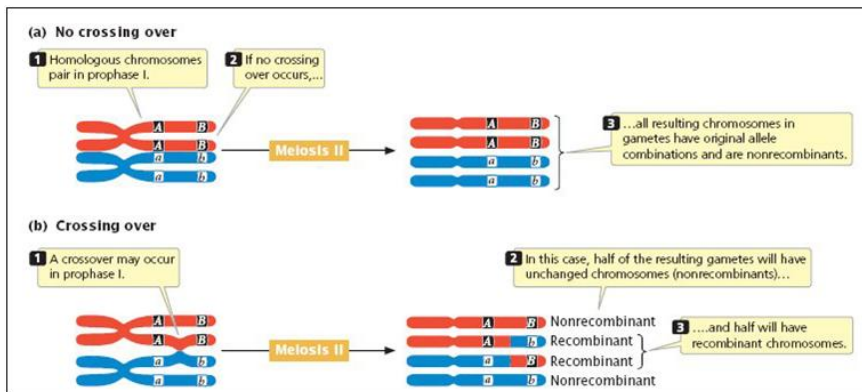


Fig. 5.3.4.b. Schema procesului de crossing-over

Fenomenul de crossing-over poate avea loc și prin mai multe puncte de contact: crossing-over dublu, triplu sau multiplu, rezultând prin recombinare numeroase variante de rearanjare a genelor.

5.3.5. Hărți cromozomiale

A. Procent de recombinare

Un pas important a fost făcut când s-au analizat valorile numerice ale abaterilor de la linkage-ul total. Aceste abateri, explicate prin fenomenul de crossing-over, reprezintă procentul în care poate avea loc procesul de recombinare, notat uzual cu RF („recombination frequency”).

B. Distanța între gene

Dacă se face ipoteza că probabilitatea crossing-over-ului este uniformă de-a lungul unui cromozom, atunci putem considera că în cazul a două gene mai depărtate între ele pe cromozom, probabilitatea de recombinare este mai mare. A apărut astfel ideea de a aprecia distanța dintre gene în funcție de frecvența recombinărilor RF.

C. Unitatea pentru distanță genetică

S-a propus numele de „morgan” pentru unitatea fundamentală de distanță genetică, cu submultiplul practic centimorgan (cM) sau m. u. (genetic map unit).

Definiție: O unitate de hartă genetică (1 m. u. = 1cM este distanța dintre două gene cu frecvența de recombinare 1/100 (1 produs din 100 meioze este recombinat).

D. Pe baza datelor experimentale culese s-au putut realiza scheme cu poziția relativă a genelor pe cromozomi. Prima hartă cromozomială a fost realizată de Sturdevant.

În figura 5.3.5 este prezentată o hartă cromozomială simplă rezultată din studii de recombinare.

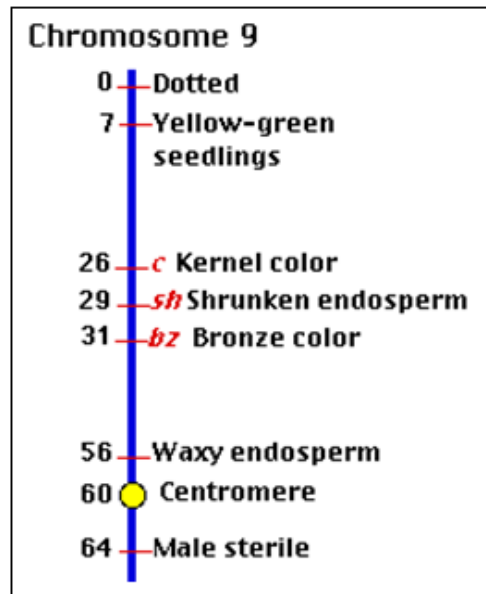


Fig. 5.3.5. Harta cromozomială a cromozomului 9 de șoarece

E. Scoruri pentru RF

Din punct de vedere formal, putem estima frecvența de recombinare printr-un scor numit LOD (Logarithm of ODDS), dat de relația:

$$LOD = Z = \log_{10} \frac{\text{probabilitatea unei secvențe pentru o valoare data a linkage-ului}}{\text{Probabilitatea secvenței fara linkage}}$$

$$LOD = \log_{10} \frac{(1-\vartheta)^{NR} \times \vartheta^R}{0.5^{(NR+R)}} \quad (5.3.5.a)$$

unde am folosit notațiile:

NR – numărul de urmași ne-recombinanți,

R – numărul de urmași recombinanți

ϑ - fracția de recombinare [= R / (NR + R)]

5.3.6. Markeri pentru hărți genetice

A. Metodele folosite pentru primele hărți au aplicabilitate limitată, bazându-se pe identificarea vizuală (eventual cu microscopul) a variațiilor fenotipice. În plus, s-a găsit că un caracter poate fi afectat de mai multe gene.

B. De aceea s-a recurs la metode biochimice pentru identificarea diferitelor fenotipuri. Moleculele utilizate pentru astfel de identificări se numesc *markeri*. Multe din rezultatele actuale au fost obținute folosind markeri ADN (aspecte mapate care nu sunt gene).

5.4. Genetica moleculară

5.4.1. Corespondențe moleculare

A. După descoperirea ADN și a structurii sale, interesul s-a mutat spre determinarea corespondențelor între genele determinate prin metode genetice și secvențele corespunzătoare din moleculele de ADN.

B. De importanță covârșitoare a fost dezvoltarea metodelor de determinare a secvențelor nucleotidelor în acizii nucleici, respectiv a secvențelor aminoacizilor în proteine. Acestea au condus în final la stabilirea unor hărți genetice detaliate, cu localizarea genelor responsabile pentru anumite caractere, inclusiv în cazul unor boli genetice.

5.4.2. Structura unui cromozom

A. Reunind datele din biologia celulară și moleculară cu cele din genetică s-a putut reprezenta schematic structura generală a unui cromozom (figura 5.4.2.).

5.4.3. Hărți fizice

A. Hărțile construite prin metode genetice arată localizările pe cromozomi a genelor pe baza evidențelor fenotipice sau utilizând markeri ADN. Ele au însă un grad limitat de rezoluție și o acuratețe de asemenea limitată.

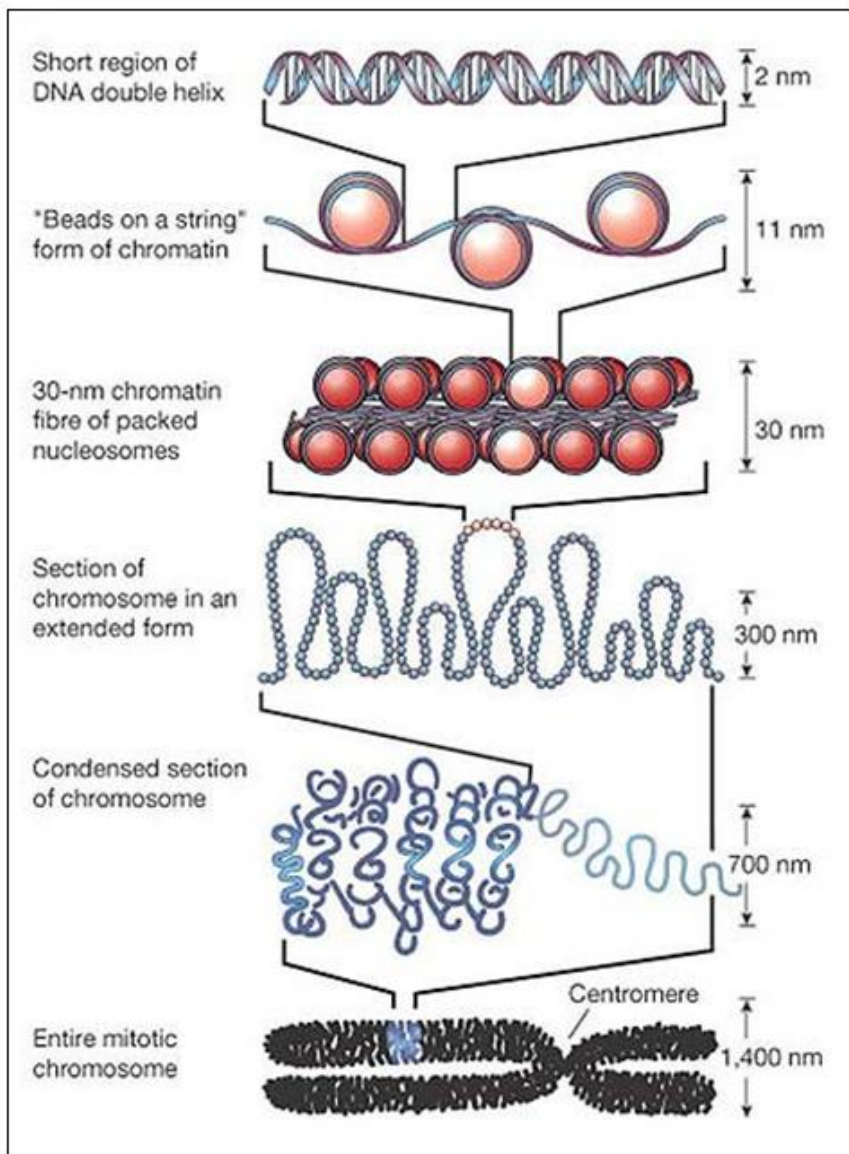


Fig. 5.4.2. Structura unui cromozom

B. O metodă alternativă este construcția unor hărți „fizice” bazate pe date de biologie moleculară, prin analiza directă a moleculelor de ADN. Menționăm existența a trei categorii de metode:

- „restriction mapping”, prin localizarea punctelor de acțiune al enzimelor de restricție (endonucleaze de restricție),
- hibridizarea fluorescentă *in situ* (FISH fluorescent *in situ* hybridization),
- „sequence tagged site mapping” (STS) - se mapează pozițiile unor secvențe scurte prin analiza fragmentelor de genom.

5.4.4. Exemplu de hartă cromozomială

A. Prezentăm ilustrativ în figura 5.4.4. un exemplu de hartă cromozomială.

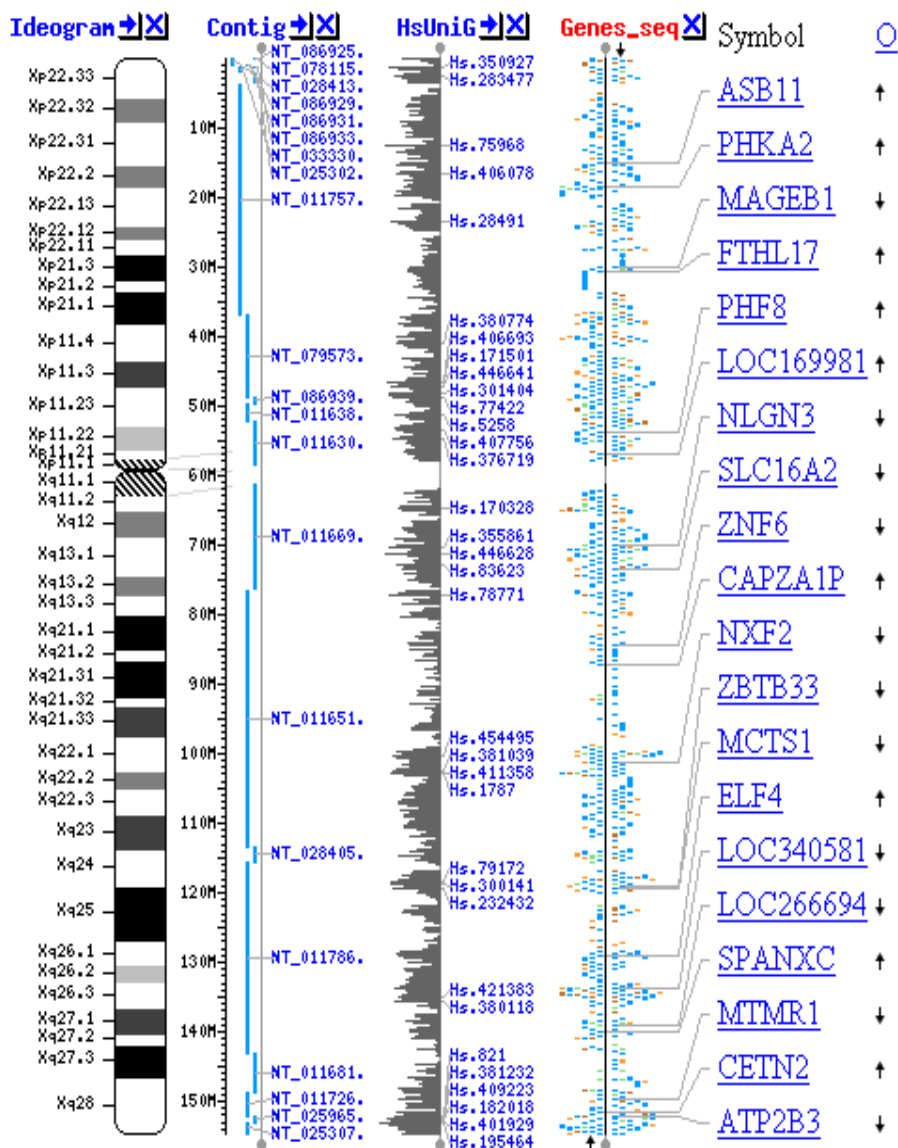


Fig. 5.4.4. Hartă cromozomială

B. În următoarele capitole vor fi prezentate metodele specifice de bioinformatică pentru analiza secvențelor și interpretarea structurilor din punct de vedere informațional.

6. Analiza secvențială individuală

6.1. Introducere

6.1.1. Obiectul bioinformaticii

A. Cercetările biologice, care au cunoscut o puternică dezvoltare în ultimele decenii, au condus la acumularea unei cantități imense de informație (sub formă de date), mai ales în domeniul biologiei moleculare. Acum marea provocare este de a da sens acestor informații, adică transformarea „datelor” în „cunoștințe”, care să ne facă capabili să interpretăm datele. Acesta este chiar obiectul bioinformaticii. Din multitudinea de definiții considerăm cea mai potrivită cea din glosarul NIH.

B. *Definiție:* „Bioinformatica este studiul structurii informației biologice și a sistemelor biologice. Ea aduce împreună datele de cercetare gnomică cu teoria și instrumentele matematicii și a științei calculatoarelor”.

6.1.2. Capitolele bioinformaticii

A. Există mai multe posibilități de a privi structural bioinformatica și a o diviza în capitole – fie în funcție de moleculele studiate, în funcție de aspectele analizate sau în funcție de metodele folosite.

B. Vom aborda aici o clasificare hibridă combinând „aspecte/metode”, conform uzanțelor:

- a) Baze de date moleculare (se vor trata la lucrări practice)
- b) Analiza secvențială (cap. 6, 7, 8 și 9 precum și la lucrări practice)
- c) Lanțuri Markov (cap. 10 și 11 precum și la lucrări practice)
- d) Analiza filogenetică (cap. 12 precum și la lucrări practice)
- e) Structuri tridimensionale (se vor trata la lucrări practice)
- f) Biologia sistemelor (nu a fost inclusă în acest manual).

6.1.3. Obiectul analizei secvențiale

A. poziție centrală în bioinformatică este ocupată de analiza secvențială. Într-adevăr, confirmarea ipotezei că informația biologică este înregistrată sub forma unor secvențe în structura acizilor nucleici și a proteinelor, a adus în prim plan metodele de analiză a acestor secvențe.

- B. Prin analiză secvențială putem aborda:
- analiza secvențială individuală,
 - compararea a două secvențe,
 - alinierea multiplă.

6.2. Analiza secvențială grafică

6.2.1. Baze teoretice

A. Proprietățile AA sunt diferite.

Evident, putem extrage o serie de informații privind proprietățile unei macromolecule – aici ne gândim în primul rând la proteine – pornind de la compoziția și structura sa, având în vedere faptul că elementele structurale – aminoacizii (AA) – au proprietăți diferite:

- reziduurile (catenele laterale) pot fi hidrophile sau hidrofobe,
- dimensiunea,
- sarcina electrică (aceasta depinde și de pH-ul mediului),
- flexibilitatea legăturilor.

Pentru reprezentări grafice se utilizează în practică reprezentarea unei proprietăți de-a lungul unei secvențe. Prezența au absența unei periodicități este vizibilă în astfel de reprezentări.

Pentru o mai ușoară evidențiere a unor posibile periodicități se folosesc și ferestre de netezire (în care un punct de grafic reprezintă media unei regiuni în jurul său). Practic s-au dovedit utile ferestre de 9 – 15 elemente. Un astfel de grafic este prezentată în figura 6.2.1.

Nishikawa face o comparație între evoluția sistemelor nevii, dictată de condițiile externe și a celor vii, dictată de informație.

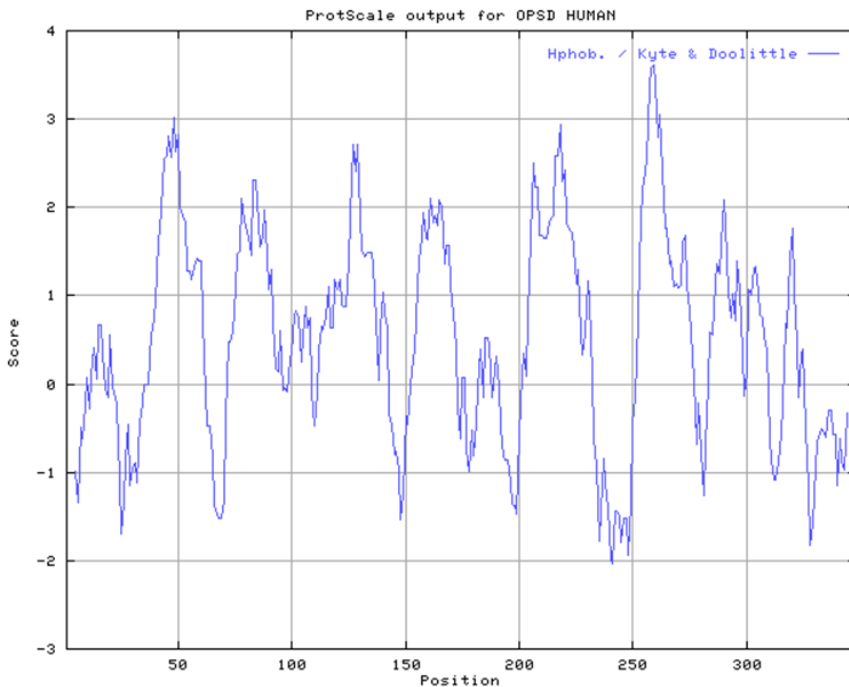


Fig. 6.2.1. Hidrofobicitatea rodopsinei umane (AC P08100) creat cu ExPASy Service ProtScale. S-a folosit o secvență de 350 AA și o fereastră de 9 unități, pe scala de hidrofobicitate Kyte-Doolittle.

B. Analiza proporțiilor AA – există o apariție preferențială a AA în anumite structuri secundare.

Așa cum am văzut în capitolul privind „Structura proteinelor”, structurile tridimensionale, începând cu structura secundară depind de structura primară, adică de *secvența* aminoacizilor. Pornind de la distribuția statistică a AA în diferite structuri secundare, s-a constatat o distribuție preferențială a AA, dictată de fapt de configurația spațială a atomilor fiecărui AA. Iată câteva exemple:

- **Glu** se găsește aproape exclusiv în α – *helix*
- **Val (Pro, Tyr)** se găsește în principal în fâșii β
- **Gly** este numită și „spărgător de elice”, iar **Pro** este de asemenea rar întâlnită în α – *helix* (**Ala, Cys, Leu** = promotori de *helix*).

6.2.2. Aspecte structurale

A. AA hidrofobi tind să apară în interiorul proteinelor globulare, iar cei hidrofilii spre exterior

- în proteinele membranare regiunile de inserare în membrane, regiunile de inserare în membrană sunt de asemenea bogate AA hidrofobi (având în vedere natura lipidică a moleculelor membranare).

B. Detectarea anumitor periodicități ale hidrofobității poate de asemenea sugera apartenența la o structură elice sau fâșie.

O aplicație pornind de aici ar fi predicția epitopilor antigeni, prin relația lor cu suprafața expusă.

6.3. Structura tridimensională

6.3.1. Roți elicoidale și elice amfipatice

A. Aranjarea aminoacizilor cu anumite proprietăți (de exemplu hidrofobitatea) în structura secundară ne poate sugera proprietăți și chiar funcții ale proteinei respective. S-au găsit astfel de aranjamente în cazul proteinelor structurale, în special pori sau canale transmembranare.

B. bună metodă de vizualizare a regiunilor elicoidale o reprezintă imaginea structurală a unei secvențe proiectată de-a lungul axei (vedere „din avion”), numită „roata elicoidală”. Un astfel de exemplu este redat în figura 6.3.1.

C. În partea de sus, (A), este reprezentată secvența de 18 AA care poate forma un α – *helix*. Luând ca punct de plecare pentru reprezentare „poziția orei 12”, plecând în sens orar, sunt trecuți toți aminoacizii din secvență. Cum o spiră α – *helix* cuprinde cca 3,6 aminoacizi, unghiul între doi aminoacizi succesivi va fi cca 100° , deci vom avea un aminoacid pe aceeași generatoare după o secvență de 18 reziduuri (adică 5 spire). Observăm că jumătate din dreapta elicei cuprinde doar aminoacizi nepolari, în timp ce cea din stânga are doar aminoacizi polari. Această aranjare permite agregarea moleculei noastre cu alte suprafețe hidrofobe și să poată funcționa ca por sau canal în membrană.

7. Compararea a două secvențe

7.1. *Introducere*

7.1.1. *Fundamentele analizei secvențiale*

A. Observațiile noastre privind lumea vie în mod frecvent bazate pe comparații, găsind adesea asemănări sau deosebiri între diferite structuri sau între diferite funcții. Pentru o abordare științifică, s-a simțit, desigur, nevoia de a introduce metode care să permită o exprimare cantitativă a concluziilor comparației, să putem să evaluăm „cât de mult” seamănă (sau se deosebesc) două (sau mai multe) structuri.

B. Încă de la început ne dăm seama că o comparație reală este multiaxială – putem compara diferite aspecte: forme, dimensiuni, proprietăți de tot felul etc. În plus, un anumit grad de arbitrar (sau convențional) va fi mereu prezent, pentru a grada diverse nivele de deosebiri. Însă, chiar dacă instrumentele pe care le construim sunt imperfecte, ele constituie un suport puternic în interpretarea cantităților imense de informații care însoțesc studiile materiei vii.

C. Metodele elaborate pentru analiza comparativă a două secvențe s-au dovedit a ocupa o poziție centrală în dezvoltarea bioinformaticii, motiv pentru care acestui capitol i se acordă o deosebită importanță în toate manualele și tratatele de bioinformatică.

7.1.2. *Evenimente*

A. Când comparăm două secvențe urmărim adesea evidențe că ar fi putut proveni dintr-un ancestor comun. Din acesta structurile diverg prin procese de mutație și selecție. „Evenimentele” sau procesele mutaționale de bază sunt:

- substituții: un element (nucleotid în cazul acizilor nucleici, respectiv aminoacid în cazul proteinelor) este înlocuit cu altul,
- inserții: unul sau mai multe elemente sunt introduse într-o secvență,
- deleții: unul sau mai multe elemente sunt eliminate dintr-o secvență.

B. În cazul comparației a două secvențe, unei inserții în secvența X îi corespunde în secvența Y un loc gol, notat „-” și numit „gap”. În cazul unei deleții din secvența X, vom plasa „gap”-urile în secvența X, ele având corespondent real Y. Vedem deci că inserțiile și delețiile sunt tratate similar, motiv pentru care s-a introdus și termenul „indel”.

7.1.3. *Termeni*

A. În analiza secvențială putem întâlni diferite variante de “asemănări” (similaritate), pentru care se folosesc termeni specifici.

Termenul ce desemnează similaritatea se folosește ca adjectiv și poate fi:

- a) *homolog* (omolog):
 - când secvențele provin din 2 organisme diferite, cu ancestor comun
- b) *ortolog*:
 - sunt diferențe datorită speciației

- se reține funcționalitatea în evoluție
 - c) *paralog*:
 - din 1 organism, evenimente apărute la duplicare
 - d) *xenolog*:
 - nu au aceeași origine prin evoluție
 - apar prin evenimente „orizontale” (simbioză, viruși etc.)
- B. Să remarcăm că „similaritatea” \neq „homologie”.

7.2. Graficele de puncte „dot plots”

7.2.1. Principiul metodei

Este o metodă simplă și poate cea mai veche [Maizel și Lenk].

A. Definiție: un grafic de puncte (dot plot) este o reprezentare vizuală a similarității între două secvențe, în care cele două secvențe se plasează pe cele două axe ale unei matrice sau ale unui grafic și se marchează prin puncte toate elementele comune.

B. În forma primară („fără filtre”), un grafic de puncte pentru compararea a două secvențe de acizi nucleici (aceste secvențe conțin doar 4 simboluri), vor apărea multe „zgomote” – coincidențe întâmplătoare, datorită repetitivității simbolurilor.

7.2.2. Filtre

A. De aceea, pentru eliminarea sau reducerea acestor zgomote, se pot aplica „filtre” de diferite lungimi. Astfel, un „filtru de 2 elemente” va permite marcarea unei coincidențe numai dacă regăsim în ambele secvențe un element împreună cu succesorul său (adică se identifică și se marchează toate grupele de câte 2 simboluri). Putem extinde fereastra la 3 sau mai multe elemente, ajungând să evidențiem într-adevăr anumite secvențe care apar în ambele structuri.

B. În cazul graficelor aplicate pentru proteine, zgomotele sunt mai rare (aceste secvențe conțin 20 de simboluri), în plus, putem introduce chiar scări cu diferențe nuanțe de gri în funcție de rezultatul comparării regiunilor din ferestrele preluate din cele două secvențe. Vom reveni la aceste reprezentări după ce vom prezenta primele „modele de scor”.

C. În figura 7.2.2.a se prezintă graficul de puncte între secvențele X și Y, fără filtru iar în 7.2.2.b este aplicat un filtru de două elemente pentru aceleași secvențe.

D. Graficele de puncte sunt foarte sugestive, fiind recomandate ca un pas inițial de analiză. Din grafic se observă dacă va fi vorba de o similaritate globală sau locală, dacă sunt repetiții sau inversiuni etc. Se folosesc în acest sens și grafice ale unei secvențe comparate cu ea însăși, în care, pe lângă diagonala care apare evident, se pot evidenția unele repetiții.

E. Există software dedicat pentru dot-plots: programul „dotter” ce poate fi descărcat de pe serverul ftp EBI.

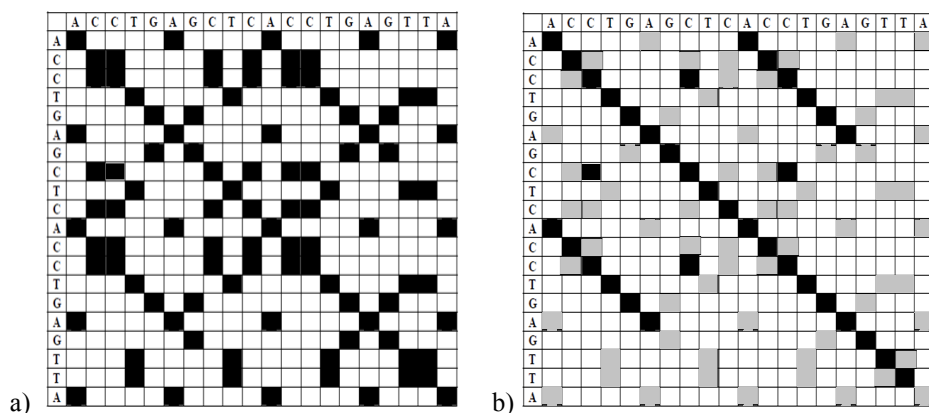


Fig. 7.2.2.: a – dot plot fără filtru, b – dot plot cu filtru de 2

7.3. „Distanțe” între secvențe

7.3.1. Noțiunea de aliniere

Vizualizarea prin „dot-plots”, deși este foarte sugestivă, nu caracterizează cantitativ gradul de similaritate sau deosebire între secvențe. În acest sens s-au propus varii metode de estimare cantitativă a asemănării între secvențe, pe care le vom studia în continuare.

A. Toate pornesc de la **alinierea secvențelor**, pe care o vom defini aici.

Definiție: Alinierea secvențelor (sequence alignment) este o metodă de comparare a două (sau mai multe secvențe) prin scrierea lor una sub alta (element) cu scopul de a evalua similaritatea între aceste secvențe.

În acest capitol ne vom limita la alinierea a două secvențe, urmând apoi a generaliza metoda pentru aliniere multiplă.

Situațiile cu care ne putem întâlni sunt următoarele:

- dacă același simbol apare în ambele secvențe vom considera că poziția s-a conservat în evoluție,
- dacă simbolurile sunt diferite, se presupune că cele două derivă dintr-un ancestor comun, care poate fi unul din cele două sau chiar altul decât cele două,
- secvențele comparate pot avea lungimi diferite datorită posibilelor inserții sau deleții.

B. Alinierea globală și locală

În cazul în care comparăm două secvențe în totalitatea lor, vorbim de **aliniere globală**.

În cursul evoluției au putut să apară diferite inserții, deleții sau chiar mutări ale unei secvențe dintr-o zonă cromozomială în alta, astfel încât avem uneori porțiuni care se pretează bine la comparații urmate de porțiuni care nu au nimic în comun. În aceste situații este preferabil a limita comparația numai pentru anumite porțiuni din secvențe – așa numita „**aliniere locală**”.

7.3.2. Distanțe

A. Distanța Hamming

Este cea mai simplă exprimare a similarității între două secvențe, fiind dată de:

$$DH = \text{numărul total de nepotriviri} \quad (7.3.2.a)$$

Se folosește pentru alinierea globală.

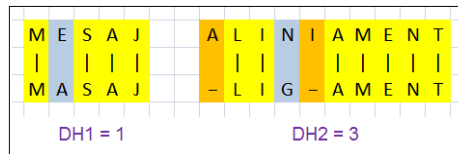


Fig. 7.3.2.a Distanța Hamming

B. Distanța Levenshtein

Se poate îmbunătăți evaluarea similarității dacă vom diferenția evenimentele:

- Conservarea o vom cota cu un punctaj pozitiv.
- Substituțiile și gap-urile le vom cota cu un punctaj negativ, penalizând mai puternic gap-ul decât substituția. Obținem astfel o schemă de scor. Totalul obținut prin aplicarea unei astfel de scheme se numește „distanță Levenshtein”.

Un model de scor utilizat frecvent este:

$$S = \begin{cases} +2 & \text{pentru potrivire (conservare)} \\ -1 & \text{pentru nepotrivire (substituție)} \\ -2 & \text{pentru gap} \end{cases} \quad (7.3.2.b)$$

$$DL = \sum S_i \quad (7.3.2.c)$$

C. Exemplu:

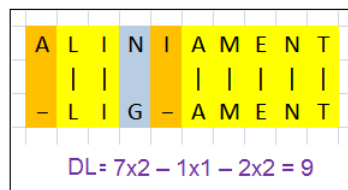


Fig. 7.3.2.b. Distanța Levenshtein

7.3.3. „Problemele” analizei secvențiale

A. Să mai luăm următorul exemplu (figura 7.3.3.).

Am scris normal prima secvență (PARAGINITUL), apoi am scris a doua secvență sub prima literă până aproape de sfârșit, când după R sub I am văzut că, dacă introducem un gap, următoarea literă, U, se va potrivi.

Calculând distanța Levenshtein obținem:

$$DL_3 = +2 \times 6 - 1 \times 2 - 2 \times 4 = +2$$

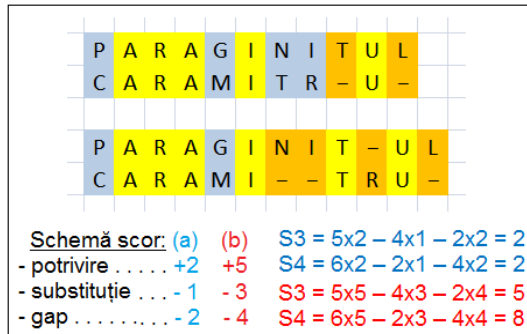


Fig. 7.3.3. Distanța Levenshtein cu 2 scheme de scor.

Observăm însă că, după potrivirea lui I sub I, în a doua secvență avem un T, pe care l-am putea forța, introducând două gap-uri, să se alinieze sub T din prima secvență; mai introducemos un gap sus să îl aliniem și pe U sub U. Obținem deci varianta de dedesubt, cu 3 gapuri în secvența de jos.

Este această aliniere mai bună? Avem într-adevăr mai multe conservări, dar și mai multe gap-uri.

$$DL_4 = +2 \times 6 - 1 \times 2 - 2 \times 4 = +2$$

Printr-o coincidență, obținem aceeași distanță Levenshtein, deci, cele două alinieri ar părea similare.

B. Putem încerca însă și o altă schemă de scor, coloana (b), folosită destul de frecvent:

$$S = \begin{cases} +5 & \text{conservare} \\ -3 & \text{substituție} \\ -4 & \text{gap} \end{cases} \quad (7.3.3.a)$$

Acum avem:

$$DL_3 = +5 \times 5 - 3 \times 4 - 4 \times 2 = +5$$

$$DL_4 = +5 \times 6 - 3 \times 2 - 4 \times 4 = +8$$

Deci cu această schemă, alinierea a doua ar fi mai bună.

C. Concluziile pe care le putem trage din acest exemplu:

i) Rezultatul comparației depinde de schema de scor; avem deci nevoie de scheme de scor cât mai realiste, din care să se înlăture, pe cât posibil, orice arbitrar.

ii) Există mai multe alinieri posibile între secvențe, fiind necesari algoritmi pentru selecția alinierii optime.

În continuare vom studia alegerea alinierii optime în ipoteza unei scheme de scor date, iar în finalul capitolului vom vedea cum putem construi o schemă de scor realistă.

7.4. Programare dinamică – Algoritmul Needleman – Wunsch

7.4.1. Alinierea globală

Așa cum am văzut mai sus, problema centrală din analiza secvențială este găsirea unui algoritm care să permită selecția alinierii optime între două secvențe, în cadrul unei scheme de scor date.

Primul și cel mai cunoscut algoritm pentru analiza globală a două secvențe a fost elaborat de Needleman și Wunsch în 1970. El poate fi aplicat pentru orice fel de secvențe (atât acizi nucleici cât și proteine).

Vom descrie în continuare algoritmul, pas cu pas, mergând în paralel cu un exemplu.

7.4.2. Principiul algoritmului

A. Plasând cele două secvențe pe cele două axe ale unei matrice, ne propunem să calculăm un element de matrice $F(i, j)$, unde i reprezintă coloana și j linia.

Considerăm cele două secvențe de comparat:

X , cu elementele $X_i, i = 1, n$

Y , cu elementele $Y_j, j = 1, m$

B. Construim o matrice $M \times N$ cu m linii și n coloane, plasând pe axe – pe axa liniilor secvența X , iar pe cea a coloanelor secvența Y :

X Y	X_1	X_2	...	X_{i-1}	X_i	...	X_n
Y_1	$F(1,1)$	$F(2,1)$...		$F(i,1)$...	$F(n,1)$
Y_2	$F(1,2)$	$F(2,2)$	
...
Y_{i-1}			...	$F(i-1, j-1)$	$F(i, j-1)$...	
Y_i	$F(1, j)$...	$F(i-1, j)$...	
...
Y_m	$F(1,m)$						$F(n,m)$

Fig. 7.4.2.a. Matricea elementelor programării dinamice.

C. Valoarea elementului de matrice $F(i,j)$ se calculează după următorul raționament:

- natural este ca elementul $F(i,j)$ să se calculeze din precedentul din secvență, adică din $F(i-1, j-1)$, acordând un „bonus” pentru similaritatea $X_i \equiv Y_j$, sau aplicând o „penalizare” pentru substituția lui X_i cu un element $Y_j \neq X_i$.

Valoarea „bonus”-ului sau a „penalizării”, notată $S(X_i, Y_j)$ este dată, în general, într-o matrice de substituție”. În exemplele anterioare noi am folosit niște valori generice: +2 (sau +5) pentru conservare, respectiv - 1 (sau - 3) pentru **orice** substituție. Ar fi desigur, preferabil ca toate aceste valori să fie diferite, ținând cont de gradul de asemănare a proprietăților diferiților substituți; vom scrie în acest caz:

$$F(i,j) = F(i-1, j-1) + S(X_i, Y_j) \quad (7.4.2.a)$$

- este însă posibil să fie preferabil un „grup” în una din secvențe, pentru a asigura o potrivire mai bună în elementele următoare din secvență. Acest lucru este greu de anticipat, însă va ieși în evidență după completarea întregii matrice. De aceea, pentru fiecare element de matrice $F(i,j)$, calculăm și valoarea care ar corespunde în cazul în care într-una din secvențe vom introduce un gap. Astfel, dacă presupunem că am introduce un gap în secvențele Y, am avea:

$$F(i,j) = F(i-1, j) - d \quad (7.4.2.b)$$

Unde d este „penalizarea” pentru gap; în exemplele anterioare noi am luat valoarea - 2 (respectiv - 4).

- o În cazul unui gap în secvența X vom utiliza relația:

$$F(i,j) = F(i, j-1) - d \quad (7.4.2.c)$$

Schematic putem reprezenta calculul elementului din matrice $F(i,j)$ ca în figura 7.4.2.b.

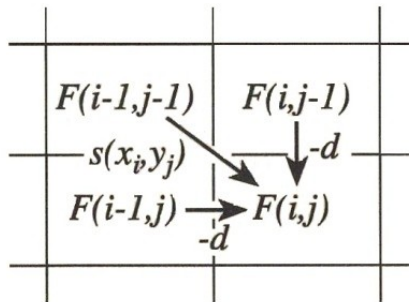


Fig. 7.4.2.b. Algoritmul Needleman-Wunsch.

Algoritmul Needleman – Wunsch pune în poziție centrală alegerea valorii maxime dintre cele trei posibilități prezentate mai sus, exprimată pe relația:

$$F(i,j) = \max \begin{cases} F(i-1, j-1) + S(X_i Y_j) \\ F(i-1, j) - d \\ F(i, j-1) - d \end{cases} \quad (7.4.2.d)$$

7.4.3. Marcarea traseului

A. Pentru a marca varianta aleasă, se obișnuiește a se plasa în căsuța $F(i, j)$ o săgeată, fie diagonală orientată spre stânga-sus „ \nwarrow ”, fie la stânga „ \leftarrow ”, fie în sus „ \uparrow ”, care să indice care a fost varianta utilizată. În cazul în care două variante de calcul dau aceeași valoare, vom plasa ambele săgeți.

B. Aplicând aceste reguli, pe rând, la toate elementele matricei vom completa matricea de scor pentru determinarea alinierii optime a celor două secvențe. Algoritmul prin care se realizează „alinieră optimă” se numește „trace-back”.

C. Se pornește de la elementul $F(n, m)$, din colțul dreapta-jos al matricei, urmărind săgețile care indică „din care element” a fost calculat. Vom marca astfel fiecare „pas”, pornind de la „coada” secvențelor spre început. În cele din urmă vom obține alinierea globală optimă a celor două secvențe.

D. Este posibil ca soluția să nu fie unică și să avem două (sau chiar mai multe) alinieri, caracterizate prin distanțe identice. Desigur matricea de substituție influențează scorul (valoarea „distanței”) pentru fiecare aliniere.

7.4.4. Exemplu

A. Să ilustrăm aplicarea algoritmului Needleman-Wunsch pe un exemplu simplu. Să considerăm secvențele ADN:

X: C G T A

Y: C T A

De asemenea vom aplica o matrice de substituție foarte simplă: $potrivire = + 1$, $nepotrivire = - 1$, $gap = - 2$, adică $S(i, j) = + / - 1$ și $d = 2$

B. Construim matricea de aliniere. Plasăm secvența X pe orizontală și Y pe verticală. Vom „inițializa” matricea adăugând în plus o linie și o coloană înainte de începerea fiecărei secvențe, corespunzătoare situației în care primul (primele) element(e) din oricare dintre secvențe ar fi un „gap”. Vom nota aceste elemente cu $F(i, 0)$ pentru linia inițială, respectiv $F(0, j)$ pentru coloana de inițializare.

	Col	(0)	(1)	(2)	(3)	(4)
Lin			C	G	T	A
(0)		0	-2	-4	-6	-8
(1)	C	-2				
(2)	T	-4				
(3)	A	-6				

Fig. 7.4.4.a. Inițializarea matricei de aliniere

Elementului $F(0,0)$ îi dăm valoare 0, apoi cazul în care gap-urile s-ar succeda, pentru fiecare element aplicăm o penalizare de gap, deci elementele liniei de inițializare vor avea valorile: $d, 2d, 3d$ etc. Similar pentru elementele primei coloane. Matricea inițializată în cazul exemplului nostru arată ca în fig. 7.4.4.a.

Calculăm primul element, $F(1,1)$; aplicăm pe rând cele 3 relații din (7.4.2.d) și obținem:

a: pornind din stânga-sus: $F(1,1) = F(0,0) + S(C, C) = 0 + (+1) = 1$

b: pornind din stânga: $F(1,1) = F(0,1) - d = -2 - 2 = -4$

c: pornind de sus: $F(1,1) = F(1,0) - d = -2 - 2 = -4$

Alegem valoarea maximă (cazul a): $F(1,1) = 1$ și marcăm săgeata corespunzătoare; rezultatul este vizibil în figura 7.4.4.b.

Continuăm cu elementul $F(2,1)$ pentru care, avem:

a: $F(2,1) = F(1,0) + S(G, C) = -2 - 1 = -3$

b: $F(2,1) = F(1,1) - d = 1 - 2 = -1$

c: $F(2,1) = F(2,0) - d = -4 - 2 = -6$

Alegem varianta b și plasăm și săgeata corespunzătoare.

	col	(0)	(1)	(2)	(3)	(4)
Lin			C	G	T	A
(0)		0	-2	-4	-6	-8
(1)	C	-2	+1	← -1	← -3	← -5
(2)	T	-4	↑ -1	↖ 0	↖ 0	← -2
(3)	A	-6	↑ -3	↖ -2	↖ -1	↖ +1

Fig. 7.4.4.b. Completarea matricei de aliniere

Procedăm similar pentru toate elementele de pe linia (1):

$S(T, C): F(3,1) = \max(-4 - 1, -1 - 2, -6 - 2) = \max(-5, -3, -8)$

$S(A, C): F(4,1) = \max(-6 - 1, -3 - 2, -8 - 2) = \max(-7, -5, -10)$

Pentru linia a doua:

$S(C, T): F(1,2) = \max(-2 - 1, -4 - 2, +1 - 2) = \max(-3, -6, -1)$

$S(G, T): F(2,2) = \max(+1 - 1, -1 - 2, -1 - 2) = \max(0, -3, -3)$

$S(T, T): F(3,2) = \max(-1 + 1, 0 - 2, -3 - 2) = \max(0, -2, -5)$

$S(A, T): F(4,2) = \max(-3 + 1, 0 - 2, -5 - 2) = \max(-2, -2, -7)$

În fine pentru ultima linie:

$F(1,3) = \max(-4 - 1, -6 - 2, -1 - 2) = \max(-5, -8, -3) = -3$

$F(2,3) = \max(-1 - 1, -3 - 2, 0 - 2) = \max(-2, -5, -2)$

$F(3,3) = \max(0 - 1, -2 - 2, 0 - 2) = \max(-1, -4, -2)$

$F(4,3) = \max(0 + 1, -1 - 2, -2 - 2) = \max(+1, -3, -4)$

Matricea de aliniere obținută este prezentată în figura 7.4.4.b

C. Pornim acum algoritmul „trace-back”.

Marcăm căsuța din colțul dreapta-jos $F(4,3)$. Acest element corespunde coincidenței la capătul secvențelor A – A. (figura 7.4.4.c).

Observăm că $F(3,4)$ avea o singură săgeată, în diagonală.

Vom alege deci în continuare elementul $F(3,2)$, corespunzător corespondenței T – T. Acest element are de asemenea o singură săgeată, indicând proveniența sa din $F(2,1)$.

	Col	(0)	(1)	(2)	(3)	(4)
Lin			C	G	T	A
(0)		0	-2	-4	-6	-8
(1)	C	-2	+1 ↖	← -1	← -3	← -5
(2)	T	-4	↑ -1	↖ 0	↖ 0	← -2
(3)	A	-6	↑ -3	↖ -2	↖ -1	↖ +1

Fig. 7.4.4.c. Pornirea "trace-back"

Dar acesta are o săgeată spre stânga, adică spre $F(1,1)$ (fig. 7.4.4.d). Cu alte cuvinte, după corespondența C – C ilustrată de $F(1,1)$, urmează în secvența Y un „gap”, având următoarea căsuță, $F(2,1)$ tot pe orizontală, adică lui „G” din secvența X nu-i corespunde nici un nucleotid în secvența Y.

	Col	(0)	(1)	(2)	(3)	(4)
Lin			C	G	T	A
(0)		0	-2	-4	-6	-8
(1)	C	-2	+1 ↖	← -1	← -3	← -5
(2)	T	-4	↑ -1	↖ 0	↖ 0	← -2
(3)	A	-6	↑ -3	↖ -2	↖ -1	↖ +1

Fig. 7.4.4.d. Rezultatul obținut prin "trace-back"

Rezultatul alinierii îl exprimăm în forma de mai jos (figura 7.4.4.e):

C	G	T	A	
C	-	T	A	
+1	-2	+1	+1	= +1

Fig. 7.4.4.e. Alinierea finală

Rezultatul obținut era de așteptat. Observăm că „gap-urile” apar vizibile în matricea de aliniere acolo unde nu se urmează traseul pe diagonală. Dacă apar căsuțe succesive selectate pe aceeași linie, atunci gap-ul se va găsi în secvența a doua (cea plasată pe axa verticală a matricii); căsuța cea mai din stânga grupului rămâne în secvență, cealaltă (celelalte, dacă sunt mai multe succesiv) vor reprezenta gap-uri. Similar se procedează dacă apar căsuțe succesive pe aceeași coloană; în acest caz gap-urile se vor găsi în secvența plasată pe axa orizontală a matricii de aliniere.

