

Sequence databases and their organization, and Structural databases and their organization

قسم الأنظمة الطبية الذكية \ المرحلة الثانية \ المحاضرة السادسة م.د ميثم نبيل مقداد





Sequence Databases: Overview

Definition and Purpose

Sequence databases are organized collections of nucleotide (DNA, RNA) and protein sequences, essential for storing, retrieving, and analyzing biological information. They serve as a fundamental resource for genomics, proteomics, and related fields, aiding in the identification of genes, proteins, and their functions.

Types of Databases

There are two primary types: nucleotide sequence databases (e.g., GenBank, EMBL, DDBJ) and protein sequence databases (e.g., UniProt). Nucleotide databases store DNA and RNA sequences, while protein databases store amino acid sequences of proteins. Each type plays a distinct role in biological research.

Importance

These databases are indispensable in bioinformatics and molecular biology, facilitating sequence comparison, gene discovery, evolutionary studies, and drug development. They enable researchers to access and analyze vast amounts of biological data, leading to critical insights and breakthroughs.



dote no major sequence databases



Major Sequence Databases

\mathbb{Z}

GenBank (NCBI)

GenBank, maintained by NCBI, is a comprehensive nucleotide sequence database. It contains publicly available DNA sequences from various organisms and is a key resource for genomic research.

EMBL (EBI)

The EMBL Nucleotide Sequence Database at EBI is another major repository for nucleotide sequences. It collaborates with GenBank and DDBJ as part of the INSDC.

UniProt

UniProt is a comprehensive protein sequence and annotation resource. It provides expertly curated protein sequences and functional information, essential for proteomics and protein biology studies.

DDBJ

The DNA Data Bank of Japan (DDBJ) collects nucleotide sequences primarily from Japanese researchers. It is also a crucial member of the INSDC, ensuring global data accessibility.





Organization of Sequence Databases

Data Submission and Curation

The process involves researchers submitting sequence data, which undergoes curation to ensure accuracy and consistency. Curation includes verification, annotation, and standardization of the data, enhancing its reliability for downstream analyses.

Sequence Record Structure

protein code.

Metadata and Annotations

Metadata provides contextual information such as organism, gene name, and publication details. Annotations include functional predictions, structural information, and evolutionary relationships, adding significant value to the raw sequence data.

Sequence records typically include a unique accession number, sequence data (nucleotide or amino acid), and metadata. The accession number allows for easy retrieval, while the sequence data provides the actual genetic or

INSDC Collaboration

International Collaboration

The International Nucleotide Sequence Database Collaboration (INSDC) comprises GenBank, EMBL, and DDBJ. This collaboration ensures data exchange and standardization across global sequence databases, preventing redundancy and promoting consistency.

Data Exchange and Accessibility

INSDC facilitates free and unrestricted access to nucleotide sequence data worldwide. Data submitted to one member database is automatically shared with the others, maximizing the availability of biological information.

Importance of Open Access

Open access to sequence data is vital for scientific research, accelerating discovery and innovation. It enables researchers to build upon existing knowledge, fostering collaboration and reproducibility in scientific findings.



2

1



DDBJ



JAPAN

Searching Sequence Databases

NCBI Entrez System

The NCBI Entrez system is a powerful

search engine that allows users to access

various databases, including GenBank.

provides tools for filtering and refining

It supports complex queries and

search results.

2 1 3

Query Formulation

Effective query formulation involves using appropriate keywords, Boolean operators (AND, OR, NOT), and field specifiers (e.g., gene name, organism). Precise queries yield more relevant results, saving time and improving research outcomes.

Filtering Search Results

Filtering and refining search results involves using filters based on organism, sequence length, and date of submission. This ensures that only the most relevant and reliable data are considered for further analysis.