7.5. Alinierea locală. Algoritmul Smith – Waterman

7.5.1. Deosebiri față de algoritmul NW

A. Algoritmul Needleman – Wunsch prezentat mai sus se aplică pentru alinierea globală a două secvențe, adică a secvențelor întregi. Este însă cunoscut faptul că, în evoluția lor, structurile biologice au suferit numeroase modificări, însă au păstrat destul de frecvent anumite porțiuni cu modificări minore, conservând în felul acesta funcționalitatea acestor molecule, în timp ce alte porțiuni din moleculă pot avea evoluții divergente. De aceea, destul de des dorim să detectăm doar domeniile comune din proteine și să apreciem gradul de similaritate între aceste porțiuni. Acest gen de analiză se numește „alinie locală”, care este aplicabilă și pentru alinierea unor secțiuni extinse de ADN, însă cel mai des se aplică pentru detecția similarității între secvențe divergente, cu origine comună, în care anumite porțiuni se păstrează. Astfel, prin alinierea locală se vor exclude porțiunile terminale din secvențe, porțiuni ce ar fi avut numeroase gap-uri și un procent înalt de substituții.

B. Pentru alinierea locală, Smith și Waterman au propus în 1980 un algoritm asemănător cu algoritmul Needleman – Wunsch, însă introducem următoarele modificări:

- la inițializare, marginile vor avea numai valoarea 0 în loc de penalizările pentru gap-uri;
- în calculul unui element de matrice, vom compara, ca în algoritmul NW valorile obținute din cele trei căsuțe vecine de origine potențială (diagonală stânga-sus, căsuța din stânga, respectiv cea de sus), alegând valoarea maximă; însă dacă toate sunt negative, vom alocă valoarea 0; în felul acesta matricea de aliniere va avea elementele numai pozitive sau 0;
- la “trace-back” vom porni nu din colțul dreapta-jos, ci de la cea mai mare valoare;
- alinierea se poate opri oriunde în matrice, ajungând eventual la marginea superioară sau la cea din stânga în cazul în care vreuna din secvențele alinate prezintă domeniul de interes chiar la început.

7.5.2. Descrierea algoritmului SW

A. În mod formal putem scrie pentru algoritmul Smith – Waterman următoarele relații și notații:

- H = matricea de aliniere, cu elementele $H(i,j)$
- m = lungimea secvenței X (pe axa orizontală - sus)
- n = lungimea secvenței Y (pe axa verticală - stânga)
- la inițializare:

$$H(i, 0) = 0, 0 \leq i \leq m \quad (i = nr. \text{coloanei}) \quad (7.5.2.a)$$

$$H(0, j) = 0, 0 \leq j \leq n \quad (j = nr. \text{liniei})$$

B. Calculul elementelor matricei H urmează regula:

conservare/substituție

$$H(i, j) \max \begin{cases} H(i-1, j-1) + s(X_i, Y_j) & \text{conservare/substituție} \\ H(i-1, j) + g(X_i, -) & \text{deletie din } X \\ H(i, j-1) + g(-, Y_i) & \text{insertie in } X \\ 0 & \end{cases} \quad (7.5.2.b)$$

Uzual se ia penalizarea pentru gap aceeași pentru oricare dintre secvențele aliniat X și Y, adică

$$g(X_i, -) = g(-, Y_i) = -d \quad (7.5.2.c)$$

- valorile $s(X_i, Y_j)$ se iau din matricea de substituție s , la fel ca în cazul alinierii globale; de obicei valorile sunt pozitive pentru conservarea elementului în secvență (potrivire) și negative pentru substituție.

C. Exemplu

Luăm secvențele:

$$X: \text{GAATTCAGTTA} \quad (7.5.2.d)$$

$$Y: \text{CGGATCGA}$$

Matricea de substituție s (figura 7.5.2.a) este construită foarte simplu: conservare = +5, substituție = -3, gap = -4 (prescurtat: +5/-3/-4).

	A	C	G	T
A	5	-3	-3	-3
C	-3	5	-3	-3
G	-3	-3	5	-3
T	-3	-3	-3	5
gap = -4				

Fig. 7.5.2.a. Matricea de substituție

Procedăm asemănător cu alinierea globală:

- inițializăm prima linie și coloană (de indici 0) cu valoarea 0
- calculăm $H(1,1) = (0-3,0-4,0-4,0) = 0$
- continuăm cu toată linia 1, apoi și celelalte linii; în cazul în care valoarea aleasă provine din vreuna din căsuțele alăturate, vom marca acest lucru printr-o săgeată ce

indică originea valorii (pot fi și două săgeți); dacă valorile calculate din căsuțele vecine au fost negative, vom lua valoarea 0 și nu inserăm nici o săgeată.

În figura 7.5.2.b este prezentată matricea de aliniere după efectuarea acestor calcule.

	Col	r̂(0)	r̂(1)	r̂(2)	r̂(3)	r̂(4)	r̂(5)	r̂(6)	r̂(7)	r̂(8)	r̂(9)	r̂(10)	r̂(11)
Lin		G	A	A	T	T	C	A	G	T	T	A	
r̂(0)		0	0	0	0	0	0	0	0	0	0	0	0
r̂(1)	C	0	0	0	0	0	0	5	← 1	0	0	0	0
r̂(2)	G	0	↖ 5	← 1	0	0	0	1	↖ 2	6	← 2	0	0
r̂(3)	G	0	↖ 5	↖ 2	0	0	0	0	0	7	← 3	0	0
r̂(4)	A	0	↑ 1	↖ 10	↖ 7	← 3	0	0	5	↖ 3	4	0	↖ 5
r̂(5)	T	0	0	↖ 6	↖ 7	↖ 12	← 8	← 4	↑ 1	2	↖ 8	↖ 9	← 5
r̂(6)	C	0	0	↑ 2	↖ 3	↑ 8	↖ 9	↖ 13	← 9	← 5	↑ 4	↖ 5	↖ 6
r̂(7)	G	0	↖ 5	← 1	0	4	5	↖ 9	↖ 10	14	← 10	← 6	← 2
r̂(8)	A	0	↑ 1	↖ 10	← 6	← 2	↑ 1	↖ 5	↖ 14	← 10	↖ 11	← 7	↖ 11

Fig. 7.5.2.b. Matricea de aliniere

- începem procedura “trace-back”; observăm că valoarea maximă din matrice este 14, corespunzătoare elementului H(7,8), care provine în diagonală din H(6,7), legat în sus de H(6,6) și apoi de H(5,5). De aici sunt două săgeți, deci putem merge fie prin H(4,4) la H(3,4), fie prin H(4,5) și în diagonală la H(3,4); deci vom avea aici două trasee echivalente; mai departe mergem la H(2,3) și H(1,2). De fapt și la pornire, puteam pleca din H(8,7) care avea tot valoarea 14, ajungând în H(6,6), urmând de acolo traseul reprezentat anterior.

Structura aliniierilor posibile rămase sunt vizibile în figura 7.5.2.c.

	col	r̂(0)	r̂(1)	r̂(2)	r̂(3)	r̂(4)	r̂(5)	r̂(6)	r̂(7)	r̂(8)	r̂(9)	r̂(10)	r̂(11)
Lin		G	A	A	T	T	C	A	G	T	T	A	
r̂(0)		0	0	0	0	0	0	0	0	0	0	0	0
r̂(1)	C	0	0	0	0	0	0	5	← 1	0	0	0	0
r̂(2)	G	0	↖ 5	← 1	0	0	0	1	↖ 2	6	← 2	0	0
r̂(3)	G	0	↖ 5	↖ 2	0	0	0	0	0	7	← 3	0	0
r̂(4)	A	0	↑ 1	↖ 10	↖ 7	← 3	0	0	5	↖ 3	4	0	↖ 5
r̂(5)	T	0	0	↖ 6	↖ 7	↖ 12	← 8	← 4	↑ 1	2	↖ 8	↖ 9	← 5
r̂(6)	C	0	0	↑ 2	↖ 3	↑ 8	↖ 9	↖ 13	← 9	← 5	↑ 4	↖ 5	↖ 6
r̂(7)	G	0	↖ 5	← 1	0	4	5	↖ 9	↖ 10	14	← 10	← 6	← 2
r̂(8)	A	0	↑ 1	↖ 10	← 6	← 2	↑ 1	↖ 5	↖ 14	← 10	↖ 11	← 7	↖ 11

Fig. 7.5.2.c. Trace-back pentru alinierea locală

Vom sintetiza rezultatele în figura 7.5.2.d., în care am reprezentat separat variantele de drum din prima porțiune a alinierei, respectiv din a doua parte – în cele din urmă vom avea 4 variante posibile echivalente.

-	G	A	A	T	T	C	C	A	G	T	T	A	
C	G	G	A	-	T	C	C	-	G	-	-	A	
-	G	A	A	T	T	C	C	-	A	G	T	T	A
C	G	G	A	-	T	C	C	G	A	-	-	-	-

Fig. 7.5.2.d. Alinieri locale posibile

7.6. Modele complexe

7.6.1. Potriviri repetate (Repeated matches)

A. Algoritmul Smith – Waterman găsește (uzual) o singură cea mai bună aliniere între două secvențe. Este cunoscut însă faptul că secvențele reale conțin uneori mai multe copii ale unui domeniu repetat („motiv”). Ar fi deci de dorit să avem un algoritm care să evidențieze toate aceste regiuni. Există mai multe metode prin care putem realiza acest lucru, una dintre ele, bazată pe „lanțuri Markov ascunse” va fi menționată în capitolul respectiv. În continuare vom prezenta aici pe scurt o extensie a algoritmului Smith – Waterman, bazat pe programarea dinamică, prin care să putem evidenția repetările.

B. Să considerăm cele două secvențe de comparat $X(x_i, i=1, \dots, n)$ și $Y(y_j, j=1, \dots, m)$ și să presupunem că secvența Y conține un domeniu care se repetă. Să mai reamintim că și în acest caz, ca și în cele discutate până acum, porțiunile aliniată conțin nu numai elemente conservate ci și substituții sau gap-uri. Intuim deci dificultatea ce apare prin semnalarea oricăror scurte porțiuni din Y ce se aseamănă cu scurte porțiuni din X . De aceea, pentru a obliga algoritmul să semnaleze doar secvențele cu relevanță, se introduce un prag, T , (threshold), astfel încât numai regiunile în care scorurile depășesc pragul vor fi scoase în evidență. Pragul se introduce în relația de recurență pentru inițializarea primei linii. Astfel, în timp ce la algoritmul de aliniere locală $S. W.$ toate valorile erau zero: $H(i,0) = 0, i = 1, \dots, n$, în extensia SW pentru potriviri repetate vom lua $H(0,0) = 0$, apoi:

$$H(i, 0) = \max \begin{cases} H(i - 1, 0) \\ H(i - 1, j) - T, j = 1, \dots, m \end{cases} \quad (7.6.1.a)$$

În felul acesta secvența X va fi partiționată în regiuni care prezintă aliniere cu porțiuni din Y și regiuni fără potriviri.

Posibila prezență a unor valori diferite de zero pe prima linie (linia de inițializare) permite terminarea unor alinieri când se ajunge la marginea superioară a tabelului prin procedura de trace-back chiar dacă valoarea din căsuță nu este zero (în timp ce în algoritmul clasic SW , alinierea se considera încheiată întotdeauna când se ajungea la o căsuță cu valoarea zero - indiferent de poziția ei în matrice).

Vom ilustra detecția unei potriviri repetate printr-un exemplu.

C. Exemplu

Fie secvențele proteice X și Y definite astfel:

X : DTWGGICDAWHGL

(7.6.1.b)

Y : CDAWGEP

Obținem alinierea unei porțiuni din Y (zona [C]DAWG) cu două regiuni din X :

D	T	W	G	G	I	C	D	A	W	H	G	L
D	A	W	G	.	.	C	D	A	W	-	G	.

Fig. 7.6.1. Reprezentarea unei alinieri repetate

Am marcat cu un punct (.) regiunile de nepotrivire.

D. Procedura “trace-back” are niște particularități în acest caz. Pentru aplicarea ei procedăm astfel:

- introducem în matricea de aliniere o căsuță în plus pe prima linie: $H(n+1,0)$ cu o valoare calculată după aceeași relație de recurență (7.6.1.a); valoarea din această căsuță va reprezenta scorul total al tuturor potrivirilor, având valoarea T scăzută pentru fiecare potrivire (deci la k potriviri s-au scăzut kT puncte). Există și posibilitatea ca valoarea să fie 0, dacă pragul T nu a fost niciodată depășit.
- din căsuța $H(n+1,0)$ vom trasa o săgeată către căsuța din care provine (fie cea din stânga ei, fie o altă căsuță din coloana anterioară); putem deci ca în procedura “trace-back” să avem „salturi”.
- în continuare procedăm similar cu procedura aplicată în celelalte cazuri.
- de menționat că, de fiecare dată când ajungem la marginea superioară, dacă valoarea nu este 0 vom avea un salt la căsuța cu valoarea maximă din coloana din stânga.
- procedura se încheie la o căsuță cu valoarea 0, oriunde în matrice, cel mai adesea fie pe marginea superioară (deci restul din secvența X va fi considerat “gap” terminal), fie pe marginea din stânga (deci “gap-ul” terminal va fi din secvența Y).

E. O problemă importantă la detecția potrivirilor repetate este alegerea pragului T

- o valoare prea mare ar putea exclude unele potriviri, rămânând nedetectate.
- o valoare prea mică divizează mult secvențele semnalând și potriviri mai slabe, uneori întâmplătoare (mai ales în secvențele de acizi nucleici, unde alfabetul secvențelor are doar 4 litere).

7.6.2. Potriviri suprapuse

A. altă situație cu care ne putem întâlni este cazul în care o porțiune mai lungă dintr-o secvență are ea însăși zone repetitive, sau este inclusă în cealaltă secvență. Vom fi într-o situație asemănătoare cu cea de la potriviri repetate, cu deosebirea că zonele de potrivire din secvența de referință (notată anterior cu X) nu sunt disjuncte ci se pot suprapune parțial.

Putem întâlni astfel situații la compararea fragmentelor genomice de ADN sau la compararea unor secvențe cromozomale mari.

B. Pentru detecția acestor situații vom face o nouă extensie la algoritmiile de programe dinamică, pornind tot de la varianta Smith – Waterman pentru alinierea locală și vom ține cont și de utilitatea pragului T introdus pentru detecția repetărilor, în cazul în care urmărim și acest aspect.

- Ca elemente specifice ale algoritmului în această extensie avem:
 - păstrăm regula alinierii locale de a nu penaliza “gap-urile” terminale (început sau sfârșit de secvență)
 - marginile se inițiază cu 0, ca la algoritmul SW; există totuși versiuni în care se urmăresc și repetițiile, în care caz se aplică relația (7.6.) cu un prag T ales corespunzător
 - pentru calculul elementelor matricei de aliniere se aplică algoritmul Meedlemann – Wunsch al alinierii globale; vom putea avea deci și elemente negative
 - se alege valoarea maximă F_{\max} ca valoare maximă de pe linia de jos sau coloana din dreapta (deci nu cea mai mare din matrice)
 - procedura “trace-back” va începe de la F_{\max} și va merge până la marginea de sus sau din stânga
 - toate celelalte reguli rămân nemodificate.

C. Să luăm din nou un exemplu:

Fie secvențele X și Y din figura 7.6.2.

X:	D	A	W	T	L	A	C	E	L	P	A	C	S
Y:		A	C	T	L	A	C	H					
Y:						A	C	T	L	-	A	C	H

Fig. 7.6.2. Alinierea suprapusă a două secvențe proteice, X și Y

Se observă că apare o suprapunere între cele două potriviri găsite de algoritmul de aliniere.

7.6.3. Potriviri hibride (Hybrid match conditions)

A. Din cele prezentate anterior rezultă că metodele programării dinamice sunt destul de flexibile putând fi introduse extensii care să evidențieze elementele pe care dorim să le extragem în procesul de comparație a două secvențe. Fiecare metodă prezentată are avantajele și dezavantajele sale. De aceea, este bine a avea la dispoziție proceduri prin care putem scoate în evidență anumite particularități.

B. Enumerăm în continuare câteva situații care necesită abordări specifice, fără a intra în detalii privind aspectele teoretice ale acestor abordări. Pentru detalii privind extensiile algoritmilor în aceste cazuri recomandăm tratatul lui Durbin și colaboratorii.

Iată câteva exemple care necesită abordări specifice:

- secvență repetitivă care tinde să fie găsită în copii tandem neseperate,
- când căutăm secvențe ce încep la startul ambelor secvențe dar se pot termina în orice punct,
- când există o probabilitate ridicată ca o secvență să fie regăsită integral în cealaltă, dar și o probabilitate mare se a găsi numai un segment – așa numitele „căutări în familie”.

C. Din punct de vedere teoretic cazurile standard sunt limitate, dar putem găsi ceva apropiat. În plus, se pot efectua așa-numitele „post-procesări”, adică o rafinare a rezultatului oferit de algoritmi prin selecția variantelor cele mai plauzibile de aliniere, pentru care există alternative justificatoare.

7.7. *Gap-uri afine*

7.7.1. *Baze teoretice*

A. Un aspect important, dar neglijat până acum, a fost cel privind penalizarea gap-urilor. Variantele descrise mai sus au aplicat o penalizare uniformă pentru fiecare gap ce era „cerut” de algoritm pentru optimizarea scorului. Totuși, când legăm aspectul formal, exprimat prin relațiile propuse, cu fenomenul real, putem lesne observa că originea gap-urilor, generate de procese de inserție sau deleție în fenomenele de copiere (indiferent de nivel – replicare, transcriere, translație) cuprind cel mai adesea o porțiune dintr-o secvență, ce poate avea lungimi variabile, de la 1 component (nucleotid sau aminoacid) la zeci de componente. Deci „gravitatea” fenomenului generator de “gap” în aliniere, este mai curând legată de apariția “per-se” a fenomenului și mai puțin de dimensiunea sa (adică de numărul de elemente implicate). De aceea ar trebui să ținem cont de prezența „gap-urilor în lanț” și să ajustăm algoritmul în acest sens.

7.7.2. *Tipuri de gap-uri*

Din punct de vedere formal vom spune că avem două tipuri principale de penalizare a gap-urilor:

a) penalizare liniară, exprimată prin relația:

$$\gamma(g) = -gd \quad (7.7.2.a)$$

unde:

- g = nr. gap-uri
- d = penalizarea pentru un grup (este variant utilizată în prezentările anterioare)

b) gap-uri afine, pentru care penalizarea se face conform relației:

$$\gamma(g) = -d - (g - 1) \times e$$

unde:

- „ d ” este penalizarea pentru primul gap
- „ e ” este penalizarea pentru un gap din lanț; de obicei $e \ll d$.

Pentru a introduce aceste elemente, va fi nevoie să modificăm corespunzător și relațiile de recurență folosite atât în algoritmi de aliniere globală cât și în cei de aliniere locală, cu toate extensiile lor.

8. Matrici de substituție

8.1. *Introducere*

A. În prezentările anterioare ne-am concentrat atenția asupra algoritmilor prin care putem optimiza alinierea a două secvențe. De fiecare dată am acceptat ideea că, în cazul în care se observă conservarea (potrivirea) unui component într-o anumită poziție în cele două secvențe, să tratăm evenimentul ca fapt pozitiv, acordându-i puncte la scorul de aliniere, în timp ce prezența unui gap sau substituția – adică înlocuirea unui element cu altul – să fie penalizate. În acest capitol ne vom ocupa de modul în care stabilim penalizările ce vor fi utilizate în calculul scorului de aliniere, pentru substituția unui element dintr-o secvență, cu un alt element.

B. Ne vom ocupa separat de matricile de substituție utilizate în studiile privind acizii nucleici și cele din compararea secvențelor proteice.

Există o serie de asemănări fundamentale între abordările utilizate în cele două categorii de molecule, dar și unele deosebiri. Vom considera aici ca esențiale asemănările, astfel încât nu vom descrie separat elementele constructive pentru fiecare categorie de molecule.

8.2. *Matrici de substituție pentru proteine*

Este evident că matricile de substituție trebuie să acorde valori în funcție de tipul de substituție. Cea mai simplă variantă ar fi „uniformizarea” substituțiilor, variantă prin care toate substituțiile sunt penalizate la fel. Din exemplele prezentate anterior, conform notației formale [+2/-1/-2] sau [+5/-3/-4], valorile din schema de scor au semnificația [conservare/substituție/gap], deci noi am avut exemple în care substituția era penalizată fie cu 1 punct (față de +2 pentru conservare), respectiv cu 3 puncte (față de +5 pentru conservare). Oricum, ne putem imagina că au fost introduse o serie de simplificări, utile din punct de vedere didactic.

Totuși, scopul de aliniere, care este obiectivul principal în oricare din studiile de aliniere secvențială, depinde în mare măsură de valorile din matricea de substituție. De aceea, ne vom concentra atenția asupra modalităților prin care putem oferi o matrice cu valori rezonabile pentru a fi utilizate în algoritmii de aliniere secvențială (programare dinamică).

8.2.1. *Matrici PAM*

A. Denumirea PAM provine de la “**P**oint **A**ccepted **M**utations”, reprezentând probabilitățile mutațiilor. Să încercăm să explicăm această denumire.

- Ipoteza de bază: fiecare substituție a unui aminoacid (AA) cu alt AA este independentă de alte schimbări anterioare (este clar că ipoteza nu corespunde realității, însă reprezintă o bună aproximație pentru început).

Pentru a găsi niște valori rezonabile, cât mai aproape de realitate, Margaret Dayhoff (1978) a pornit de la analiza unor date privind substituțiile întâlnite evolutiv între familii

similare de proteine, extrapolându-le pentru distanțe evolutive mari. S-au utilizat numai date confirmate privind secvențele din 71 familii de proteine cu grad de similaritate de cel puțin 85%, totalizând 1572 schimbări; înrudirea secvențelor a fost analizată prin construcția unui arbore filogenetic, folosind metoda parsimoniei (pe care le vom prezenta în capitolul de analiză filogenetică). Prin acest studiu s-a putut construi o matrice conținând frecvențele tuturor perechilor de reziduuri între secvențe și ancesorii lor imediați din arbore. Din această matrice, printr-o serie de operații de normalizare și scalare în timp s-a obținut o matrice generică (figura 8.2.1.a), numită PAM1 conținând, pentru fiecare aminoacid probabilitatea de a fi substituit de un alt aminoacid într-un anumit interval de timp, când numărul așteptat de substituții a fost 1% (de unde denumirea PAM1 = Point Accepted Mutation 1%).

Example PAM1 matrix (normalized probabilities multiplied by 10000)

	Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	Ile	Leu	Lys	Met	Phe	Pro	Ser	Thr	Trp	Tyr	Val
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
Ala A	9867	2	9	10	3	8	17	21	2	6	4	2	6	2	22	35	32	0	2	18
Arg R	1	9913	1	0	1	10	0	0	10	3	1	19	4	1	4	6	1	8	0	1
Asn N	4	1	9822	36	0	4	6	6	21	3	1	13	0	1	2	20	9	1	4	1
Asp D	6	0	42	9859	0	6	53	6	4	1	0	3	0	0	1	5	3	0	0	1
Cys C	1	1	0	0	9973	0	0	0	1	1	0	0	0	0	1	5	1	0	3	2
Gln Q	3	9	4	5	0	9876	27	1	23	1	3	6	4	0	6	2	2	0	0	1
Glu E	10	0	7	56	0	35	9865	4	2	3	1	4	1	0	3	4	2	0	1	2
Gly G	21	1	12	11	1	3	7	9935	1	0	1	2	1	1	3	21	3	0	0	5
His H	1	8	18	3	1	20	1	0	9912	0	1	1	0	2	3	1	1	1	4	1
Ile I	2	2	3	1	2	1	2	0	0	9872	9	2	12	7	0	1	7	0	1	33
Leu L	3	1	3	0	0	6	1	1	4	22	9947	2	45	13	3	1	3	4	2	15
Lys K	2	37	25	6	0	12	7	2	2	4	1	9926	20	0	3	8	11	0	1	1
Met M	1	1	0	0	0	2	0	0	0	5	8	4	9874	1	0	1	2	0	0	4
Phe F	1	1	1	0	0	0	0	1	2	8	6	0	4	9946	0	2	1	3	28	0
Pro P	13	5	2	1	1	8	3	2	5	1	2	2	1	1	9926	12	4	0	0	2
Ser S	28	11	34	7	11	4	6	16	2	2	1	7	4	3	17	9840	38	5	2	2
Thr T	22	2	13	4	1	3	2	2	1	11	2	8	6	1	5	32	9871	0	2	9
Trp W	0	2	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	9976	1	0
Tyr Y	1	0	3	0	3	0	1	0	4	1	1	0	0	21	0	1	1	2	9945	1
Val V	13	2	1	1	3	2	2	3	3	57	11	1	17	1	3	2	10	0	2	9901

Fig. 8.2.1.a. Matrice generică PAM1

B. Pentru intervalele mai îndelungate procentul de mutații crește. Într-o accepțiune a liniarității fenomenului, se poate calcula o matrice pentru un procent n% de mutații acceptate prin ridicarea matricei PAM1 la puterea n:

$$PAM_n = (PAM1)^n \tag{8.2.1.a}$$

C. Pentru utilizarea matricilor de substituție în analiza secvențială este necesară încă o normalizare, prin care să ținem cont de frecvența cu care apare aminoacidul substituit în lanțurile proteice.

Să notăm un element din matricea PAM1 cu $P(b|a,t=1)$, reprezentând probabilitatea ca aminoacidul **a** să fie substituit cu **b** într-o unitate de timp arbitrară ($t=1$) și cu f_b

probabilitatea de apariție a aminoacidului **b**. Atunci un element din matricea de substituție, $S(a,b)$, reprezentând numărul de puncte ce se atribuie în cazul substituției aminoacidului **a** cu aminoacidul **b** este dat de relația (pentru orice timp t , pe care îl omitem):

$$S(a, b) = k \times \log \frac{P(b|a)}{f_b} \quad (8.2.1.b)$$

unde k este o constantă.

Observăm că s-a introdus o scară logaritmică, pentru a obține o scară aditivă.

D. Matricile PAM sunt foarte utilizate, fiind posibil a construi variante pentru diferite procente de similaritate. Astfel, pentru un procent de similaritate de 50% se folosește matricea PAM80, pentru 40% - PAM120 etc. Cea mai des folosită matrice este matricea PAM250, prezentată în figura 8.2.1.b (Să menționăm că indicele matricei PAM este legat puternic de intervalul de timp în care se acceptă mutația și nu reprezintă procentul de mutații!)

The PAM250 scoring matrix

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

Fig. 8.2.1.b. Matricea de substituție PAM250

Se observă că valorile au fost rotunjite la numere întregi și matricea a devenit simetrică. Valorile de pe diagonală reprezintă scorul acordat la conservarea aminoacidului pe poziția respectivă în secvență.

8.2.2. Matrici BLOSUM

A. Matricile PAM au fost construite pornind de la probabilitățile de substituție pentru intervale de timp relativ scurte. S-a constatat că extensia pentru intervale îndelungate (PAM250) se depărtează de realitate. Destul de frecvent ne întâlnim cu substituții succesive în aceeași poziție sau cu substituții în bloc, deci o îndepărtare de la ipoteza că fiecare substituție este independentă de alte substituții.

B. Pentru a înlătura aceste neajunsuri, Stephan și Georgia Henikoff (1992) au propus o nouă variantă de matrice de substituție numită BLOSUM “BLOcks (Amino Acid) SUBstitution Matrix”. Ei au analizat 2000 de patternuri de aminoacizi organizați în blocuri, pornind de la bazele de date conținând alinierea multiplă a unor proteine înrudite mai îndepărtat, luând regiunile fără gap-uri și incluzând în bloc porțiunile (clustere) pentru care procentul de reziduuri identice depășea un prag L. Calculând frecvențele cu care un reziduu **a** dintr-un cluster se alinia cu un reziduu **b** din alt cluster, s-au putut estima probabilitățile de substituție.

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

Fig. 8.2.2. Matricea de substituție BLOSUM62

C. Din punct de vedere formal, un element al matricei de substituție este dat de relația:

$$S(a, b) = (1/\lambda) \times \log[p(a, b)/f_a f_b] \tag{8.2.2.a}$$

unde λ este o constantă,

$p(a, b)$ – probabilitatea alinierii reziduuului a cu b

$f_a f_b$ – probabilitățile (frecvențele) aminoacizilor a și b

Relația (8.2.2.a) este similară cu (8.2.1.b) ținând cont de faptul că probabilitatea unei substituții a lui a cu b poate fi exprimată prin relația:

$$P(b|a) = p(a, b)/f_a \tag{8.2.2.b}$$

D. Matricile BLOSUM au valori destul de apropiate de matricile PAM echivalente, deși au pornit de la baze de date diferite și au fost construite în mod diferit.

În practică se folosesc uzual matricile BLOSUM62 și BLOSUM50. Varianta BLOSUM50, cu procent de similaritate de 50% corespunde unui timp evoluționar mai lung, deci este mai potrivită pentru secvențele mai divergente, conținând eventual și

gap-uri, în timp ce BLOSUM62, care este folosită și mai frecvent, este considerată standard pentru porțiuni fără gap-uri.

În figura 8.2.2 este prezentată matricea BLOSUM62.

8.3. Matrici de substituție pentru acizi nucleici

La fel ca și în cazul proteinelor, punctul de plecare pentru analiza substituțiilor pornește de la analiza evolutivă a structurilor moleculare din materia vie.

Matricile de substituție pentru acizii nucleici s-au construit înaintea celor pentru proteine, pornind chiar de la o abordare teoretică a modelelor evolutive.

8.3.1. Matricea Jukes – Cantor

A. Cea mai simplă versiune a unei matrici de substituție a fost abordată de Jukes și Cantor (1969). În modelul lor ei au considerat că:

- fiecare site dintr-o frecvență polinucleotidică poate fi tratat ca independent (adică substituția într-o poziție este independentă de alte substituții),
- toate substituțiile au aceeași probabilitate,
- nu sunt luate în considerare inserții au deleții,
- toate nucleotidele au aceeași probabilitate de apariție într-o frecvență,
- probabilitatea unei substituții este constantă în timp.

B. Din punct de vedere formal, dacă notăm cu α probabilitatea unei substituții într-un interval de timp t (numită rata substituției), atunci matricea probabilităților acestor substituții vor fi cele din figura 8.3.1.a, unde A,C,G,T sunt simbolurile celor 4 tipuri de nucleotide ale ADN.

	A	C	G	T
A	$1 - 3\alpha$	α	α	α
C	α	$1 - 3\alpha$	α	α
G	α	α	$1 - 3\alpha$	α
T	α	α	α	$1 - 3\alpha$

Fig. 8.3.1.a. Matricea probabilităților substituțiilor în modelul Jukes – Cantor

C. Dacă am lua probabilitatea unei substituții de 1% am obține matricea generică PAM1 pentru ADN conform modelului Jukes – Cantor prezentată în figura 8.3.1.b.

	A	C	G	T
A	0.99			
C	0.0033	0.99		
G	0.0033	0.0033	0.99	
T	0.0033	0.0033	0.0033	0.99

Fig. 8.3.1.b. Matricea generică PAM1 în modelul Jukes – Cantor

D. Aplicând și în acest caz o relație logaritmică, cu logaritm în baza 2, obținem matricea de substituție PAM1 prezentată în figura 8.3.1.c.

	A	C	G	T
A	2			
C	-6	2		
G	-6	-6	2	
T	-6	-6	-6	2

Fig. 8.3.1.c. Matricea de substituție PAM1 pentru ADN în modelul Jukes – Cantor

S-au luat valori întregi (s-a aproximat $1/300 \approx 2^{-8}$), și probabilitățile la echilibru $p_A = p_C = p_G = p_T = 1/4$.

8.3.2. Matricea Kimura

A. Condițiile simplificatoare ale modelului Jukes – Cantor au permis o abordare mai facilă a substituțiilor în acizii nucleici, dar cu riscul unor aproximări mai grosolane.

O primă îmbunătățire a fost realizată de Kimura (1980) care a luat în considerare marea diferență între probabilitățile de tranziție și cele de transversie. (Reamintim aici că tranzițiile reprezintă substituții în cadrul aceleiași clase – purinică sau pirimidinică, adică substituții între A și G, respectiv între C și T, iar transversiile sunt substituții între clase, adică A sau G cu C sau T). Probabilitățile transversiiilor sunt mai mici decât ale tranzițiilor, fiind vorba de substituții între molecule cu configurație, dimensiune și proprietăți diferite.

O schemă a acestor substituții, împreună cu probabilitățile lor în unitatea de timp este prezentată în figura 8.3.2.a.

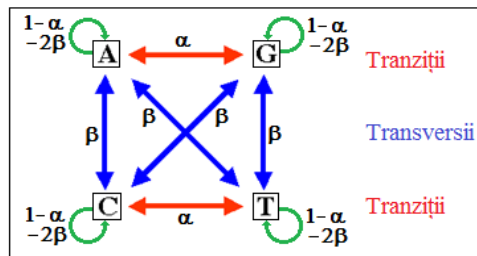


Fig. 8.3.2.a. Tranziții și transversii

B. Din punct de vedere formal matricea probabilităților tuturor substituțiilor posibile se scrie asemănător ca mai sus. Vom nota cu α rata tranzițiilor și cu β rata transversiiilor ($\beta \ll \alpha$). Atunci matricea probabilităților substituțiilor conform modelului Kimura arată ca cea din figura 8.3.2.b.

	A	C	G	T
A	$1 - \alpha - 2\beta$	β	α	β
C	β	$1 - \alpha - 2\beta$	β	α
G	α	β	$1 - \alpha - 2\beta$	β
T	β	α	β	$1 - \alpha - 2\beta$

Fig. 8.3.2.b. Matricea probabilităților substituțiilor în modelul Kimura

C. Putem și în acest caz să construim o matrice PAM1. Vom lua între α și β raportul $\alpha = 3\beta$. În acest caz matricea generică PAM1 va fi cea din figura 8.3.2.c.

	A	C	G	T
A	2			
C	-6	2		
G	-6	-6	2	
T	-6	-6	-6	2

Fig. 8.3.2.c. Matricea generică PAM1 pentru ADN conform modelului Kimura

D. Matricea *log odds* corespunzătoare celei din figura 8.3.2.c. este cea din figura 8.3.2.d., aplicând din nou logaritmi în baza 2 și rotunjind la valori întregi.

	A	C	G	T
A	2			
C	-7	2		
G	-5	-7	2	
T	-7	-5	-7	2

Fig. 8.3.2.d. Matricea de substituție PAM1 pentru acizi nucleici conform modelului Kimura

8.4. Testarea semnificației alinierii

8.4.1. Scoruri

A. După cum am observat din abordările anterioare, abordarea probabilistică ocupă o poziție importantă în studiile de bioinformatică. De aceea, va trebui să luăm în considerare și faptul că anumite alinieri între frecvențe au o anumită probabilitate de a apărea absolut din întâmplare. Dorim deci să facem o evaluare a acestei probabilități pentru a oferi un anumit nivel de confidență concluziilor noastre.

B. Se poate demonstra că, în cazul în care încercăm să aliniem două secvențe, X și Y, de lungime n , și respectiv m , și acordăm tuturor alinierilor posibile un scor, atunci numărul E al alinierilor cu un scor minim S este dat de relația:

$$E = k \times n \times m \times e^{-\lambda \times S} \quad (8.4.1.a)$$

unde k și λ sunt constante, fiind parametrii statistici cu scorul S. Putem face și o reprezentare grafică, luând

$$S = 10 \times \log x \quad (8.4.1.b)$$

și vom obține o distribuție Poisson, reprezentată în graficul din figura 8.4.1.

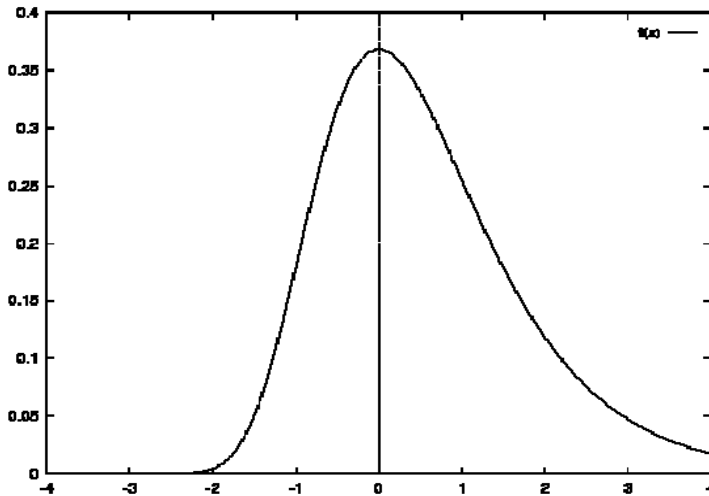


Fig. 8.4.1. Distribuția scorurilor normalizate după lungime

8.4.2. Semnificația scorurilor

A. Pentru semnificația alinierii să reamintim că pentru exprimarea cantității de informație dintr-un scor folosim ca unități de măsură:

- bit-ul, când se lucrează cu logaritmi în baza 2
- nat-ul, când se lucrează cu logaritmi naturali.

B. Conversia la biți se face astfel:

$$S' = (\lambda S - \ln k) / \ln 2 \quad (8.4.2.a)$$

și atunci vom avea pentru E relația:

$$E = n \times m \times 2^{-S'} \quad (8.4.2.b)$$

C. Exprimăm acum valoarea probabilității de a obține un scor minim S la întâmplare:

$$P = 1 - e^{-E} \quad (8.4.2.c)$$

D. Deseori nici nu mai calculăm probabilitatea din relația (8.4.2.c), ci doar impunem condiția ca exponentul să fie cât mai mic ($E \ll 1$) adică E să fie semnificativ mai mic decât unitatea. Prin logaritmarea relației (8.4.1.a.) putem scrie condiția

$$S > T + \frac{\log(mn)}{\lambda} \quad (8.4.2.d)$$

T fiind o constantă dependentă de k.

Valorile care se obțin în mod curent pentru P din relația (8.4.2.c.) sunt foarte mici, de ordinul 2^{-10} până la 2^{-100} . Pentru valori $P > 0.1\%$ nu se mai consideră asemănări între secvențe, eventualele coincidențe se vor atribui întâmplării.

9. Alinierea multiplă

Analiza secvențială cuprinde, gradat, analiza secvențelor individuale, compararea a două secvențe sau compararea mai multor secvențe, numită „alinieră multiplă”, prescurtat MSA (Multiple Sequence Alignment).

9.1. *Semnificația alinierii multiple*

Pornind de la datele experimentale biologiei au izbutit să realizeze manual alinierea de bună calitate a mai multor secvențe, bazându-se pe o serie de cunoștințe acumulate din experiență.

9.1.1. *Factorii luați în considerare pentru MSA*

Biologii includ pe lista factorilor importanți în alinierea, multiplă atât cunoștințe privind evoluția secvențelor proteice cât și elemente structurale. Să trecem în revistă acești factori:

- alinieri specifice „pe coloane” ale unor componente ce prezintă un înalt grad de conservare sau reziduuri hidrofobe mascate,
- influența structurii secundare și terțiare, cum ar fi alternanța componentelor hidrofile și hidrofobe pe coloanele fâșiilor beta,
- pattern-uri de inserții și deleții care tind să alterneze cu blocuri de secvențe conservate,
- constrângeri impuse de relații filogenetice.

9.1.2. *Obiective în alinierea multiplă*

A. Dezvoltarea metodelor automate pentru alinierea multiplă ocupă o poziție centrală în bioinformatică. În principiu ele își propun să atribuie alinierii un scor care să reflecte calitatea alinierii. Totuși trebuie să facem distincție între alinierea „optimă” după un anumit criteriu de scor și alinierea cea mai bună ținând cont și de criterii structurale și de evoluție. De aceea se insistă actualmente pe modalitățile de a transpune în algoritmi de scor criteriile biologice.

B. În principiu, în MSA scriem secvențele una sub alta aliniind reziduurile „omoloage” în coloane. Termenul „omolog” (*homologous*) este luat atât în sens structural cât și de evoluție. Se apreciază că elementele plasate pe o coloană ocupă poziții similare în structura tridimensională a moleculelor și toate provin cel mai probabil de la un ancestor comun. Identificarea structurală și de evoluție a pozițiilor omoloage este dificilă și adesea găsim mai multe alinieri potențial corecte (deci avem unele ambiguități), datorită imposibilității ca două secvențe cu unele diferențe să fie superpozabile în întregime. Unul dintre cele mai bune exemple, în care structura se conservă foarte bine chiar dacă apar diferențe secvențiale este oferit de „familia globinelor”.

C. Chiar dacă – în principiu – admitem că există o singură aliniere corectă evolutiv, în practică este chiar mai dificilă inferența evolutivă decât cea structurală. Aceasta datorită faptului că, pentru alinierea structurală avem criterii clare (cum ar fi

coordonatele spațiale din structura tridimensională determinate prin studii cristalografice sau RMN), în timp ce pentru alinierea evolutivă istoria unei familii de secvențe nu se determină prin metode independente. În plus, secvențele diverg mai rapid decât structurile.

D. În cazul alinierii secvențelor foarte asemănătoare, adesea se poate găsi o aliniere neambiguă, însă, în general, pentru cazurile de interes trebuie să reținem că nu există vreun mijloc obiectiv pentru a defini o aliniere corectă neambiguă. Adesea avem în familiile analizate proteine ce au în comun doar vreo 30% secvențe cu alinieri locale în perechi de bună calitate.

Ne ajută însă deseori câte un mic subset de reziduuri cheie identificabile, ce pot fi aliniate fără ambiguități.

E. Testarea calității alinierii trebuie să țină cont de aceste elemente. În loc să forțăm algoritmi să producă alinieri identice cu cele ce pot fi obținute manual, atenția pe subseturile din coloanele ce corespund reziduurilor cheie.

9.2. Scop și motivație pentru alinierea multiplă

Vom prezenta în continuare sintetic punerea problemei în alinierea multiplă și o serie de aplicații în care MSA și-a dovedit utilitatea.

9.2.1. Punerea problemei

A. Fiind dat un set de mai multe secvențe și o metodă de scor pentru aliniere, să se determine corespondențele între secvențe astfel încât scorul de aliniere să fie maxim.

B. Revenim aici la comentariul anterior privind criteriile incluse în scor; avem deseori posibilitatea de a aborda aceleași secvențe prin mai multe scheme de scor.

9.2.2. Motivație

Rezultatele MSA pot fi folosite pentru:

- stabilirea datelor de intrare pentru analiza filogenetică,
- determinarea istoriei evolutive a unui set de secvențe, precizând în ce punct au apărut anumite mutații,
- detecția unor „motive” comune într-un set de secvențe (de exemplu, secvențe de AND care leagă aceeași proteină),
- caracterizarea unui set de secvențe (de exemplu o familie de proteine),
- construcția profilelor pentru căutarea bazelor de date de secvențe (de exemplu, Psl-BLAST).

9.3. Scoruri pentru alinierea multiplă

9.3.1. Privire generală

A. Algoritmi pentru stabilirea unor scoruri de MSA pornesc de la o serie de ipoteze simplificatoare. O ipoteză folosită în majoritatea metodelor este independența coloanelor individuale.

B. Formula generală a unui scor MSA este:

$$Score(m) = G + \sum S(m_i) \quad (9.3.1.a)$$

unde: G este o funcție de gap, iar i este indicele care precizează coloana, deci $S(m_i)$ este un scor al coloanei „ i ”.

Dintre metodele uzuale ne vom opri pe scurt la:

- suma perechilor,
- entropia minimă.

9.3.2. Suma perechilor

Prin această metodă se calculează suma scorurilor din alinierea perechilor, folosind relația:

$$S(m_i) = \sum_{k < l} S(m_i^k, m_i^l) \quad (9.3.2.a)$$

unde am notat cu m_i^k caracterul (simbolul) din secvența „ k ” în coloana „ i ”. Elementele s sunt elementele matricei de substituție.

9.3.3. Entropia minimă

A. În această metodă ideea de bază este încercarea de minimizare a entropiei informaționale a fiecărei coloane. Sunt „bune” coloanele ce pot fi comunicate cu puțini biți (care au de fapt o bună conservare a unui element).

Reamintim, din teoria informației, faptul că un cod optim folosește un număr de $-\log_2 p$ biți pentru a codifica un mesaj/simbol cu probabilitatea p .

B. În MSA, mesajul este considerat pe coloană. Entropia informațională a unei coloane va fi:

$$S(m_i) = - \sum_a c_{ia} \times \log_2 p_{ia} \quad (9.3.3.a)$$

în care:

- m_i = coloana „ i ” din secvența/alinierea „ m ”
- c_{ia} = de câte ori apare caracterul „ a ” în coloana „ i ”
- p_{ia} = probabilitatea caracterului „ a ” în coloana „ i ”.

9.3.4. Reprezentări grafice

A. Pentru alinierea multiplă se folosește uzual o reprezentare intuitivă: plasarea secvențelor una sub alta, astfel încât să fie vizibile coincidențele, substituțiile și gap-urile. Deseori se marchează cu un chenar sau culori diferite coloanele importante (cu conservare puternică) (figura 9.3.4.a.).