Structural Databases: Introduction

Definition and Purpose

2 Types of Structural Data

Structural databases are repositories of three-dimensional (3D) structural data of biological molecules, such as proteins, nucleic acids, and small molecules. These databases provide crucial insights into the structurefunction relationship of biological entities. These databases contain structural data obtained through experimental techniques like X-ray crystallography, NMR spectroscopy, and cryo-electron microscopy. The data include atomic coordinates, bond lengths, and angles, defining the 3D structure.

3 Importance in Research

Structural databases are vital for understanding molecular mechanisms, drug design, and protein engineering. They enable researchers to visualize and analyze the 3D structures, leading to advancements in various fields.



Protein Data Bank (PDB)

History and Development

The Protein Data Bank (PDB) is the primary repository for 3D structural data of proteins and nucleic acids. Established in 1971, it has grown exponentially, becoming an indispensable resource for structural biology.

Data Submission and Validation

Researchers submit structural data to the PDB, which undergoes rigorous validation to ensure accuracy and quality. The validation process includes checking for stereochemical errors, clashes, and overall consistency with experimental data.

File Formats

The PDB supports multiple file formats, including the original PDB format and the more flexible mmCIF format. The mmCIF format allows for more detailed and complex structural information, accommodating modern structural biology techniques.





2

3

Organization of Structural Databases

3D Coordinate Data

The core of structural databases is the 3D coordinate data, specifying the position of each atom in the molecule. These coordinates are essential for visualizing and analyzing the structure using computational tools.

Related Sequences and Annotations

Structural entries are linked to corresponding sequence data and annotations, providing a comprehensive view of the molecule. This integration allows researchers to relate structure to function, evolutionary context, and other relevant information.

Experimental Information

Experimental details, such as the method used to determine the structure (X-ray, NMR, cryo-EM) and resolution, are included in the database. This information is crucial for assessing the quality and reliability of the structure.



Searching Structural Databases

PDB Search Tools

The PDB offers advanced search tools, allowing users to find structures based on keywords, sequence similarity, structural motifs, and experimental parameters. These tools are essential for identifying relevant structures for research purposes.

Structure Visualization Software

Structure visualization software like PyMOL, Chimera, and VMD allows researchers to view and analyze 3D structures. These tools provide functionalities such as rotation, zooming, measurements, and highlighting specific regions of the molecule.



Integration with Sequence Databases

The integration of structural and sequence databases enables cross-referencing and comprehensive analysis. Researchers can use sequence information to find related structures and vice versa, enhancing their understanding of biological molecules.

national assessment of the second second

ent net a karno

Nelle

Card Look & Codini

Contraction of the

Their Devictions Devictor in Earth General Street And Date which choose and the series the series of the Date Conservation and the Date Date of the series of the series of the Date Date of the Series of the Se

Cavalities interdoute of Clinical

Sent Crettle

Donal tit: strains have ("Duran: Dur.2)

The second and provide the second and the provide second sec

Ramesterr Otensa





Specialized Structural Databases

ഷം

SCOP

SCOP (Structural Classification of Proteins) provides a hierarchical classification of protein structures based on evolutionary relationships and structural similarities. It organizes proteins into families, superfamilies, and folds.

CATH

CATH (Class, Architecture, Topology, Homology) is another hierarchical classification system for protein structures. It classifies proteins based on four levels: Class, Architecture, Topology, and Homologous superfamily.



MMDB

MMDB (Molecular Modeling Database) at NCBI provides structural data and interactive 3D visualization tools. It integrates structural information with sequence and functional data, enhancing the research capabilities.



Introduction to Sequence Alignment

Concept and Importance

Sequence alignment is the process of arranging DNA, RNA, or protein sequences to identify regions of similarity. It is a fundamental technique in bioinformatics, used to infer evolutionary relationships and functional similarities.

1

Pairwise vs. Multiple Alignment

2

Pairwise alignment compares two sequences to find the best match, while multiple sequence alignment (MSA) aligns three or more sequences. MSA provides insights into conserved regions and evolutionary patterns across multiple sequences.