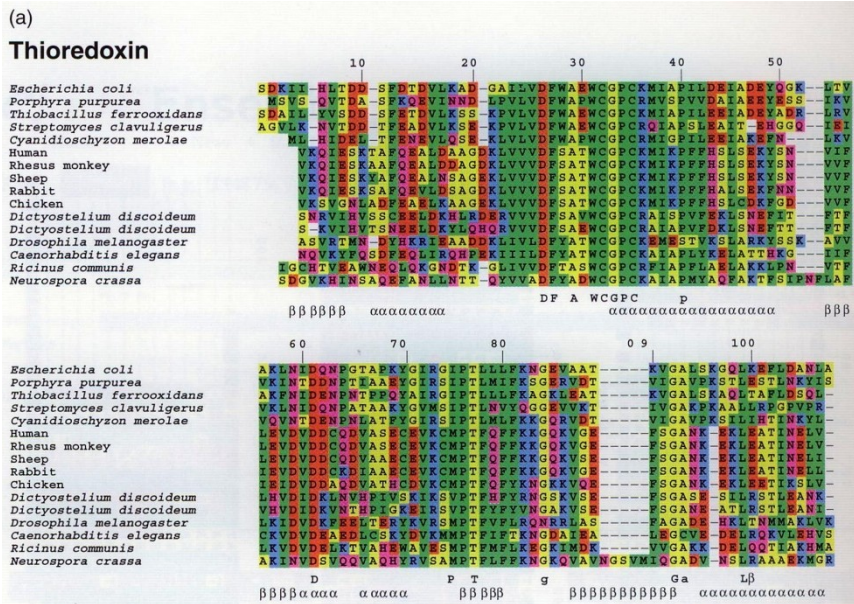


Fig. 9.3.4.a. Reprezentarea alinierii multiple cu coloane

B. În practică se mai folosește și o prezentare sugestivă, în care pentru fiecare poziție din secvență se scriu deasupra simbolurile găsite în acea poziție cu litere de diverse dimensiuni - cu cât simbolul respectiv a fost mai frecvent în acea coloană, cu atât litera este mai mare. În plus, literele pot avea și culori diferite (în cazul aminoacizilor) reprezentând caracterul polar/nepolar, acid/bazic. O astfel de reprezentare este redată în figura 9.3.4.b.

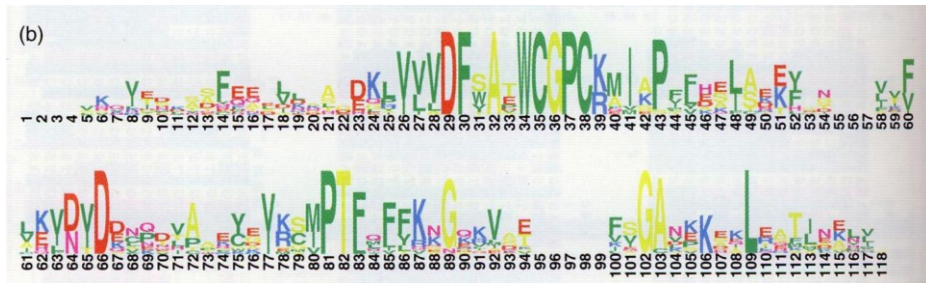


Fig. 9.3.4.b. Reprezentarea alinierii multiple cu simboluri

9.4. Algoritmi pentru alinierea multiplă

S-a dezvoltat o largă varietate de algoritmi pentru MSA, pornind adesea de la algoritmi dezvoltati pentru compararea a două secvențe.

Vom trece în revistă cei mai cunoscuți și utilizați algoritmi.

9.4.1. Programare dinamică

Succesul programării dinamice în cazul comparării a două secvențe a făcut ca algoritmi să fie generalizați pentru SA.

În cazul comparării a k secvențe ar fi necesară construcția unei matrici de dimensiune k , în care fiecare element să reprezinte scorul pentru k secvențe în loc de două. Numărul de combinații cu gap în oricare secvență (exceptând gap în toate) crește rapid cu k , fiind $2^k - 1$. Observăm deci că programarea dinamică devine complet nepractică.

9.4.2. Metode de aliniere progresivă

A. Principiul alinierii progresive

O abordare simplă a alinierii multiple este oferită de alinierea progresivă. Aceasta se construiește printr-o succesiune de alinieri pereche (compararea a două secvențe). Se iau două secvențe, se aliniază printr-una din metodele clasice (cel mai adesea prin programare dinamică, cu algoritmul Needleman – Wunsch, Smith – Waterman sau alte variante) din care se obține o nouă secvență care le sintetizează. Se ia apoi a treia secvență, care se compară cu sinteza primelor două obținând o nouă aliniere/secvență sintetică a celor 3. Procesul urmează iterativ până la epuizarea celor k secvențe.

B. Variante de aliniere progresivă

Au fost propuse mai multe variante de aliniere progresivă în funcție de:

- ordinea în care sunt introduse secvențele la aliniere, element care poate influența puternic rezultatul final,
- maniera de parcurgere a progresiei – adăugarea liniară a câte unei secvențe sau gruparea secvențelor întâi pe familii și reunirea ulterioară,
- procedurile de scor ale aliniierilor.

C. Caracterul euristic al alinierii progresive

Este evident faptul că în abordarea descrisă aici, procedura de optimizare nu este separată de calcularea scorului, neexistând o optimizare a unui scor global. Totuși, alinierea progresivă este rapidă, simplă și, dacă se urmează o procedură de ierarhizare în succesiunea secvențelor la aliniere, rezultatele sunt rezonabile.

Se recomandă alinierea la început a perechilor celor mai apropiate, existând și aici mai multe variante, care vor fi prezentate sumar.

9.4.3. Algoritmul MSA

Vom prezenta pe scurt un algoritm de construcție a unei secvențe sintetice. Vom considera cazul unor secvențe de acizi nucleici.

A. Reprezentarea unei secvențe prin matricea coeficienților

Fie o secvență $X(x_1x_2\dots x_i\dots x_n)$. Elementele din secvență au valori dintr-un alfabet al tipului de secvență. Pentru a acoperi cazul general din analiza secvențial, când într-o anumită poziție putem avea un gap, alfabetul are $m+1$ valori posibile, adică în cazul acizilor nucleici avem 5 valori posibile (a, c, g, t, -), iar în cazul proteinelor avem 21 valori posibile (cu „-” am notat un gap).

Vom nota vectorul valorilor posibile cu V . În cazul acizilor nucleici vom avea:

$$V(j) = \{a, c, g, t, \text{“-”}\} \quad (9.4.3.a)$$

adică $v_1 = V(1) = a$, $v_2 = V(2) = c$, etc.

Pentru orice secvență primară putem construi o matrice a coeficienților T astfel:

$$Tx(i, j) = \begin{cases} 1 & \text{daca } x_i = v_j \\ 0 & \text{daca } x_i \neq v_j \end{cases} \quad (9.4.3.b)$$

De exemplu, pentru secvența $X = (cgga-tg)$, matricea T(X) va avea valorile din tabelul (9.4.3.a)

Tabelul 9.4.3.a. Matricea coeficienților pentru secvența X

	c	g	g	a	-	t	g
a	0	0	0	1	0	0	0
c	1	0	0	0	0	0	0
g	0	1	1	0	0	0	1
t	0	0	0	0	0	1	0
-	0	0	0	0	1	0	0

B. Sinteza a două secvențe

Fie secvențele X și Y aliniate printr-unul din algoritmi clasici pentru perechi. Dacă nu au lungimi egale, vom completa secvența mai scurtă cu gap-uri.

Deci putem nota $X(x_1x_2...x_i...x_n)$ și $Y(y_1y_2...y_i...y_n)$. Construim matricea T_z a coeficienților secvenței sintetice z ce va rezulta în urma combinării X cu Y.

Calculăm elementele matricei T_z cu relațiile:

$$Tz(i, j) = [T_x(i, j) + T_y(i, j)]/2 \quad (9.4.3.c)$$

Să completăm exemplul început mai sus al secvenței $X = (cgga-tg)$, cu $Y = (tga-tt)$. Matricea coeficienților va fi (tabelul 9.4.3.b):

Tabelul 9.4.3.b Matricea coeficienților sintezei Z = (XY)

X	c	g	g	a	-	t	g
Y	t	g	a	-	-	t	t
a	0	0	0.5	0.5	0	0	0
c	0.5	0	0	0	0	0	0
g	0	1	0.5	0	0	0	0.5
t	0.5	0	0	0	0	1	0.5
-	0	0	0	0.5	1	0	0

C. Sinteza a k secvențe

Relația (9.4.3.c) se va generaliza devenind:

$$T_z(i, j) = [\sum_k T_k(i, j)]/k \quad (9.4.3.d)$$

D. Matricea de substituție

Față de matricile de substituție clasice, în cazul alinierii multiple vom introduce și gap-ul ca unul din simbolurile posibile. Să mai facem și observația că vom avea și situația în care un gap se conservă peste mai multe secvențe din set, având o penalizare mai mică decât un grup nou. Un model de matrice utilizabilă în alinierea multiplă este redat în tabelul 9.4.3.c.

Tabelul 9.4.3.c. Model de matrice de substituție pentru alinierea multiplă a secvențelor de ADN

	a	c	g	t	-
a	4	-2	-1	-2	-3
c	-2	4	-2	-1	-3
g	-1	-2	4	-2	-3
t	-2	-1	-2	4	-3
-	-3	-3	-3	-3	1

E. Programarea dinamică cu secvențe sintetice

Procedura de programare dinamică poate fi aplicată în cazul în care una sau ambele secvențe sunt sintetice, deosebirea că fiecare poziție din secvență intră în program doar cu coeficientul său de pondere.

9.5. Modele de ordonare

În metodele euristice de aliniere multiplă, cum este alinierea progresivă, ordinea în care sunt luate secvențele devine foarte importantă. În acest sens s-au propus mai multe variante. Vom prezenta în continuare două dintre ele, cel mai des folosite.

9.5.1. Modelul „stea”

- A. Se dau k secvențe ce urmează a fi aliniate: x_1, \dots, x_k .
- Se alege o secvență x_c ca și „centru”.
 - Pentru fiecare secvență $x_i \neq x_c$ se determină între x_i și x_c o aliniere optimală.
 - Se reunesc alinierea perechii.
 - Rezultatul de aliniere multiplă se obține prin agregarea alinierea perechii.
- B. Pentru alegerea centrului se încearcă fiecare secvență ca centru și se ia cea mai bună aliniere multiplă rezultată. Procedura este laborioasă însă oferă o alegere rezonabilă a centrului.
- C. Aplicarea modelului stea este ilustrată în figurile 9.5.1.a și 9.5.1.b

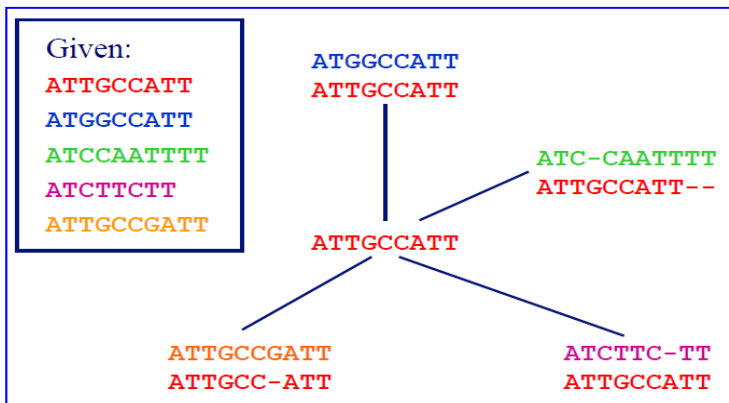


Fig. 9.5.1.a. Aplicarea modelului „stea” pentru un set de 5 secvențe

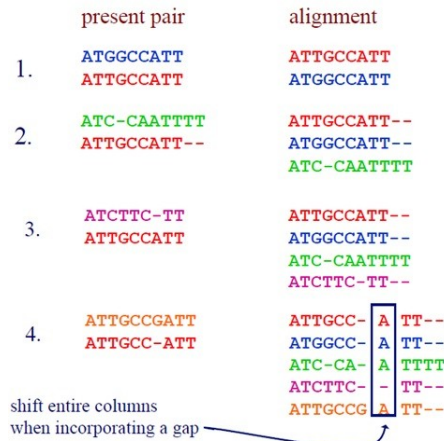


Fig. 9.5.1.b. Introducerea unui gap în modelul „stea”

9.5.2. Modelul „arbore”

Varianta care stă la baza unui soft folosit frecvent este modelul „arbore”.

A. Ideea de bază este organizarea alinierii folosind un „arbore ghid” (guide tree), în care:

- frunzele reprezintă frecvențele,
- nodurile interne reprezintă alinieri.

(Noțiunile fundamentale despre arbori vor fi introduse în capitolul de analiză filogenetică).

B. Alinierea de determină pornind de la bază în sus, iar alinierea multiplă rezultată la rădăcina arborelui.

C. Varianta uzuală este aplicată în programul CLUSTAL W (Thompson, 1994).

În funcție de nodul intern din arbore, putem avea de aliniat:

- o secvență cu o secvență
- o secvență cu un profil (alinieare parțială)
- un profil cu un profil.

Pentru cazul profilelor se recomandă scorul SP.

În figura 9.5.2 este prezentat un model de arbore folosit pentru alinierea multiplă.

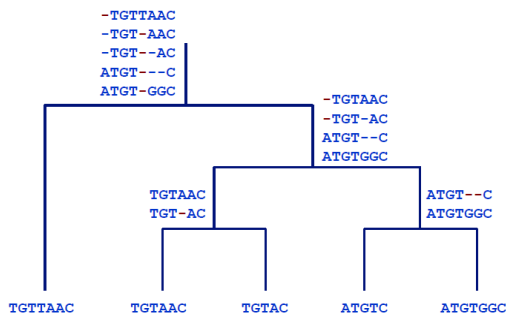


Fig. 9.5.2. Modelul „arbore” construit pentru alinierea multiplă a 5 frecvențe

10. Lanțuri Markov

În abordările anterioare am considerat că probabilitatea de a găsi un anumit element/reziduu (aminoacid în cazul unei secvențe de acid nucleic), este independentă de elementele ce ocupă pozițiile învecinate. Există însă numeroase date experimentale care arată că întâlnim adesea regiuni în care apar cu predilecție anumite perechi sau chiar microsecvențe. În cele ce urmează vom aborda analiza secvențială în acest context.

10.1. Lanțuri Markov simple

10.1.1. Descrierea unui lanț Markov

A. Pentru a lua în considerare dependența probabilității de a găsi un element într-o secvență de elementul anterior, vom construi o schemă aplicabilă în cazul secvențelor de acizi nucleici, în care alfabetul nostru este format din numai 4 litere: A, C, G și T. (figura 10.1.1)

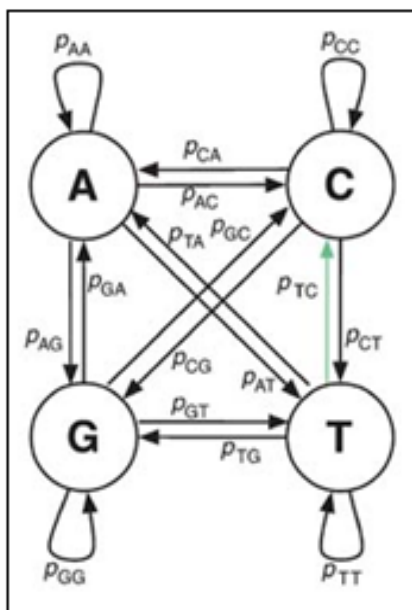


Fig. 10.1.1. Probabilitatea de succesiune în acizii nucleici

B. Am notat indicii probabilităților în ordinea în care apar elementele în secvență, de exemplu p_{AC} este probabilitatea ca C să fie precedat de A, cu alte cuvinte probabilitatea tranziției $A \rightarrow C$. Să mai menționăm aici că termenul „tranziție” în analiza secvențelor individuale, cum este cazul aici, are simpla semnificație ca un element să (de ex.: A) să fie urmat de un alt element (de ex.: C) și nu semnificația din matricile de substituție!

10.1.2. Probabilitatea unei secvențe

A. Din punct de vedere formal putem nota probabilitatea unei tranziții de la elementul „s” la „t” astfel:

$$a_{st} = P(x_i = t | x_{i-1} = s) = P(x_t | x_s) \quad (10.1.2.a)$$

Dacă notăm $P(x_s)$ probabilitatea de a găsi elementul „s”, atunci, probabilitatea de a găsi perechea „st”:

$$P(x_s, x_t) = P(x_t | x_s) \times P(x_s) \quad (10.1.2.b)$$

B. Proprietatea fundamentală a unui lanț Markov este că probabilitatea fiecărui element x_i depinde numai de simbolul precedent x_{i-1} , nu de întreaga secvență din față, adică:

$$P(x_i | x_{i-1}, \dots, x_1) = P(x_i | x_{i-1} = a_{x_{i-1}x_i}) \quad (10.1.2.c)$$

Putem aplica recurent (10.1.2.c) pentru o frecvență de lungime L , notată:

$X = [x_1 x_2 \dots x_L]$ și obținem:

$$P(X) = P(x_L | x_{L-1}) \times P(x_{L-1} | x_{L-2}) \dots P(x_2 | x_1) \times P(x_1) \quad (10.1.2.d)$$

$$P(x) = P(x_1) \prod_{i=2}^L a_{x_{i-1}x_i} P(x_i | x_{i-1}) \quad (10.1.2.e)$$

Deci probabilitatea unei secvențe se calculează din produsul tuturor tranzițiilor de la probabilitatea primului simbol din frecvență.

C. Pentru omogenizarea relației (10.1.2.e), adică introducerea primului simbol tot sub forma unei „tranziții”, s-a convenit completarea schemei din figura 10.1.1 cu încă o stare denumită „Început” (sau “Begin”) din care pot exista tranziții către oricare din simbolurile ce pot fi întâlnite în frecvență (figura 10.1.2). Observăm că starea „B” de început nu are decât săgeți emergente. În acest context vom putea nota probabilitatea ca secvența X să înceapă cu un anumit simbol, de exemplu „s”, ca pe o tranziție din starea „B” în „s”:

$$P(x_1 = s | B) = a_{Bs} \quad (10.1.2.f)$$

Relațiile (10.1.2.d.) și (10.1.2.e) devin:

$$P(X) = P(x_L | x_{L-1}) \times P(x_{L-1} | x_{L-2}) \dots P(x_2 | x_1) \times P(x_1 | B) \quad (10.1.2.d')$$

$$P(X) = \prod_{i=1}^L a_{x_{i-1}x_i} \quad (10.1.2.e')$$

Putem proceda similar și pentru încheierea secvenței, introducând încă o stare specială, denumită „Sfârșit” (sau “End”), care va avea numai săgeți incidente. De exemplu, dacă secvența X se va termina cu simbolul „t” vom putea marca acest lucru ca o tranziție din starea „t” în „E”:

$$P(E | x_L = t) = a_{tE} \quad (10.1.2.g)$$

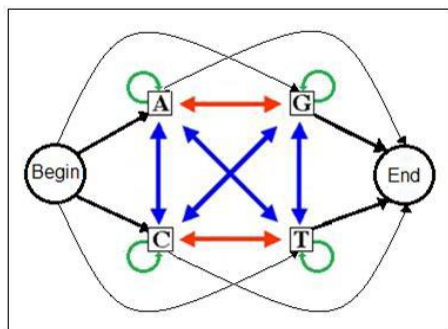


Fig. 10.1.2. Introducerea stărilor “Begin” și “End” la lanțul Markov pentru analiza secvențelor ADN

D. Exemplu

Să calculăm probabilitatea de a găsi secvența „cggg” la începutul unui lanț ADN.

Înlocuind în formula (10.1.2.d’) obținem:

$$P(cggt) = P(c|B) \times P(g|c) \times P(g|g) \times P(t|g) \times P(E|t) \quad (10.1.2.h)$$

Aceste probabilități sunt cunoscute; pentru porțiuni care nu sunt „insule CpG” (pe care le vom discuta în acest capitol) aceste valori sunt:

$$P(c|B) = 0.13, P(g|c) = 0.08, P(g|g) = 0.30, P(t|g) = 0.21$$

De obicei sfârșitul secvenței nu este introdus în model presupunându-se că secvența se poate termina oriunde; sfârșitul se introduce doar dacă trebuie considerată explicit lungimea frecvenței.

10.2. Estimarea parametrilor modelului

10.2.1. Punerea problemei

În cazul modelării secvențelor prin lanțuri Markov, este evident că probabilitățile tranzițiilor ocupă o poziție centrală.

A. Enunțul problemei: Fiind dat un set de date, D (de exemplu, o succesiune de secvențe), să se determine probabilitățile tranzițiilor.

Formal problema este determinarea parametrilor ϑ care maximizează $P(D|\vartheta)$, adică face setul de date D cel mai probabil.

Există mai multe abordări posibile ale problemei: fie prin metoda lui Bayes a probabilităților condiționate, fie prin metoda asemănării maxime. În cele ce urmează vom prezenta metoda asemănării maxime (maximum likelihood).

După cum se remarcă ușor din formulele prezentate mai sus, vom avea de folosit în mod frecvent produsul unor probabilități, iar aceste probabilități se vor obține din analiza unor secvențe care constituie setul de date D . Ce ne facem însă în cazul în care vreun simbol nu apare? În acest caz probabilitatea sa va fi zero și întregul produs se anulează! Desigur, astfel de situații ar fi rare, însă nu imposibile, mai ales când setul de date este mai restrâns. Pentru astfel de situații vom apela la proceduri prin care să aproximăm valorile probabilităților, evitând însă valoarea zero, aproximații numite „estimate Laplace”.

Să pornim însă cu abordarea clasică a exprimării probabilității de apariție a unui simbol.

B. Probabilități individuale

(i) Pornim de la cazul simplu în care setul de date este compus dintr-un anumit număr de secvențe.

Pentru moment să facem următoarele asumții:

- fiecare poziție este independentă de celelalte
- fiecare poziție este generată de aceeași distribuție multinomială.

(ii) În cazul în care ne referim la acizi nucleici, vom avea 4 componente (nucleotide) ale căror probabilități dorim să le exprimăm. În cazul ADN, vom reprezenta nucleotidele prin simbolurile „a, c, g, t”, deci dorim să exprimăm parametrii $P(a), P(c), P(g), P(t)$.

(iii) Să considerăm un set de secvențe și să notăm cu n_a numărul de apariții a simbolului „a” în aceste secvențe. Atunci:

$$P(a) = n_a / \sum n_i \quad (10.2.1.a)$$

unde indicele i ia valorile „a, c, g, t”, adică numitorul reprezintă de fapt numărul total al elementelor din toate secvențele.

$$\sum n_i = n_a + n_c + n_g + n_t \quad (10.2.1.b)$$

În mod similar vom avea relații pentru probabilitățile celorlalte tipuri de nucleotide: „c, g și t”.

(iv) Să luăm un exemplu.

Fie secvențele:

X: accgcgctta

Y: gcttagtgac

Z: tagccgttac.

Am luat aici secvențe egale ca lungime, însă la modul general ele au lungimi diferite. Observăm că putem calcula ușor probabilitățile celor 4 tipuri de nucleotide:

$$P(a) = 6/30 = 0.2 \quad P(g) = 7/30 = 0.233$$

$$P(c) = 9/30 = 0.3 \quad P(t) = 8/30 = 0.267$$

unde numitorul satisface relația (10.2.1.b), adică $30 = 6+9+7+8$.

10.2.2. Estimarea Laplace

A. Să continuăm exemplul anterior, însă într-o versiune în care un simbol dispare, de exemplu, să considerăm că în aceste secvențe toate nucleotidele „a” au fost substituie cu „g”. Secvențele ar arăta acum astfel:

X': gccgcgcttg

Y': gcttggtggc

Z': tggccgttgc.

(10.2.2.a)

Acum probabilitățile au devenit:

$$P(a) = 0/30 = 0$$

$$P(g) = 13/30 = 0.433$$

$$P(c) = 9/30 = 0.3$$

$$P(t) = 8/30 = 0.267$$

B. Estimarea Laplace se face provenind de la cunoștințele anterioare (“beliefs”) prin care presupunem că totuși fiecare component are o probabilitate nenulă de apariție. În acest scop vom introduce la numărătoare o unitate, care în realitate nu apare în secvențe, numită “pseudocount”. Dar, pentru a nu favoriza numai un component, vom acorda acest “pseudocount” pentru toate componentele.

C. Formal acest lucru se exprimă astfel:

$$P(a) = (n_1 + 1) / \sum(n_i + 1) \quad (10.2.2.b)$$

D. Reluând calculele de mai sus vom obține:

$$P(a) = (0+1)/(30+4) = 1/34 = 0.029$$

$$P(c) = (9+1)/34 = 10/34 = 0.294$$

$$P(g) = 14/34 = 0.412$$

$$P(t) = 9/34 = 0.264$$

Vedem că valorile de aici sunt o bună aproximație a celor reale, evitând valoare 0.

10.2.3. Estimarea Laplace generalizată

A. În varianta descrisă mai sus am introdus doar câte „1” pseudocount pentru fiecare simbol. Putem generaliza formula (10.2.m) introducând un număr „m” de pseudocounts, adică instanțe virtuale:

$$P(a) = (n_a + p_a m) / [(\sum n_i) + m] \quad (10.2.3.a)$$

unde p_a reprezintă „probabilitatea anterioară” (“prior probability”) a lui „a”.

B. De exemplu, putem pentru $P(c)$, în cazul $m = 8$ și $p_c = 0.25$

$$P(c) = (9 + 0.5 \times 8) / (30 + 8) = 11/38 = 0.289.$$

10.2.4. Estimarea parametrilor de ordin I

A. După ce am prezentat estimarea Laplace a probabilităților individuale, putem reveni la punctul de vedere prezentat la începutul descrierii lanțurilor Markov, adică să ținem cont de faptul că probabilitatea apariției într-o secvență a unui simbol depinde de pozițiile anterioare. În cazul în care vom lua în considerare doar 1 poziție anterioară, probabilitățile calculate le vom numi probabilități de ordin 1. Cu această convenție, probabilitățile individuale independente s-ar numi probabilități de ordin 0.

B. Probabilități Laplace de ordin 1.

Definiție: Probabilitatea de ordin 1 este probabilitatea unui element, „t”, de a fi precedat de un anumit element, „s”, adică probabilitatea secvenței „st”.

$$P(t|s) = (n_{st} + 1) / (\sum n_{si} + 1) \quad (10.2.4.a)$$

unde i ia toate valorile posibile ale simbolurilor.

C. Exemplu

Fie din nou secvențele din (10.2.2.a)

X': gccgcgcttg

Y': gcttggtggc

Z': tggccgttgc.

Iată un model de calcul aplicând (10.2.4.a):

$$P(c|g) = (7 + 1)/(12 + 4), n_{gc} = 7$$

$$n_{ga} + n_{gc} + n_{gg} + n_{gt} = 0 + 7 + 3 + 2 = 12$$

Valorile care se obțin pentru exemplul nostru sunt:

$$P(a|g) = (0+1)/(12+4)$$

$$P(a|c) = (0+1)/(7+4)$$

$$P(c|g) = (7+1)/(12+4)$$

$$P(c|c) = (2+1)/(7+4)$$

$$P(g|g) = (3+1)/(12+4)$$

$$P(g|c) = (3+1)/(7+4)$$

$$P(t|g) = (2+1)/(12+4)$$

$$P(t|c) = (2+1)/(7+4)$$

$$P(a|a) = (0+1)/(0+4)$$

$$P(a|t) = (0+1)/(8+4)$$

$$P(c|a) = (0+1)/(0+4)$$

$$P(c|t) = (0+1)/(8+4)$$

$$P(g|a) = (0+1)/(0+4)$$

$$P(g|t) = (5+1)/(8+4)$$

$$P(t|a) = (0+1)/(0+4)$$

$$P(t|t) = (3+1)/(8+4)$$

10.2.5. Lanțuri Markov de ordin superior

A. Generalizând definiția lanțului Markov de ordin 1 se pot defini lanțuri Markov de ordin superior. Deoarece un lanț Markov este rezultatul unui proces prin care se realizează secvența analizată, se mai folosește și termenul de „proces Markov”.

B. Definiție: Un proces Markov de ordin n este un proces stohastic în care fiecare eveniment depinde de n evenimente precedente, adică un simbol depinde de n simboluri anterioare.

Formal putem scrie acest lucru prin relația:

$$P(x_i | x_{i-1}, x_{i-2}, \dots, x_1) = P(x_i | x_{i-1}, \dots, x_{i-n}) \quad (10.2.5.a)$$

adică din întreaga secvență din față ($x_1 x_2 \dots x_{i-1}$) s-a reținut doar porțiunea de n elemente care determină următorul (al n+1 -lea).

C. Numărul de parametri ce trebuie determinați crește exponențial cu ordinul lanțului. Am văzut că pentru ADN, unde alfabetul are doar 4 simboluri, numărul de parametri (probabilități) pentru ordinul 0 era 4, iar pentru ordinul 1 era 4². În general, pentru ordinul „n” numărul de probabilități va fi 4ⁿ⁺¹.

D. Aplicație – lanțuri Markov de ordinul 5

O problemă importantă în utilizarea lanțurilor Markov este alegerea ordinului. Practic s-a găsit că ordinul 5 este cel mai potrivit pentru găsirea genelor, procedura folosită de către programul **GeneMark**, realizat de către Borodovski.

Să remarcăm în încheiere că se pot efectua analize secvențiale și cu lanțuri Markov neomogene, în care se consideră diferite distribuții în diferite regiuni de secvență.

10.3. “Open Reading Frames” (ORF)

10.3.1. Noțiunea de „cadru de citire”

A. Una dintre aplicațiile majore ale proceselor Markov este identificarea regiunilor care codifică sinteza unei protein, adică găsirea „genelor”.

Definiție: o secvență ADN care ar putea codifica o proteină se numește “Open Reading Frame” – ORF (cadru de citire).

B. Un ORF are următoarele proprietăți:

- începe cu un codon de start (ATG),
- se termină cu un codon de stop (TAG, TAA, TGA),
- nu are codon de stop intern,
- satisface anumite cerințe minime de lungime.



Fig. 10.3.1.a. Secvențe potențial ORF

C. Este cunoscut faptul că un codon este format din 3 nucleotide, deci unei secvențe de 3 nucleotide îi corespunde un aminoacid.

O problemă esențială în înțelegerea procesului de sinteză a proteinelor este modul în care este inițializat procesul. Dintr-o secvență de ADN, formată din două lanțuri complementare, putem – în principiu – să efectuăm 6 variante de citire (figura 10.3.1.b)

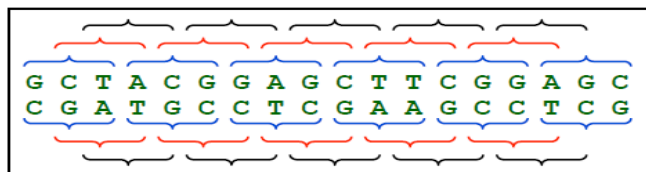


Fig. 10.3.1.b. Codificări posibile ale unei proteine dintr-un lanț de ADN

D. Exemplu

Să găsim un ORF în secvența de mai jos și să stabilim structura polipeptidului ce poate fi sintetizat.

Fie secvența:

CTGGATGATTCAGCTGGTCTAGTGACTAGC

Vom proceda după cum urmează:

- aplicăm o fereastră de căutare cu lățimea unui codon, adică 3 nucleotide, pe care o deplasăm poziție cu poziție până când găsim un codon de start (care corespunde metioninei)

CTGGATGATTCAGCTGGTCTAGTGACTAGC

- împărțim restul lanțului în secvențe de câte 3 elemente (în codoni)

- urmărim codonii de-a lungul lanțului până când întâlnim un codon de stop

CTGG ATG ATT CAG CTG GTC TAG TGACTAGC

- în final asociem fiecărui codon aminoacidul corespunzător conform codului genetic, adică avem pentapeptidul:

Met-Ile-Gln-Leu-Ala

sau, utilizând simbolurile din Bioinformatică:

MIQLA.

10.3.2. Metode pentru găsirea genelor

A. Un obiectiv important în bioinformatică este găsirea genelor ce determină anumite caracteristici inclusiv gene implicate în procese patologice.

Prin metodele prezente până acum putem vedea că avem o paletă mai largă, care cuprinde:

- căutarea prin similaritate de secvență (sss – Search by Sequence Similarity), adică se caută prin analiză secvențială potriviri cu secvențe cunoscute a fi legate de gene
- căutare prin semnal (Search by Signal), metodă bazată pe identificarea „semnalelor” implicate în expresia genei (v. semnalizare celulară)
- căutare prin conținut (Search by Content) – găsirea genelor prin proprietățile statistice ce disting ADN-ul ce codifică proteinele (exoni) de cel ce nu codifică (introni).

B. Căutarea prin metode markoviene, cu identificare de ORF-uri este inclusă în metoda căutării prin conținut. Să mai adăugăm următoarele observații:

- anumiți AA apar mai frecvent în exoni decât în introni (exemplu: Leu mai frecvent ca Trp)
- diferiți AA au un număr diferit de codoni (de exemplu, pentru Trp avem doar 1 codon, în timp ce pentru Leu sunt 6 codoni)
- pentru AA care sunt codificați prin mai mulți codoni, unii codoni apar mai frecvent decât alții, această „preferință” pentru codoni variază cu specia.

C. În tabelul 10.3.2 sunt trecute comparative câteva frecvențe de codificare ale unor AA.

Tabelul 10.3.2. Frecvențe de codificare a unor aminoacizi prin diferiți codoni

AA	codon	freq per 1000			
Gly	GGG	1.89	Glu	GAG	15.68
Gly	GGA	0.44	Glu	GAA	57.20
Gly	GGU	52.99	Asp	GAU	21.63
Gly	GGC	34.55	Asp	GAC	43.26

11. Modele Markov Ascunse

11.1. Suportul biologic al abordării

11.1.1. Insulele CpG

A. Un aspect interesant găsit în studiile privind frecvența unor microsecvențe de câte două nucleotide se referă la așa-numitele „insulele CpG”.

B. În cursul proceselor biologice se poate întâmpla ca citozina (C) să fie metilată și substituția sa cu timină (T) devine mult mai probabilă, acest fenomen fiind mai frecvent când citozina este urmată secvențial de guanină. Această succesiune citozină-guanină într-un lanț ADN se notează „CpG”, pentru a face distincție între această succesiune și notația CG folosită pentru a desemna perechea citozină-guanină din cele două lanțuri ale dublului helix ADN, menținut prin punțile de hidrogen.

C. Iată însă că natura a găsit o cale de a se apăra împotriva acestei probabilități ridicate de mutație, prin scăderea frecvenței microsecvențelor CpG în lanțurile ADN. În tabelul 11.1.1, secțiunea din dreapta, marcată cu semnul „-” observăm că frecvența reală a perechii „CG” este de 0.078, de cca. 4 ori mai mică decât o frecvență aleatoare (0.25)!

Tabelul 11.1.1. Secțiunea din stânga, notată cu „+” reprezintă frecvența perechilor de 2 nucleotide în „insulele CpG”, iar secțiunea din dreapta, notată „-” conține secvențele în regiunile neinsulare.

Pe axa verticală s-a notat primul nucleotid din pereche, iar pe cea orizontală, al doilea.

		P(c a)									
		a	c	g	t						
+	a	.18	.27	.43	.12	-	a	.30	.21	.28	.21
	c	.17	.37	.27	.19		c	.32	.30	.08	.30
	g	.16	.34	.38	.12		g	.25	.24	.30	.21
	t	.08	.36	.38	.18		t	.18	.24	.29	.29
		CpG						null			

D. Un alt aspect interesant este că în alte regiuni ale lanțurilor ADN, această „protecție” nu mai este menținută - este vorba mai ales de regiunile care au alte mecanisme de protecție, cum ar fi zonele „promoter”, care conțin secvențele de inițializare a proceselor de sinteză, zone care sunt de obicei controlate prin cuplarea unor proteine cu rol de control al mecanismelor de sinteză. Aceste regiuni în care frecvența perechilor microsecvențiale CG este apropiată de cea naturală, poartă denumirea de „insulele CpG”. Studiul a 48 astfel de insule a permis stabilizarea probabilității secvențelor din tabelul 11.1.1.

11.1.2. Distincția „stare” - „simbol”

A. Remarcăm deci, un fapt esențial: în secvențele reale ADN întâlnim regiuni cu probabilități diferite ale microsecvențelor – perechi de nucleotide succesive. Vom numi aceste regiuni - „stări”, iar apariția/generarea unui simbol va fi numită „emisie”. Este

evident că probabilitatea de „emisie” a unui simbol va fi dependentă de „starea” în care se găsește sistemul.

Din punct de vedere formal, emisia unui simbol este echivalentă cu „tranziția” definită în descrierea unui lanț Markov.

B. Vom introduce deci două feluri de probabilități:

a) probabilitatea ca sistemul să treacă dintr-o stare în alta; să notăm cele două stări „ l ” și „ k ”, starea inițială fiind k și cea după „tranziție” fiind l . Probabilitatea de tranziție $k \rightarrow l$ va fi dată de relația:

$$a_{kl} = P(\pi_i = l | \pi_{i-1} = k) \quad (11.1.2.a)$$

unde π_i și π_{i-1} reprezintă probabilitatea stării i , respectiv a stării anterioare $i-1$; să mai menționăm că probabilitatea de a rămâne în aceeași stare este $1 - a_{kl}$.

b) probabilitatea de „emisie” a unui simbol „ b ” când sistemul este în starea „ k ”:

$$e_k(b) = P(x_i = b | \pi_i = k) \quad (11.1.2.b)$$

11.2. Enunțul problemei în HMM

11.2.1. Denumirea HMM (Hidden Markov Model)

A. În modelele Markov clasice există o singură stare, cunoscută, și problema era să se determine probabilitatea emisiei unui simbol în funcție de simbolul emis anterior.

B. În cazul nostru problema are un aspect nou: simbolul („ b ”) poate fi emis fie în starea „ k ”, fie în starea „ l ”, cu probabilități diferite. Problema este trivială în cazul în care cunoaștem starea. Devine însă foarte interesantă problema de a determina starea în care a fost sistemul când a emis simbolul. Datorită faptului că tocmai starea este necunoscută când se generează o succesiune de simboluri fiind posibilă și trecerea sistemului – tranziția – dintr-o stare în alta, a generat denumirea de „model Markov ascuns” (Hidden Markov Model), în sensul că starea în care s-a emis un element a fost necunoscută (ascunsă).

11.2.2. Reprezentarea schematică a unui HMM

A. Să considerăm un exemplu în care sistemul poate intra fie în starea „1”, fie în „2”, fiecare cu probabilități specifice de emisie a celor 4 simboluri dintr-un lanț ADN.

După emiterea unui simbol, sistemul poate fie să rămână în aceeași stare, (1, respectiv 2), cu probabilități cunoscute, fie trecerea în noi stări: 3 din 1, sau 4 dacă starea anterioară a fost 2.

Din nou, sistemul emite către un simbol, după care poate rămâne în aceeași stare (3 sau 4), sau poate încheia emisia, trecând în starea 5 (end) – figura 11.2.2.

B. Problema principală este: determinarea succesiunii stărilor, pornind de la o secvență emisă.

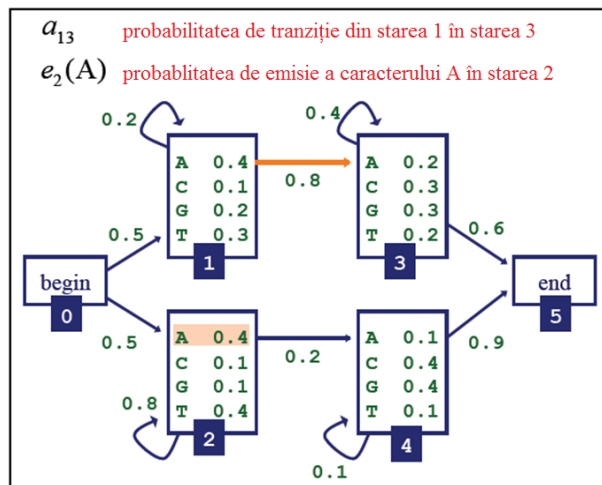


Fig. 11.2.2. Exemplu de HMM

11.3. Algoritmi de calcul pentru HMM

Există mai multe metode prin care putem aborda problema detecției succesiunii stărilor care au generat o frecvență dată, cunoscute fiind atât probabilitățile de tranziție între stări, precum și probabilitățile de emisie ale tuturor simbolurilor, pentru toate stările.

- Algoritmul “Forward” – prin care se calculează cât de probabilă este o secvență dată;
- Algoritmul “Viterbi” prin care se determină care este cea mai probabilă „cale” (secvență de stări ascunse = path) pentru a genera secvența dată;
- Algoritmul Baum – Welch (“Forward - Backward”), care încearcă să răspundă la întrebarea: „cum putem afla parametrii modelului HMM dintr-un set de secvențe date?”

În cele ce urmează vom prezenta în detaliu doar algoritmul Viterbi – cel mai folosit.

11.3.1. Algoritmul Viterbi

A. Algoritmul propus de Andrew Viterbi (1967) pornește de la niște presupuneri simplificatoare:

- evenimentele observate precum și cele ascunse sunt într-o secvență, ce corespunde – cel mai adesea – cu o succesiune temporală,
- cele două secvențe – a evenimentelor observate, respectiv a stărilor – trebuie să fie aliniate și unei instanțe de eveniment observat îi corespunde doar unei instanțe de stare (eveniment ascuns),
- calculul celei mai probabile căi a secvenței stărilor până la un moment „t” depinde numai de evenimentul observat la momentul t și de cea mai probabilă secvență până la momentul anterior t-1.

Aceste condiții sunt în general îndeplinite în modelele HMM de ordinul 1.

B. Din punct de vedere formal, păstrăm notațiile din formulele (11.1.2.a și b) și mai notând:

- $x_0, \dots, x_t =$ șirul de observații
- $y_0, \dots, y_t =$ secvența stărilor - cea mai probabilă să fi produs observațiile.

Cu aceste notații putem scrie o relație de recurență pentru probabilitatea $V_{t,k}$ a celei mai probabile secvența de stări responsabile pentru primele $t + 1$ observații (am început indexarea de la 0), având k drept stare finală:

$$V_{t,k} = P(x_t|k) \times \max_{y \in Y} (a_{y,k} \times V_{t-1,y}), \text{ cu}$$

$$V_{0,k} = P(x_0|k) \times \pi_k \quad (11.3.1.a)$$

Observăm că produsele succesive de probabilități ; totdeauna subunitare, vor genera numere din ce in ce mai mici, astfel încât s-a propus exprimarea într-un spațiu logaritm, calculându-se $\log [V_{t,k}]$.

11.3.2. Exemplu

A. Considerăm două persoane, A și B (Ann și Bob) care stau în localități diferite dar comunică zilnic. Bob îi transmite activitățile dominante din fiecare zi, care pot fi: (plimbare/walk, cumpărături/shopping sau curățenie/clean) prescurtat (w, s, c). Probabilitatea fiecărei activități depinde de vremea de afară, adică „starea” sistemului, care poate fi {ploioasă/rainy sau însorită/sunny (R, S)}. Din comunicările telefonice, Ann știe că Bob a avut în trei zile succesive acțiunile: (w, s, c).

B. Problema ar fi: să se determine „starea” cea mai probabilă în fiecare din cele trei zile. Se consideră cunoscute probabilitățile stărilor R („Ploaie”), respective S („Soare”) în regiunea respectivă. Totodată se mai consideră cunoscute și probabilitățile de inițiere a oricărei acțiuni în diferite condiții meteorologice – figura 11.3.2.a.

C. Problema: determinarea succesiunii stărilor cunoscând o secvență output („emisiu”). Notații: Stări (R,S), Emisii (w, s, c).

Rezolvare:

$$\text{Ziua 1: } p'_1(w) = p(R) \times p(w|R) = 0.6 \times 0.1 = 0.06$$

$$p''_1(w) = p(S) \times p(w|S) = 0.4 \times 0.6 = 0.24$$

$$\text{Ziua 2: } p'_2(s|w) = p'_1(w) \times p(R|R) \times p(s|R) = 0.06 \times 0.7 \times 0.4 = 0.0168$$

$$p''_2(s|w) = p'_1(w) \times p(S|R) \times p(s|S) = 0.06 \times 0.3 \times 0.3 = 0.0054$$

$$p'''_2(s|w) = p''_1(w) \times p(R|S) \times p(s|R) = 0.24 \times 0.4 \times 0.4 = 0.0384$$

$$p''''_2(s|w) = p''_1(w) \times p(S|S) \times p(s|S) = 0.24 \times 0.6 \times 0.3 = 0.0432$$

$$\text{Ziua 3: } p(c|ws) = p''''_2(s|w) \times p(R|R) \times p(c|R) = 0.0384 \times 0.7 \times 0.5 = 0.01344$$

$$p(c|ws) = p''''_2(s|w) \times p(S|R) \times p(c|S) = 0.0384 \times 0.3 \times 0.1 = 0.001152$$

$$p(c|ws) = p''''_2(s|w) \times p(R|S) \times p(c|R) = 0.0432 \times 0.4 \times 0.5 = 0.00864$$

$$p(c|ws) = p''''_2(s|w) \times p(S|S) \times p(c|S) = 0.0432 \times 0.6 \times 0.1 = 0.002592$$

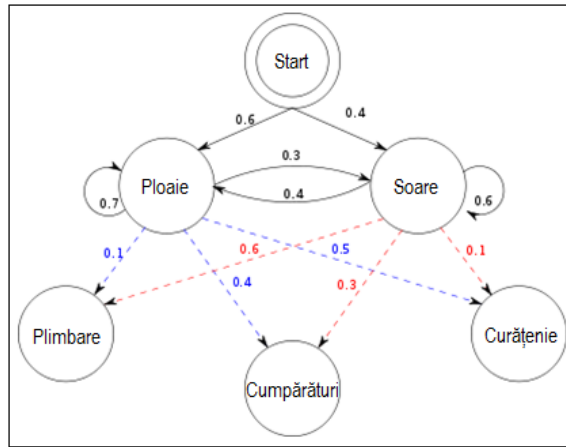


Fig. 11.3.2.a. Schema HMM pentru exemplul 11.3.2.