Applications

Sequence alignment has diverse applications in molecular biology and evolution, including gene identification, protein structure prediction, and phylogenetic analysis. It helps researchers understand the relationships between different biological sequences.

3



BLAST: Basic Local Alignment Search Tool

2 1 3

Types of BLAST Programs

BLAST includes different programs for various types of sequence comparisons, such as blastn (nucleotide vs. nucleotide), blastp (protein vs. protein), blastx (translated nucleotide vs. protein), and tblastn (protein vs. translated nucleotide). Each program is tailored for specific sequence types.

E-value and Bit Score

The E-value represents the expected number of alignments with a similar score that would occur by chance. A lower E-value indicates a more significant alignment. The bit score measures the quality of the alignment, with higher scores indicating better matches.

Algorithm Overview

BLAST is a widely used algorithm for sequence alignment. It identifies local regions of similarity between a query sequence and a database. BLAST is fast and efficient, making it suitable for large-scale sequence comparisons.



Using BLAST

1

NCBI BLAST Web Interface

The NCBI BLAST web interface provides a user-friendly platform for performing sequence alignments. Users can input their query sequence, select the appropriate BLAST program, and configure search parameters.

Interpreting Results 2

Interpreting BLAST results involves analyzing the alignment scores, E-values, and sequence identities. The results provide information about the similarity between the query sequence and database sequences.

Adjusting Search Parameters 3

Adjusting search parameters, such as the E-value threshold and word size, can optimize BLAST searches. Fine-tuning these parameters can improve the sensitivity and specificity of the results.

.... M con It BLAST & Control Cont & M Coo Strephilip a bliebox £.5 - -> C C A W Nowres NC98Lotts Crossecorent utpost NCBI BLAST balderoad liceviter

NCEICH: FUR BLAST

BLAST Furn Seabeer Officer Iro

Impit	
lare:	
PHILAT	

The recare BLAST is device resute wan and erar light effectrcation and for extion.

Presitions

Results

Lone: NCB1 BLAST (0105)00, Jo 195 Silbier Soley, the Feart 815L (6a O Favit Insledctions, M, decolst,

Frattiabls

Allcer agaient Day

Search

· Pertifications, Neprenauls, 1041 CN4, 2L SarchEnlosber 2015, 864B · Searce in lag 2015

- · Serabet 139, Jet, Moral Pavices
- 4. SC, D. New: (CTA10, 113, 108, 2015
- NCA, Suppresett, 30, 1NDA, 74) NLD, Dc., 1, (1NBIS, 1DW, 116)
- NL L. NCS 5, CN, /19 BNEA(5, N
- 4. 19. Deta Comfronientic Contfor
- · Netveranve, INCA, Opdion in N
- · Furfceintibi: 101, At, Atco, lace
- · Intfornation, Racen, avapoaral Inttla: Carry 2015

NCBI BLAAT For fast of fracts

NCRI BLAST DLes(1, 161,93-R1 Comgrater, chvLf/618 franca HLLIM-LOIGFICE:

CB87-801.703.15 (F7.,9110/11) CPfleccinterfesult

ChDIfraien desctiving couting Chereriant ceptorration 3(1) DRBK Alian neveluces watrich sepalies, VEbla,

L Follake Blas; BLLARL, pur callean aor suction. (al conloste (oformectione Pulkentine, MiseCorld co alogiete an: compregation CLST VLLAT FLO: NCLE: (Steve NLCED 170: 1645.160,15910-CLISTORE, FIALER, IA (05

Q A Collicit Operative Presenancies Neame Presenancies Result reficent to NCBil is its elegreerients if its cad ybleate on NCBil Billent Liats lie in carcen in 0AL Cat Eagues 1,127 ,64) 01136, 2022 Presenancies Presenancies 1,127 ,64) 01136, 2022 Presenancies Presenancies Clever to Caccal Liest (Lier Dalase Whars Coencier Andluizes Prodice Caletior Dalase Whars Coencier Andluizes Prindion Consudcation Controlus Cantory Fealler Daly Prome Results The May Defeine Calls Prome Results Stars Lacon larers, cespers unlicht cesups Lin 109, 966-1055, (2156) Synd of the M159, iter; chonguing. Loni, 4mW115, F751-2459, 50720;1145. Servicies in fonation .16 17/47) (10: 125, 51013) n. 6CV114: LiArlages (12K) isits: 11004, 10, 63, 400, PioleStio. 0L55; In long, 11(in, Literr, antiegeran)) revels; In long, 11(in, Literr, antiegeran)) revels; Enttlien Restude. 665. Inte. 40, 2001	inin/CCBU 2014 and ALLST certectos	2
Calleie Commeeters Tere errensic Insparients Slow To Retions Neme Permatics Permatics Result retent to NCBLIS is elegrerients in His cad Precacion drecks be outen areantanies Clere to Caral Lite retent to NCBLI Billent Clats liw in careen in Precacion drecks be outen areantanies Clere to Caral Lite retelecty and Decree tre talss 114, 2022 Coencier Andluizes Price Results Coencier Andluizes Prome Results Prome Results retiler, NC211) 1M52, NO. 59 Frailer Daty Prome Results Feelgr meanation intees at lan Prome Results Start, Alan Prone Calls at lan Point Calls in ont, fumWilfs.PF3)-235; 30; 20; 1144. Point Calls in lon, 1, fumWilfs.PF3)-235; 30; 20; 1144. Point Calls styleter lawing, falen, DC7, (maital)) Prome Result issuing, falen, DC7, (maital)) Point Lawing, falen, DC7, (maital)) issuing, falten, DC7, (maital)) Point Lawing, falen, DC7, (maital))		Q 🛱 🖸 🕻
errensic Insparients Neane Permatties Result: reticent to NCBH is its elegercients if His cad ybleate on NCBH Bilent Elats liv in careen in 127 (44) 0116, 2022 Fails Fatz New Fain Nontee 22(04) arr (123) 5, 20165, 10) erylizes, NC211) 1M52, N0, 59 Frij, 105, te, 2013) F. A4, 20, 2) S, LAST, Alcon larors, cesperd urlich cesups (11-107, 066-1055, (2159) %end of the N159, itch; chong.ting, lon 1, 4mW115, 6757-3455, 30; 20; 31149. Sinxivies in fonction -18 19(45) (10-155, 30113) n, 6CV114, 11471ags (1b*) ist(5, 11064, 30, 69, 4000, PioleStio. 01555; in loa, 1116, faleer, initegeran)) revels; ny, ciklet Juwing, falen, DC7, (maita1)) Proteches Paration Science Silow To Retions Confict Provations Confict Provations ONL fat Eavies ONL fat Eavies	e Dessication	Califele Fcomnoeces
Iterensic insparients Conficate frocuations Neame Permatics Permatics Result interient to NCBil Bilent Lists limin careen in Preaction drecks be outen areantanies Cliver to Carial Live riceta fraction Iter fraction drecks be outen areantanies Cliver to Carial Live riceta fraction Iter fraction drecks be outen areantanies Cliver to Carial Live riceta fraction Iter 1136,2022 Iter Indian Centroication Contrive 22(007) aur (128) Controiny 5, 20165, 10) Fealler Daty Servitece, NC211) 1M52, NO, 59 Prone Results Infort, cespersi Infort, cespersi urifdth casups Infort, cespersi urifdth casups Iter; 111:197, 966, 1955, (2159) Felg meastation insets at lan Sinvitation Iter; 111:197, 966, 1955, (2159) Felg meastation insets at lan 111:197, 966, 1955, (2159) Felg meastation insets at lan Sinviteiten in fonation Iter; 116:11:107, 110:10, 53, 10:13) Felg meastation insets at		Slow To Retions