D. Vom alege acum valoarea cea mai mare pentru probabilitatea ultimei stări (este 0.1344 – corespunzătoare acțiunii de curățenie într-o zi ploioasă, efectuată după o zi de cumpărături într-o zi ploioasă, după ce anterior a fost la plimbare într-o zi însorită). Deci parcurgem drumul înapoi, exact ca la orice “trace back” din algoritmi de programare dinamică.

Acest lucru poate fi ilustrat prin niște diagrame specifice, numite “trellis”: o astfel de diagramă este prezentată în figura 11.3.2.b.

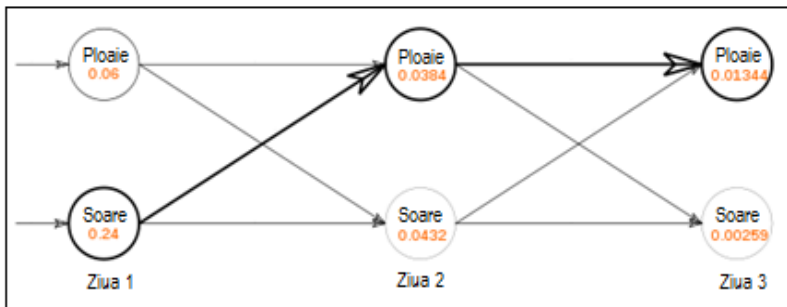


Fig. 11.3.2.b. Diagrama “trellis” pentru exemplul 11.3.2

11.4. Aplicarea HMM pentru discriminare

11.4.1. Notății

A. Modelele Markov Ascunse, HMM, permit identificarea succesiunii stărilor posibile care au generat o anumită secvență. Această capacitate poate fi utilizată pentru recunoașterea/discriminarea diferitelor regiuni din secvențele analizate.

Am menționat anterior (§ 11.1.1) existența „insulelor CpG”, iar în tabelul 11.1.1 am prezentat și probabilitățile tuturor perechilor de nucleotide în cele două tipuri de secvențe de ADN – din insule CpG (modelul notat cu „+”), respectiv din afara insulelor (modelul numit “null” și notat cu „-”).

- B. Notațiile folosite uzual în aceste analize sunt:
- probabilitatea de „tranziție în interiorul fiecărui model”, notată a_{st}^+ pentru insulele CpG (modelul „+”), respectiv a_{st}^- înafara insulelor:

$$a_{st}^t = c_{st}^t / \sum_{t'} c_{st}^t \quad (11.4.1.a)$$

unde c_{st}^+ arată de câte ori simbolul t a fost urmat de simbolul s (t și s au valori din alfabetul de 4 litere: A, C, G, T).

O relație similară se scrie pentru a_{st}^- .

De fapt, valorile din tabelul 11.1.1 sunt chiar acestea.

- C. Observăm asimetria tablourilor: A urmat de G este mult mai frecvent decât G urmat de A, mai ales în insule, iar T urmat de A apare mult mai rar decât GA și în insule și în afara lor.

11.4.2. Scoruri “log-odd”

- A. Pentru a utiliza modelele Markov pentru discriminarea diferitelor regiuni se introduce un scor, numit “raport log-odd”, $S(x)$, dat de relația:

$$S(x) = \log \frac{P(x|model^+)}{P(x|model^-)} = \sum_{i=1}^L \log \frac{a_{x_{i-1}x_i}^+}{a_{x_{i-1}x_i}^-} = \sum_{i=1}^L \beta_{x_{i-1}x_i} \quad (11.4.1.b)$$

în care x este o secvență de două simboluri, iar β este logaritmul raportului probabilităților corespunzătoare tranzițiilor în cele două modele, „+” și „-”.

Folosind logaritmi în baza 2, raportul β se va măsura în biți.

- B. Pentru discriminarea insulelor CpG față de celelalte regiuni, valorile lui β sunt cele din tabelul 11.4.2. prezentat mai jos:

Tabelul 11.4.2. Valorile lui β pentru discriminarea CpG

β	A	C	G	T
A	-0.740	0.419	0.580	-0.803
C	-0.913	0.302	1.812	-0.685
G	-0.624	0.461	0.331	-0.730
T	-1.169	0.573	0.393	-0.679

- C. Făcând o reprezentare grafică a distribuției scorurilor (nominalizate în funcție de lungimea secvențelor), se observă o discriminare rezonabilă între regiunile CpG față de celelalte regiuni (figura 11.4.2)

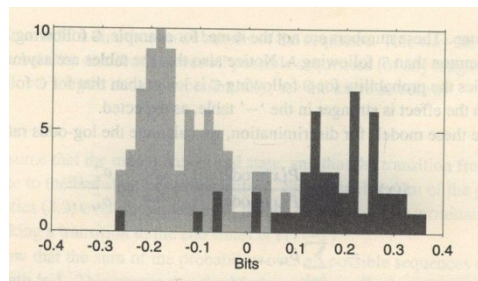


Fig. 11.4.2. Histograma scorurilor normalizate după lungime. Insulele CpG sunt reprezentate cu negru, iar celelalte regiuni cu gri

11.5. Modele Markov ascunse cu inserții și deleții

11.5.1. Stări silențioase

Într-un model HMM am întâlnit până acum două stări în care sistemul nu emitea nici un simbol erau Begin și End. Modelul poate fi generalizat prin introducerea unor stări de acest fel oriunde în lanț.

11.5.2. Schema HMM generală

Să notăm cu M_j stările succesive prin care trece sistemul, reprezentate prin pătrate, apoi I_j inserțiile, reprezentate în romburi și cu D_j delețiile, plasate în cercuri, vom obține o schemă generală de model Markov ascuns, aplicabil și în analiza secvențială multiplă modelul propus de Haussler, 1993 (figura 11.5.2).

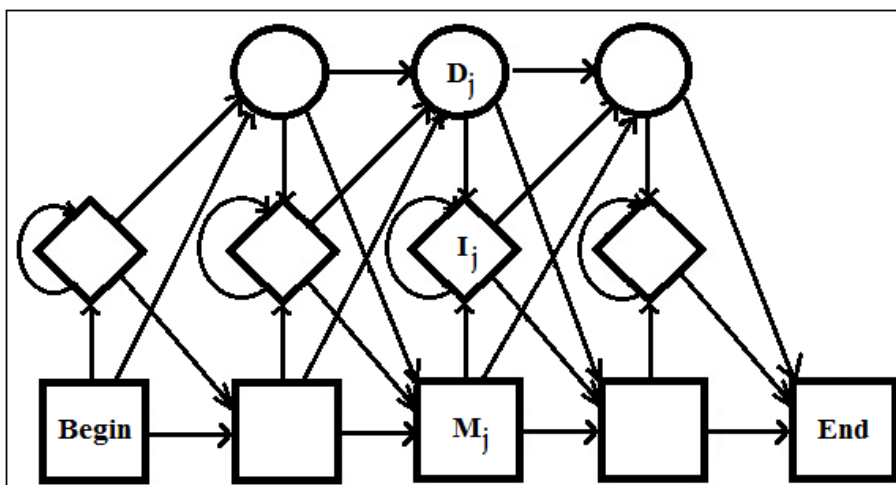


Fig. 11.5.2. Structura generală a unui profil HMM

12. Analiza filogenetică

12.1. Inferență filogenetică

12.1.1. Obiectivul analizei filogenetice

A. Unul dintre cele mai importante capitole ale bioinformaticii este reprezentat de analiza filogenetică. Am văzut și în capitolele anterioare că în cursul evoluției au loc diferite evenimente - deleții și substituții, însoțite de diverse modificări ale proprietăților moleculelor implicate. Uneori aceste modificări sunt minore, se conservă proprietățile importante precum și funcția în mecanismele celulare. Alteori aceste modificări sunt mai mari, pot fi reflectate și la nivel fenotipic și/sau funcțional, fiind considerate mutații. Analiza originii, evoluției și relațiilor între diferitele molecule reprezintă obiectul analizei filogenetice.

B. Să mai menționăm că relațiile evolutive sunt importante nu numai pentru căutarea ancestorilor comuni, ci permit o mai bună abordare a alinierii multiple. În acest sens am și pomenit de metoda arborilor pentru ierarhizarea moleculelor la introducerea în alinierea progresivă.

C. În cele ce urmează vom aborda analiza filogenetică prin construcția arborilor filogenetici, pornind de la câteva noțiuni introductive despre arbori și apoi ne vom concentra asupra principalelor două metode de construcție a arborilor: metoda distanțelor și metoda parsimoniei, trecând sumar și prin metoda asemănării maxime.

12.2. Noțiuni generale despre arbori

12.2.1. Terminologie

A. Biologia moleculară a adus argumente incontestabile privind similaritatea mecanismelor moleculare din materia vie, care sugerează că toate organismele vii de pe pământ, din toată istoria lumii, trebuie să fi avut un ancestor comun. Deci între oricare două organisme/specii se pot găsi diverse grade de rudenie. Relațiile de înrudire între specii se numește filogenie. Aceste relații pot fi reprezentate printr-un arbore filogenetic.

B. Din punct de vedere istoric, prima lucrare care sugerează extragerea informației din secvențe moleculare pentru a studia filogenia speciilor a fost publicată de Zuckenkandl și Pauling [1962]. Apoi Langley și Fitch [1974] arată că proteinele se modifică cu rate diferite, iar Doolittle [1996] demonstrează o bună corelație între lungimi și perioadele de timp de evoluție.

12.2.2. Utilitate, motivați

- A. Studiul arborilor filogenetici are o mare însemnătate teoretică și practică:
- putem înțelege mai bine relațiile evolutive ale speciilor,
 - putem înțelege cum au evoluat diverse funcții,
 - arborii filogenetici aduc informații pentru alinierea multiplă,

- se pot identifica cele mai importante elemente, porțiunile care se conservă în unele clase de secvențe.

B. Exemple

În figura 12.2.2.a este prezentat un arbore filogenetic *al genelor* ce codifică informații privind globinele.

În figura 12.2.2.b este o schemă simplificată a arborelui *speciilor* de babuini.

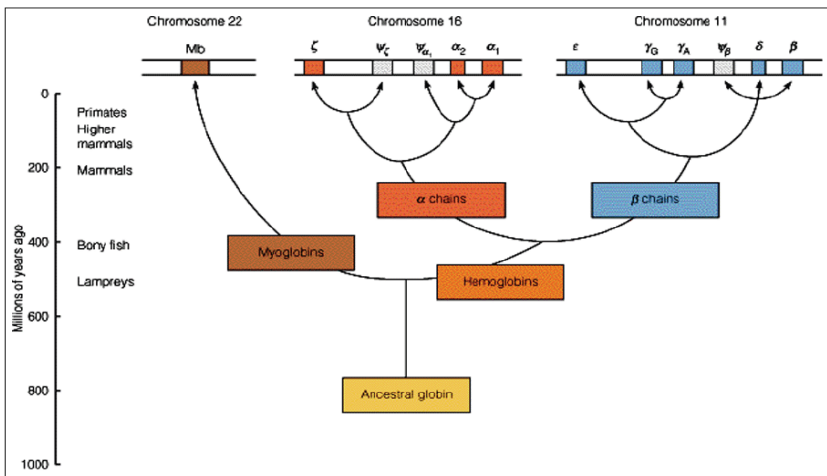


Fig. 12.2.2.a. Arbore filogenetic pentru gene: globinele

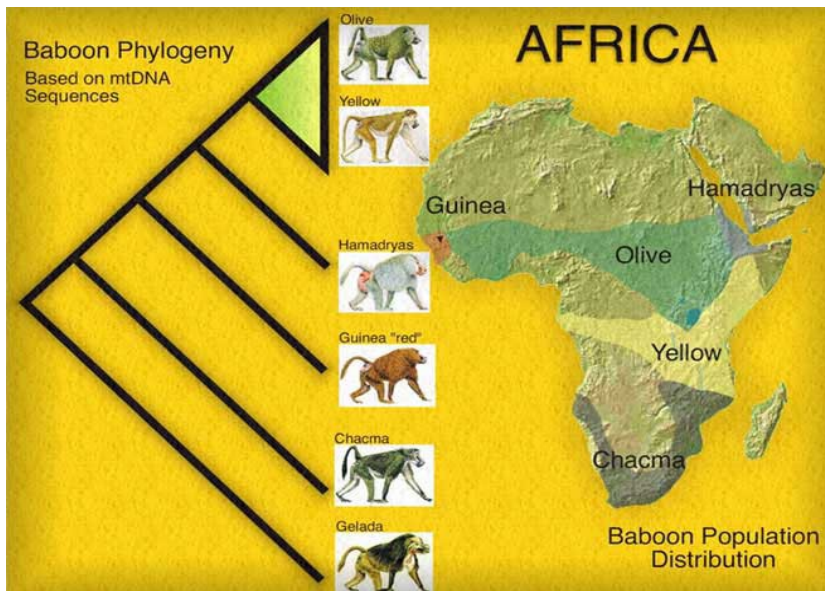


Fig. 12.2.2.b. Arbore filogenetic pentru specii: babuini

12.3. Proprietățile arborilor

12.3.1. Definiție, structură

A. Definiție: un arbore (“tree”) este un graf aciclic, nedirecționat.

B. Elementele componente ale unui arbore sunt:

- frunze (“leaves”) - elemente primare, obiecte (de exemplu: specii, gene, secvențe de proteine, chiar componente ale unei secvențe) cu poziție terminală în arbore; se mai numesc și noduri exterioare sau noduri de grad „1”; de obicei sunt notate cu numere (uneori cu litere)
- noduri (“nodes”) - reprezintă intersecții de ramuri; se numerotează de la $n+1$ în sus (în cazul unui arbore cu n frunze)
- ramuri (“edges”) - sunt legăturile dintre noduri; deseori au asociată o „lungime” calculată după diverse criterii.

Observație: în cazul arborilor în care se prezintă relații între specii, o frunză se mai numește “taxon” (pl: taxa)!

12.3.2. Tipuri de arbori

A. Clasificare Arborii filogenetici pot fi de două tipuri:

- a) arbori fără rădăcină (“unrooted tree”) - în care se specifică doar relațiile dintre noduri (figura 12.3.2.a)
- b) arbori cu rădăcină (“rooted tree”) - rădăcina este ultima ramură ce ajunge la ultimul nod, presupus a reprezenta ancestorul comun (figura 12.3.2.b).

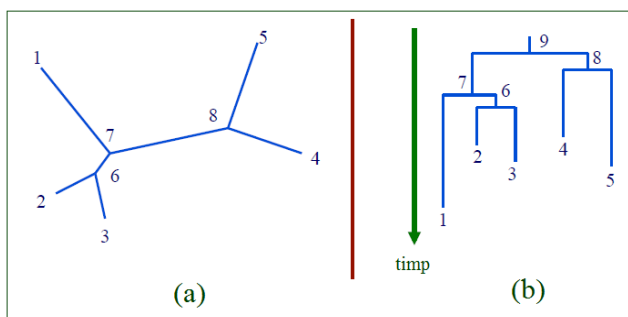


Fig. 12.3.2. Tipuri de arbori

B. Pentru arborii cu rădăcină se consideră convențional un sens al evoluției în timp, uzual plasând frunzele în partea de jos iar ancestorul comun (rădăcina) în partea superioară. De asemenea, lungimile ramurilor ar trebui să fie proporționale cu „distanța” dintre noduri (similară distanțelor dintre secvențe). Însă în arborii fără rădăcină nu avem un sens de desfășurare (nu avem o axă a timpului), și nici lungimile nu trebuie să fie dependente de gradul de similaritate (deși uneori se practică o reprezentare aproximativ proporțională). De aceea spunem despre arborii fără rădăcină că reprezintă topologia arborelui.

C. Inferență filogenetică

Operațiunea principală (concluzia) studiilor filogenetice, ce reprezintă stabilirea unui arbore filogenetic care caracterizează linia evolutivă între specii sau gene se numește „inferență filogenetică”.

12.3.3. Numărul de arbori

A. Numărul de noduri și ramuri

Pentru arborii cu rădăcină

Considerăm un arbore cu n frunze de la care pleacă în sus n ramuri. Unirea a fiecare două frunze va genera un nod de ordin superior și va reduce numărul ramurilor cu o unitate. Vom avea deci $2n-1$ noduri și $2n-2$ ramuri (dacă nu mai numărăm ultima ramură - cea de deasupra nodului rădăcină).

Pentru arborii fără rădăcină cu n frunze vom avea $2n-2$ noduri și $2n-3$ ramuri.

B. Calculul numărului de arbori fără rădăcină N_{AFR}

Pornind de la cel mai simplu arbore fără rădăcină – cel cu 3 frunze (figura 12.3.3.a), vedem că putem genera din el 3 variante de arbori cu 4 frunze asociind a 4-a frunză pe oricare dintre cele $2 \times 3 - 3 = 3$ ramuri (figura 12.3.3.b). Acești arbori cu 4 frunze au $2 \times 4 - 3 = 5$ ramuri.

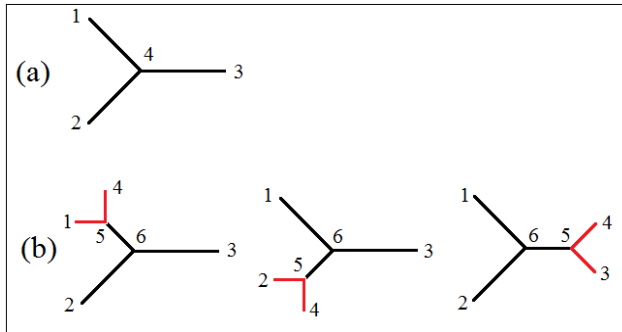


Fig. 12.3.3.a si b. Formarea incrementală a arborilor fără rădăcină

Deci pentru $n = 4$ frunze, pe cei 3 arbori fiecare cu $2n-3$ ramuri putem adăuga a 5-a frunză în $3 \times 5 = 15$ poziții diferite. Putem stabili o relație de recurență cu datele din tabelul 12.3.3.a.

Tabel 12.3.3.a. Numărul arborilor fără rădăcină

Nr. frunze	3	4	5	...	n
Nr. noduri	4	6	8	...	$2n - 2$
Nr. ramuri	3	5	7	...	$2n - 3$
Nr. arbori	1	3×5	$3 \times 5 \times 7$...	$(2n - 5)!!$

Putem sintetiza aceste date în relația (12.3.3.a):

$$N_{AFR} = \prod_{i=3}^n (2i - 5) = (2n - 5)!! \tag{12.3.3.a}$$

C. Calculul numărului de arbori cu rădăcină N_{AR}

Putem face un raționament similar cu cel de mai sus, pornind tot de la cel mai simplu arbore fără rădăcină, cel cu 3 frunze (figura 12.3.3.c). Vom putea din nou obține 3 variante de arbori cu rădăcină, plasând-o pe fiecare dintre cele 3 ramuri (figura 12.3.3.d).

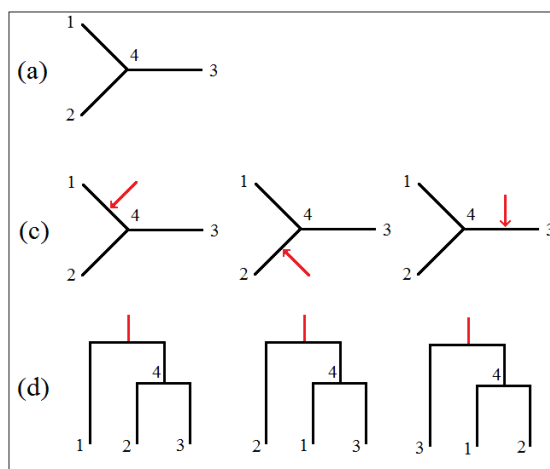


Fig. 12.3.3.c și d. Arbori cu rădăcină derivați din arbori fără rădăcină

La modul general, dintr-un arbore fără rădăcină cu n frunze care are $2n-3$ ramuri putem plasa rădăcina pe oricare din ramuri, deci vom avea relațiile:

$$N_{AR} = N_{AFR} \times (2n - 3) = (2n - 3)!! \quad (12.3.3.b)$$

12.4. Construcția arborilor filogenetici - generalități

12.4.1. Date pentru construcția arborilor

În construcția arborilor filogenetici pornim de la o serie de date experimentale:

- Distanțe – măsuri sau estimări ale distanțelor între specii sau între gene,
- Caractere – colecția de aspecte morfologice (ex.: forma boabelor, culoarea ochilor, mărimea aripilor etc.), secvențe ADN sau proteine,
- Ordinea genelor – din hărțile genetice după ordinea lineară a genelor ortoloage în genomurile date.

12.4.2. Metode de construcție a arborilor filogenetici

Au fost elaborate mai multe metode de construcție a arborilor filogenetici ce pot fi grupate în mai multe clase.

- metode bazate pe distanțe – arborele explică distanțele evolutive estimate,
- metoda parsimoniei – se alege arborele care necesită numărul minim de schimbări pentru a explica datele,
- metoda asemănării maxime – pornind de la modele probabiliste de evoluție.

În cele ce urmează vom prezenta metodele bazate pe distanțe și metoda parsimoniei.

12.4.3. Comparatie între metodele bazate pe distanțe și cele bazate pe parsimonie

A. Să considerăm secvențele frunzele de mai jos (cu câte 7 elemente, prezentate în tabelul 12.4.3.a).

Tabel 12.4.3.a. Exemplu pentru compararea metodelor de construcție a arborilor filogenetici

		pozitii						
		1	2	3	4	5	6	7
secvente	1	t	t	a	t	t	a	a
	2	a	a	t	t	t	a	a
	3	a	a	a	a	a	t	a
	4	a	a	a	a	a	a	t

B. În metodele bazate pe distanțe construim întâi o matrice a distanțelor între toate perechile. Luăm cea mai simplă distanță – distanța Hamming, ce exprimă numărul de nepotriviri, prezentată în tabelul 12.4.3.b.

Tabel 12.4.3.b. Matricea distanțelor pentru secvențele din tabelul 12.4.3.a

		Distanțe			
		1	2	3	4
secvente	1	0			
	2	3	0		
	3	5	4	0	
	4	5	4	2	0
		1	2	3	4
		secvente			

Fără a intra acum în detalii privind modul de construcție, prezentăm arborele din figura 12.4.3.a, care satisface distanțele între frunze – pe fiecare ramură a fost precizată „lungimea/ponderea” sa. Distanța între oricare două frunze este suma lungimilor ramurilor între ele.

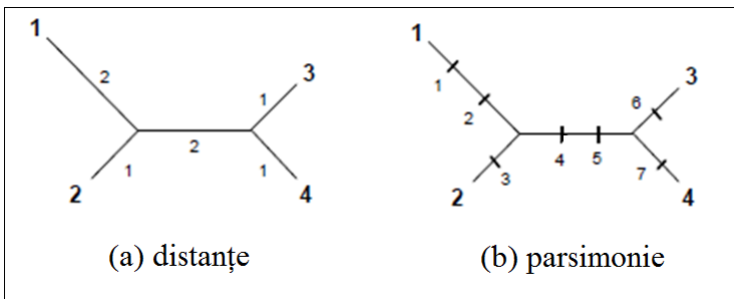


Fig. 12.4.3. Arborii corespunzătorii secvențelor din tabelul 12.4.3.a. Varianta (a) are pe ramuri lungimile corespunzătoare conform metodei distanțelor. Varianta (b) corespunde metodei parsimoniei, având pe ramuri precizate pozițiile care trebuie schimbate.

C. În metodele bazate pe parsimonie se precizează schimbările necesare pentru a alinia secvențele și se alege arborele care necesită numărul minim de schimbări. În figura 12.4.3.b sunt prezentate chiar și pozițiile din secvență care trebuie modificate la alinierea secvențelor. Numărul total de modificări între două secvențe este totalul de pe ramurile ce unesc cele două frunze.

12.5. Metode bazate pe distanțe

12.5.1. Proprietățile distanțelor, formularea problemei

A. Metodele bazate pe distanțe sunt simple și au un caracter intuitiv. În capitolul 7 am descris o serie de abordări privind exprimarea „distanțelor” între diferite secvențe, distanțele fiind invers proporționale cu gradul de similaritate.

B. Proprietățile distanțelor

Din punct de vedere formal, orice metrică ce introduce o distanță între două secvențe, x_i și x_j , notate d_{ij} sau $\text{dist}(x_i, x_j)$, trebuie să acorde distanțelor următoarele proprietăți:

- (a) $\text{dist}(x_i, x_j) \geq 0$ (distanțe pozitiv definite)
- (b) $\text{dist}(x_i, x_j) = 0$ (distanța punctuală)
- (c) $\text{dist}(x_i, x_j) = \text{dist}(x_j, x_i)$ (simetria)
- (d) $\text{dist}(x_i, x_j) \leq \text{dist}(x_i, x_k) + \text{dist}(x_k, x_j)$ (regula triunghiului)

C. Fiind dată o matrice M între taxonii i și j de dimensiune $n \times n$ (n = număr de taxoni /frunze), să se construiască un arbore cu ramuri ponderate (“edge weighted tree”).

D. În figura 12.5.1.a este prezentat un exemplu de matrice a distanțelor iar alături este arborele corespunzător în varianta arbore cu rădăcină, având și o scară pentru proporționalitatea lungimii ramurilor.

În cele ce urmează vom descrie un algoritm utilizat frecvent pentru construcția arborilor filogenetici prin metoda distanțelor.

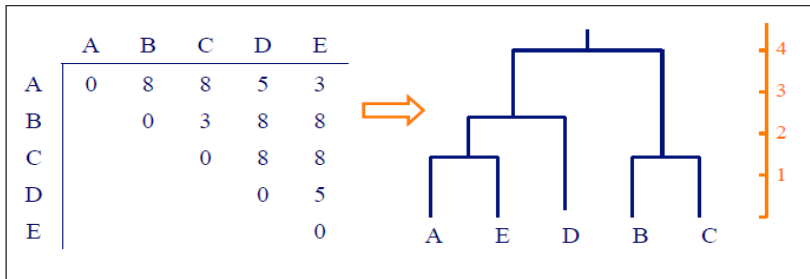


Fig. 12.5.1.a. Exemplu de arbore prin metoda distanțelor

12.5.2. Algoritmul UPGMA

A. Algoritmul UPGMA a fost propus de Sokal&Michner (1958) cu numele de “Unweighted Pair Group Method using arithmetic Averages”. Este o metoda simplă și intuitivă.

B. Ideea de bază este:

- compunerea succesivă a câte două componente (taxoni, clustere) formând un (nou) cluster
- se creează un nou nod pentru noul cluster

distanța între două clustere (între taxonii sau perechi de taxoni din fiecare cluster) este o distanță ponderată definită prin (12.5.2.a):

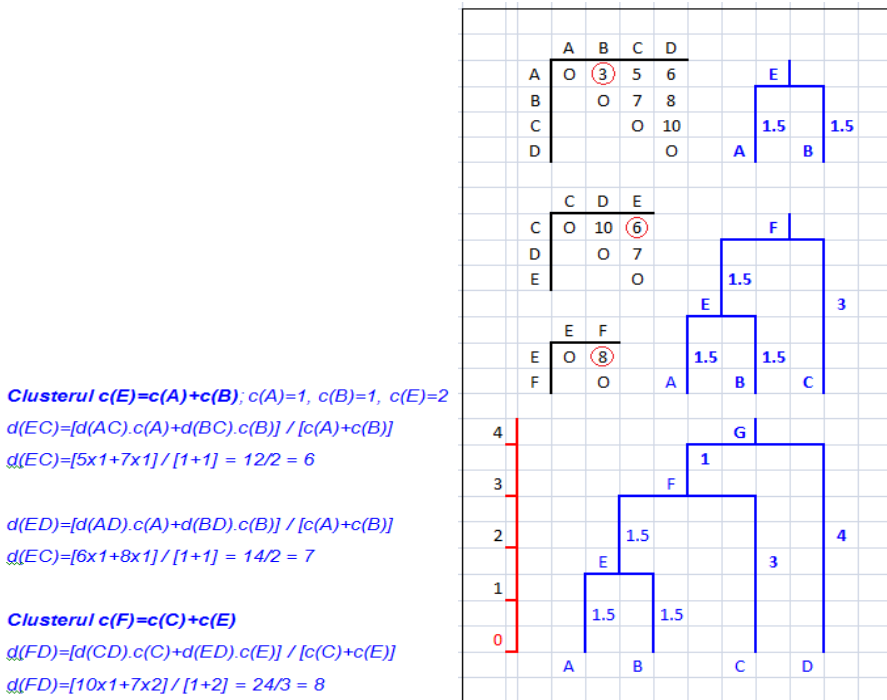
$$d_{ij} = \frac{1}{|C_i| + |C_j|} \sum_{p \in C_i, q \in C_j} d_{pq} \quad (12.5.2.a)$$

C. Descrierea algoritmului UPGMA

- Se consideră fiecare taxon ca un cluster
- Se definește o frunză pentru fiecare taxon; se plasează la înălțimea “0” pe scara distanțelor
- Când sunt mai mult de două clusterse:
 - Se aleg două clusterse, i și j , pentru care distanța d_{ij} este minimă
 - Se definește un nou cluster $C_k = C_i \cup C_j$
 - Se definește un nod k părinte al i și j ; se plasează la înălțimea $d_{ij} / 2$
 - Se înlocuiesc clustersele i și j cu k
 - Se calculează distanța între k și celelalte clusterse.

Ultimele două clusterse i și j se unesc cu o rădăcină la înălțimea $d_{ij} / 2$.

D. Să luăm un exemplu și să construim arborele urmând pas cu pas algoritmul UPGMA.



12.5.3. Distanțe ultrametrice

A. Algoritmul UPGMA descris mai sus presupune implicit o ipoteză numită „ipoteza ceasului molecular”, prin care se consideră că divergența secvențelor apare cu aceeași rată în orice punct din arbore. Datele care respectă această proporționalitate a divergenței cu timpul se numesc „date ultrametrice”, iar distanțele corespunzătoare sunt „distanțe ultrametrice”.

B. Totuși, ipoteza nu este în general valabilă. Procesul de selecție variază în diverse perioade de timp, variază cu organismul, cu genele unui organism sau cu regiunile unei gene. Putem astfel avea situații în care două secvențe apropiate ca grad de rudenie să prezinte distanțe mai mari decât distanțele reale față de secvențe mai îndepărtate filogenetic.

C. Cu alte cuvinte, putem avea situații în care prin metoda UPGMA să construim un arbore nereal. O astfel de situație este exemplificată în figura 12.5.3, în care în partea stângă - (a), este reprezentat arborele real, iar în dreapta - (b), o reconstrucție prin algoritmul UPGMA - reconstrucție de fapt greșită.

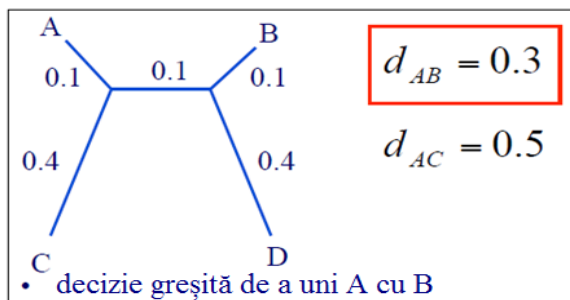


Fig. 12.5.3 Metoda Neighbor Joining

12.5.4. Metoda Neighbor Joining

A. Neajunsurile algoritmului UPGMA în cazul datelor care nu sunt ultrametrice, sunt înlăturate printr-o metodă derivată, numită “neighbor joining”.

B. Într-un algoritm propus de Saitou & Nei [1987] în care distanțele reale d_{ij} între două frunze i și j se corectează prin scăderea mediei distanțelor către toate celelalte frunze, ceea ce compensează ramurile lungi.

C. Din punct de vedere formal se definesc distanțele corectate D_{ij} prin relația

$$D_{ij} = d_{ij} - (r_i + r_j) \quad (12.5.4.a)$$

unde media distanțelor de la frunza i la celelalte este

$$r_i = \frac{1}{|L|-2} \sum d_{ik} \quad (12.5.4.b)$$

în care L este numărul de frunze.

În algoritmul Saitou se pornește de la ideea că frunzele i și j pentru care D_{ij} este minim sunt frunze vecine (adică provin din același nod proxim).

D. Algoritmul Saitou-Nei este oarecum asemănător cu UPGMA cu următoarele deosebiri:

- nu aplică ipoteza ceasului molecular
- se creează un arbore fără rădăcină
- presupune „aditivitate”, adică distanța între perechile de frunze este suma lungimilor ramurilor care le conectează.

Nu prezentăm aici algoritmul în detaliu.

12.6. Metode bazate pe „parsimonie”

12.6.1. Enunțul problemei

A. Una dintre clasele de metode ce și-a câștigat mulți adepți, atractivă prin simplitate și logică, pornește de la ipoteza că, pe linie filogenetică au apărut în mod natural secvențe divergente, însă nu cu multe modificări simultan, deci cea mai probabilă cale de căutare a unui posibil „părinte” este prin căutarea secvenței ce necesită un număr minim de schimbări.

Însăși termenul are ca origine cuvântul „parsimonios”, cu semnificația de „econom, zgârcit”.

B. În metodele bazate pe „parsimonie” se caută arborele care explică datele cu un număr minim de schimbări.

Să subliniem faptul că nu mai urmărim lungimile ramurilor ci numai topologia.

C. Să luăm ca exemplu 4 secvențe simple AAG, AAA, GGA, AGA. Putem prezenta înrudirea lor și evoluția către setul dat în mai multe variante, deci soluția nu este unică. Două din aceste soluții sunt prezentate în figura 12.6.1.

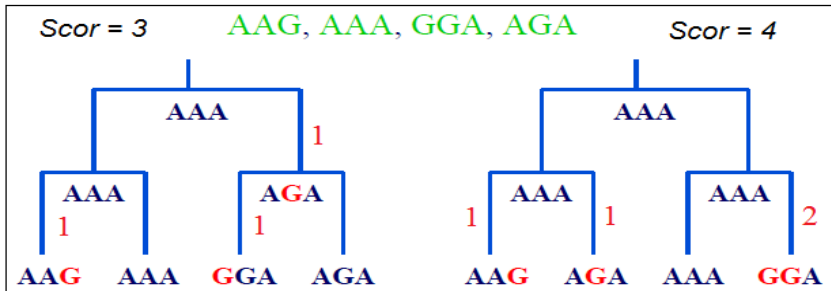


Fig. 12.6.1. Comparație între arbori pentru metoda parsimoniei

Se observă că în arborele din stânga evoluția a necesitat 3 schimbări (deci vom acorda scorul „3”), iar în arborele din dreapta scorul este mai slab, necesitând 4 schimbări.

D. Principiile algoritmilor din metoda parsimoniei

În algoritmii din clasa metodelor bazate pe parsimonie se pleacă de la următoarele ipoteze:

- orice element (stare, “state”) [nucleotid sau aminoacid] poate fi convertit în orice alt element,
- costurile schimbărilor sunt uniforme (Fitch) sau ponderate (Sankoff),
- pozițiile sunt independente (se poate calcula numărul minim separat pentru fiecare poziție).

12.6.2. Algoritmii lui Fitch

A. Cel mai popular algoritm este cel propus de Fitch (1971), în care arborele se traversează de două ori:

- în primul pas se traversează arborele de la frunze la rădăcină, determinând setul de stări posibile (de exemplu: nucleotide) pentru fiecare nod intern,

- în al doilea pas se traversează arborele de la rădăcină spre frunze alegând stările ancestrale pentru nodurile interne (un fel de "trace-back").

B. Vom descrie cei doi pași, ilustrându-i și cu un exemplu.

- a) pasul 1: traversarea "post-order" sau "bottom-up", adică de la frunze spre rădăcină.

Să stabilim stările posibile pentru nodurile interne.

- Se determină stările R_i ale nodului intern „i” care are copiii (urmașii) „j” și „k” conform regulii:
 - reuniunea R_j și R_k dacă nu au elemente comune,
 - intersecția R_j cu R_k dacă au elemente comune.

Formal scriem

$$R_i = \begin{cases} R_j \cup R_k, & \text{if } R_j \cap R_k = \emptyset \\ R_j \cap R_k, & \text{otherwise} \end{cases} \quad (12.6.2.a)$$

- Numărul de schimbări (scorul inițial) este dat de numărul de reuniuni!
- Avem în figura 12.6.2.a. un exemplu.

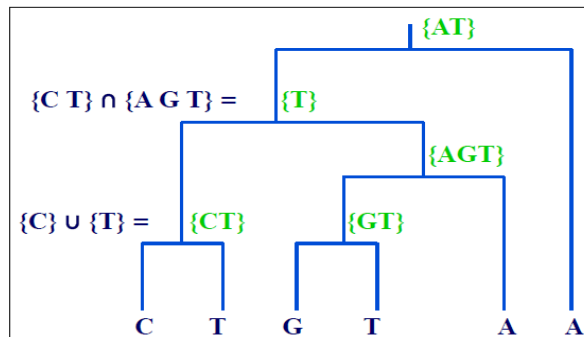


Fig. 12.6.2.a. Algoritmul Fitch, pasul „post-order”

Observăm că am efectuat 3 reuniuni și două intersecții, deci scorul stabilit este $s=3$.

- b) pasul 2: traversarea "pre-order" sau "top-down" de la rădăcină la frunze (echivalent "trace-back").

Se alege starea r_j a unui nod intern „j” în funcție de starea R_i a părintelui „i” astfel:

- egală cu cea a părintelui, dacă este inclusă
- stare arbitrară în caz contrar

Formal scriem:

$$r_j = \begin{cases} r_i, & \text{if } r_i \in R_j \\ \text{arbitrary state} \in R_j, & \text{otherwise} \end{cases} \quad (12.6.2.b)$$

Observație: în nod vom păstra starea după parcurgerea pasului 2.

Continuând exemplul de mai sus, marcăm cu literă închisă la culoare în figura 12.6.2.b starea din setul anterior care satisface relația (12.6.2.b).

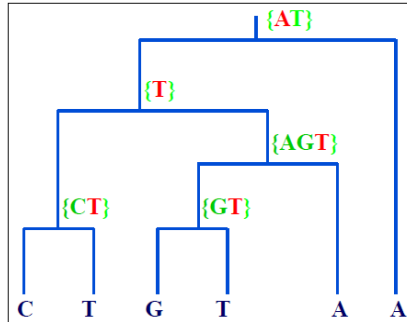


Fig. 12.6.2.b. Algoritmul Fitch, pasul „pre-order”

C. În final vom reține în noduri numai elementele selectate în pasul 2 și vom marca cu „X” ramurile pe care s-a înregistrat substituție, pentru a evidenția scorul final (figura 12.6.2.c).

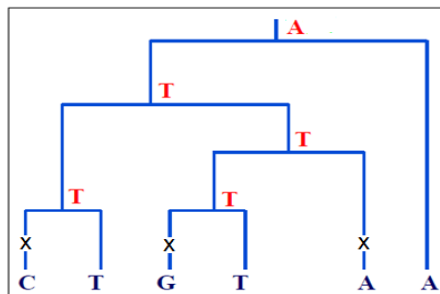


Fig. 12.6.2.c Arbore filogenetic construit prin metoda parsimoniei, algoritmul Fitch

12.6.3. Parsimonie ponderată

A. Varianta propusă de Sankoff & Cedergren [1983] aduce metodei parsimoniei aceleași ajustări ca modelul Kimura față de Jukes-Cantor, adică, în loc să considere toate schimbările echiprobabile, introduce costuri diferite $S(a,b)$ pentru diverse substituții ale lui „a” cu „b” (tabelul 12.6.3).

B. Vom avea deci o matrice de substituție S , conținând costurile tuturor substitutelor posibile. De exemplu, în cazul secvențelor ADN vom acorda penalizări mai mici pentru tranziții ($a \leftrightarrow g, c \leftrightarrow t$) decât pentru transversii ($\{a,g\} \leftrightarrow \{c,t\}$). Conservarea va avea cost nul (nu se penalizează). În exemplul ce va fi prezentat mai jos apare o astfel de matrice.

Tabel 12.6.3. Matrice de substituție pentru metoda parsimoniei

	a	c	g	t
a	0	0.8	0.2	0.9
c	0.8	0	0.7	0.5
g	0.2	0.7	0	0.1
t	0.9	0.5	0.1	0

C. Primul pas al algoritmului Sankoff urmărește propagarea costurilor în sus pe arbore. Algoritmul propus este similar programării dinamice, având ca subproblemă să se determine costul $R_i(a)$ pentru subarborile cu rădăcina în „i”, plasând caracterul „a” în nodul „i”.

D. Formal vom putea descrie algoritmul prin următoarele relații, folosind notațiile din figura 12.6.3.a

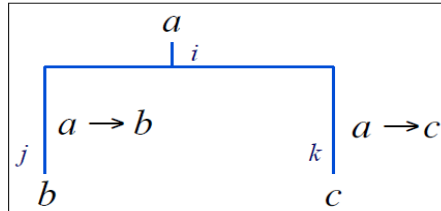


Fig. 12.6.3.a. Introducere de ponderi în parsimonia ponderată

- pentru frunze

$$R_i(a) = \begin{cases} 0 & \text{daca "a" este un caracter in frunza} \\ \infty & \text{in caz contrar} \end{cases} \quad (12.6.3.a)$$

- pentru noduri interne

$$R_i(a) = \min_b [R_j(b) + s(a, b)] + \min_c [R_k(c) + s(a, c)] \quad (12.6.3.b)$$

E. Să luăm un arbore (fig. 12.6.3.b) pe care să urmărim formulele costurilor pentru noduri.

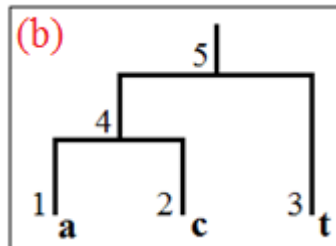


Fig. 12.6.3.b. Aplicație pentru parsimonia ponderată

Aplicăm relațiile (12.6.3.a) în cele 3 frunze, notate (1), (2) și (3) apoi formulele (12.6.3.b) în nodurile (2) și (1) și obținem:

$$R_1[a] = 0, R_1[c] = \infty, R_1[g] = \infty, R_1[t] = \infty$$

$$R_2[a] = \infty, R_2[c] = 0, R_2[g] = \infty, R_2[t] = \infty$$

$$R_3[a] = \infty, R_3[c] = \infty, R_3[g] = \infty, R_3[t] = 0$$

$$R_4[a] = \{R_1[a] + S(a, a)\} + \{R_2[c] + S(a, c)\} = \{0 + 0\} + \{0 + 0.8\} = 0.8$$

$$R_4[c] = \{R_1[a] + S(c, a)\} + \{R_2[c] + S(c, c)\} = \{0 + 0.8\} + \{0 + 0\} = 0.8$$

$$R_4[g] = \{R_1[a] + S(g, a)\} + \{R_2[c] + S(g, c)\} = \{0 + 0.2\} + \{0 + 0.7\} = 0.9$$

$$R_4[t] = \{R_1[a] + S(t, a)\} + \{R_2[c] + S(t, c)\} = \{0 + 0.9\} + \{0 + 0.5\} = 1.4$$

$$\begin{aligned}
 R_5[a] &= \min\{R_4[a] + S(a, a), R_4[c] + S(a, c), R_4[g] + S(a, g), R_4[t] + S(a, t)\} \\
 &\quad + \{R_3[t] + S(a, t)\} \\
 &= \min\{0.8 + 0, 0.8 + 0.8, 0.9 + 0.2, 1.4 + 0.9\} + \{0 + 0.9\} \\
 &= \min\{0.8, 1.6, 1.1, 2.3\} + \{0.9\} = 0.8 + 0.9 = 1.7
 \end{aligned}$$

$$\begin{aligned}
 R_5[c] &= \min\{R_4[a] + S(c, a), R_4[c] + S(c, c), R_4[g] + S(c, g), R_4[t] + S(c, t)\} \\
 &\quad + \{R_3[t] + S(c, t)\} \\
 &= \min\{0.8 + 0.8, 0.8 + 0, 0.9 + 0.7, 1.4 + 0.5\} + \{0 + 0.5\} \\
 &= \min\{1.6, 0.8, 1.6, 1.9\} + \{0.5\} = 0.8 + 0.5 = 1.3
 \end{aligned}$$

$$\begin{aligned}
 R_5[g] &= \min\{R_4[a] + S(g, a), R_4[c] + S(g, c), R_4[g] + S(g, g), R_4[t] + S(g, t)\} \\
 &\quad + \{R_3[t] + S(g, t)\} \\
 &= \min\{0.8 + 0.2, 0.8 + 0.7, 0.9 + 0, 1.4 + 0.1\} + \{0 + 0.1\} \\
 &= \min\{1.0, 1.5, 0.9, 1.5\} + \{0.1\} = 0.9 + 0.1 = 1.0
 \end{aligned}$$

$$\begin{aligned}
 R_5[t] &= \min\{R_4[a] + S(t, a), R_4[c] + S(t, c), R_4[g] + S(t, g), R_4[t] + S(t, t)\} \\
 &\quad + \{R_3[t] + S(t, t)\} \\
 &= \min\{0.8 + 0.9, 0.8 + 0.5, 0.9 + 0.1, 1.4 + 0\} + \{0 + 0\} \\
 &= \min\{1.7, 1.3, 1.0, 1.4\} + \{0\} = 1.0 + 0 = 1.0
 \end{aligned}$$

Sumar, pentru rădăcină:

$$R_5[a] = 1.7, R_5[c] = 1.3, R_5[g] = 1.0, R_5[t] = 1.0$$

F. Pasul 2: parcurgerea arborelui de sus în jos.

Observăm următoarele:

- caracterul de cost minim pentru nodul (5) este fie „g” fie „t” cost (1. 0)
- caracterul de cost minim pentru nodul (4) este „g”!

Acest lucru este impus de calea costului minim pentru nodul 1. Într-adevăr, observăm că minimul în $R_5[g]$ sau $R_5[t]$ se obține din termenul $R_4[g]$, chiar dacă valorile lui R_4 calculate din R_1 și R_2 sunt mai mici pentru „a” sau „c”.