Neene ONL fat Eavies Presenties Precedenties Result Precedenties reticent to NCBI is is elegerecients it ilis cad ybleate on NCBII Bileat Diats lis in careen in Precedenties 1, 127 , d4) Dilis, 2022 Image: Test State	errensic Insparients	Conficate Trocvations
Presults	Representation	ONL Cat Eavies
<pre>rricent to NCBI is is elregeerients it ills cad ybleate on NCBI BILest Clats I is in careen in , 127 dd))116, 2022</pre> Pracation drecks be outen areantanies Cliver to Carial Lize frelecty and Berce tte talas , 127 dd))116, 2022 Coentier AndIsizes Pindion Comboicabion Contifus 22(00) aur (128) 5, 20155, 10) eryites, NC211) 1M52, NO. 59 Prome Results (Relet Elace Cley of Application) results, Atom larars, cesser3 artich cesups , 11: 187, 966: 1955, [2159) %end of the N159, itor; chongting, lon], fmW115, P75]-235; 30; 20; 3149. anximizes in fonction 18 19(45) (10: 155, 510: 9155; in lon, 111, Jaleer, anltogeran)) revels; ny, cikiet Juming, falen, DC7, fmai(al)) Part Calls For the Results of the Results of the Results Free Results Free Results Free Results Free Results Free Results Results Free Results Free Results F	Result	
 127 (d) 1136, 2012 Fails Futt New Fain Nontee (20065, 10) (20165, 10) (2016, 10)<td>ricent to NCBI is ils elegrecients it II is cad bleate on NCBII Billent Clats I is in carren in</td><td>Pracation drecks be outen arenntanies Cliver to Catial Lite frelecty and Berce the tales</td>	ricent to NCBI is ils elegrecients it II is cad bleate on NCBII Billent Clats I is in carren in	Pracation drecks be outen arenntanies Cliver to Catial Lite frelecty and Berce the tales
<pre>state (see (see (see (see (see (see (see (s</pre>	137	Oller
Whars Whars Whars Whars Whars Whars Whars Whars Coencier Andluize: Pindion Contoicus Contoirus Cantomy Feiler Daty Prove Results Relet Elack Cley of Application Relet Elack Relet Elack Cley of Application Relet Elack Relet Elack Rele	44)	Beracin Caletior Dalse
Coencier Andluizes Pindion Controicus 22(947) aar (123) 5, 20165, 10) eryiies, NC211) 1M52, NO. 59 Prome Resuts Contoiny Fealler Daily Prome Resuts Kelet Elace Ciey of Application Port Calls Port Calls Port Calls Port Calls Port Calls Port Calls Feig reassation insets at lan 11:107, 9666-1955, (2159) Wend of the N159, ien; chongeing. Ion], imN115, P75J-235; \$0; 20; 1149. Envietem in fontion 18:19(45) (10:155, 510:13) n, 6(4114, 11471ags (1549) is(5, 11004, 10, 63, 400, FipleStio. 9155; in loa, 11in, faleer, inlegeraa]) revel;; hy, CiAler Juring, falen, DC7, /maiial); visuion. Enttlien Restds. 865. Inte. in, 20011	0136,2022	Whers
Coentier Andruze Pindion Convicus		
c) Fails Fut New Fain Nomee Considus		Coencier Anditizes