G. Rezultatul final: soluția este dublă, avem doi arbori cu cost minim prin metoda parsimoniei ponderate figura 12.6.3.c și d.

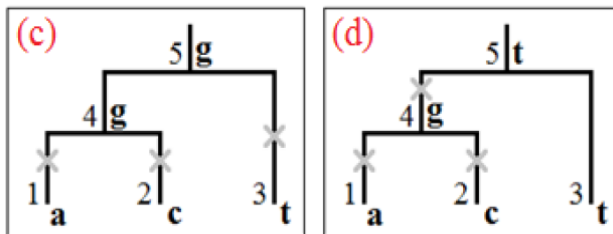


Fig. 12.6.3.c și d. Arbori de cost minim prin metoda parsimoniei ponderate

Să verificăm costul în cele două cazuri:

- $S_{41} + S_{42} + S_{53} = S(g, a) + S(g, c) + S(g, t) = 0.2 + 0.7 + 0.1 = 1.0$
- $S_{41} + S_{42} + S_{54} = S(g, a) + S(g, c) + S(t, g) = 0.2 + 0.7 + 0.1 = 1.0$

Putem verifica și că în (4) este mai potrivit să avem „g” decât „a” sau „c”, reprezentați în figura 12.6.3.e și f.

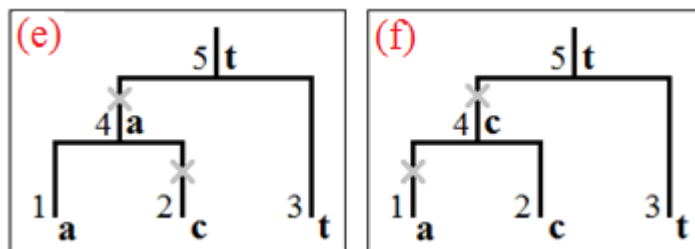


Fig. 12.6.3. e și f. Calcul comparativ de costuri în parsimonia ponderată

Se observă costul mai ridicat, deși ar fi fost efectuate doar două substituții în ambele cazuri!

La fel se poate arăta că plasarea lui „a” și „c” în nodul 4 nu aduce cost minim nici în cazul în care în 5 luăm „g”.

12.7. Testarea arborilor - metoda „bootstrap”

A. Algoritmii de construcție a arborilor descriși mai sus ne permit selectarea unui arbore „optim”, uneori, ca în cazul parsimoniei mai mulți, fără însă a da o măsură a calității arborelui ales. Desigur, este nevoie de o testare a semnificației unui aspect filogenetic, cum ar fi segregarea unui caracter pentru un grup de specii, etc.

B. Felsenstein [1985] a propus, o metodă, numită ”bootstrap”, dezvoltată ulterior de Efron & Tibshirani [1993], pe care o vom prezenta foarte pe scurt.

Fiind dat un set de secvențe aliniate, se generează un set artificial de aceeași dimensiune, preluând la întâmplare coloane din aliniere. Se aplică un algoritm de construcție a arborelui noului set de date. Se repetă procedeul de foarte multe ori (de ordinul miilor). Frecvența cu care apare un aspect filogenetic este luată ca măsură de încredere pe care o putem avea în aspectul respectiv.

Să menționăm în încheierea capitolului că pentru arborii filogenetici s-au dezvoltat și metode probabiliste bazate în special pe metoda asemănării maxime.

Partea a II-a

1. BAZE DE DATE – BD integrată NCBI (I). Utilizarea motorului de căutare ENTREZ și modulul PUBMED

1.1. Obiectivele lucrării de laborator

- căutări în bazele de date după cuvinte-cheie
- deprinderea creării interogărilor
- căutări de articole medicale în bazele de date după autor
- căutări de articole medicale în bazele de date după autor și cuvinte-cheie.

1.2. Baze de date și instrumente de căutare folosite în Bioinformatică

Bioinformatica tinde să dezvolte baze de date și algoritmi din ce în ce mai puternici în scopul accelerării și intensificării cercetării biologice. Disponibilitatea bazelor de date moleculare și a instrumentelor bioinformaticice au schimbat natura biologiei. Experimente care în trecut puteau fi efectuate doar într-un laborator, în urma efortului depus de-a lungul anilor, pot fi realizate cu ajutorul unui calculator personal și o conexiune la internet, prin analiza și interpretarea informației rezultate.

Pentru a îndeplini multitudinea de obiective la nivel global, a fost încurajată creșterea numărului bazelor de date din domeniul public. În prezent se poate obține o cantitate inimaginabilă de informații despre orice aspect al biologiei celulare prin intermediul Internetului, informații care pot fi bibliografice, genomice, structurale sau funcționale. Domeniul public, în afară de afișarea unor informații fundamentale, oferă utilitățile necesare analizei și interpretării datelor. Aceste instrumente variază de la aliniamentul multiplu al secvențelor la studiile de expresie virtuală a genelor, PCR electronic și altele. Câteva servere majore de biologie moleculară care funcționează în lume sunt:

- NCBI (National Centre for Biotechnology Information):
<http://www.ncbi.nlm.nih.gov>
- EBI (European Bioinformatics Institute):
<http://www.ebi.ac.uk>
- DDBJ (DNA DataBank of Japan):
<http://www.ddbj.nig.ac.jp>
- PDB (The Protein Data Bank):
<http://www.rcsb.org/pdb>

Centrul Național de Informație Biotehologică (NCBI) din Statele Unite ale Americii și Institutul European de Bioinformatică (EBI) din Marea Britanie sunt principalele servere științifice care mențin aceste imense baze de date precum și instrumentele de software analitic necesare analizei datelor conținute.

Serviciile oferite de aceste servere sunt posibile datorită computerelor ultraperformante cu viteză de procesare mare care pot realiza prelucrările analitice ale datelor și datorită Internetului care facilitează eforturile de comunicare electronică.

“Entrez” permite reconstituirea datelor din biologia moleculară și vizualizarea citărilor bibliografice din bazele de date integrate ale NCBI (National Center for Biotechnology Information).

Acest motor de căutare se accesează la adresa: <http://www.ncbi.nlm.nih.gov/Entrez/>

Printre modulele incluse în baza de date integrată NCBI se găsesc:

- PubMed cuprinde citări și articole biomedicale din BD Medline
- BD de secvențe nucleotidice (GenBank)
- BD de secvențe proteice
- Structuri macromoleculare 3D.

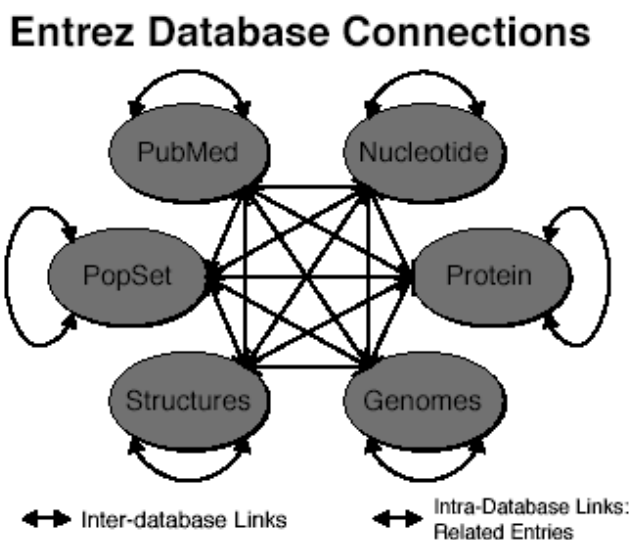


Figura 1.2.a. Legăturile intra și inter baze de date

1.3. Utilizarea Entrez

Ne propunem să utilizăm motorul de căutare Entrez folosind următoarele exemple de căutări:

- **camp activated protein kinase** (figurile 1.3.a, 1.3.b)
- **“camp activated protein kinase”**
- **“camp activated” “protein kinase”**
- o **“camp dependent protein kinase”**
- o **immunoglob*.**

The screenshot shows the NCBI Entrez Protein search interface. At the top, there are navigation tabs for 'All Databases', 'PubMed', 'Nucleotide', 'Protein', 'Genome', 'Structure', 'OMIM', 'PMC', and 'Journals'. The search bar contains the query 'camp activated protein kinase'. Below the search bar are buttons for 'Go', 'Clear', and 'Save'. There are also buttons for 'Limits', 'Preview/Index', 'History', 'Clipboard', and 'Details'. The 'Display' section shows 'Summary' selected, 'Show 20' items, and 'Sort By' options. A summary bar indicates 'All: 279' results, with 'Bacteria: 3', 'RefSeq: 128', and 'Related Structures: 251'. A 'Gene Information' box lists results for [Pka-C3](#) (*Drosophila melanogaster*): cAMP-dependent protein kinase 3, [Prkag2](#) (*Rattus norvegicus*): protein kinase, AMP-activated, gamma 2 non-catalytic subunit, and [Prkaca](#) (*Rattus norvegicus*): protein kinase, cAMP-dependent, catalytic, alpha. A 'Top Organ' dropdown menu is visible on the right.

Figura 1.3.a. Interfața motorului de căutare

The screenshot shows the 'Details' view of the search results. On the left, there is a navigation menu with links for 'About Entrez', 'Entrez Protein', 'Help | FAQ', 'Entrez Tools', 'Check sequence revision history', 'LinkOut', 'My NCBI', 'Related resources', and 'BLAST'. The main content area shows the 'Query Translation' section with the query: 'camp[All Fields] AND activated[All Fields] AND protein kinase[All Fields]'. Below this is a 'Search' button and a 'URL' field. The 'Result' section shows '279' results.

Figura 1.3.b. Exemplu de interogare a bazei de date după cuvinte-cheie

În acest scop urmăm pașii:

- se alege baza de date
- se introduce interogarea în fereastra “Details”.

Intrați pe fiecare opțiune din fereastra (“Limits”, “History” și “Details”) și observați cum puteți aplica diverse limitări ale vizualizării informației rezultat, cum puteți opera cu interogările anterioare (figura 1.3.c.) și care este expresia logică a interogării (figura 1.3.b).

Restul căutărilor le faceți individual și observați diferențele ce survin în cadrul opțiunii “Details”.

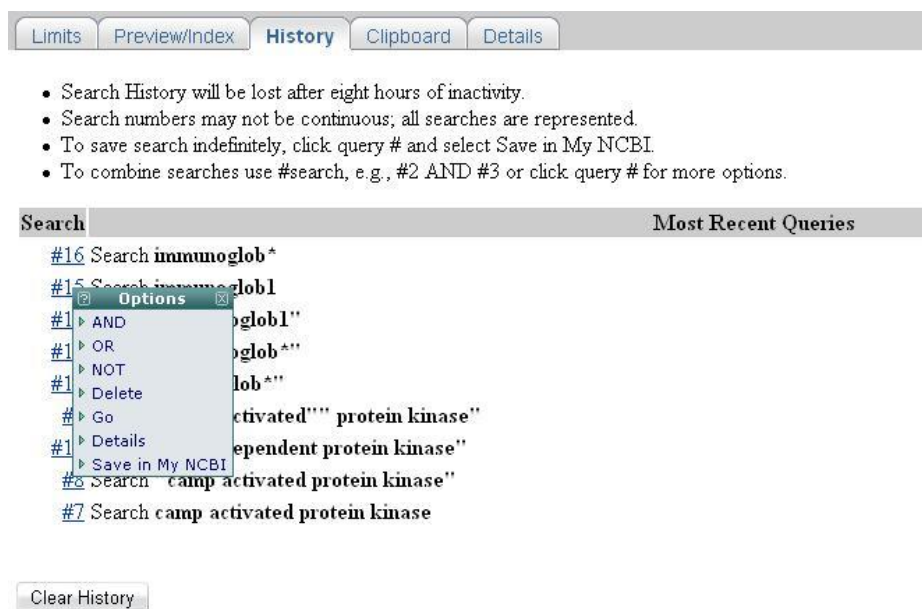


Figura 1.3.c. Prezentarea interogărilor anterioare

1.3.1. Calificatorii câmpului Entrez

Formatul general al calificatorilor “Entrez” sunt prezentați sub forma:

term [qualifier]

unde:

- **term** este șirul căutat
- **qualifier** reprezintă câmpul de căutat în BD specifică și poate fi:
 - ACCN număr de acces
 - DATE anul publicației
 - JOUR numele jurnalului
 - MAJR topica MeSH
 - PAGE prima pagină
 - PROP proprietăți
 - PTYP tipul publicației
 - TITL cuvânt de titlu
 - SLEN lungimea secvenței
 - WORD cuvânt text
 - AUTH numele autorului
 - GENE numele genei
 - KYWD cuvânt cheie
 - MDAT data modificării
 - ORGN organismul
 - PDAT data publicării/creării
 - PROT numele proteinei
 - SQID id-ul secvenței
 - SUBS substanța
 - VOL volumul

1.3.2. Exemple

a. Căutați articole după **Bourne PE** ca autor.

NCBI Entrez, The Life Sciences Search Engine

Search across databases: Bourne PE [AUTH] [GO] [Clear] Help

Result counts displayed in gray indicate one or more terms not found

134	PubMed: biomedical literature citations and abstracts	none	Books: online books
61	PubMed Central: free, full text journal articles	none	OMIM: online Mendelian Inheritance in Man
none	Site Search: NCBI web and FTP sites	none	OMIA: online Mendelian Inheritance in Animals
1	Nucleotide: Core subset of nucleotide sequence records	none	dbGaP: genotype and phenotype
none	EST: Expressed Sequence Tag records	none	UniGene: gene-oriented clusters of transcript sequences
none	GSS: Genome Survey Sequence records	none	CDD: conserved protein domain database
14	Protein: sequence database	20	3D Domains: domains from Entrez Structure
none	Genome: whole genome sequences	none	UniSTS: markers and mapping data
6	Structure: three-dimensional macromolecular structures	none	PopSet: population study data sets
none	Taxonomy: organisms in GenBank	none	GEO Profiles: expression and molecular abundance profiles
none	SNP: single nucleotide polymorphism	none	GEO DataSets: experimental sets of GEO data
none	Gene: gene-centered information	none	Cancer Chromosomes: cytogenetic databases
none	SRA: Short Read Archive	none	PubChem BioAssay: bioactivity screens of chemical substances
none	BioSystems: Pathways and systems of interacting molecules	none	PubChem Compound: unique small molecule chemical structures

Figura 1.3.2.a. Exemplu de căutare după Bourne PE [AUTH]

Dacă se dorește o căutare avansată, formatul interogărilor booleene este:

term [field] operator term [field]

Operatorii pot fi:

- **AND** (intersecție)
- **OR** (reuniune)
- **BUTNOT** (diferență)

b. Căutați articole ale autorului **Taylor P** care conțin informații despre “**protein kinase**” (figurile 1.3.2.b și 1.3.2.c).

NCBI Entrez, The Life Sciences Search Engine

Search across databases: Taylor P [AUTH] AND "protein kinase" [GO] [Clear] Help

Result counts displayed in gray indicate one or more terms not found

11	PubMed: biomedical literature citations and abstracts	2	Books: online books
21	PubMed Central: free, full text journal articles	none	Go to Books Results Page: Inheritance in Man
none	Site Search: NCBI web and FTP sites	none	OMIA: online Mendelian Inheritance in Animals
132	Nucleotide: Core subset of nucleotide sequence records	none	dbGaP: genotype and phenotype
none	EST: Expressed Sequence Tag records	none	UniGene: gene-oriented clusters of transcript sequences
none	GSS: Genome Survey Sequence records	none	CDD: conserved protein domain database
167	Protein: sequence database	32	3D Domains: domains from Entrez Structure
none	Genome: whole genome sequences	none	UniSTS: markers and mapping data
9	Structure: three-dimensional macromolecular structures	none	PopSet: population study data sets
none	Taxonomy: organisms in GenBank	none	GEO Profiles: expression and molecular abundance profiles
none	SNP: single nucleotide polymorphism	none	GEO DataSets: experimental sets of GEO data
none	Gene: gene-centered information	none	Cancer Chromosomes: cytogenetic databases
none	SRA: Short Read Archive	none	PubChem BioAssay: bioactivity screens of chemical substances
none	BioSystems: Pathways and systems of interacting molecules	none	PubChem Compound: unique small molecule chemical structures
none	HomoloGene: eukaryotic homology groups	none	PubChem Substance: deposited chemical substance records

Figura 1.3.2.b. Exemplu de căutare după Taylor P [AUTH] AND “protein kinase”

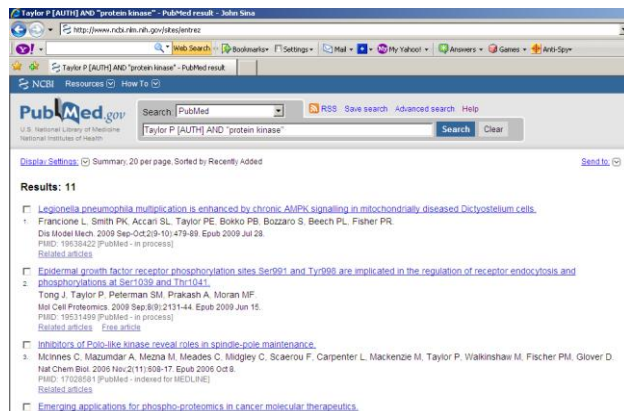


Figura 1.3.2.c. Exemplu de căutare cu ajutorul PubMed după Taylor P [AUTH] AND “protein kinase”

Notă:

Observați deosebirile între cele două motoare de căutare!

Rețineți că putem utiliza următoarele legături:

- PubMed Neighbours
- Text & MeSH
- Protein and Nucleotide Neighbours
- BLAST
- Macromolecular Structure Neighbours
- VAST

pe care le vom exemplifica în lucrarea de laborator următoare.

1.4. Utilizarea PubMed

Pentru a efectua căutări cu ajutorul secțiunii “PubMed” alegem opțiunea “PubMed: biomedical literature citations and abstracts” (figura 1.4.a).

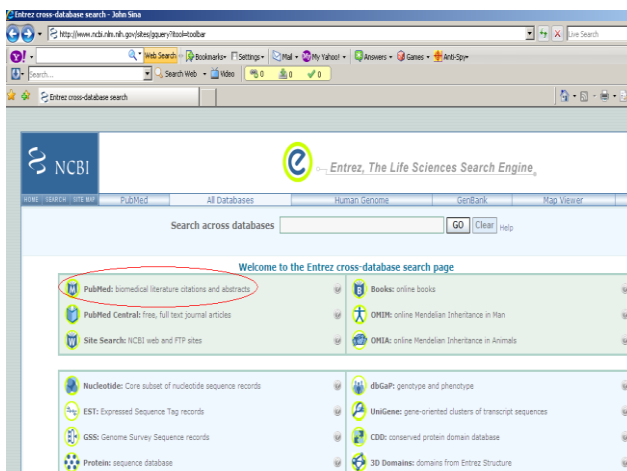


Figura 1.4.a. Opțiunea PubMed de căutare

Prin căutarea cu PubMed putem să:

- identificăm conceptele cheie din căutarea noastră
- introducem termenii în căsuța de search
- primim sugestii automate pe măsura completării termenilor în căsuța de căutare.

1.4.1. Exemple

a) Dacă dorim să aflăm ce rol are durerea în tulburările de somn introducem în căsuța de căutare cele două concepte-cheie: **pain și sleeping disorders**;

b) Dacă dorim să căutăm după autor, introducem numele de familie al autorului plus inițialele, fără punctuație, în căsuță și apăsăm butonul *search* (Ex. Watson JD, Lederberg J);

c) Apăsăm *Advanced search* pentru a căuta după *Autor, Jurnal, Data Publicării*. Pentru a căuta articole citate scrise de Bonnie W. Ramsey despre terapia genetică în fibroza chistică (cystic fibrosis), introducem următorii termeni de căutare în căsuță: **cystic fibrosis gene therapy ramsey bw**;

d) Publicații cu numele complet al autorilor pot fi căutate numai din 2002 (Ex. Joshua Lederberg, Garcia Algar, Oscar);

e) Dacă știm numai numele de familie al autorului vom introduce în căsuța de căutare după autor următoarele: **numeautor[au]**;

f) Pentru a căuta citate despre drosophila (musculița de oțet) în jurnalul “Molecular Biology of the Cell” introducem următoarele în căsuța de căutare: **molecular biology of the cell drosophila**;

g) Dacă căutăm studii sistematice despre terapia inhalatorie în pneumonie, din pagina “Clinical queries” apăsăm *Systematic Reviews* și introducem următorii termeni în căsuța de căutare: **inhalation therapy pneumonia** (figura 1.4.1.a);

NCBI Resources How To

PubMed Clinical Queries

This page provides the following specialized PubMed searches for clinicians:

- Search by Clinical Study Category
- Find Systematic Reviews
- Medical Genetics Searches

Results of searches on these pages are limited to specific clinical research areas

Search by Clinical Study Category

This search finds citations that correspond to a specific clinical study category are based on the work of Haynes RB et al. See the filter table for details.

Search Go

Category	Scope
<input type="radio"/> etiology	<input checked="" type="radio"/> narrow, specific search
<input type="radio"/> diagnosis	<input type="radio"/> broad, sensitive search
<input checked="" type="radio"/> therapy	
<input type="radio"/> prognosis	
<input type="radio"/> clinical prediction guides	

Find Systematic Reviews

For your topic(s) of interest, this search finds citations for systematic reviews development conferences, and guidelines.

For more information, see Help. See also related sources for systematic review

Search Go

Figura 1.4.1.a. Exemplu de interogare clinică prin intermediul PubMed

h) Pentru a găsi informații despre anemia falciformă (siclemia) din “Clinical Queries” apășăm *Medical Genetic Search*, apoi facem click pe *All check box* pentru a deselecta toate categoriile și *Genetic Counseling*, după care introducem următorii termeni în căsuța de căutare: **sickle cell anemia**.

1.5. Exerciții propuse

- Lucrați cu tutorialul Entrez de la adresa:
<http://www.ncbi.nlm.nih.gov/Database/tut1.html>
- Găsiți referințe despre paralizia facială. Vizualizați înregistrările în formatul: rezumat și autor.
- Găsiți secvența acidului nucleic murinic implicat în apoptoză și cancer. Descrieți toate înregistrările pe o pagină web.
- Găsiți secvența proteică cauzatoare de SIDA apărută în revista Nature. Salvați strategia de căutare pentru o căutare ulterioară.

2. BAZE DE DATE – BD integrată NCBI (II). Programul BLAST

2.1. Obiectivele lucrării de laborator

- căutări în bazele de date după secvențe de gene
- determinarea cromozomului pe care se află o secvență de nucleotide
- căutări în bazele de date după secvențe de nucleotide
- determinarea secvențelor de nucleotide responsabile de apariția unor diagnostice
- căutare în BD după secvența de aminoacizi
- găsirea gradului de potrivire al secvenței de aminoacizi introduse.

2.2. Introducere

Bioinformatica, ca știință a organizării și analizei datelor biologice complexe (reprezentate de proteine și secvențe de ADN), îmbină biologia moleculară și genetica cu tehnologia de calcul pentru a înțelege rețeaua complexă de interacțiuni dintre componentele individuale ale celulei vii și pentru a le integra în comportamentul întregului organism fiind utile în procesul de diagnosticare a bolilor și în stabilirea de noi strategii terapeutice.

Utilitatea bioinformaticii se observă mai ales în Proiectul Genomului Uman, care a avut drept scop identificarea celor 30.000 de gene din ADN-ul uman.

Bioinformatica progresaază cu o rată uimitoare, iar ariile de implicare majoră sunt achiziția de noi secvențe, încorporarea lor sub forma bazelor de date clasificate, integrarea informațiilor oferite de secvențe cu cele oferite de structuri, dezvoltarea instrumentelor pentru data mining și dezvoltarea unei platforme comune pentru folosirea resurselor. În acest scop s-a creat *International Nucleotide Sequence Database Collaboration (INSDC)*, care include cele trei baze de date genomice majore din lume: **GenBank**, **EMBL** și **DDBJ**.

Colaborarea dintre informaticieni, medici cercetători, biologi, matematicieni și biochimisti le-a permis acestora studiul bazei moleculare a unei boli cu ajutorul instrumentelor matematice și a tehnicii de calcul prin:

- analiza secvenței unei gene sau a produsului genei de interes;
- înțelegerea mai bună a organizării genelor analizate;
- predicția structurii moleculelor analizate (proteine).

Disponibilitatea informației genomice oferă bioinformaticianului un nou set de provocări. În prezent, ariile predominante ale analizei datelor bioinformaticice includ:

- aliniamentul secvențelor
- predicția structurii proteinelor.

Studiile de aliniament a secvențelor sunt (în general) de două tipuri:

- de *aliniament al secvențelor în pereche* realizat cu ajutorul programului **BLAST**

- de *aliniament multiplu al secvențelor* realizat cu programe de tipul **CLUSTAL**.

În ambele cazuri ideea este de a găsi similaritățile sau diferențele dintre seturile de secvențe. Analiza secvenței reprezintă un instrument foarte important în studii relațiilor de evoluție în genoame, duplicarea genelor, îmbinarea genelor etc.

Datele generate de Proiectul Genomul Uman sunt depozitate în bănci de date a genelor, care stochează secvențe de ADN. În prezent sunt disponibile bănci de date pentru secvențele și structurile proteice. Una din operațiile de bază din bioinformatică constă în căutarea similarității (omologiei) dintre un fragment de ADN nou secvențializat și secvențe de ADN provenite de la diferite organisme. Găsirea unei potriviri apropiate permite predicția tipului de proteină codată de noua secvență. Deși nu este posibilă deocamdată predicția completă a funcției sau structurii unei proteine *de novo* pornind de la secvența sa, pot fi trase niște concluzii folositoare în legătură cu structura și funcția proteinei, în special prin compararea secvenței proteinei cu structură și funcție necunoscută cu secvențe proteice a căror structură și funcții se cunosc. Prin compararea secvențelor proteice echivalente de la diferite specii animale, se pot trage concluzii asupra evoluției acestor specii dintr-un strămoș comun.

2.3. Programul **BLAST**

Un program popular de comparare a secvențelor de ADN este **BLAST** (*Basic Local Alignment Search Tool*). BLAST face parte dintr-un pachet de programe destinat căutării de secvențe proteice, accesibil în diverse forme la diferiți furnizori, sau prin intermediul NCBI, care mai oferă și *Entrez*, un instrument de meta-căutare care acoperă mare parte a bazelor de date de la NCBI, inclusiv cele care găzduiesc structuri tridimensionale a proteinelor, genoamele complete ale organismelor și trimiteri la jurnale științifice care însoțesc intrările din bazele de date.

Asocierea dezvoltărilor tehnologiei de calcul și moleculare deschide noi oportunități cercetărilor genetice. Folosirea combinată a informației oferită de secvențe, a instrumentelor de calcul, a bazelor de date și a biologiei tradiționale crește speranța înțelegerii funcției și reglajelor tuturor genelor și proteinelor, precum și a descifrării funcțiilor celulei.

BLAST reprezintă instrumentul de căutare a aliniamentului local de bază, fiind un set de programe de căutare a similarităților, creat pentru identificarea clasificării și a omologilor potențiali pentru o secvență dată.

Pentru a înțelege mai bine programele BLAST, trebuie cunoscute aspectele de bază ale aliniamentelor secvențelor. Acestea sunt folosite în special pentru găsirea potențialilor omologi ce vor fi folosiți ulterior pentru prezicerea posibilelor funcții ale secvenței necunoscute sau pentru modelarea structurii sale tridimensionale.

Aliniamentul global este cel mai bun aliniament, pe întreaga lungime a secvențelor specificate. Introducerea spațiilor (gaps) în secvențele respective permite alinierea lor pe întreaga lungime. Principalul avantaj al aliniamentului global este optimizarea sa pentru secvențele care au un grad înalt de similaritate, fiind astfel folositor în etapa de aliniere a secvențelor din procesul de modelare a structurii tridimensionale (bazat pe secvențele omologe cu structură tridimensională cunoscută).

Metodele de căutare ale aliniamentului local găsesc aliniamentul optim între subregiuni sau regiuni locale ale secvențelor specificate. Aliniamentul local este cel mai

potrivit pentru secvențe care au regiuni localizate de similarități. Un program de căutare a aliniamentului local este folosit de exemplu pentru găsirea motivelor, domeniilor și altor unități repetitive din secvențele respective, precum și pentru găsirea secvențelor similare pentru secvența necunoscută într-o bază de date. Pe scurt, un program de căutare al aliniamentului local este cel mai bine folosit pentru identificarea unor regiuni secvențiale mai scurte, cu un grad foarte mare de similaritate.

Toți algoritmi de comparare a secvențelor se bazează pe anumite *scheme de calcul a scorului aliniamentului*. Scorul aliniamentului este suma scorurilor mai mici, atribuite pentru fiecare din perechile sale de aminoacizi sau nucleotide. Majoritatea acestor algoritmi folosesc o *matrice de scor* pentru calcularea unui scor total fiecărui aliniament.

Teoria statistică folosită în programele BLAST a fost creată de Samuel Karlin și Steven Altschul.

Toate programele BLAST folosesc o matrice de substituție, atât în etapa de scanare a bazelor de date cât și în procesul de aliniere a secvențelor.

Schemele de substituție sunt considerate a fi cele mai bune metode de calcul al scorului aliniamentelor și se bazează pe analiza frecvenței cu care un aminoacid observat este înlocuit de un alt aminoacid în proteinele ale căror secvențe sunt aliniate.

Criteriile care diferențiază matricele de scor depind de tipul scorului pe care se bazează, astfel avem:

a) *Schemă a scorului bazată pe „identitate”*

Conform acestei scheme, perechile de aminoacizi identici sau nucleotide identice primesc un scor pozitiv, în timp ce perechile non-identice primesc scorul 0. În general scorul pozitiv atribuit perechilor identice este egal cu 1. Scorul identității globale este apoi convertit simplu (identitate procentuală).

Avantaje: această schemă de calcul este simplă și non-heuristică. Este bună în cazul secvențelor cu grad înalt de similaritate.

Dezavantaje: schema este în general inferioară aceluia care încorporează cunoștințele suplimentare, datorită în special inegalităților perechilor non-identice. De exemplu o pereche alanină-valină este mai acceptată din punct de vedere biologic decât o pereche alanină-acid aspartic. Această schemă este mai puțin efektivă în detectarea secvențelor sau a regiunilor secvențiale cu un grad redus de similaritate. Procentul identității raportat de acest aliniament nu este întotdeauna un indicator de acuratețe a gradului de omologie prezent, în special datorită dependenței acestui scor de lungime a secvenței.

b) *Schemă de calcul a scorului bazată pe „similaritate chimică”*

Această schemă a fost concepută pentru a depăși limitările asociate cu schema bazată pe „identități” și evaluează perechile de aminoacizi în funcție de caracteristicile lor chimice și structurale.

Schemele folosite de McLachlan și Feng încorporează în calcularea scorului proprietățile aminoacizilor cum ar fi polaritate, sarcină, mărime și caracteristici structurale.

Avantaje: introduce proprietățile aminoacizilor în calcularea scorului, lucru important deoarece anumite mutații care realizează o schimbare drastică în caracteristicile AA implicați au un impact mult mai mare asupra funcțiilor proteinelor decât altele. Aceste mutații, de exemplu schimbarea unui aminoacid polar cu unul non-polar, alterează mult mai mult structura și funcția proteinei respective decât o mutație implicând aminoacizi cu proprietăți similare.

Dezavantaje: mutațiile observate în natură nu sunt întotdeauna explicate prin schemele simple de calculare a scorului.

c) *Schema de calcul bazată pe „codul genetic”*

Această metodă ia în considerare numărul minim de schimbări de baze la nivel genomic, necesar pentru convertirea unui aminoacid în altul.

d) *Schema de calcul bazată pe „mutații observate”*

Această metodă de calcul a scorului unui aliniament, se bazează pe frecvența mutațiilor observate în secvențele aliniate.

Schemele bazate pe mutațiile observate reprezintă mai bine fenomenele naturale decât acelea care încearcă să explice relațiile dintre secvențe folosind matrice de calcul bazate pe similaritate chimică, identitate și cod genetic.

Algoritmii de căutare a similarităților secvențelor aliniat se bazează pe cele 210 perechi posibile de aminoacizi care sunt reprezentate de o matrice 20x20 de calcul a scorului. Numărul total de perechi posibile de aminoacizi este egal cu 210, „alfabetul” proteinelor fiind alcătuit din 20 AA. Perechile de aminoacizi identici primesc cel mai înalt scor în matrice, urmate de perechile de aminoacizi care au un anumit grad de similaritate (de ex.: Leucină și Izoleucină) și în final de acei aminoacizi care nu prezintă similarități (de ex.: Leucină și Arginină).

Programele BLAST folosesc un algoritm heuristic care identifică aliniamentele locale, găsind omologii cu secvențele cele mai apropiate, într-un timp eficient.

Serverul BLAST suportă o varietate de programe analitice care sunt fie accesate prin rețeaua Internet, fie instalate în rețele locale pentru a mări viteza de analiză. Programul BLAST bazal nu permite introducerea gap-urilor în aliniamentele sale ceea ce va reduce sensibilitatea căutării. Cu toate acestea, datele de ieșire din program oferă aliniamente regionale multiple, care pot fi folosite pentru a anticipa gap-urile din secvența de interes și cea din baza de date. În continuare sunt enumerate programele BLAST și utilizarea lor.

a) **BLASTp:** acest program permite utilizatorului să caute similaritățile dintre secvența unei proteine necunoscute și secvențele proteinelor dintr-o bază de date.

b) **BLASTx:** permite compararea secvențelor traduse în aminoacizi ale nucleotidelor cu secvențele proteinelor din bazele de date.

Secvența nucleotidică de interes este tradusă inițial în toate cele 6 catene de citire **ORF (Open Reading Frame)** posibile. Acest program este folosit în special pentru găsirea erorilor de secvențializare a nucleotidelor, prin compararea secvenței de nucleotide tradusă în aminoacizii săi proteici potențiali dintr-o bază de date cu secvențe proteice.

c) **BLASTn**: cu ajutorul acestui program se compară o secvență nucleotidică de interes cu secvențele din bazele de date nucleotidice.

d) **tBLASTn**: permite căutarea similarităților dintre o secvență proteică și secvențele traduse (translatate) ale nucleotidelor dintr-o bază de date.

Secvențele nucleotidice dintr-o bază de date sunt traduse inițial în fiecare din cele 6 catene de citire posibile și sunt apoi comparate cu secvența proteinei de interes. Acest program este util pentru găsirea erorilor de secvențializare în proteine prin compararea secvenței proteinei respective cu omologii săi potențiali obținuți prin traducerea secvențelor nucleotidice dintr-o bază de date.

e) **tBLASTx**: se compară cele 6 traduceri ale catenelor de citire ale secvenței nucleotidice chestionabile cu cele 6 catene de citire traduse ale secvențelor nucleotidice dintr-o bază.

Noul pachet de programe BLAST este menținut pe serverul BLAST 2.0 capabil să optimizeze viteza de procesare și sensibilitatea metodelor, adăugând pe de altă parte noi capacități ce permit rularea noilor programe PSI-BLAST și GAPPED-BLAST.

GAPPED BLAST – algoritmul Gapped-BLAST permite introducerea gap-urilor în aliniamentele obținute cu ajutorul programului BLAST simplu.

Introducerea gap-urilor (input) previne segmentarea regiunilor similare ale secvențelor.

Datele de intrare ale algoritmului heuristic permit reflectarea relațiilor biologice asociate aliniamentului, în special situsurile active și situsurile de legătură care au tendințe să fie mai conservate de-a lungul evoluției. Introducerea gap-urilor previne scindarea acestor regiuni în fragmente de secvențe mai puțin semnificative.

PSI-BLAST (position – specific iterated BLAST) rulează inițial programul Gapped-BLAST și folosește aliniamentul de ieșire din acesta ca input pentru PSI-BLAST. Programul construiește o matrice de calculare a scorului care înlocuiește secvența originală și este folosită pentru găsirea profiilor (secvențelor omologe) în următoarele iterații de căutare în baza de date.

Utilizatorul ar trebui să efectueze următoarele etape generale pentru rularea cu succes a programelor BLAST:

- Secvența de interes trebuie introdusă în format corect (de exemplu formatul FASTA – similarul formatului BLAST, de pe serverul EBI);
- Secvența astfel formulată va fi apoi copiată în fereastra “input sequence” a interfeței programului BLAST;
- În funcție de tipul secvenței analizate se selectează programul BLAST potrivit (de exemplu BLASTp pentru secvențele de proteine);
- În final trebuie selectată baza de date corespunzătoare. De exemplu, dacă utilizatorul este interesat numai în găsirea secvențelor omologe cu structură cunoscută trebuie selectată o bază de date ce conține structuri tridimensionale, cum ar fi PDB. Secvența de interes este transmisă serverului BLAST, iar rezultatele căutării în baza de date sunt obținute fie prin e-mail, fie văzute interactiv pe interfața Internet a programului BLAST.

Valoarea așteptată, **E**, din datele de ieșire ale programului BLAST reprezintă numărul de potriviri, „perechi” găsite aleatoriu într-o bază de date. O valoare $E=0$ semnifică faptul că pentru anumite baze de date probabilitatea de a găsi o pereche în

mod aleatoriu este 0. Această valoare descrește exponențial cu creșterea valorilor scorului S. O valoare E egală cu 1 indică probabilitatea ca identificarea similarităților secvenței să fie aleatoare.

Formatul **FASTA** constă în reprezentarea fiecărui aminoacid din structura proteinei printr-un cod format dintr-un singur caracter. Codurile acceptate de programele **BLAST** sunt:

A	alanină	P	prolină
B	aspartat sau asparagină	Q	glutamină
C	cysteină	R	arginină
D	aspartat	S	serină
E	glutamat	T	threonină
F	fenilalanină	U	selenocisteină
G	glicină	V	valină
H	histidină	W	triptofan
I	izoleucină	Y	tirozină
K	lizină	Z	glutamat sau glutamină
L	leucină	X	orice aminoacid
M	methionină	*	oprirea translației
N	asparagină nedeterminată	-	spațiu de lungime

2.4. Exemple

a. Alegem opțiunea “Nucleotide: core subset of nucleotide sequence records” (figura 2.4.a).

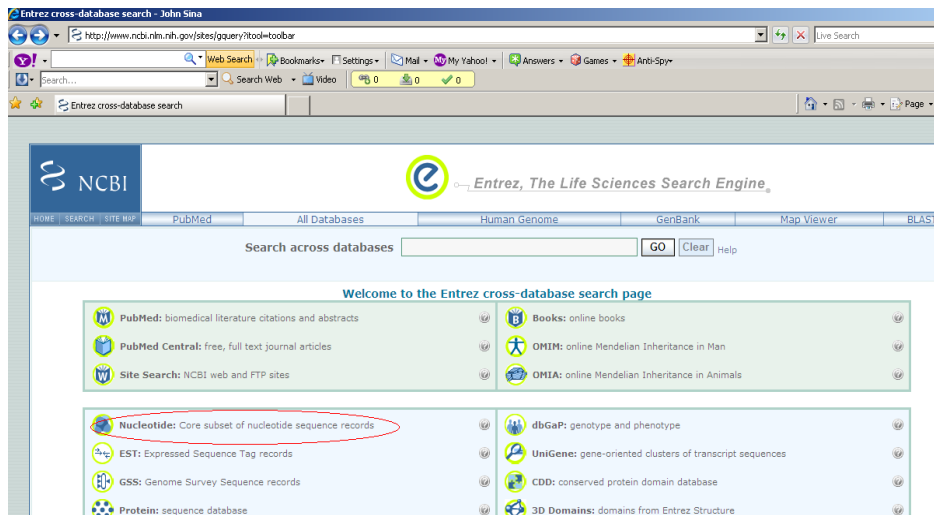


Figura 2.4.a. Interfața NCBI – opțiunea **Nucleotide**

Dacă dorim să căutăm secvența de gene care codează receptorul pentru endotelină (figura 2.4.b) introducem în căsuța search următoarele: **Endothelin receptor**.

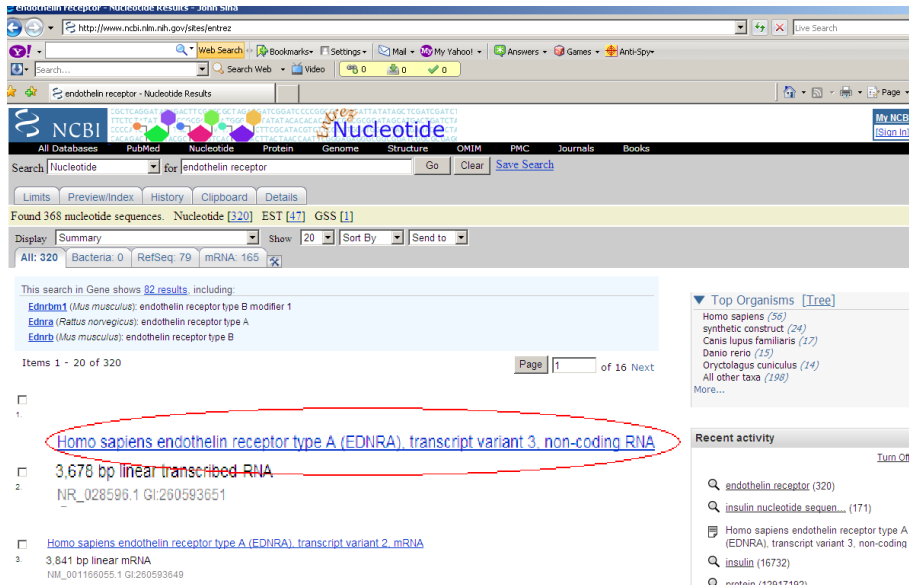


Figura 2.4.b. Exemplu de căutare a unei secvențe de nucleotide

Din mulțimea de răspunsuri vom selecta varianta receptorului pentru specia umană (*homo sapiens*)!

Va fi afișat cromozomul pe care se află secvența de nucleotide (cromozomul 40) – figura 2.4.c, localizarea genei pe cromozom, secvența de nucleotide și secvența codantă.



Figura 2.4.c. Informație despre secvența de nucleotide găsită

b. Pentru a căuta o secvență de nucleotide intrăm pe adresa: <http://blast.ncbi.nlm.nih.gov/Blast.cgi> (figura 2.4.d).

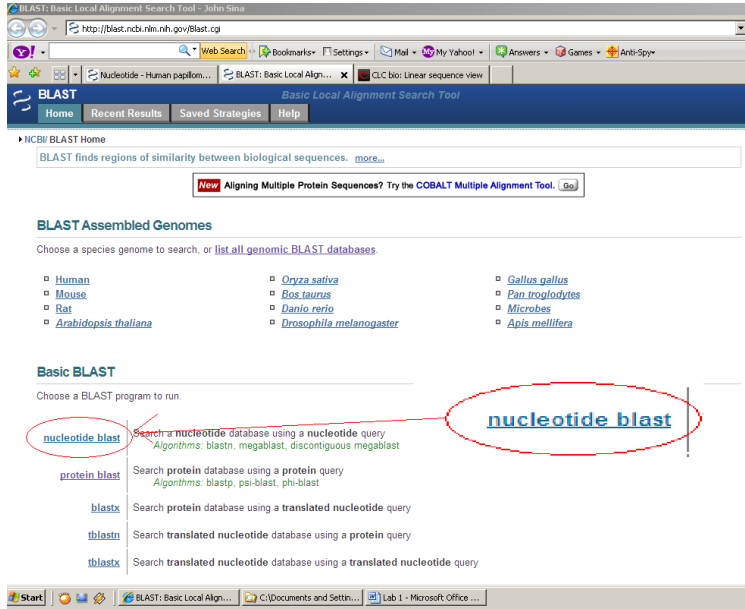


Figura 2.4.d. Opțiunile ferestrei de căutare BLAST

În căsuța de interogare (query) vom introduce următoarea secvență:

```

“1 ctgaaaccg tatgctatat aattatgtac tataaagtaa taatgtatac agtghtaagg
61 atcatgggcc atgtgctttt caaactaatt gtacataaaa caagcatcta ttgaaaatat
121 ctgacaact catcttttat tttgatgtg tgtgtgtgtg tgtgtgtgtg ttttttaac
181 agggatttgg gg”
    
```

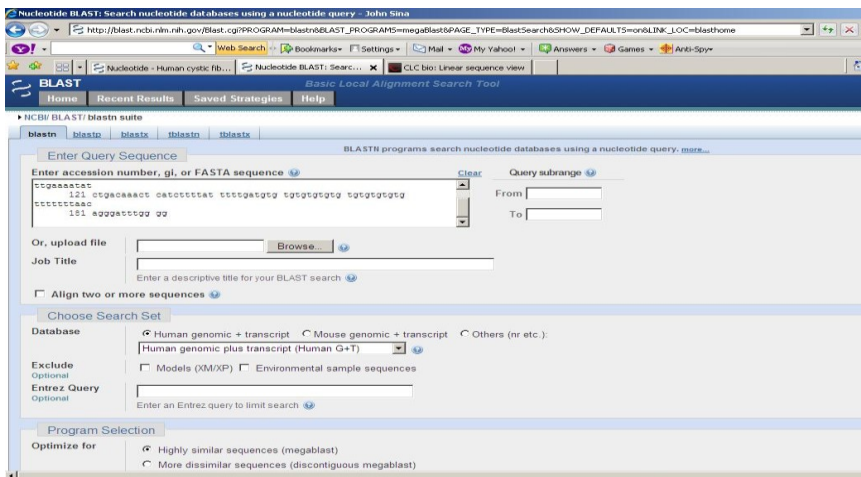


Figura 2.4.e. Interfața cu un exemplu de interogare BLASTn

Apoi vom da click pe butonul BLAST din josul paginii, ca în figura 2.4.f.

Descriptions

Legend for links to other resources: UniGene GEO Gene Structure Map Viewer

Sequences producing significant alignments:
(Click headers to sort columns)

Accession	Description	Max score	Total score
NT_007933.151	Homo sapiens chromosome 7 genomic contig, GRCh37 reference primary assembly	342	342
NW_001839071.2	Homo sapiens chromosome 7 genomic contig, alternate assembly (1)	342	342

Alignments Select All [Get selected sequences](#) [Distance tree of results](#) **NEW**

> [ref|NT_007933.151](#) Homo sapiens chromosome 7 genomic contig, GRCh37 reference primary assembly
Length=77412220

Features in this part of subject sequence:
[cystic fibrosis transmembrane conductance regulator](#)

Score = 342 bits (185), Expect = 6e-92
Identities = 190/192 (98%), Gaps = 2/192 (1%)
Strand=Plus/Plus

```

Query 1      CTAGAAACCGTATGCTATATAAATTATGTACTATAAAGTAATAATGTATACAGTGAATGG 60
             |||..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|
Sbjct 55221358 CTAGAAACCGTATGCTATATAAATTATGTACTATAAAGTAATAATGTATACAGTGAATGG 55221417

Query 61     ATCATGGGCCATGTGCTTTTCAAACTAATTGTACATAAAACAAGCATCTATTGAAAAATAT 120
             |||..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|
Sbjct 55221418 ATCATGGGCCATGTGCTTTTCAAACTAATTGTACATAAAACAAGCATCTATTGAAAAATAT 55221477

Query 121    CTGACAAACTCATCTTTTATTTTTGATGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTAAAC 180
             |||..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|..|
Sbjct 55221478 CTGACAAACTCATCTTTTATTTTTGATGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTAAAC 55221535

```

Figura 2.4.f. Rezultatul interogării

REZULTAT: Am găsit secvența de nucleotide responsabilă de apariția fibrozei chistice!

2.5. Utilizarea Bazelor de date cu secvențe proteice

a. Din pagina principală Entrez vom selecta “Protein: sequence database” (figura 2.5.a).

Entrez cross-database search - John Sina

http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=pubmed

Search across databases [Help](#)

Welcome to the Entrez cross-database search page

- PubMed: biomedical literature citations and abstracts
- PubMed Central: free, full text journal articles
- Site Search: NCBI web and FTP sites
- Nucleotide: Core subset of nucleotide sequence records
- EST: Expressed Sequence Tag records
- dbSS: Genome Survey Sequence records
- Genome: whole genome sequences
- Structure: three-dimensional macromolecular structures
- Taxonomy: organisms in GenBank
- Books: online books
- OMIM: online Mendelian Inheritance in Man
- OMIA: online Mendelian Inheritance in Animals
- dbGaP: genotype and phenotype
- UniGene: gene-oriented clusters of transcript sequences
- CDD: conserved protein domain database
- 3D Domains: domains from Entrez Structure
- UniSTS: markers and mapping data
- PopSet: population study data sets
- GED Profiles: expression and molecular abundance profiles

Figura 2.5.a. Opțiunea BD cu secvențe proteice

Dacă ne propunem să căutăm proteina insulină introducem în căsuța de search **insulin** (figura 2.5.b).

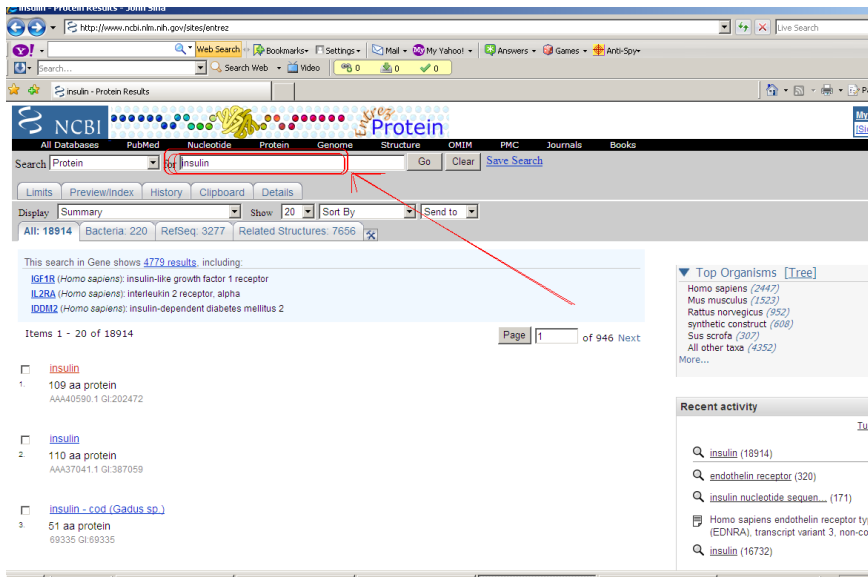


Figura 2.5.b. Exemplu de căutare a insulinei

Vom alege insulina corespunzătoare speciei umane (*homo sapiens*) și vom găsi informația potrivit căreia insulina umană este o proteină formată din 110 aminoacizi a căror reprezentare o puteți vedea în figura 2.5.c.

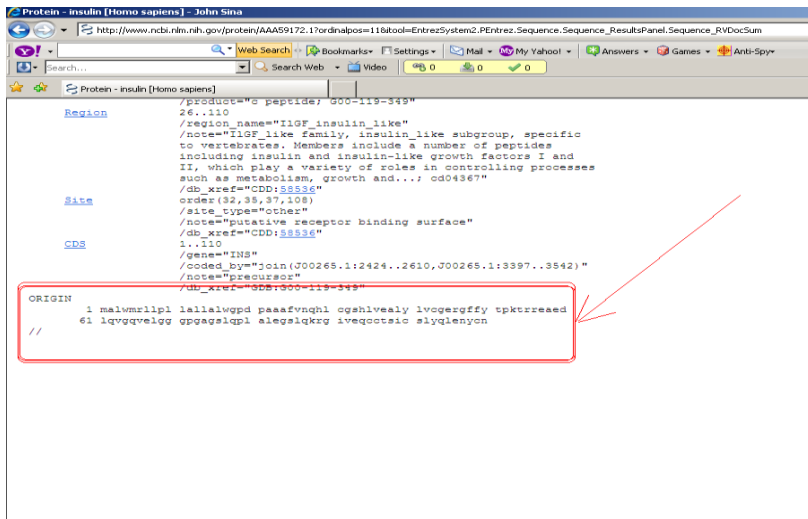


Figura 2.5.c. Secvența de aminoacizi ai proteinei insulina umană

În josul paginii va fi afișată secvența de aminoacizi din care e formată insulina, începând de la origine. După cum știți de la curs, fiecare aminoacid e codat cu o singură literă.

b. Introduceți secvența de aminoacizi **malwmrlpl** (figura 2.5.d).

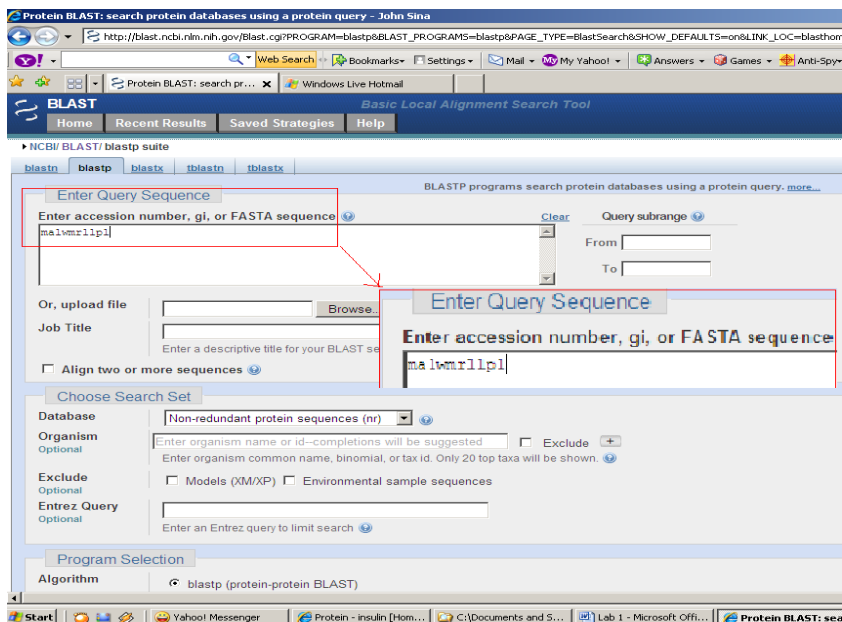


Figura 2.5.d. Intergarea BD după o secvență de aminoacizi

Apoi alegeți opțiunea cu butonul BLAST din josul paginii. Ca rezultat veți avea un grafic care va afișa gradul de potrivire al secvenței de aminoacizi introdusă (figura 2.5.e). Sub grafic obțineți și rezultate text în ordinea potrivirii lor cu secvența introdusă inițial.

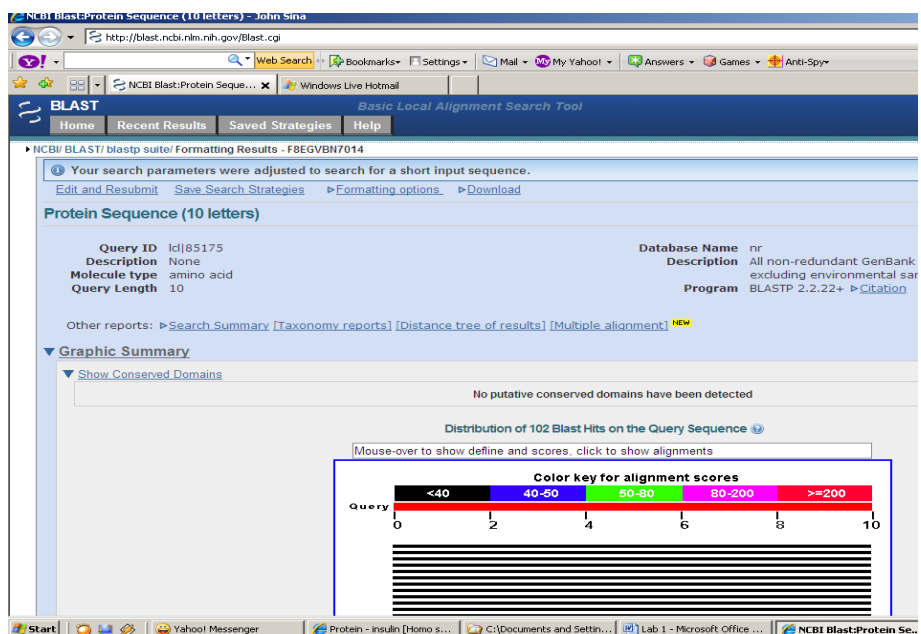


Figura 2.5.e. Rezultatul căutării gradului de potrivire al secvenței de aminoacizi

2.6. *Exerciții propuse*

Găsiți toate secvențele nucleotide de la șoarece și om adăugate în BD Entrez în anul 1997. Câte structuri corespund rezultatului?

Vizualizați structura primară a hemoglobinei.

Copiați o secvență de aminoacid și dați o căutare în BLAST (pentru a găsi asemănarea cu proteinele cunoscute).

Căutați în BLAST următoarea secvență de aminoacizi:

```
"1 mkwvtfisll flfssaysrg vfrdahlkse vahrfkdige enfkalvlia faqylqqcpf
61 edhkvlnvev tefaktcvad esaencdksl htlfgdklct vatlretyge madccakqep
121 ernecflqhk ddnplnprlv rpevdvmcta fhdneetflk kylyeiarrh pyfyapellf
181 fakrykaaft eccqaadkaa clpkldelr degkassakq rikcaslqkf gerafkawav
241 arlsqrfpka efaevsklvt dltkvhtecc hgdllccadd radlakyice nqdsisskkl
301 ecceppllek shciaevend empadlpsla adfveskdvc knyaeakdvf lgmflyeyar
361 rhpdyssvll lrlaktyett lekccaaadp hecyakvfde fkplveepqn likqncelfe
421 qlgeyfkfna llvrytkvp qvstptllev srnlgkvgsk cckhpeakrm pcaedylsvv
481 lnqlcvlhek tpsvdrvtkc cteslvnrrp cfsalevdet yvpkefnaet ftfhadictl
541 sekerqikkq talvelvkhk pkatkeqlka vmddfaafve kockaddket cfaeegkklv
601 aasqaalgl"
```

Căutați în BLAST următoarea secvență de nucleotide:

```
"1 gacacggctg tatatactg tgcctctctt gacagcgcgt ttcgggggca ctggggccat
61 ggaacctcgg tcagcgtctc ctcagcatcc cgcaccagcc ccaaggtctt cccgctgagc
121 ctctcagca cccagccaga tgggaacgtg gtcacgcct gcctgttcca gggcttcttc
181 cccagcagc cactcagtg gacctggagc gaaagcggac agggcgtgac cgccagaaac
241 ttcccacca gccagatgc ctccggggac ctgtacacca cgagcagcca gctgacctg
301 cggccacac agtgctagc cggcaagtc gtagatgcc acgtgaagca ctacacgaat"
```

Căutați secvența de nucleotide pentru gena RHD, RHCE (cromozomul 1). Căutați secvența de aminoacizi a imunoglobulinelor g (regiunea variabilă).

Căutați secvența de nucleotide care codifică sinteza fibrinogenului (factorul XIII al coagularii). Care este gena? (răspuns: F13A1).

Căutați următoarea secvență de aminoacizi:

```
1 mssvavltqe sfaehrglv pqqikvatln seesdppty kdafpplpek aaclesaqep
61 agawgnkirp ikasvitqvf hvplerkyk dmnqfgegeq akicleimqr tgahlelsla
121 kdqqlsimvs gklavamkar kdivarlqtq asatvaipke hhrfvigkng eklqdllekt
```

Căutați secvența de aminoacizi a mioglobinei.

3. Structura secundară, terțiară și cuaternară a proteinelor. Programele PDB, Cn3D și Rasmol

3.1. Obiectivele lucrării de laborator

- recunoașterea structurii secundare, terțiare și cuaternare a proteinelor
- exploatarea structurilor proteice cu PDB
- vizualizarea structurilor cu programul Cn3D
- vizualizarea structurilor cu programul Rasmol.

3.2. Noțiuni introductive

Proteinele sunt substanțe organice macromoleculare formate din lanțuri simple sau complexe de aminoacizi; ele sunt prezente în celulele tuturor organismelor vii în proporție de peste 50% din greutatea uscată. Toate proteinele sunt polimeri ai aminoacizilor în care secvența acestora este codificată de către o genă. Fiecare proteină are secvența ei unică de aminoacizi, determinată de secvența nucleotidică a genei.

3.2.1. Proprietăți fizico-chimice

a) Masa moleculară

Datorită formării aproape în exclusivitate din aminoacizi, putem considera proteinele ca fiind de fapt niște polipeptide, cu masă moleculară foarte mare între 10.000 și 60.000.000. Masa moleculară se determină prin diferite metode, mai ales în cazul proteinelor cu masa moleculară foarte mare ca de exemplu proteina C reactiva. În următorul tabel sunt enumerate diferite proteine împreună cu masa moleculară corespunzătoare:

Tabel 3.2.1.a. Exemple de proteine, sursa și masa moleculară

Denumirea proteinei	Sursa proteinei/Izolată din	Masa moleculară
Lactalbumină	lapte	17000
Gliadina	grâu	27.500
Insulina	pancreas	12.000
Hordeina	orz	27.500
Hemoglobina	globule roșii	68.000
Hemocianina	Moluște (sânge), artropode (sânge)	2.800.00
Miozina	mușchi	850.000
Pepsină	stomac	36.000
Peroxidaza	rinichi	44.000
Virusul mozaicului tutunului (capsida)	tutun	17.000.000

Deoarece la multe proteine masa moleculară apare ca un multiplu de 17,500, multă vreme s-a mers pe ipoteza că particulele proteice sunt formate prin unirea mai multor molecule de bază ce au masa moleculară în jurul valorii de 17,500. Aceste molecule de bază s-ar putea uni între ele prin așa numitele valențe reziduale, ducând la formarea de *agregate moleculare*. Atunci când are loc ruperea acestor valențe reziduale ar avea loc doar modificarea proprietăților fizice ale proteinelor, în timp ce dacă are loc ruperea legăturilor principale (legăturile peptidice), proteina își modifică proprietățile fizico-chimice.

b) Solubilitatea proteinelor

Proteinele sunt substanțe solide, macromoleculare, solubile în general în apă și insolubile în solvenți organici nepolari. Unele proteine sunt solubile în apă dar insolubile în alcool, altele sunt solubile în soluții apoase de electroliți, acizi organici. Datorită gradului diferit de solubilitate în diferiți solvenți, proteinele se pot izola, identifica și separa. Solubilitatea lor depinde foarte mult de legăturile care se stabilesc între grupările libere de la suprafața macromoleculilor și moleculele solventului. La suprafața macromoleculilor proteice se găsesc grupări libere de tip polar: $-\text{COOH}$, $-\text{NH}_2$, $-\text{OH}$, $-\text{SH}$, $-\text{NH}$, grupări cu caracter hidrofil care favorizează dizolvarea proteinelor în apă. De asemenea există grupări de tip apolar, hidrofobe, de regulă radicali de hidrocarburi: $-\text{CH}_3$, $-\text{C}_6\text{H}_5$, $-\text{C}_2\text{H}_5$, care favorizează dizolvarea proteinelor în alcool. Însă în marea lor majoritate predomină grupările polare, determinante pentru caracterul hidrofil. În contact cu apa proteinele greu solubile manifestă fenomenul de gonflare, datorită tendinței de hidratare datorată grupărilor polare. Gelatina de exemplu se îmbibă foarte puternic cu apa dând naștere prin răcire la geluri. La dizolvarea proteinelor în apă are loc fenomenul de formare a coloizilor hidrofilii. S-a constatat că în soluții diluate se găsesc macromolecule proteice izolate, iar în cazul soluțiilor concentrate se formează agregate de macromolecule proteice. Soluțiile coloidale ale proteinelor coagulează prin încălzire și prezintă efectul Tyndall (dispersia fasciculului de lumină).

c) Structura secundară

Imaginea alfa helixurilor mioglobinei, a cărei structură a fost determinată de către Max Perutz și Sir John Cowdery Kendrew în 1958 folosind crystalografia cu raze X, este prezentată în figura 3.2.1.a.



Figura 3.2.1.a. Structura secundară a mioglobinei

Structura secundară se referă la forma și la lungimea lanțurilor polipeptidice, proprietăți induse de legăturile de hidrogen. Cele mai întâlnite tipuri de structura

secundară sunt alpha helixul și lanțurile beta (figura 3.2.1.b). Elicea alpha se formează prin rotația unui lanț polipeptidic în jurul propriei axe.

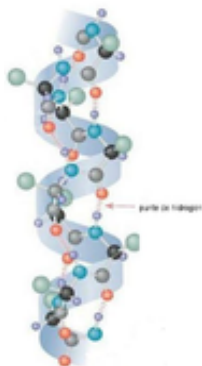


Figura 3.2.1.b. Structura secundară sub formă de lanțuri beta

Alte helix-uri cum ar fi helixul 3_{10} și helixul π sunt din punct de vedere energetic favorabile formării legăturilor de hidrogen, dar sunt rareori observate în proteinele naturale exceptând părțile terminale ale helixului α în timpul formării scheletului proteic (de obicei centrul helixului). Aminoacizii au un comportament diferit vis-a-vis de posibilitatea formării structurii secundare. Prolina și glicina sunt cunoscuți ca așa numiții „helix breakers” (spărgători de helix), deoarece afectează configurația scheletului proteic; ambii aminoacizi au abilități conformaționale neobișnuite și de regulă se găsesc în colțurile scheletului proteic. Aminoacizii care preferă să adopte configurația helixului proteic fac parte din așa numita serie MALEK (codurile formate din 1 literă a aminoacizilor: metionină, alanină, leucină, acid glutamic și lizina); prin contrast aminoacizii aromatici (triptofanul, tirozina și fenilalanina) dar și aminoacizii cu legare prin carbonul beta (izoleucina, valina și treonina) adoptă configurația β .

Structura secundară cunoaște câteva ipoteze privind formarea ei:

- Teoria polipeptidică formulată de către E. Hoffmeister în 1902 și dezvoltată ulterior de către E. Fischer, are la bază conceptul conform căruia moleculele proteice sunt formate din lanțuri polipeptidice foarte lungi. Teoria are câteva dezavantaje:
 - nu explica diferențierea biologică a anumitor proteine
 - unele proteine sunt rezistente la acțiunea enzimelor proteolitice (deși datorită lungimii lanțului nu ar trebui).
- Teoria plierii și răsucirii lanțului polipeptidice a fost elaborată de către Corey și Pauling în 1943 și a fost confirmată prin spectrele de difracție cu raze X, microscopului electronic, prin măsurarea unghiurilor de valență, a distanțelor interatomice, au confirmat faptul că lanțul polipeptidic se găsește sub formă pliată.
 - Plierea catenei are loc prin formarea legăturilor de hidrogen între gruparea carboxilică a unui aminoacid și gruparea aminică a aminoacidului vecin. Lanțul polipeptidic pliat se prezintă ca o panglică îndoită alternativ la dreapta și la stânga, plierea având loc în dreptul carbonilor metinici. Mai multe lanțuri pliate polipeptidice dau naștere unei rețele, între aceste lanțuri pliate putându-se de asemenea forma legături de hidrogen, acestea fiind în număr mai mare când grupările terminale a două lanțuri sunt aranjate diferit ($-\text{NH}_2$ și COOH , sau HOOC -și $-\text{NH}_2$). Catenele polipeptidice pliate predomină în proteinele fibrilare și mai puțin în cele globulare. După valoarea perioadei de identitate se cunosc mai multe tipuri de proteine cu structură pliată. Prin

perioada de identitate se înțelege distanța cea mai mică la care se repetă aminoacizii identici din moleculă.

- Structura α elicoidală, ipoteză lansată de Corey și Pauling e o ipoteză conform căreia lanțul polipeptidic se poate prezenta și înfășura sub formă de spirală. În acest model, fiecare spiră conține de obicei 27 aminoacizi, iar distanța între spire este de $5,44 \text{ \AA}$. Fiecare aminoacid mărește spira cu $1,47 \text{ \AA}$. În fața fiecărei grupări $-\text{CO}-$ va apare la o distanță de $2,8 \text{ \AA}$ o grupare NH de la al treilea aminoacid. Între aceste grupări se stabilesc punțile de hidrogen care asigură stabilitatea α helix-ului. În acest model lanțul polipeptidic se prezintă sub forma unui șurub cu pasul fie spre dreapta, fie spre stânga. În cazul proteinelor naturale, acestea datorită conținutului în L-aminoacizi, pasul helixului va fi spre dreapta, catenele laterale ies în afara corpului propriu-zis putând reacționa fie cu moleculele solventului fie cu alte catene polipeptidice. Canalul format în interiorul helixului este foarte îngust, în el nu poate pătrunde molecula solventului. Legăturile peptidice sunt plane, iar 2 planuri consecutive $-\text{CO}-\text{NH}-$ formează un unghi de 180° , rotirea lanțului făcându-se la carbonul α (metinic).

d) *Structura terțiară*

Prin intermediul cristalografiei cu raze X s-a dovedit faptul că macromoleculele proteice au o conformație tridimensională, realizată de obicei prin intermediul cuplării mai multor lanțuri polipeptidice scurte între ele, cuplare care duce la formarea fibrelor proteice; legăturile intercatenare pot fi principale sau secundare:

- Legături de hidrogen, sunt legături coordinativ heteropolare care se stabilesc cu ușurință între gruparea carbonil $\text{C}=\text{O}$ (electronegativă) și gruparea NH (electropozitivă), din 2 lanțuri polipeptidice alăturate, sau în cazul formelor lactam-lactimă între gruparea $-\text{OH}$ și azotul iminic $=\text{NH}$ (figura 3.2.1.c).

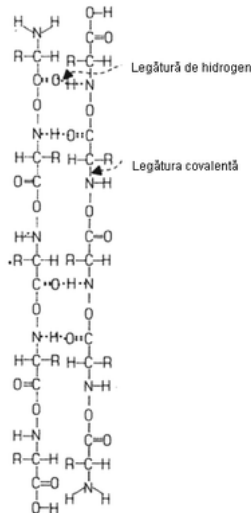


Figura 3.2.1.c. Exemplu de structură terțiară

Legăturile de hidrogen au lungimea cuprinsă între $2,7-3,1 \text{ \AA}$ și energia de $3-7 \text{ Kcal/mol}$ la peptide, iar la apă $2-3 \text{ Kcal/mol}$. Legăturile de hidrogen se pot stabili și între catenele laterale care au grupări carboxil, hidroxil, amino sau tiolice. Din punct de

vedere energetic legătura de hidrogen nu este puternică dar datorită răspândirii relativ uniforme de-a lungul scheletului proteic oferă proteinei stabilitatea necesară.

– Legături disulfidice

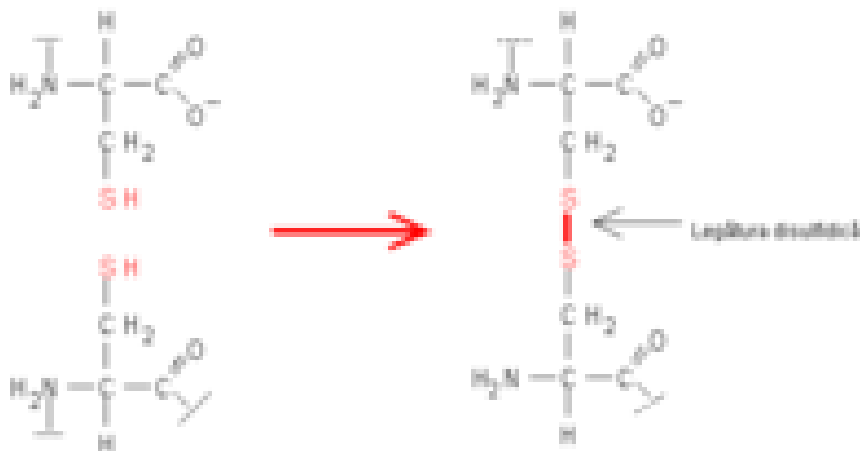


Figura 3.2.1.d. Exemplu de legătură disulfidică

Legătura disulfidică (figura 3.2.1.d) este foarte puternică, 50-100kcal/mol și are un rol foarte important în stabilizarea arhitecturii spațiale a moleculei proteice. Legătura este rezistentă la hidroliză, însă se poate desface iar prin reducere formează tioli(SH), iar prin oxidare formează acizi. În general legătura sulfidică se întâlnește la proteinele transformate care au o rezistență mecanică mare. În afară de aceste legături se mai pot stabili alte tipuri de legături: legături ionice (stabilite de obicei între grupările aminice și cele carboxilice ionizate), legături de tip *van der Waals* (legături electrostatice slabe care se stabilesc între radicalii hidrofobi), legături fosfodiesterice (între două resturi de serină și acid fosforic) și legături eterice (stabilite la nivelul aminoacizilor cu grupări hidroxilice).

e) *Structura cuaternară*

Structura cuaternară se referă la modul cum se unesc subunitățile proteice. Enzimele care catalizează asamblarea acestor subunități poartă denumirea de holoenzime, în care o parte poartă denumirea de subunități reglatoare și subunități catalitice.



Figura 3.2.1.e. Vizualizarea 3D a hemoglobinei

Vederea 3 D a hemoglobinei (figura 3.2.1.e) prezintă 4 subunități roșu și galben, iar unitatea hemică e verde. Numele de hemoglobină este format din hem și globină, denumire ce denotă faptul că hemoglobina are la bază proteine globulare cuplate cu o grupare hem. Proteinele care au structura cuaternară sunt hemoglobina, ADN polimeraza și canalele ionice, dar și nucleozomi și nanotubuli, care sunt complexe multiproteice. Fragmentele proteice pot suferi transformări în structura cuaternară, transformări care se reflectă fie în structurile individuale fie în reorientările fiecărei subunități proteice. Numerele subunităților din oligomerice sunt denumite prin adăugarea sufix-ului *-mer* (grecescul pentru subunitate), precedat de numele subunității, astfel:

Tabel 3.2.1.b. Numerotarea subunităților oligomerice

• 1 = monomer	• 7 = heptamer	• 13 = tridecimer	• 19 = nonadecamer
• 2 = dimer	• 8 = octamer	• 14 = tetradecamer	• 20 = eicosamer
• 3 = trimer	• 9 = nonamer	• 15 = pentadecamer*	• 21-mer
• 4 = tetramer	• 10 = decamer	• 16 = hexadecamer	• 22-mer
• 5 = pentamer	• 11 = undecamer	• 17 = heptadecamer*	• 23-mer
• 6 = hexamer	• 12 = dodecimer	• 18 = octadecamer	• etc.

3.3. *Determinarea structurii proteinelor prin tehnici experimentale și metode de modelare comparativă*

Modelarea moleculară este ansamblu de activități întreprinse pentru calcularea unor proprietăți moleculare, statice și/sau dinamice, din care se pot estima proprietăți microscopice. Modelarea moleculară caută să simuleze interacțiuni intra- și inter-moleculare pentru a înțelege procesele fizico-chimice cu importanță în biologie și medicină, pentru a testa ipoteze și pentru a prezice evenimente noi.

Un rol central în modelare îl joacă interacția cu calculatorul, fie prin mijloace grafice fie prin tastatură-folosind un grup de comenzi care permit monitorizarea și influențarea stimulărilor.

Modelarea macromoleculilor biologice a fost inițiată în anii 1940, cu mult înaintea apariției uneltelor computaționale avansate. Modelarea structurilor macromoleculilor biologice ne permit studierea în profunzime a caracteristicilor funcționale moleculare.

Primele macromolecule modelate au fost dublu-helixul de ADN și mioglobina. Absența uneltelor computaționale avansate a îngreunat mult aceste eforturi de modelare.

Linus Pauling a fost părintele modelării moleculare, structura alpha-helixului modelată de el permițând identificarea unor „modele” (pattern) structurale numite structuri secundare (alpha-helix, beta, strand) care există la nivelul molecular al proteinelor în toate organismele.

Sir Francis Crick și James Watson au primit Premiul Nobel pentru descoperirea structurii dublu catenare a moleculei de ADN, rezolvând astfel multe ambiguități legate de ereditate.

În zilele de astăzi, cristalografia, rezonanța magnetică nucleară (NMR) și microscopia crio-electronică ce permite defracție electronică sunt tehnicile folosite pentru obținerea structurilor de înaltă rezoluție ale macromoleculilor biologice.

În cristalografie proteina cristalizată este bombardată cu electroni, modelul de difracție al electronilor fiind folosit pentru determinarea structurii atomice a moleculei.

Modelul de difracție este folosit pentru calcularea coordonatelor atomice bazat pe măsurarea densității electronice asociată cu fiecare dintre atomii detectați. Întrucât această metodă nu este limitată de mărimea moleculelor, poate fi folosită în principiu pentru determinarea structurii oricărei macromolecule și în cristalografia cu dezavantaje: dificultatea de a obține cristale regulate a milioane de unități identice, multe din proteinele noastre esențiale (de exemplu proteinele membranare) sunt incapabile să cristalizeze când sunt îndepărtate din mediul lor; calcule extensive asociate cu analiza datelor. O altă limitare a metodei este lipsa informației legată de dinamica și flexibilitatea macromoleculelor hologice în mediul lor, cristalografia oferind doar o imagine rapidă a unei molecule deformate expunerii de lungă durată necesară colectării modelelor de defracție.

O altă metodă folosită pentru determinarea structurii macromoleculare, este rezonanța magnetică nucleară. Chiar dacă metoda pare să fie avantajoasă și rezonabilă, datorată folosirii semnăturilor protonice care, se suprapun în cazul moleculelor mai mari, impun folosirea NMR-ului în cazul macromoleculelor mai mici de 30 KD, de aceea sunt necesare tehnici alternative (custalograph, difracția electronică) pentru determinare structurii moleculelor mai mari.

Calculare extensive din analizele cristalografice și RMN sunt efectuate cu acuratețe de programe automate și semi automate ce reduc mai mult timpul necesar obținerii structurilor macromoleculare.

Cu ajutorul computerelor performante din ziua de azi, analiza unei molecule proteice tipică poate fi realizată în câteva zile în timp ce în absența acestor tehnologii ar fi necesari mult mai mulți ani pentru cercetarea respectivă. Cunoașterea aspectelor structurale ale proteinei de interes va genera o imensă cantitate de informație despre funcția sa potențială și despre relațiile cu alte macromolecule esențiale.

Predicția structurii de înaltă rezoluție a proteinelor și al mecanismelor corecte de pliere al lor constituie însă o problemă majoră în biologie, deoarece proteinele sunt compuse din 20 de aminoacizi diferiți care introduc variabilitatea esențială a secvențelor observate la aceste molecule.

Fiecare din acești aminoacizi adopta variate conformații în raport cu celelalte reziduuri din proteină. Teoretic, numărul posibilităților structurale la care proteina se poate conforma în funcție de secvența sa de aminoacizi este fenomenal însă în realitate doar una sa câteva structuri active sunt formate în soluție.

Numeroase tehnici și metodologii au fost utilizate pentru descoperirea acestor structuri active ale proteinelor din numărul imens de posibilități conformaționale, cele mai folosite fiind cele care folosesc relația proteine cu omologi cunoscuți și cele de minimizare energetică axate pe termodinamică.

Așa cum am amintit anterior, rata de obținere a structurilor macromoleculelor prin cristalografie sau rezonanță magnetică nucleară este mult mai mică decât numărul foarte mare de secvențe noi ADN și proteice introdus zilnic, fiind astfel necesară o abordare automatizată a designului structural. Grupurile de informaticieni și cercetători își direcționează eforturile spre identificarea macromoleculelor biologice-cheie (de ex. proteina) responsabile de procesele patologice și propunerea inhibitorilor potențiali ai acestor molecule. Cel mai mare obstacol întâmpinat este lipsa datelor structurale ele

moleculii de interes. În majoritatea cazurilor macromoleculele studiate sunt proteine structurale sau reglatorii. În cazul moleculelor pentru care nu există date structurale obișnuite cristalografic sau prin NMR, pot fi folosite date ale proteinelor omologe care au o structură cunoscută. Modelarea comparativă este bazată pe similaritățile structurale dintre proteina necunoscută și omologă sau cu structura tridimensională cunoscută și unul din programele performante care realizează această modelare este programul „Homology” din pachetul software Insight II al companiei MSI.

În continuare vor fi prezentate detalii legate de modelarea comparativă cu avantajele și limbajele acesteia.

Multe dintre proteinele esențiale prezente în organismul uman se găsesc și în alte organisme vii. Aceste proteine au un rol cheie în menținerea vieții. De exemplu proteinele implicate în catalizarea unor procese esențiale, cum ar fi replicarea ADN-ului, trebuie să-și conserve funcționalitate de-a lungul evoluției și sunt comune tuturor organismelor în timp ce alte proteine mai puțin esențiale sunt mai caracteristice pentru anumite organisme.

Proteinele omologe sunt de obicei proteine care au o compoziție similară de aminoacizi și o origine evoluționară comună. Schimbări în secvența aminoacizilor din proteină poate conduce la schimbarea structurii sale tridimensionale. Această relație între secvența aminoacizilor din proteină și structura tridimensională a proteinei ne permite compararea proteinelor care nu au date structurale oferite de cristalografie sau NMR cu secvențele omologe având structuri tridimensionale cunoscute.

În modelarea proteică bazată pe omologie, structurile determinate experimental reprezintă modelele „template” pentru secvențele omologe „țintă” care nu au coordonate structurale.

Protocolul standard al modelării comparative cuprinde 5 etape (figura 3.3.a):

1. izolarea secvenței ADN
2. găsirea omologilor cu structura cunoscută pentru proteina „țintă”, această căutare putând fi realizată cu programul BLAST;
3. alinierea secvențelor „țintă” cu aceste secvențe „model” pentru identificarea structurii model cea mai apropiată;
4. construirea modelului tridimensional al secvenței „țintă” pornind de la modelele cele mai corespunzătoare găsite în etapa a 2-a;
5. evaluarea calitativă a modelului secvenței țintă, folosind criterii diverse (de ex. energetice, stereochemice).

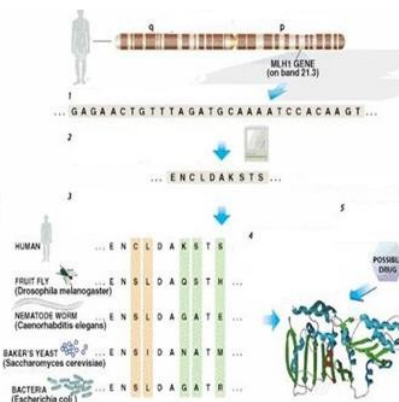
În mod ideal compoziția în aminoacizi a unei proteine necunoscute și cea cunoscută a omologilor săi structurali sunt oarecum similare și în baza de date proteică (PDB) sunt prezente mai multe structuri omologe provenind din specii diferite.

Programul de căutare BLAST folosește ca și date de intrare pentru căutare secvenței de aminoacizi a proteinei țintă, pentru căutarea secvențelor aminoacizilor omologe în baza de date specializată (de ex.: PDB, SWISS-Prot). La sfârșitul căutării sunt afișate structurile omologe cele mai apropiate de structura țintă cu atât mai înrudite cu cât scorul BLAST este mai mare.

Coordonatele structurale ale proteinelor celor mai similare cu proteina țintă sunt salvate și folosite ca date de intrare pentru programul de modelare comparativă (de ex. software-ul „Homology” din Insight II).

Metode de găsim a moleculelor țintă în terapia antitumorală

1. Izolarea secvenței de ADN
2. Translația secvenței de ADN în secvența de aminoacizi folosind programe de calculator
3. Căutarea secvențelor similare în bazele de date proteice (ariile verzi indică diferențe mari, ariile portocalii diferențe mici)
4. Modelarea proteinei țintă, bazată pe structura cunoscută a unor proteine omologe
5. Găsirea unei molecule-medicament care să se lege de proteina modelată



Exemplu folosind gena MLH1, de pe cromozomul 3, asociată cu cancerul de colon

Sursă bibliografică: Howard K., Scientific American, Iulie 2000

Figura 3.3.a. Protocolul standard al modelării comparative

Coordonatele atomice ale structurilor „matriță” (template) sunt aliniate structural pentru a reprezenta regiunile structurale conservate (RSC_s) ale proteinei „țintă”. Folosirea SCR_s – urilor oferă posibilitatea construirii caracteristicilor structurale conservate evoluționist pentru proteina țintă.

Multe regiuni de buclă ale proteinelor omologe au o mare similaritate a secvențelor de aminoacizi dar au caracteristici structurale tridimensionale variabile. Cu alte cuvinte, relația între regiunile structurale conservate și secvențele de aminoacizi aliniate nu este așa de evidentă cum ar fi de așteptat. De aceea în procesul modelării comparative este foarte importantă stabilirea acestor RSC_s, majoritatea acestei regiuni conservate fiind reprezentate de structurile secundare (α - helix, structură β , etc.). Structurile regiunilor din afara SCR-urilor sunt de obicei buclele structurale (loops) și capetele terminale ale moleculelor.

După stabilirea coordonatelor aminoacizilor din secvențele suprapuse la nivelul SCR-urilor trebuie găsite cele mai corespunzătoare coordonate pentru regiunile de buclă ale proteinei țintă, precum și ale capetelor terminale. Aceste deziderate pot fi realizate cu ajutorul programului „Homology” din pachetul de programe comerciale „Insight II” al companiei MSI.

Uneltele software de predicție a structurilor proteice sunt în general clasificate în două domenii: domeniul public accesibil și cel comercial, privat.

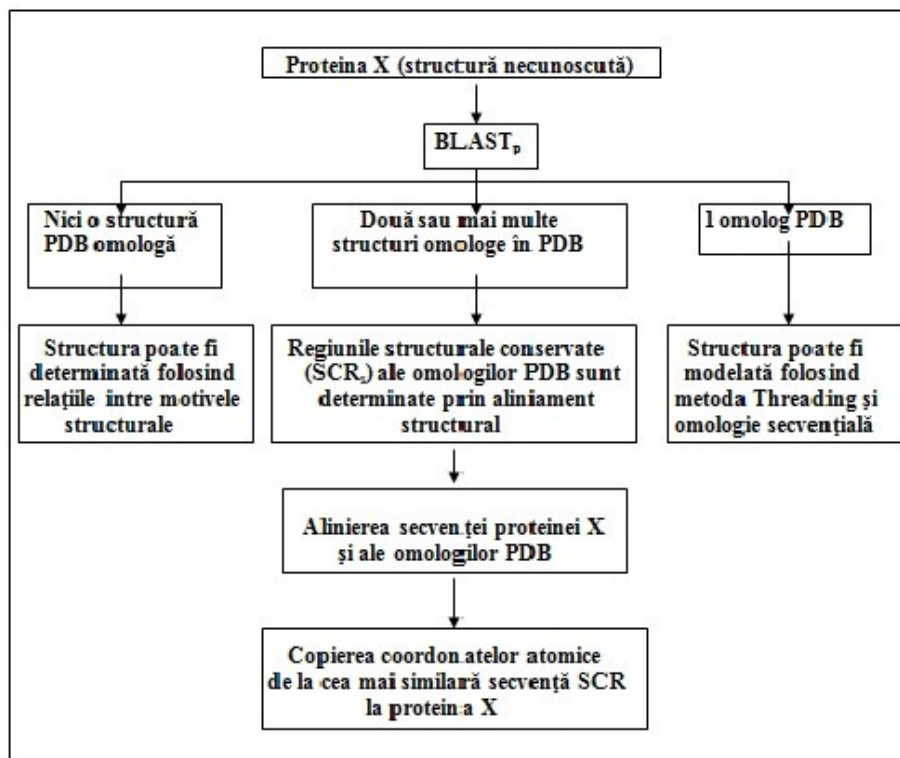


Figura 3.3.b. Modelarea structurilor proteice bazată pe omologie

Unul din programele comerciale foarte des folosit este mediul Insight II al companiei biotehnologiei și bioinformaticii MSI (Molecular Simulation Inc.) care a revoluționat studiile de design medicamentos prin adăugarea elementelor de predicție a structurii proteice la abordarea empirică învechită.

La ora actuală programele Insight II rulează numai pe computere Silicon Graphics și IBM RISC system/6000 workstation. Majoritatea modulelor software ale companiei MSI necesită programul grafic Insight II 3-D graphics, un mediu prietenos și ușor de folosit de către utilizatori.

3.4. Grafica moleculară

Reprezentarea computerizată a modelelor biologice a fost realizată pentru prima oară de către Cyrus Levinthal și colegii săi care au creat un sistem ce a afișat pe un osciloscop reprezentări ale structurilor macromoleculare.

Primul sistem pentru afișajul interactiv al structurilor moleculare a fost creat la MIT în 1966 (Eric Francoeur, „Early Interactive Molecules Graphics at MIT, 1998, <http://www.umass.edu/molvis/francoeur/levinthal/lev-index.html>).

Folosind unul din primele calculatoare, Project MAC, C. Levinthal și colegii săi au conceput un program „model Building” care să lucreze cu structurile proteice.

Programul a permis studiul interacțiunilor de scurtă durată între atomi precum și „manipularea” structurilor moleculare. Terminalul de afișare a structurii moleculare a

fost un osciloscop monocromic arătând structura în formă atomică cu reprezentare grafică de „sârmă”. Efectul tridimensional a fost obținut prin rotația constantă a structurii pe ecran.

La începutul anilor 1970, pentru prima dată, structura unei proteine a fost obținută cristalografic și vizualizată în întregime cu ajutorul calculatoarelor de către Jane și David Richardson și colegii. La sfârșitul anilor 1970, Thomas Porter a dezvoltat algoritmi pentru reprezentările spațiale care au revoluționat vizualizarea macromoleculilor dar au avut acces limitat, numai pentru specialiști.

În anii 1980 David și Jane Richardson au realizat reprezentări grafice ale structurilor moleculare folosind programul CHAOS iar în 1992 au descris „Kinemage” (imaginea cinetică) și programele create de ei, MAGE și PREKIN.

- MAGE este accesibil la adresa <http://kinemage.biochem.duke.edu>, „The Richardsons 3D Protein Structure Laboratory and Kinemage Home Page”. Cuprinde următoarele:
- MAGE– pentru afisajul Kinemage-lor
- PREKIN prepară imaginile cinetice (Kinemage’s) moleculare necesare ca input pentru programul MAGE, pornind de la coordonatele în format PDB
- REDUCE – pentru adăugarea și optimizarea hidrogenilor
- PROBE – calculul contactelor interatomice.

MAGE poate arăta diferite conformații ale unei molecule.

Această „capacitate de animație moleculară” a programului MAGE nu este valabilă și pentru programul Rasmol.

În ultimii ani au fost create programe de vizualizare grafică relativ necostisitoare, accesibile populației interesată de cercetarea moleculară cum ar fi RASMOL, Kinemage și Chemscape Chime Viewer accesibil la adresa <http://www.mdli.com/chemscape/chime>.

Grafica moleculară este pasul spre dimensiunea a treia: analiza 3D, manipularea în 3D, precum și compararea diferitelor conformații.

Animația joacă un rol esențial în vizualizarea simulărilor de dinamică moleculară.

Echipamentul hardware, parte componentă a tehnologiei de grafică moleculară, este completat de software-ul adecvat, necesar unor tehnici curente de vizualizare: zooming (alterarea scalei de vizualizare), clipping (simularea mișcării în planul Z al monitorului, care permite vizualizarea în adâncime), precum și algoritmi de rotație-translație, acuratețea geometrică de reprezentare a atomilor și lungimilor de legătură, manipularea precisă a obiectelor în procesul de superpoziție atomică, vizualizarea spațiului torsional (diagrama Ramachandran).

Alte suprafețe și volume vizualizate curent sunt orbitalele moleculare, potențialul electrofilic sau lipofilic proiectate pe suprafața van der Waals sau pe suprafața accesibilă solventului (de tip “Connolly”). Suprafața poate fi la rândul ei opacă, translucență sau reprezentată sub formă de puncte. Suprafețele Connolly sunt adesea vizualizate în combinație cu potențialul electrostatic.

Pentru proteine există convenții de vizualizare a structurii secundare folosind obiecte grafice de tip „panglică” sau „cilindru”.

Programele de modelare moleculară sunt capabile să execute asemenea operații fie pe baza structurii predefinite (în format PDB), fie pe baza unor reguli care se referă la unghiurile de torsiune phi-psi, putând fi astfel vizualizată cu claritate structura secundară, terțiară sau cuaternară a unei proteine.

3.5. Exemplu de explorare a structurilor proteice cu PDB (Protein Data Bank)

- Accesați adresa de internet <http://www.rcsb.org/pdb/home/home.do>
- căutați “Protein Kinase C Interacting Protein with PDB-ID 1AV5” (figura 3.5.a).

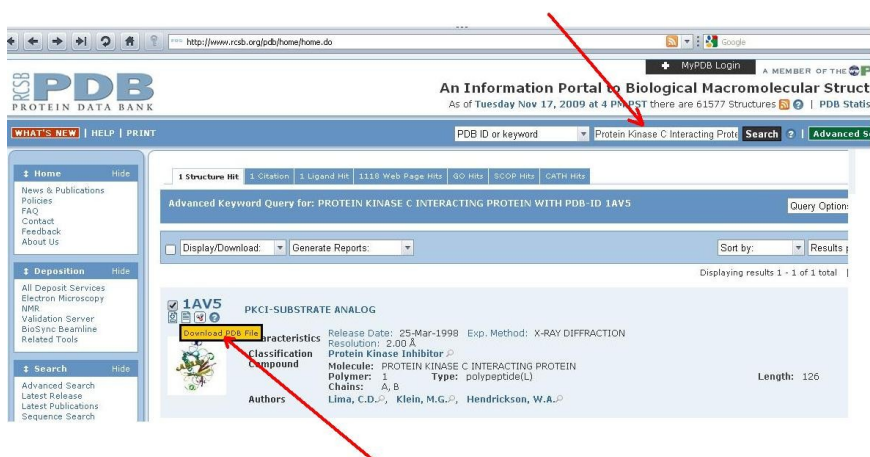


Figura 3.5.a. Interogare cu PDB

- Salvați structura în format.pdb
- Structura terțiară vizualizată cu jmol o alegeți ca în figura 3.5.a, iar rezultatul e prezentat în figura 3.5.b.

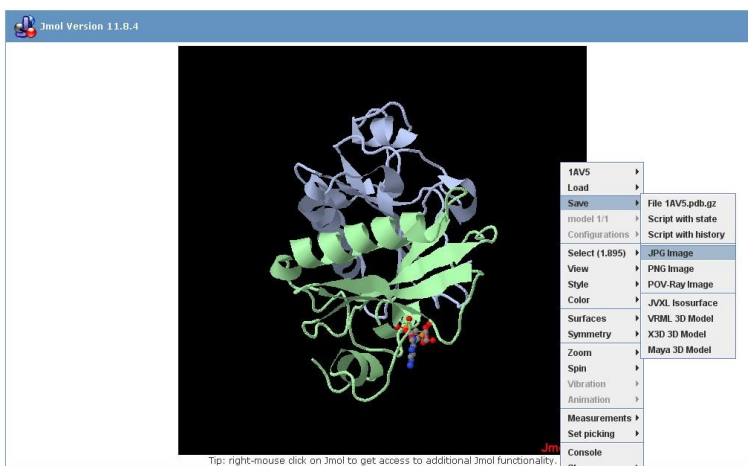


Figura 3.5.b. Exemplu de structură terțiară vizualizată cu jmol

3.6. Utilizarea Cn3D

Acest modul de program îl puteți accesa de pe adresa <http://www.ncbi.nlm.nih.gov/sites/gquery> (figura 3.6.a).

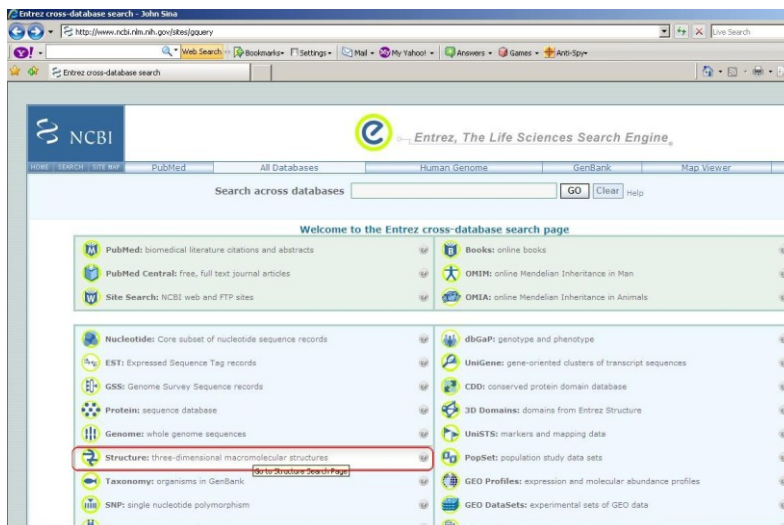


Figura 3.6.a. Opțiunea de vizualizare a structurii 3D moleculare

În căsuța “search” vom scrie *fibrinogen* și vom obține rezultatul căutării în figura 3.6.b:

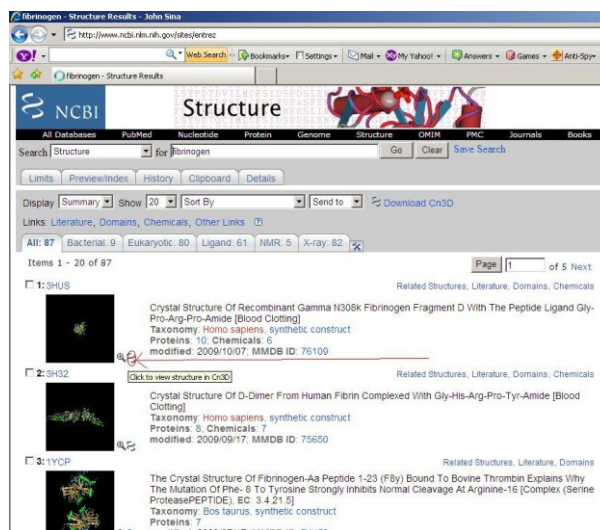


Figura 3.6.b. Rezultatele interogării

Prin accesarea butonului indicat de săgeată (figura 3.6.b) se va deschide programul Cn3D (figura 3.6.c).

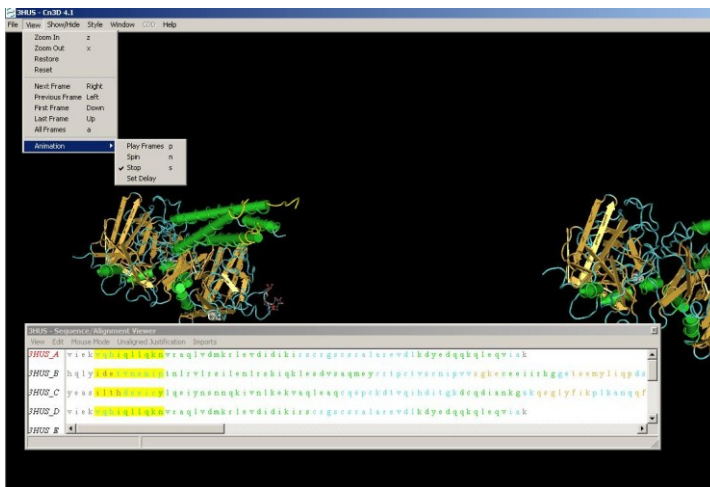


Figura 3.6.c. Vizualizarea celor 3 structuri ale fibrinogenului

În Cn3D vom avea afișată structura secundară și terțiară în ecranul principal, precum și structura primară (lanțul de aminoacizi) în ecranul “Sequence alignment viewer” (figura 3.6.c). În acest ecran aminoacizii sunt organizați pe domenii care corespund vederii 3D a proteinei din background. Tot în Cn3D avem la dispoziție meniul “View”, din care putem alege *zoom in* (pentru a vedea anumite domenii ascunse), precum și opțiunea *animantion* prin care putem crea o animație în care structura 3D a proteinei se va roti. De asemenea avem la dispoziție meniul “Styles” unde putem selecta opțiunea *Rendering shortcuts*. Aceasta va afecta modul de vizualizare al proteinei, fiind disponibile următoarele variante (figura 3.6.d): **worms** – vizualizarea cu spirale și panglici (cea mai sugestivă); **tubes**; **wire**; **ball and stick** - vizualizează atomii și legăturile dintre molecule.

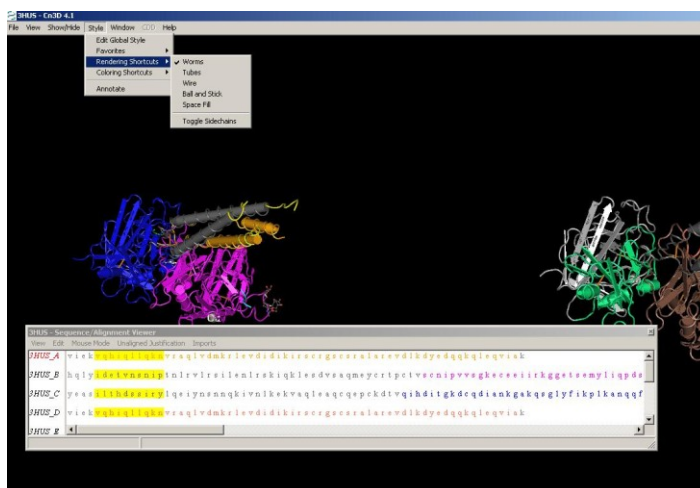


Figura 3.6.d. Moduri diferite de vizualizare ale structurilor terțiare și cuaternare

Tot din meniul „Styles” avem opțiunea *Coloring Shortcuts* pentru a selecta culoarea de reprezentare a structurii 3D:

- **Secondary structure** - va colora structura secundară
- **Domain** - va colora un domeniu (partea funcțională a proteinei)

- **Molecule** - va colora o anumită moleculă
- **Temperature** - va colora în funcție de temperatură
- **Charge** - va colora în funcție de încărcătura electrică
- **Hidrophobicity** - va colora în funcție de hidrofobicitate.

3.7. Exemplu de utilizare a programului RasMol

Denumirea **RasMol** provine de la *raster display of molecules*. **Raster** este un tip de afișaj computerizat folosit în special pentru a afișa suprafețe solide. RasMol deschide fișiere de tipul **.pdb** (protein database bank).

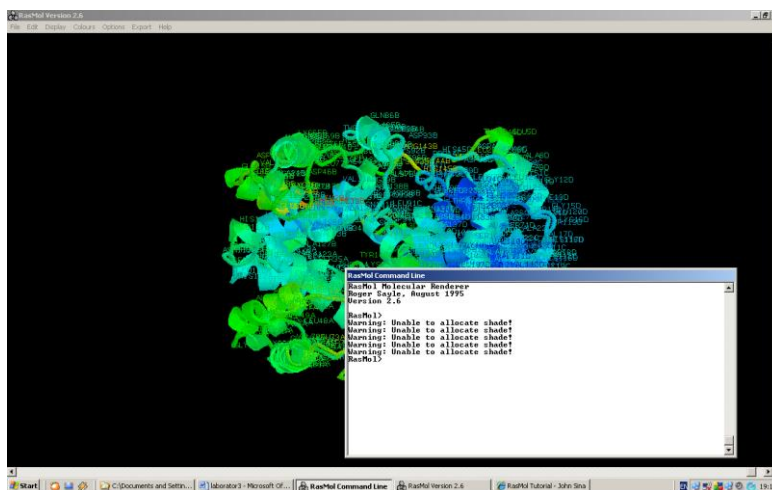


Figura 3.7.a. Vizualizarea proteinei cu programul RasMol

Fereastra “Linia de Comandă” apare în fața ecranului principal (figura 3.6.a). În continuare ne propunem să introducem câteva comenzi în liniile de comandă:

RasMol > **help**

Citiți introducerea din fereastra de *help*.

RasMol > **help commands**

Windows: Main Window

Ecranul principal este ecranul pe care va fi afișată structura proteinei. Se recomandă ca ecranul cu linia de comandă să fie așezat în josul ecranului principal. Cu acest aranjament făcut se pot vedea și ultimele linii de comandă și mesajele trimise de RasMol. În ecranul principal se prezintă o structură sub formă de cadru de sârmă a unui model de proteină (figura 3.7.a).

Display: Backbone

Așa se pot vedea doar atomii de carbon din poziția alpha. Liniile care conectează atomii de *Carbon alpha* sunt virtuale de obicei. Se pot vizualiza elemente de structură secundară ca elice alpha sau ca foițe beta pliate.

Încercați și alte comenzi din meniul “Display”. Fiecare comandă din acest meniu modifică felul în care RasMol reprezintă sau afișează modelul. La final afișați modelul *backbone*. Astfel structura secundară va fi și mai evidentă.

Colours: Structure

Această comandă colorează elicea alpha în roz, foițele beta în galben și alte părți de moleculă în gri deschis.

RasMol > select hetero and not hoh

Această comandă va selecta numai grupările hetero (non-proteice) din acest fișier, excluzând moleculele de apă care sunt frecvent incluse în modelele cristal. Nimic nu se întâmplă până nu este introdusă o altă comandă. **Atenție! Comanda afectează numai atomii selectați în mod curent.** De asemenea nu este nevoie de a aduce fereastra cu linia de comandă în primul plan pentru a introduce comenzi.

Display: Ball & Stick

RasMol afișează un model de “bețe” și “mingi” ale atomilor selectați. Observați că această comandă nu afectează lanțul proteic pentru că nu e selectat.

Colours: CPK

RasMol colorează atomii selectați conform unor convenții chimice larg răspândite: Carbonul (C) este gri deschis, Azotul (N) se colorează în albastru deschis, iar oxigenul se colorează în roșu.

Dați click pe Carbonul din poziția alpha a proteinei și urmăriți linia de comandă. Această manevră se numește „extragerea unui atom”. Când veți alege un atom din fereastra principală RasMol îl va identifica (figura 3.5.b). O să observați o expresie de genul: *Atom: CA 652 Group: PRO 81* care ne spune că am selectat carbonul din poziția alpha (CA) din restul 81 în citocromul b5, care este o prolină. RasMol definește restul “ca un “grup” (de aminoacizi), iar porțiunea non-proteică dintr-o moleculă ca “hetero”. Dacă vedeți un atom colorat în galben după modelul “ball and stick” dați click pe el. Va fi afișat ceva de genul: *Atom: FE 754 Hetero: HEM 201*. Atomul galben numit FE este Fierul (Fe³⁺, de fapt) în centrul grupului hem al citocromului b5. Folosind comanda “select hetero and not hoh” și afișând totul într-un stil contrast (similar cu displayul Display:Ball&Stick) putem găsi foarte rapid grupările nonproteice din fișierul nostru PDB. Selectând un atom dintr-un grup putem afla numele în PDB al celui grup, fapt care ne va folosi mai târziu.

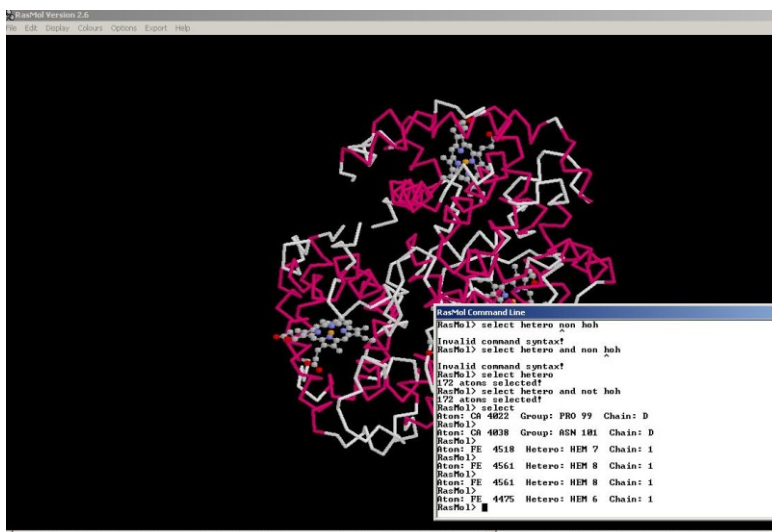


Figura 3.7.b. Tipuri de vizualizări ale atomilor și fereastra de comandă

Atenție! Selectarea unui atom funcționează cel mai bine cu modul de vizualizare **Wireframe** și **Sticks** și nu funcționează deloc cu modul **Ribbons**, **Strands** și **Cartoons** (fiind inexact prin vizualizarea cu **Spacefill** și **Ball&Stick**). Dedesubt este afișată lista de comenzi utilizată până acum (figura 3.5.b). Puteți folosi această listă de comenzi pentru a va crea cât mai repede o idee asupra oricărei macromolecule pentru care aveți fișierul PDB cât și pentru a afișa eventualele grupări hetero (cofactori, inhibitori, apa):

Display: Backbone

Colours: Structure

Rasmol > **select hetero and not hoh**

Display: Ball & Stick

Colours: CPK.

3.8. *Exerciții propuse*

1. Descărcați RasMol de la adresa <http://www.umass.edu/microbio/rasmol/>
2. Din ce clasă face parte proteina “Protein Kinase C Interacting Protein with PDB-ID 1AV5”?

R: Face parte din familia inhibitorilor de protein kinaza.

3. Căutați și vizualizați hexokinaza umană I (Swissprot cod HXK1_HUMAN, număr P19367, PDB-id 1HKB).

4. Vizualizați domeniul de legare cu beta-D –glucoza cu ajutorul programului Ligand explorer.

R:Dupa ce ați căutat proteina 1HKB la secțiunea Ligand Chemical Component=> beta-D- glucose accesați Ligand explorer.

5. Folosiți GOR IV pentru a prezice structura secundară a proteinei cu structura primară următoare (lanțul alpha al hemoglobinei):

vlspadktnv kaawgkvqah ageygaeale rmflsfpttk tyfphfdlsh gsaqvkgghk kvadaltnav
ahvddmpnal salsdlhahk lrvdpvnfkl lshellvtla ahlpaeftpa vhasldkfla svstvltsky r

(http://npsa-pbil.ibcp.fr/cgi/in/npsa_automat.pl?page=npsa_gor4.html)

6. Comparați rezultatele cu structura vizualizată în Cn3D a proteinei de mai sus (ID 2HCO).

Exerciții cu RasMol

Mergeți la link-ul PDB și descărcați:

1ZNI (insulina porcina)

1LPH (insulina umana creata artificial)

Dați dublu click pe icoana *Raswin* de pe desktop. Va apărea fereastra *RasMol* și fereastra cu linia de comandă *RasMol* minimalizată. Accesați fereastra cu linia de comandă *RasMol* pentru a afișa linia de comandă. Scrieți **background white**, apăsați “enter”. Se va schimba fundalul ecranului principal.

Duceți fișierul PDB 1ZNI în fereastra principală *RasMol*. Acum puteți vedea structura insulinei afișată cu displayul wireframe în ecranul principal.

Câți aminoacizi are proteina (insulina)?

Câte lanțuri polipeptidice are?

Care este aminoacidul din capătul N-terminal și cel din capătul C-terminal al lanțurilor A și B?

Care este structura secundară?

Care este cel mai lung Helix?

R: Câte lanțuri polipeptidice are insulina?

Introduceți următoarele comenzi în linia de comandă:

```
reset  
restrict backbone  
ribbons on  
wireframe off  
color chain
```

Puteți salva un script pe hard-disk introducând comanda: *write script C:\zni.sp*

R: Câți aminoacizi are lanțul A?

```
select *.ca and *a (selectează toți atomii CA din lanțul A)
```

(Răspuns: “21 atoms selected!”), deci există 21 de aminoacizi în lanțul A)

R. Care este capătul N-terminal și C-terminal?

```
select 1,21 and *.ca and *a label
```

R: Care este structura secundară? Care este cel mai lung Helix ?

Există pantru lanțuri polipeptidice în structura insulinei, fiecare afișat cu o culoare diferită. Un fișier PDB folosește o literă unică pentru fiecare lanț polipeptidic, începând cu litera “A”.

De exemplu introduceți comanda: *restrict *a*

Se va afișa numai lanțul „A” în ecranul principal RasMol.

R: Cum puteți folosi RasMol pentru a afla capătul N-terminal și C-terminal al lanțului A?

Accesați ambele capete ale lanțului și veți vedea afișat în linia de comandă RasMol, aminoacizii corespunzători.

4. Aplicația Vector NTI – programul Align X

4.1. Obiectivele lucrării de laborator

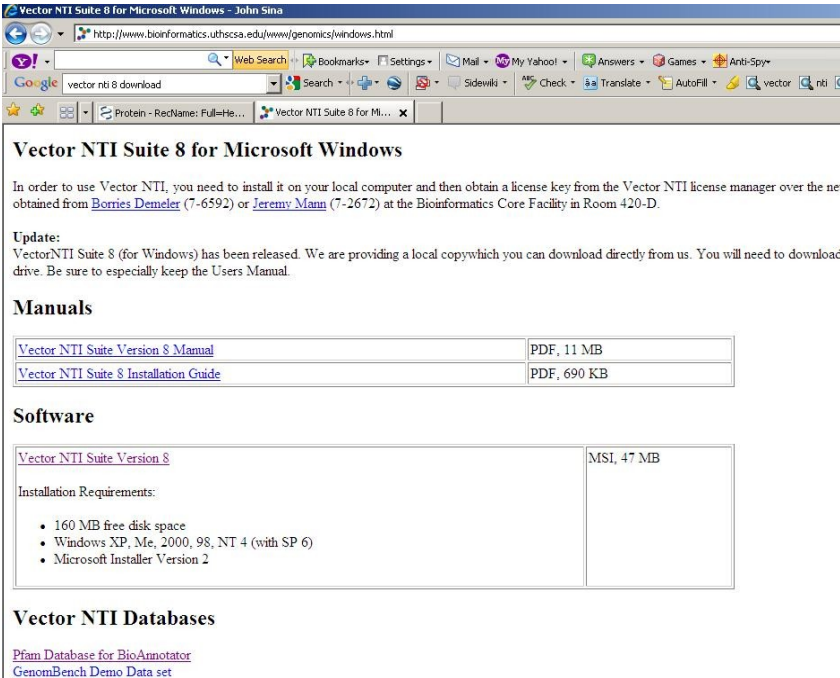
- compararea secvențelor ADN și a secvențelor proteice
- descrierea generală a unei proteine, determinarea numărului de aminoacizi din care e formată, funcția ei, greutatea moleculară și încărcarea electrică
- alinierea proteinelor
- vizualizarea arborelui filogenetic
- afișarea aliniamentului folosind matricele de scor.

Align X este un program, parte componentă din Vector NTI care poate fi folosit pentru a compara secvențe de ADN sau de proteine.

Pentru a instala programul Vector NTI suite 8 vom accesa următoarea adresă:

<http://www.bioinformatics.uthscsa.edu/www/genomics/windows.html>

De asemenea se pot descărca și tutorialele. După ce lansăm în execuție programul VectorNTI, dăm click pe butonul “Local Database” (figura 4.1.a). Vom parcurge această lucrare prin intermediul a 3 exemple.



Vector NTI Suite 8 for Microsoft Windows

In order to use Vector NTI, you need to install it on your local computer and then obtain a license key from the Vector NTI license manager over the netw obtained from [Borries Demeler](#) (7-6592) or [Jeremy Mann](#) (7-2672) at the Bioinformatics Core Facility in Room 420-D.

Update:
VectorNTI Suite 8 (for Windows) has been released. We are providing a local copy which you can download directly from us. You will need to download a drive. Be sure to especially keep the Users Manual.

Manuals

Vector NTI Suite Version 8 Manual	PDF, 11 MB
Vector NTI Suite 8 Installation Guide	PDF, 690 KB

Software

Vector NTI Suite Version 8 Installation Requirements: <ul style="list-style-type: none">• 160 MB free disk space• Windows XP, Me, 2000, 98, NT 4 (with SP 6)• Microsoft Installer Version 2	MSI, 47 MB
--	------------

Vector NTI Databases

[Pfam Database for BioAnnotator](#)
[GenomBench Demo Data set](#)

Figura 4.1.a. Modul de descărcare al programului Vector NTI

4.2. Exemplul 1 – Ne propunem să creăm o secvență de nucleotide pentru care să vizualizăm matricea DotMatrix

Dacă dorim să vizualizăm componentele bazei de date trebuie să alegem opțiunea “Exploring” din meniul principal al programului (figura 4.2.a).

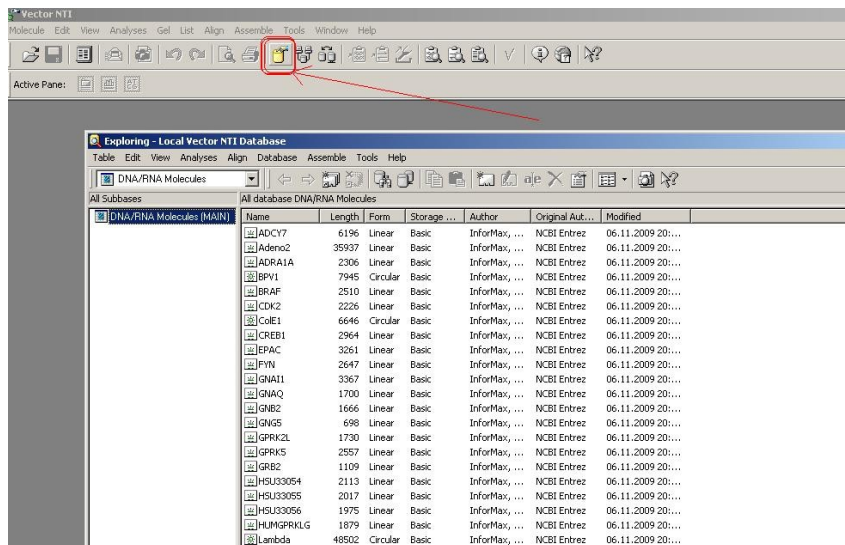


Figura 4.2.a. Afășarea fișierelor cu molecule ADN/ARN din BD a programului

Selectăm DNA/RNA molecules, facem click dreapta și selectăm “New Item” pentru a crea o nouă moleculă, căreia îi dăm numele “test”.

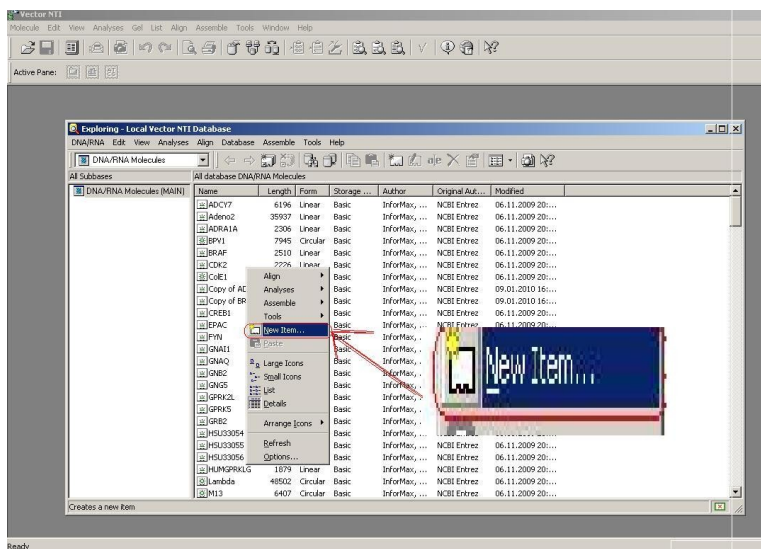


Figura 4.2.b. Exemplu de creare a unei noi molecule

Se va deschide o fereastră în care vom putea crea o moleculă test pentru care vom afișa alinierea vizuală a nucleotidelor din exemplul de la curs (Dot Matrix) - vom numi noua moleculă creată “Test” (figura 4.2.c).

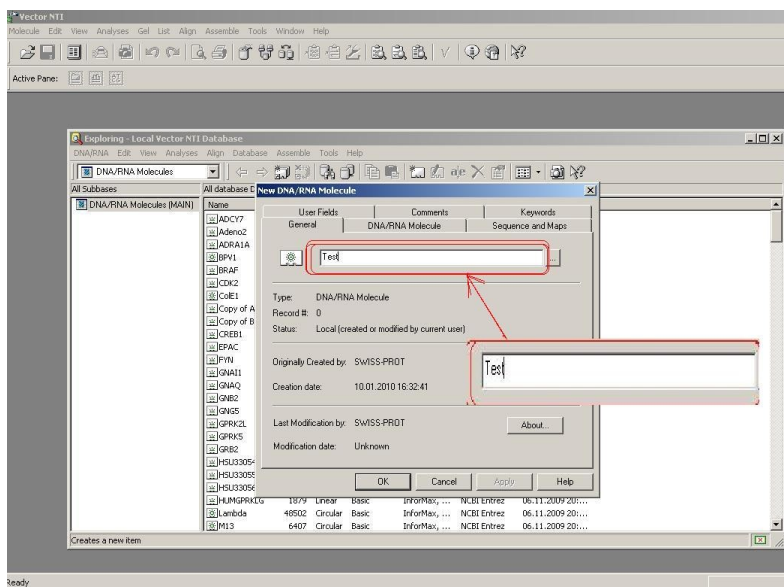


Figura 4.2.c. Denumirea noii molecule

Pentru a completa cu informație noua moleculă dați click pe tab-ul *Sequence and Maps*, iar după aceea click pe butonul *Edit Sequence* și introduceți următoarea secvență de nucleotide **ACCTGAGCTCACCTGAGTTA** (figura 4.2.d).

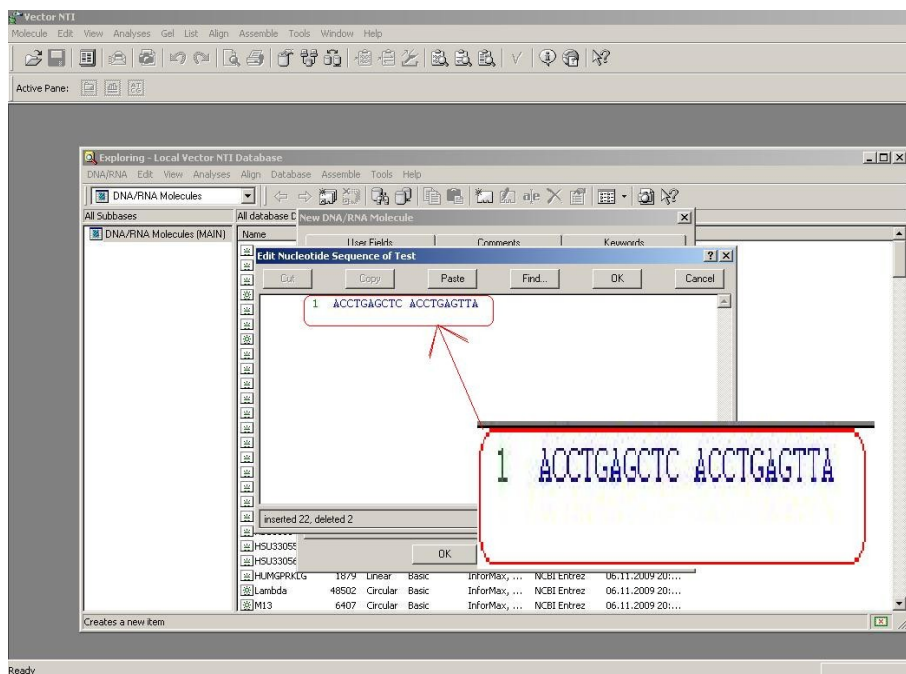


Figura 4.2.d. Introducerea secvenței de nucleotide în fereastra de editare a programului

Apoi apăsați tasta OK, după care faceți click dreapta pe molecula **Test** apoi alegeți opțiunea **Align X** și **Align Selected Molecules**.

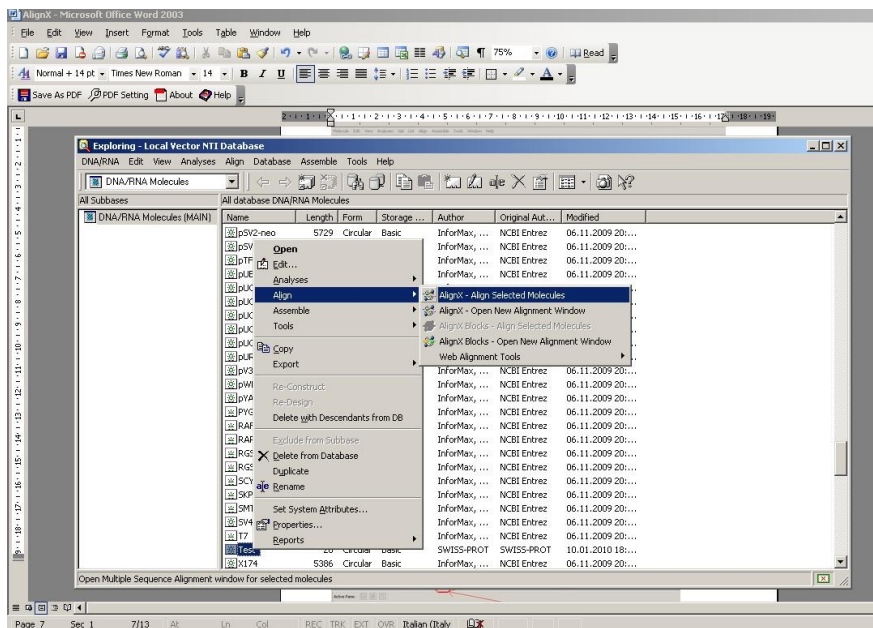


Figura 4.2.e. Opțiunile de aliniere a secvențelor

Se va deschide fereastra **Align X** iar pentru a afișa alinierea vizuala (Dot Matrix) pentru secvența “test” vom alege această opțiune (figura 4.2.e).

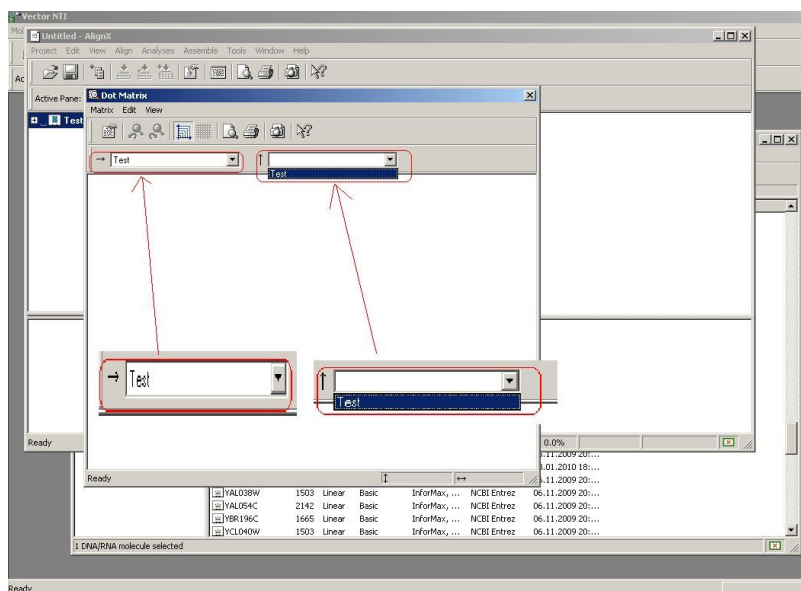


Figura 4.2.e. Alinierea secvenței “test” cu ea însăși

Se va afișa fereastra cu Dot Matrix. Observăm că pe orizontală vor fi afișate nucleotidele moleculei **Test**. Vom alege astfel încât și pe verticală să fie afișate nucleotidele moleculei **Test**.

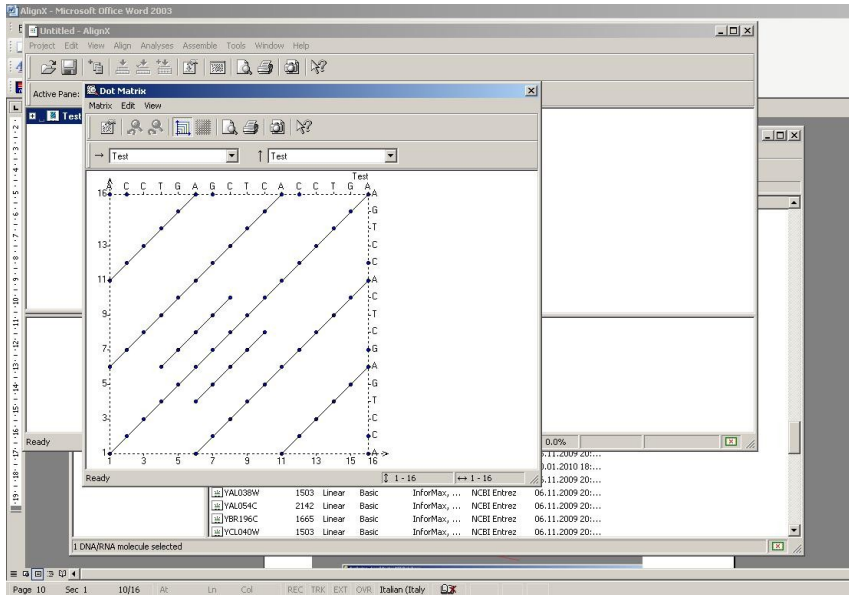


Figura 4.2.f. Conținutul matricei Dot Matrix

Observăm că atât pe orizontală cât și pe verticală se află doar 16 nucleotide din cele 20 introduse. Pentru aceasta vom alege din meniul *Matrix =>Matrix setup => window* iar în caseta *window* vom introduce valoarea 1.

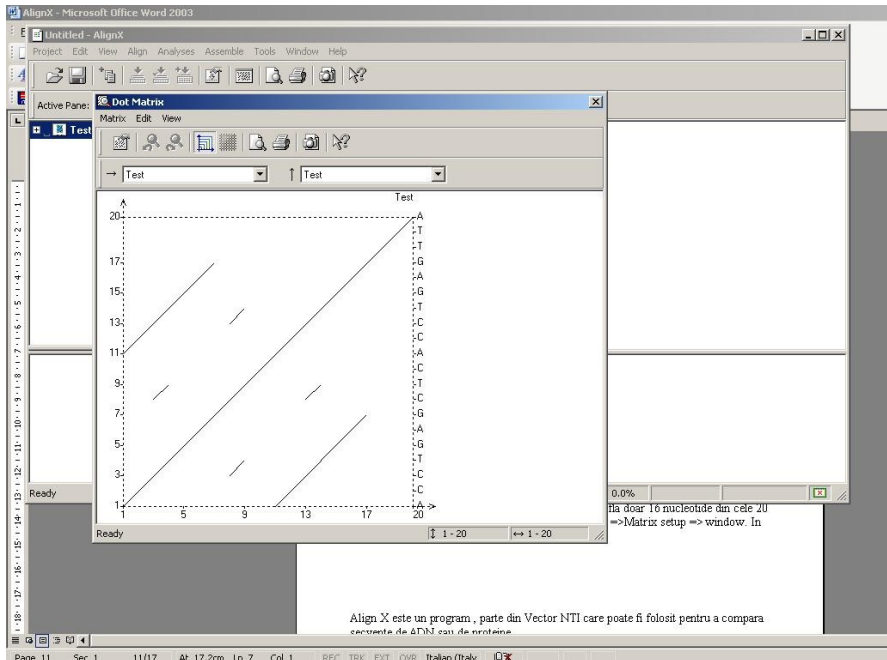


Figura 4.2.g. Reprezentare finală a Dot Matrix

4.3. Exemplu 2 – Compararea a două proteine cu programul Vector Align X

Din secțiunea *Database explorer* selectăm protein.

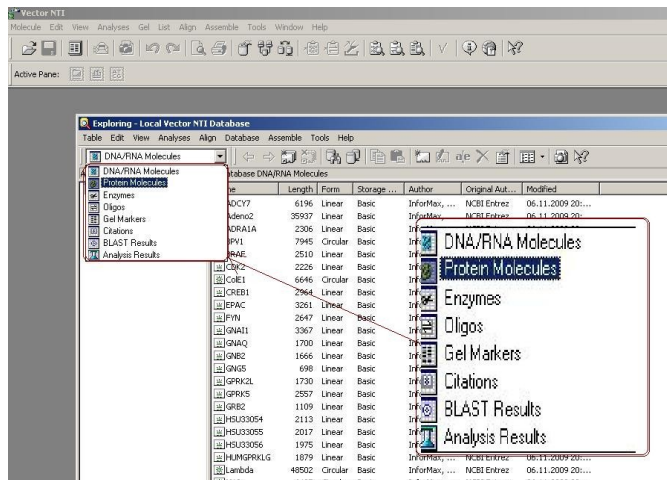


Figura 4.3.a. Selecția proteinelor

Apoi selectăm 2 molecule: **41BB_HUMAN** și **4F2_HUMAN** (figura 4.3.b).

Pentru alinierea lor facem click dreapta și apoi tastăm *Align => Align selected molecules*.

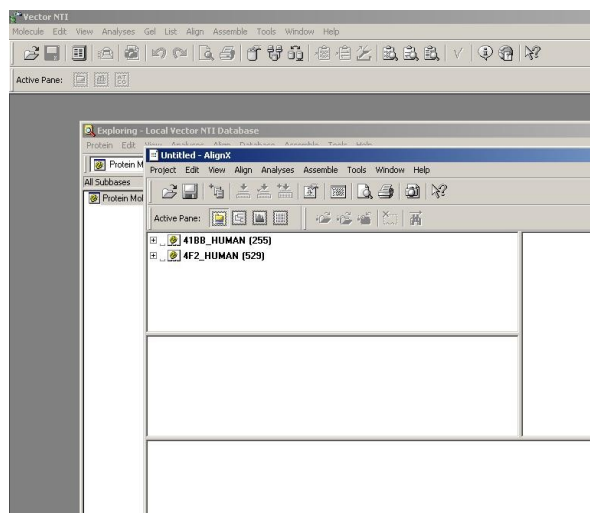


Figura 4.3.b. Cele două molecule alese spre comparare, în fereastra AlignX

Se va deschide subprogramul alignX. Dacă vom da click pe „+” vom putea obține următoarele informații despre proteine: descrierea generală a proteinei, numărul de aminoacizi din care e formată, funcția, greutatea moleculară și încărcarea electrică.

Vom selecta apoi DotMatrix, unde vom vedea secvențele de proteine comparate (figura 4.3.c).

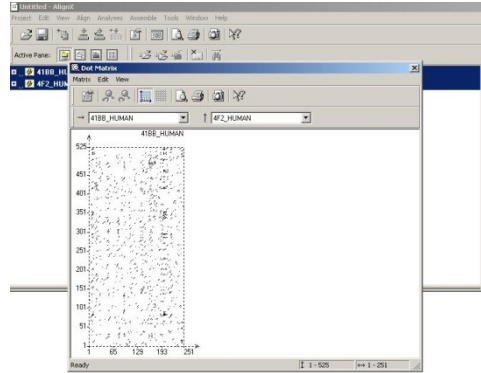


Figura 4.3.c. DotMatrix pentru alinierea celor două secvențe

4.4. Exemplul 3

La fel ca în exemplul 2, după ce deschidem programul Vector NTI, din secțiunea *Database explorer* alegem opțiunea “protein” după care vom selecta următoarele proteine: **41BB_HUMAN**, **5H1A_HUMAN**, **5H1A_MOUSE** și **5H1A_RAT**.

Executăm click dreapta => *Align X* => *Align selected molecules* pentru a alinia secvențele proteinelor de tip 5H1A. Va apărea următorul ecran care conține următoarele panouri ca în figura 4.4.a:

- Panoul cu text (Text Pane)
- Panoul care conține arborele filogenetic (phylogenetic Pane)
- Panoul care conține aliniamentul (alignment Pane)
- Panoul care conține analiza (Analysis Pane) sau panoul grafic.

Selectăm moleculele și apăsăm butonul *Align*.

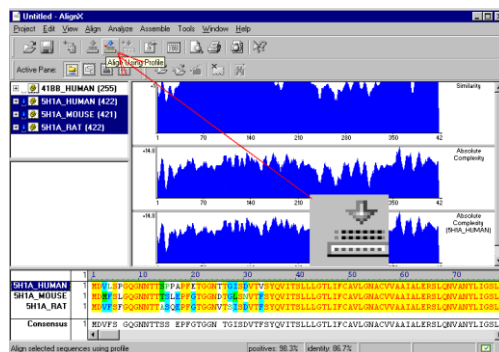


Figura 4.4.a. Panouri cu vizualizarea secvențelor aliniate

În panoul “Text” dați dublu click pe o moleculă și vizualizați datele despre aceasta. Aceste informații conțin tipul moleculei, forma moleculei, comentarii, referințe și date în format GenBank.

Panoul implicit de analiză conține trei reprezentări grafice ale rezultatelor aliniamentului.

- Primul grafic conține profilul calității aliniamentului. Se analizează dacă secvența este identică, similară sau slab similară cu secvența de comparat.

- Al doilea grafic conține semnificația statistică a profilului (complexitatea absolută a unui aliniament). Este suma tuturor scorurilor calculate conform matricei de substituție.
- Al treilea grafic prezintă semnificația statistică a aliniamentului pentru o moleculă selectată în raport cu molecula de comparat (figura 4.4.a).

Putem adăuga analize adiționale prin tastarea butonului “*Graphics Pane*”. Apoi vom apăsa pe “*View*” => *list of analyses*.

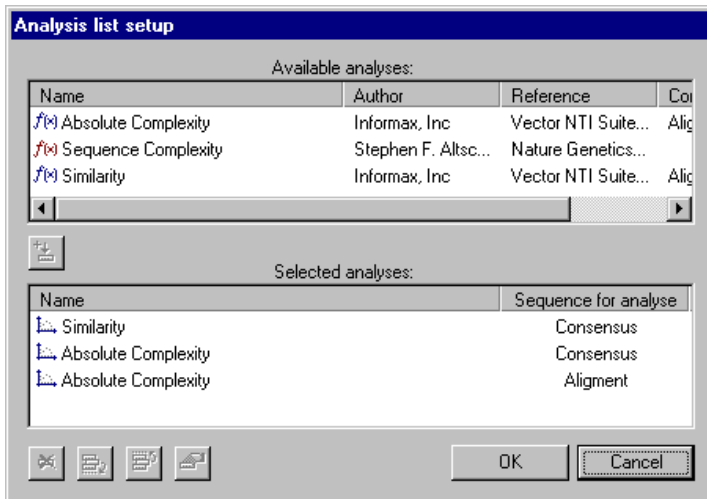


Figura 4.4.b. Lista analizelor adiționale

Sunt afișate toate analizele disponibile pentru ADN sau proteine, depinzând de tipul moleculelor de aliniat.

Analiza filogenetică este o metodă de studiere a presupuselor relații evolutive ale proteinelor. Panoul care conține arborele filogenetic îl va afișa doar dacă sunt comparate cel puțin 4 molecule (proteine) – figura 4.4.c.

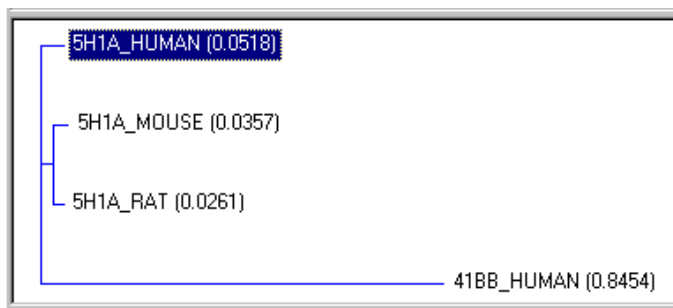


Figura 4.4.c. Arborele filogenetic pentru cele 4 molecule selectate inițial

Arborele filogenetic este creat folosind metoda *Neighbor Joining* (NJ) of Saitou and Nei.

Setarea parametrilor aliniamentului se obține prin intermediul butonului “*alignment setup*”. Aici vom putea seta parametrii pentru *Pairwise Alignment*, *Multiple sequence alignment*, precum și *Score Matrix* (matricea de scor) – figura 4.4.d.

The screenshot shows the 'Alignment Setup' dialog box with the 'Score Matrix' tab selected. The matrix name is 'blosum62mt2' and the order is 'ABCDEFGHIKLMNPQRSTVWXYZ'. The matrix is a 20x20 grid of scores for amino acid comparisons.

	A	B	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	X	Y	Z
A	8	-4	0	-4	-2	-4	0	-4	-2	-2	-2	-4	-2	-2	-2	2	0	0	-6	0	-4	-2	
B	-4	8	-6	8	2	-6	-2	0	-6	0	-8	-6	-6	-4	0	-2	0	-2	-6	-8	-2	-6	2
C	0	-6	18	-6	-8	-4	-6	-6	-2	-6	-2	-2	-6	-6	-6	-6	-2	-2	-2	-4	-4	-4	-6
D	-4	8	-6	12	4	-6	-2	-2	-6	-2	-8	-6	-2	-2	0	-4	0	-2	-6	-8	-2	-6	2
E	-2	2	-8	4	10	-6	-4	0	-6	2	-6	-4	0	-2	4	0	0	-2	-4	-6	-2	-4	8
F	-4	-6	-4	-6	-6	12	-6	-2	0	-6	0	0	-6	-8	-6	-6	-4	-4	-2	-2	-2	-6	-6
G	0	-2	-6	-2	-4	-6	12	-4	-8	-4	-8	-6	0	-4	-4	-4	0	-4	-6	-4	-2	-6	-4
H	-4	0	-6	-2	0	-2	-4	16	-6	-2	-6	-4	2	-4	0	0	-2	-4	-6	-4	-2	-4	0

Figura 4.4.d. Afișarea elementelor matricei de scoruri BLOSUM62

Din tab-ul *score matrix* vom apăsa pe butonul *select matrix* și vom selecta matricea de scor pe care vrem să o folosim. Matricea de scor se prezintă ca un fișier text (figura 4.4.d).

Putem selecta matricele de tipul BLOSUM, PAM, DAYHOFF sau cu ajutorul utilitarului *matrix editor*, putem crea propria noastră matrice de scor.

Editarea aliniamentului se realizează prin butonul *Edit alignment*. Aici vom putea selecta o secvență de aminoacizi și o vom putea muta în cadrul aliniamentului (figura 4.4.e).

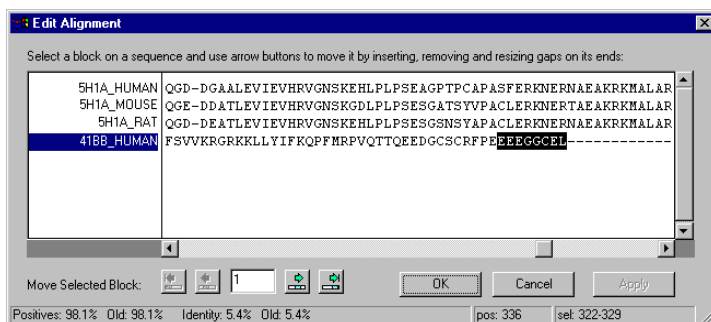


Figura 4.4.e. Exemplu de editare în cadrul aliniamentului secvențelor selectate

4.5. Exerciții propuse

1. Să se selecteze cele 3 proteine de tipul 5H1D și proteina 5H1E, să se alinieze aceste proteine, după care să se vizualizeze arborele filogenetic. Afișați aliniamentul folosind matricele de scor PAM 120, PAM 50 și BLOSSUM 55.
2. Selectați proteinele de tipul _ECOLI (esicheria coli), aliniați-le și vizualizați arborele filogenetic. Să se afișeze aliniamentul folosind matricele de scor PAM10, PAM20, PAM40, BLOSUM 100, BLOSUM 85.
3. Să se selecteze proteinele de tipul YLR și să se alinieze aceste proteine. Vizualizați arborele filogenetic și afișați aliniamentul folosind matricele de scor PAM 120, PAM 50 și BLOSSUM 55.
4. Afișați DOT MATRIX pentru proteinele din exercițiile de mai sus (2 câte 2).

5. Cu ajutorul algoritmului *Needleman-Wunsch* din fișierul excel calculați matricea de scoruri asociată alinierii celor două secvențe de nucleotide: GAATTCAGTTA și GGATCGA.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	
1	Potriviri:	5				Key:	Red: move diagonally up and left, sequences are aligned												
2	Nepotriviri:	-3					Green: move up, gap in top (query) sequence												
3	Gap (InDel):	-4					Blue: move left, gap in left-side (subject) sequence												
4							White: move up or left or diagonally, alternate optimal alignments												
5	Algoritmul Needleman-Wunsch																		
6			G	A	A	T	T	C	A	G	T	T	A						
7	0	-4	-8	-12	-16	-20	-24	-28	-32	-36	-40	-44	-48	-52	-56	-60			
8	G	-4	5R	-1G	-2G	-7G	-11G	-15G	-19G	-23	27G	-11G	-35G	-39G	-43G	-47G	-51G		
9	G	-8	1	2R	-2	-6	-10	-14	-18	-22G	-26G	-30G	-34G	-38G	-42G	-46G			
10	A	-12	-1	6R	2R	4R	-1G	-5G	-9	-13G	-17	-21	-25G	-29G	-33G				
11	T	-16	-7	-1	3	12R	8	5G	0G	-4G	-8	-12	-16G	-20G	-24G	-28G			
12	C	-20	-11	-2	-1	9	11R	9G	5G	1G	-3G	-7G	-11G	-15G	-19G	-23G			
13	G	-24	-15	-6	-5	6	5	10R	14R	10G	6G	2G	-2G	-6G	-10G	-14G			
14	A	-28	-19	-10	-11	3	1	7	14R	10	11R	7	11R	7G	3G	-1G	-5G		
15		-32	-23	-14	-15	-4	-3	1	10	11G	7	8R	4	10R	12	8	4		
16		-36	-27	-18	-19	-8	-7	-1	5	7	8R	4	8R	12	15R	17	13		
17		-40	-31	-22	-23	-12	-11	-1	3	4	5R	1	10R	17	20R	22			
18		-44	-35	-26	-27	-16	-15	-1	-1	0	1	2R	6	15R	22	11R			
19		-48	-39	-30	-31	-20	-19	-1	-5	-4	-3	-2	3R	11	20R	27			
20		-52	-43	-34	-35	-24	-23	-1	-9	-8	-7	-6	3	12R	16	23R			
21		-56	-47	-38	-39	-28	-27	-1	-13	-12	-11	-10	-1	8	17R	21			
22		-60	-51	-42	-43	-32	-31	-1	-17	-16	-15	-14	-5	4	13	22R			

R:

6. Aceeași aliniere realizați-o cu algoritmul *Smith-Waterman*.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	
1	Potriviri:	5				Key:	Red: move diagonally up and left, sequences are aligned												
2	Nepotriviri:	-3					Green: move up, gap in top (query) sequence												
3	Gap (InDel):	-4					Blue: move left, gap in left-side (subject) sequence												
4							White: move up or left or diagonally, alternate optimal alignments												
5	Algoritmul Smith-Waterman																		
6			G	A	A	T	T	C	A	G	T	T	A						
7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
8	G	0	5R	1G	0	0R	0R	0R	0R	5R	1G	0	0R	0R	0R	0R	0R	0R	0R
9	G	0	5R	2R	0	0R	0R	0R	0R	5R	2R	0	0R	0R	0R	0R	0R	0R	0R
10	A	0	-1	10R	7R	5G	0G	0R	3R	1	7R	0R	5R	1G	0	0R	0R	0R	0R
11	T	0	0	-1	7R	12R	8	-4G	-1	2R	0R	7R	9G	2R	0	0R	0R	0R	0R
12	C	0	0R	-1	3	-1	5R	13R	9G	5G	-1	3	0R	0	0R	0R	0R	0R	0R
13	G	0	5R	1G	0	1	5	9	10R	14R	10G	6G	2G	1R	0	0R	0R	0R	0R
14	A	0	-1	10R	6	2G	1	-1	10R	10	11R	7	11R	7G	3G	0G	0R	0R	0R
15		0	0	5	7	3	0	1	10	11R	7	0R	7	10R	12	8	1R		
16		0	0R	0	3	4R	0	0	7	7	0R	4	5R	12	21R	17	13		
17		0	0R	0	0	0	0	0	3	4	5R	1	10R	17	26R	22			
18		0	0R	0R	0R	0R	0	0R	0	0	1	2R	6	15R	22	15R			
19		0	0R	0R	0R	0R	0R	0R	0R	0R	0	0	7R	11	20R	27			
20		0	0R	0R	0R	0R	0R	0R	0R	0R	0R	0R	0R	5R	12R	16	25R		
21		0	0R	0R	0R	0R	0R	0R	0R	0R	0R	0R	0R	5R	10R	17R	21		
22		0	0R	0R	0R	0R	0R	0R	0R	0R	0R	0R	0R	5R	10R	15R	22R		

R:

5. Clustal X

5.1. Obiectivele lucrării de laborator

- alinierea multiplă a secvențelor de proteine cu programul **ClustalX**
- vizualizarea arborelui filogenetic cu același program
- setarea parametrilor aliniamentului
- scriere aliniamentului ca fișier Postscript.

5.2. Introducere

Programul CLUSTALX oferă o interfață grafică de utilizator pentru algoritmul CLUSTALW de aliniere multiplă a secvențelor (MSA – Multiple Sequences Alignment).

Publicațiile originale care descriu CLUSTALX și CLUSTALW sunt unele dintre cele mai citate documente - atunci când au fost publicate au prevăzut o creștere majoră în acuratețea și rapiditatea cu care utilizatorii ar putea construi MSAs, cu accent pe secvențele de proteine - în cazul CLUSTALX a fost adăugată în plus o interfață ușor de folosit (GUI) care să conducă la o mare popularitate a programelor.

Există în prezent o serie de alte produse software MSA care sunt, datorita multipleror seturi de secvențe și analize, o alegere mai bună decât CLUSTALX dacă țelul vostru este de a obține o aliniere (MSA) a unui set de secvențe de proteine - și, desigur, dacă sunteți interesați de alinierea secvențelor ADN (ceva ce CLUSTALW/X nu au avut ca obiectiv principal). Cu toate acestea CLUSTALX rămâne util deoarece în unele circumstanțe face o bună aliniere a secvențelor (judecate după standardele de azi) și poate cel mai important din cauza mai multor caracteristici ale GUI care oferă interfețe ușor de utilizat pentru mai multe sarcini cheie care sunt necesare, în multe cazuri, pentru a construi o bună aliniere multiplă a secvențelor (MSA).

5.3. Utilizarea Clustal X

ClustalX poate fi descărcat de la adresa: <http://www.clustal.org/download/current/>.
Observație! Dacă obțineți eroare de tipul "lipsă fișier mingw10.dll", la rularea programului ClustalX, o puteți corecta prin descărcarea acestui fișier de la adresa: <http://www.dll-files.com/dllindex/dll-files.shtml?mingwm10> după care îl salvați în locația C:\Windows\System.

O vedere generală asupra programului ClustalX o oferă imaginea din figura 5.3.a.

Bara de meniu-oferă acces la o gamă largă de caracteristici și acțiuni care pot fi aplicate la aliniere: de exemplu, salvarea unei alinieri, ajustarea sistemului de colorat, etc.

Căsuța de alegere a modului - este folosit pentru a comuta între mai multe tipuri de alinieri: alinierea multiplă și alinierea după profil (așa cum este arătat mai sus).

Rețineți că unele opțiuni de meniu sunt disponibile numai în alinierea multiplă sau numai în alinierea după profil.



Figura 5.3.a. Interfața cu utilizatorul a programului ClustalX

Căsuța cu mărimea fontului este folosită pentru a specifica dimensiunea fontului utilizat pentru a reprezenta secvențele în **zona de afișare a aliniamentului**.

Numele de identificare al secvențelor folosite în aliniament sunt afișate în **Panoul cu numele secvențelor**.

Linia de consens oferă un rezumat al gradului de conservare a reziduurilor din coloana de aliniere corespunzătoare:

- „*” (reziduuri identice în toate secvențele)
- „.” (conservate bine în coloană)
- „-” (conservate slab în coloană).

Liniarul de aliniere indică poziția în aliniere (prima coloană are atribuită poziția 1, etc.) a unei coloane date.

Zona de conservare a alinierii oferă o diagramă bară care indică gradul de conservare al fiecărei coloane în aliniament.

5.4. Crearea fișierului de intrare pentru alinierea multiplă

Acest program, ca oricare program de calculator are nevoie ca datele pe care le utilizează (fișier de intrare) să fie într-un format pe care să îl recunoască. Pentru a crea fișierul de intrare puteți folosi un editor de texte, de preferabil Notepad. Clustal citește fișiere într-unul din cele 7 formate de fișiere de intrare, înlocuind orice secvență care este deja încărcată. Toate secvențele trebuie să se găsească într-un singur fișier, secvență după secvență. Formatele de fișier care sunt recunoscute automat sunt:

NBRF/PIR, EMBL/SWISSPROT, Pearson (Fasta), Clustal (*.aln), GCG/MSF (Pileup), GCG9/RSF și GDE. Toate caracterele non-alfabetice (spații, cifre, semne de punctuație) sunt ignorate, exceptând „-” care e folosită pentru a indica un GAP („-” în MSF/RSF).

Programul încearcă să recunoască automat diferitele tipuri de fișiere folosite și să stabilească dacă secvențele sunt de aminoacizi sau de nucleotide.

FASTA și NBRF/PIR sunt recunoscute pentru că au caracterul „>” ca primul caracter din fișier.

EMBL/Swiss Prot sunt recunoscute după literele “ID” de la începutul fișierului.

CLUSTAL este recunoscut după cuvântul “CLUSTAL” de la începutul fișierului.

GCG/MSF este recunoscut după unul din următoarele caractere:

- cuvântul “PileUp” la începutul fișierului
- cuvântul “!!AA_MULTIPLE_ALIGNMENT” sau “!!NA_MULTIPLE_ALIGNMENT” la începutul fișierului.
- cuvântul “MSF” la începutul primei linii din fișier și caracterele “..” la sfârșitul acestei linii.
- cuvântul “!RICH SEQUENCE” la începutul fișierului.

Dacă 85% sau mai mult dintre caracterele secvenței sunt din A,C,G,T,U sau N atunci secvența va fi tratată ca o secvență de nucleotide.

5.4.1. Exemplu

Copiați secvența următoare în fișierul text:

```
gi|15146064|gb|AY040893.1| Homo sapiens individual VP37 mitochondrial control region
```

```
GGTCTATCACCTATTAACCACTCACGGGAGCTCTCCATGCATTTGGTATTT  
TCGTCTGGGGGGTGTGCA
```

```
CGCGATAGCATTGCGAGACGCTGGAGCCGGAGCACCTATGTGCGAGTATC  
TGTCTTTGATTCTGCCTC
```

```
ATCCTGTTATTTATCGCACCTACGTTCAATATTACAGGCGAACATACTACT  
AAAGTGTGTTAATTAATT
```

```
AATGCTTGTAGGACATAATAATAACAATTG
```

```
gi|15146065|gb|AY040894.1| Homo sapiens individual VP5 mitochondrial control region
```

```
GGTCTATCACCTATTAACCACTCACGGGAGCTCTCCATGCATTTGGTATTT  
TCGTCTGGGGGGTATGCA
```

```
CGCGATAGCATTGCGAGACGCTGGAGCCGGAGCACCTATGTGCGAGTATC  
TGTCTTTGATTCTGCCTC
```

```
ATCCTATTATTTATCGCACCTACGTTCAATATTACAGGCGAACATACTACT  
AAAGTGTGTTAATTAATT
```

ClustalX poate recunoaște mai multe formate de secvențe, dar noi vom utiliza în exemplul următor formatul FASTA. Formatul FASTA poate fi recunoscut ușor deoarece prima linie începe cu caracterul „>”. Această linie conține titlul secvenței. Secvența va începe de la linia următoare. Caracterul „>” va fi urmat de un singur cuvânt, pe care ClustalX îl va folosi ca nume al secvenței în alinierea multiplă pe care o creează.

Numele fișierului este bine ca să aibă salvat ID-ul de tip Genbank pentru a putea fi ușor de recunoscut ulterior.

5.4.2. Alinierea multiplă

Programarea dinamică poate fi folosită pentru a alinia mai multe secvențe. Poate fi creat un aliniament optim, dar nu pot fi folosite mai mult de 5 secvențe, datorită timpului de calculare. De aceea este aplicată metoda progresivă a alinierii multiple.

Clustal efectuează o aliniere multiplă-globală după metoda progresivă. Etapele sunt următoarele:

- a) Efectuează alinierea a câte două secvențe folosind programarea dinamică,
- b) Folosește scorul alinierii pentru a produce un arbore filogenetic cu metoda neighbour-joining (NJ),
- c) Aliniază mai multe secvențe folosind arborele filogenetic.

Deși cele mai apropiate secvențe sunt alinate primele, iar apoi se adaugă secvențe sau grupuri de secvențe, utilizând aliniamentul inițial pentru a produce o aliniere multiplă, arătând în fiecare coloană variația secvenței în cadrul aliniamentului.

Cu cât sunt adăugate mai multe secvențe la profil, gap-urile se acumulează și influențează alinierea secvențelor următoare. Clustal calculează gap într-un mod nou, proiectat să plaseze gap-urile între domeniile conservate. Gapurile găsite în aliniamentul inițial rămân fixate. Pot fi adăugate gap-uri noi, când sunt adăugate secvențe noi, dar gap-urile nu pot fi șterse, numai adăugate. Clustal de asemenea implementează metode care încearcă să compenseze matricea scor (de ex. PAM), gap-urile care pot fi prevăzute și diferențele în lungimea secvenței.

Clustal are opțiuni avansate:

- Adaugă secvențe cu greutate
- Adaugă greutate la poziții diferite în secvență
- Adaugă o secvență sau o aliniere la alinierea existent
- Folosește arborele definit de utilizator pentru aliniere.

Unele dintre acestea vor fi discutate în capitolele următoare.

Problema cu alinierea progresivă este reprezentată de dependența ultimei secvențe din aliniamentul multiplu de aliniamentele inițiale. Primele secvențe care vor fi alinate sunt cele mai strâns legate din arborele secvențial. Dacă aceste secvențe se vor alinia bine, vor exista câteva erori în aliniamentele inițiale. Cu toate acestea, cu cât sunt mai departe înrudite secvențele, cu atât mai multe erori vor fi făcute, și aceste erori vor fi propagate la aliniamentul secvențial multiplu.

O a doua problemă a metodei alinierii progresive este alegerea potrivită a matricelor scor, a penalizării gap-urilor care se aplică la setul de secvențe.

5.4.3. Introducerea datelor în programul ClustalX

După ce porniți ClustalX, ecranul va arăta în felul următor (figura 5.4.3.a):

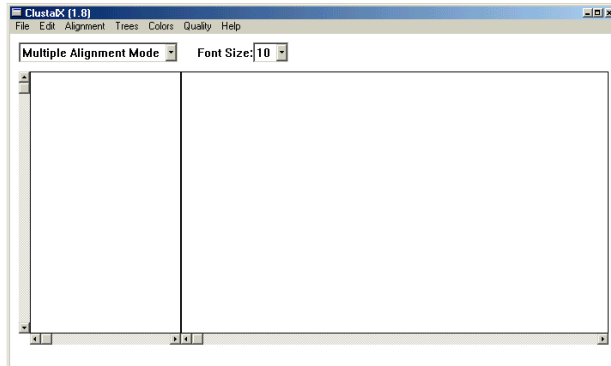


Figura 5.4.3.a. Fereastra programului Clustal

Alegeți din meniul *File* opțiunea *Load Sequences*. Pe urmă alegeți folderul care conține fișierul de intrare (fișierul text care conține secvențele pentru aliniere în diferite formate FASTA).

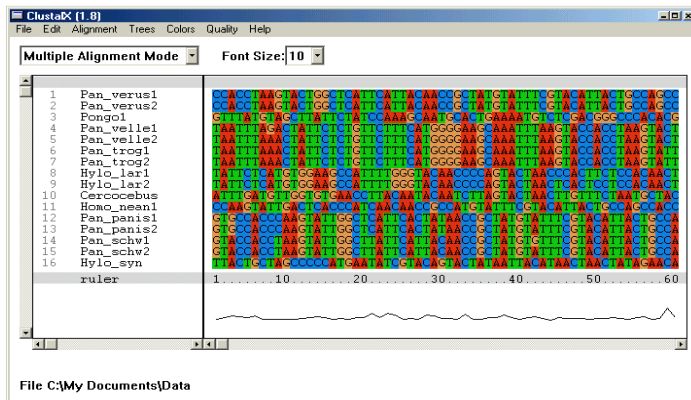


Figura 5.4.3.b.Exemplu de secvențe de nucleotide sub formă de fișier text

Panoul din stânga (figura 5.4.3.b) afișează secvențele conform cu numele care urmează după caracterul „>” din fișierul de intrare. Panoul din dreapta afișează începutul fiecărei secvențe. Puteți derula mai în dreapta pentru a vedea rezultatul fiecărei secvențe folosind bara de scroll din josul panoului.

5.5. Setarea parametrilor aliniamentului

Aliniamentul este realizat după câteva etape succesive:

- Resetarea tuturor Gap-urilor (*Alignment->Alignment parameters, Edit->Remove all Gaps*)
- Rafinarea parametrilor aliniamentului pereche (*Alignment->Alignment parameters*)
- Rafinarea parametrilor aliniamentului multiplu (*Alignment->Alignment parameters*)
- Rafinarea formatului fișierului de ieșire (*Alignment->Output Format Options*)

- Scrierea aliniamentului ca un fișier Postscript (*File->Write Alignment as Postscript*)
- Evaluarea calității aliniamentului:
 - a. Dacă nu sunteți mulțumiți -> mergeți la pasul 1
 - b. Dacă sunteți mulțumiți -> rafinați aliniamentul manual.

5.5.1. Parametrii alinierii pereche

Pentru a crea alinierea pereche, ClustalX trebuie să știe ce penalități să aplice pentru crearea fiecărui gap și pentru extensia aceluia gap. Alegeți *Pairwise Alignment Parameters* din meniul *Alignment*. Veți vedea apoi o casuță dialog asemănătoare cu figura următoare.

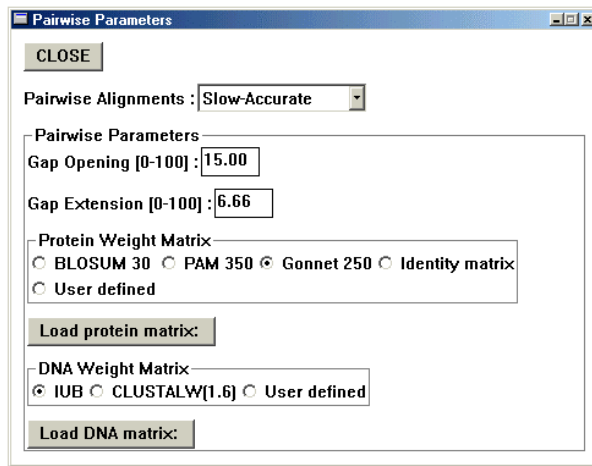


Figura 5.5.1.a. Afișarea parametrilor alinierii pereche

Prima opțiune în cadrul alinierii pereche vă dă posibilitatea să alegeți dintre metodele **Slow-accurate** (mai lentă, dar precisă) și **Fast-approximate** (rapidă, dar aproximativă). Metoda Slow este preferată, dar dacă aliniați mai multe secvențe, sau secvențele de aliniat sunt foarte lungi, programul va rula într-un timp mai lung, așadar veți dori să alegeți metoda Fast. Metoda Fast folosește o metodă K-tuple pentru alinierea pereche, spre deosebire de metoda Slow care folosește algoritmi din programarea dinamică. Casuța afișează valorile implicite pentru penalitățile Gap Opening și Gap Extension (deschiderea și întinderea gap-urilor). Micșorând penalitățile gap veți putea să introduceți mai multe gap-uri și mai puține nepotriviri. Acest lucru va avea ca rezultat potriviri care nu reflectă omologia. Creșterea penalităților pentru gap va avea un efect invers, dar s-ar putea să omitem potriviri care sunt omoloage.

Parametrii matricelor scor pot fi modificați și ei. Parametrii matricei IUB DNA are ca scor pentru potrivire 1.9 și pentru nepotrivire 0. Matricele de scor pentru proteine sunt echivalente aceluiași matrice folosite ca modele evoluționare în producția dendrogramelor. Toate matricele scor au avantajele și dezavantajele lor: PAM a fost folosită mult timp, dar este oarecum depășită, iar Gonnet poate fi mai potrivită să o folosiți. Matricea BLOSUM pare a fi cea mai bună pentru căutarea în cadrul bazelor de date. Puteți crea și încărca propria dumneavoastră matrice în ClustalX. Pentru descrierea formatului matricei puteți să studiați fișierul *matrices.h* care arată în felul următor:

- Pentru aminoacizi:

```
char *amino_acid_order = "ABCDEFGHIKLMNPQRSTVWXYZ";
short blosum30mt[]={
4,
0, 5,
-3, -2, 17,
0, 5, -3, 9,
0, 0, 1, 1, 6,
-2, -3, -3, -5, -4, 10,
0, 0, -4, -1, -2, -3, 8,
-2, -2, -5, -2, 0, -3, -3, 14,
0, -2, -2, -4, -3, 0, -1, -2, 6,
0, 0, -3, 0, 2, -1, -1, -2, -2, 4,
-1, -1, 0, -1, -1, 2, -2, -1, 2, -2, 4,
1, -2, -2, -3, -1, -2, -2, 2, 1, 2, 2, 6,
0, 4, -1, 1, -1, -1, 0, -1, 0, 0, -2, 0, 8,
-1, -2, -3, -1, 1, -4, -1, 1, -3, 1, -3, -4, -3, 11,
1, -1, -2, -1, 2, -3, -2, 0, -2, 0, -2, -1, -1, 0, 8,
-1, -2, -2, -1, -1, -1, -2, -1, -3, 1, -2, 0, -2, -1, 3, 8,
1, 0, -2, 0, 0, -1, 0, -1, -1, 0, -2, -2, 0, -1, -1, -1, 4,
1, 0, -2, -1, -2, -2, -2, -2, 0, -1, 0, 0, 1, 0, 0, -3, 2, 5,
1, -2, -2, -2, -3, 1, -3, -3, 4, -2, 1, 0, -2, -4, -3, -1, -1, 1, 5,
-5, -5, -2, -4, -1, 1, 1, -5, -3, -2, -2, -3, -7, -3, -1, 0, -3, -5, -3, 20,
0, -1, -2, -1, -1, -1, -1, -1, 0, 0, 0, 0, 0, -1, 0, -1, 0, 0, 0, -2, -1,
-4, -3, -6, -1, -2, 3, -3, 0, -1, -1, 3, -1, -4, -2, -1, 0, -2, -1, 1, 5, -1, 9,
0, 0, 0, 0, 5, -4, -2, 0, -3, 1, -1, -1, -1, 0, 4, 0, -1, -1, -3, -1, 0,2,4};
```

- Pentru ADN:

```
char *nucleic_acid_order = "ABCDGHKMNRSUVWXY";
short clustalvdnamt[]={
10,
0, 0,
0, 0, 10,
0, 0, 0, 0,
0, 0, 0, 0, 10,
0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 10,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 10,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0};
```

5.5.2. Parametrii alinierii multiple

Alegeți opțiunea *Multiple Alignment Parameters* din meniul *Alignment* (figura 5.5.2.a).

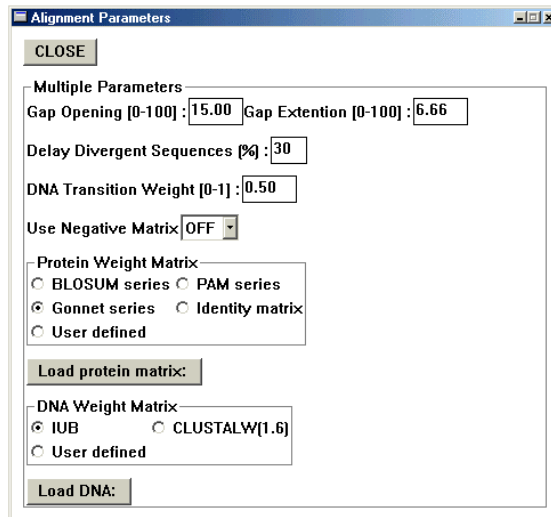


Figura 5.5.2.a. Fereastra de dialog pentru alinierea multiplă

Parametrii alinierii pereche și multiple sunt configurați independent, deoarece Clustal are nevoie de amândoi. Cum a fost precizat anterior, o matrice a alinierii pereche este calculată inițial, iar apoi pe baza acelor distanțe este formată alinierea multiplă. Folosind diferite setări pentru acești pași ai alinierii vom avea mai multă flexibilitate asupra modului în care este efectuată alinierea.

În comparație cu alinierea pereche vom avea câțiva parametrii în plus. *Delay Divergent Sequences* determină cum trebuie să fie două secvențe pentru ca aliniamentul lor să fie decalat. Această opțiune încearcă să compenseze pentru bias introdusă în metoda alinierii progresive. Scorul tranziției ADN poate fi modificat. Weight 0 înseamnă că tranzițiile sunt punctate ca și nepotriviri, iar weight 1 înseamnă că tranzițiile au același scor ca transversile. Pentru secvențe de ADN slab înrudite, scorul ar trebui să fie aproape de 0; pentru secvențe apropiate poate fi util să atribuim un scor mai ridicat. Nu mai trebuie să alegeți individual matricea scor, pentru că știm deja că de similare sunt secvențele. Clustal X va alege automat cea mai potrivită matrice scor dintr-o serie de matrice. Astfel, dacă ați schimbat matricea scor în cadrul parametrilor alinierii pereche, să efectuați și aici aceeași schimbare.

5.5.3. Formatul de ieșire al alinierii

Ultimul lucru care trebuie modificat înaintea efectuării aliniamentului este formatul de ieșire.

Formatul fișierului de ieșire se poate modifica selectând *Output Format Options* din meniul *Alignment*.

Când ClustalX creează o aliniere, scrie secvențele aliniate într-un fișier. Există mai multe tipuri de formate de ieșire, care vor fi necesare în funcție de programul care vreți să îl folosiți pentru analiza ulterioară.

Pentru construcția arborelui filogenetic alegeți *Phylip* ca format al fișierului de ieșire. Să nu uitați să scrieți secvențele aliniate și într-un fișier de tip Clustal.

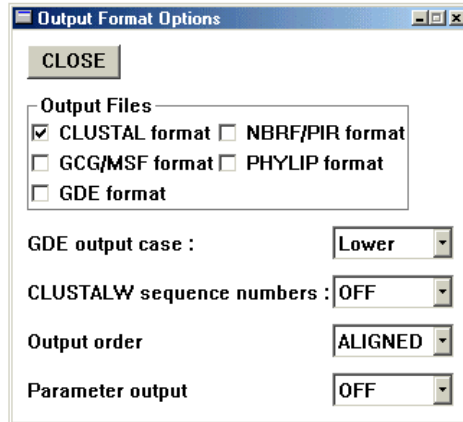


Figura 5.5.3.a. Opțiunile formatului de ieșire al alinierii

5.6. Crearea alinierii

Alegeți opțiunea *Do Complete Alignment* din meniul *Alignment*. ClustalX va spune ce face în fiecare moment, dar nu ar trebui să folosiți alt program în timp ce clustal efectuează alinierea. După ce alinierea a fost efectuată, ecranul principal va fi înprospătat cu secvențele aliniate.

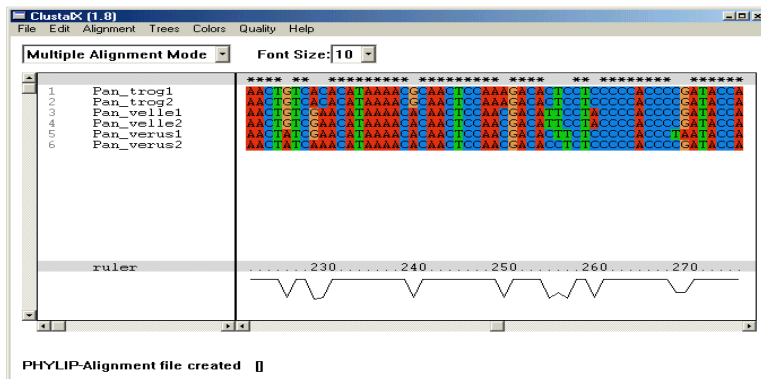


Figura 5.6.a. Exemplu de aliniere de tip Phylip

Bazele sunt colorate, ceea ce face ca evaluarea alinierii să fie mult mai ușoară. Histograma de sub linia indică gradul de similaritate (figura 5.6.a). Vârfulurile indică pozițiile cu cea mai mare similaritate, iar văile indică pozițiile cu cea mai mică similaritate. Linia gri aflată deasupra secvențelor este folosită pentru a marca porțiunile conservate. Caracterul „*” indică poziții care au fost conservate în totalitate (secvențe identice).

5.7. Scrierea aliniamentului ca fișier Postscript

Este posibil să folosim fișierele pe care deja le-am creat pentru a construi un arbore filogenetic, dar calitatea și valoarea acelu arbore nu va fi mai bună decât valoarea aliniamentului. Trebuie reținut că oricât de diferite ar fi două secvențe ClustalX va produce întotdeauna o aliniere. Doar faptul că există o aliniere nu înseamnă că secvențele sunt înrudite. Este la atitudinea utilizatorului să determine dacă secvențele din setul de date sunt omoloage, adică pot fi alinate. Calitatea alinierii este mai ușor de verificat, dacă există copii tipărite ale secvențelor. Există două posibilități:

- puteți tipări alinierea în format ClustalX și folosi aceasta, dar veți pierde informațiile colorate;
- puteți să instalați **Ghostscript** și **Ghostview**, care vă vor permite să manipulați documentele Postscript și să tipăriți alinierea ca postscript incluzând culorile.

Mergeți la meniul *File* și selectați opțiunea *Write Alignment as Postscript* - figura 5.7.a.

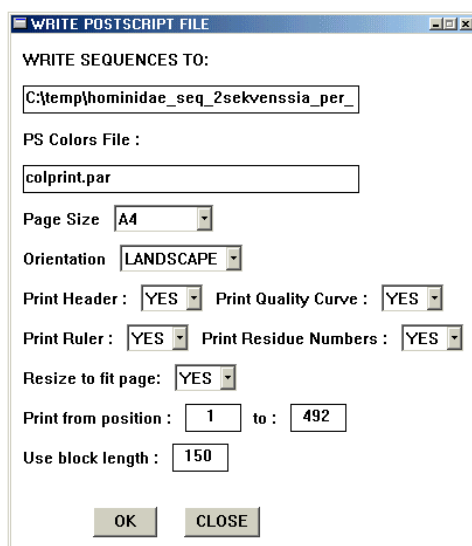


Figura 5.7.a. Exemplu de scriere a aliniamentului ca fișier Postscript

Ghostscript poate fi descărcat de la adresa:

<http://sourceforge.net/projects/ghostscript/files/GPL%20Ghostscript/9.00/gs900w32.exe/download>

și **ghostview** de la adresa:

<http://www.seas.ucla.edu/~ee5cta/ghostView/>.

Vor trebui schimbate anumite opțiuni, iar după ce le-ați setat, apăsați “OK”; apoi deschideți fișierul postscript în Ghostview și tipăriți alinierea folosind meniul *File=>Print*.

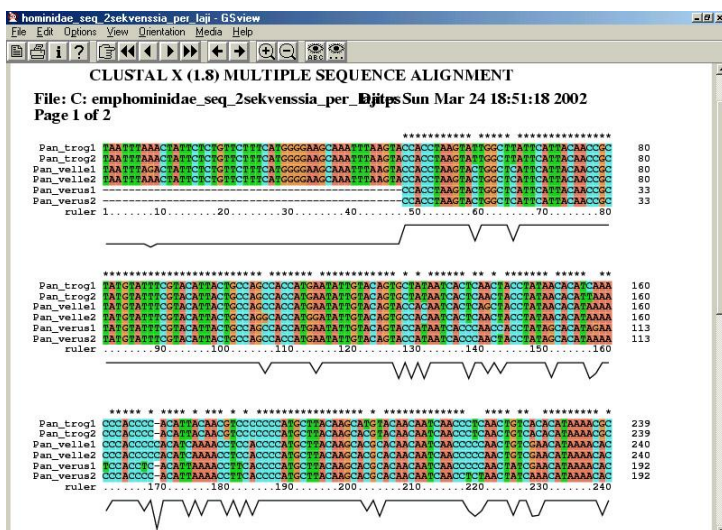


Figura 5.7.b Vizualizarea fişierului Postscript ce conţine alinierea secvenţelor, cu programul Ghostview

5.8. *Arbori filogenetici folosind ClustalX*

Pentru o analiză avansată a arborelui filogenetic se poate folosi pachetul **Phylip**, de la adresa:

<http://bioweb.pasteur.fr/seqanal/phylogeny/phylip-uk.html>.

NJ plot- este programul folosit de obicei cu ClustalX, de la adresa:

<http://pbil.univ-lyon1.fr/software/njplot.html>

Treeweb este o altă metodă software de vizualizare, găsită la adresa:

<http://taxonomy.zoology.gla.ac.uk/rod/treeweb.html>

ClustalX include implementarea algoritmului Neighbour-Joining (NJ) care ne permite să construim arbori filogenetici pornind de la o aliniere multiplă.

Atenție: există o diferență între arborele ghid (care a fost construit pe baza distanțelor între perechi înainte efectuării aliniamentului) și arborele NJ (construit după aliniament). Arborele NJ este construit calculând distanțele dintre fiecare pereche de secvențe în cadrul alinierei multiple. Alinierea dintre o pereche de secvențe poate fi diferită în cadrul aliniamentului multiplu.

5.9. *Exerciții propuse*

1. Vom porni de la un aliniament de secvențe peptidice al proteinei *Homoserine O-succinyltransferases* care poate fi descărcată de la următorul link:

http://www.bigre.ulb.ac.be/Users/jvanheld/bioinformatics_introduutory_course/web_course/data/sequence_analysis/metA_family.aln

2. Deschideți acest fișier cu ClustalX. Observați că nu este nevoie să aliniați această secvență deoarece sunt deja aliniate.

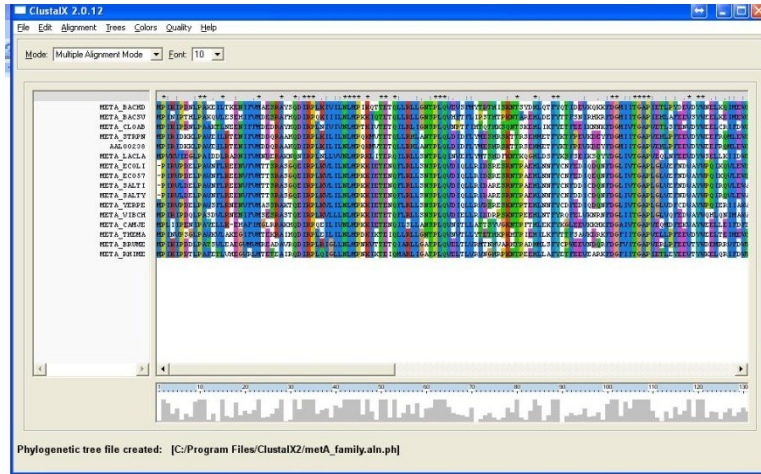


Figura 5.9.a. Aliniamentul de secvențe peptidice al proteinei *Homoserine O-succinyltransferases*

3. În meniul *Trees* selectați comanda *Draw Tree*. Aceasta va crea un fișier cu extensia *.ph* în directorul Clustal.

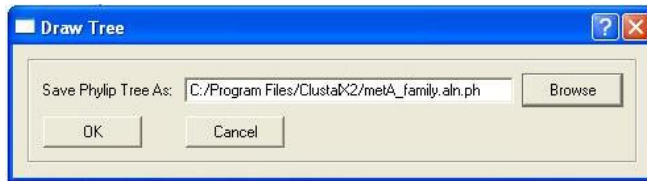
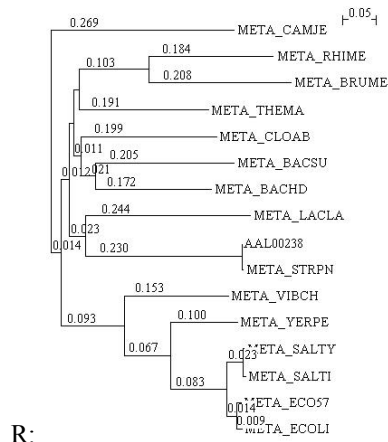


Figura 5.9.b. Exemplu de creare a fișierului cu descrierea arborelui filogenetic

4. Vom aplica și procedura *bootstrap* pentru a estima gradul de înrudire dintre diferite ramuri ale arborelui. În meniul *Trees* apăsați *Bootstrap N-J tree*. Aceasta va crea un fișier cu extensia *.php* în același director cu aliniamentul Clustal.

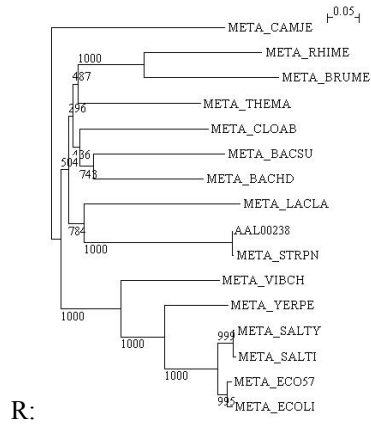
5. Deschideți programul **Njplot**.

6. Folosind Njplot deschideți fișierul care conține arborele NJ (extensia *.ph*). Verificați opțiunea *Branch Lengths*.



R:

7. Folosind Njplot deschideți fișierul care conține arborele NJ bootstrapped. Verificați opțiunea *Bootstrap Values*.



8. Deschideți diferitele tipuri de fișiere care conțin arborii filogenetic folosind sau Njplot sau Treeview.

- .dnd ClustalX guide tree
- .ph ClustalX neighbor-joining tree
- .phb ClustalX neighbor-joining bootstrap tree
- .tree the result of Phylip

Bibliografie

1. Lesk Arthur M. Introduction to Bioinformatics (3rd edition), Oxford Univ Press, Oxford UK, 2008
2. Durbin R., Eddy S., Krogh A., Mitchison G., Biological Sequence Analysis, Cambridge Univ Press, Cambridge Uk, 1998
3. Claverie J.M., Notredam C., Bioinformatics for dummies, Wiley, Hoboken N.J., 2007
4. Mihalaş G.I., Lungeanu Diana, Informatică medicală și biostatistică, Ed. Victor Babeş, Timișoara, 2009
5. Brown T. A. Genomes, John Wiley & Sons, New York USA, 1999
6. Mihalaş G.I., Neagu M., Neagu A., Textbook of Biophysics, Ed. Eurobit, Timișoara, 2001
7. Swanson T. A., Kim S.I., Glucksman M.J., Biochemistry Molecular Biology and Genetics (4th edition), Wolters Kluwer, Philadelphia, 2007
8. Simon-Gruita A., Saitan T., Biologie, CD Press, București, 2009
9. www.clustal.org –Jarno Tuimala-ClustalX
10. <http://www.embl.de/~seqanal/courses/commonCourseContent/usingClustalx.html>
11. http://www.bigre.ulb.ac.be/Users/jvanheld/bioinformatics_introduitory_course/web_course/practicals/phylogeny.html#contents Jacques van Helden-bioinformatics_introduitory_course