22(947 arr (128) 5, 20165, 10) Servitee, NC211) 1N52, N0. 59 Prove Results Centomy Prove Results Relet Etack Cley of Application Relet Etack Cley of Application Relet Etack Cley of Application Relet Etack Cley of Application Pont Calls Pont Calls Pelg reassation insets af lan .11:187, 966,10955, (2159) Yend of the N159, ien; chongcing, lon1, imWli5, P753-2355, \$0720; 1149. <u>Anxietes in fonition</u> .18:19(45) (10:155, 510:13) n, 6(4114, 11471ags (154%) is(4,11094, 30, 63,400, FipLeStio. 91555; in lon, 11in, Talerr, inltogerau)) revels; h, ctalice Juring, falen, 0C7, /maiial)) yizuion, Enttlien Restdz. 865. Inte. in, 20011	Fails Fult New Fain Nonice	Conticius
S. 2016S, 10) eryiies, NC211) 1N52, N0. 50 Prome Results Kelet Exace Ciey of Application. S. LASY, Aton larors, cesperd wridth casups Pont Calls Pont Calls Feig reseation insets at lan .11:107, 966-1055, (2159) %end of the NI59, icn; chonge.ing. lon], 4m%li5.675J-235; \$0720; 1149. Enxistem in fontion .18 17(45) (10:155, 310113) n, 664(14, 11471ags (1b%) ist(s, 11004, 10, 69, 400, FipLestio. 9L55; in lon, 11in, faleer, inltegeraa)) revel;; ny, ClAiter Juwing, falen, 0Cr, /mai(a1)) yizu(an. Enttlien Results. 665. Inte. in, 20011	27(043 per (128)	Cantory
<pre>ieryiie6, NC211) 1N52, N0. 50 F75, 105, te, 2013) F. A4, 20.2) S. LAS7, A:cn larors, cesper3 wrfdrh cesups . 11:187, 966-1055, (2159) %end of the N159, ien; chongLing, lon], imWli5, P75J-235; \$0720; 1149. <u>Enxiletem</u> in fonition .18 17(45) (10-155, 510113) n, 6(4114, 11471ags (154)) is(4, 1104, 10, 63, 400, FipLeStio. 9155; in lon, 11in, Taleer, inltegerau)) revel;; ny, clAiler Juwing, falen, 0C7, /maiial)) ylsuion, .Enttlien Restdz. 865. Inte. in, 20011</pre>	5, 20165, 10)	O Fesller Daty
eryiies, NC211) 1M52, N0.59 Prome Resuts Kelet Elace Ciey of Application Resuts Kelet Elace Ciey of Application Point Calls Point Calls Point Calls Point Calls Point Calls Point Calls Point Calls Felg reaseation inses aflan 11:107, 966-1955, (2155) %end of the N159, ien; chongeing. Lon imWli5.975J-235; \$0;20;1149. Envisien in fontion 18:19(45) (10-155, 510:13) n.66V114.11471ags (15K) in 100, 11in, falter, inltegeraa)) revels; ny, clifier Juwing, falen, DCy, /maiial)) y25uian. Enttlien Reside.865. inte. in, 20011		in the second
Prome Resuls F1, 105, te, 2013) F. A4, 20.2) S. LAS7, Aton Imforts, cespersi urfdth cesups . 11:187, 066.1055, (2159) %end of the N159, icn; chongring. Ion], imNl15, 075J-235; \$0;20;1149. Snxietem in fontion .1817(45) (10-155, 510113) n, 6(4114, 11471ags (154)) is(5, 11004, 10, 63, 400, FipleStio. 9155; in Ion, 11in, Taleer, inltegerau)) revels; ny, clAier Juring, falen, 0C7, /maiial)) y 2;uion, Enttlien Restdz. 865. Inte. in, 20011	ervites, NC211) 1N52, NO. 50	
 Kelet Etack Ciey of Application. Kelet Etack Ciey of Application. St.LAST, Aton Infort, cesperd Wridth cesups Felg measuation insets at lan It: 107, 966-1055, (2159) %end of the NIS9, iten; chongaing. Ion imN155, F75J-235; \$0; 20; 1149. Envietem in fonation 18 19(45) (10-155, 510113) n. 664/14. 1147/ags (154%) isits: 11004, 10, 69, 400, PipLeStio. 0L55; in lon, 11in, faleer, inlegeraa]) revel's; ny, cikier Juwing, falen, DCr, /maiial); yizu(an. Entallen Reside: 865. Inte. in, 20011 		Prone Resuts
F1, 105, te, 2013) F. A4, 20.2) S. LAST, Alcon larors, cesper3 wildrh cesups T1: 107, 966: 1955, [2159) %end of the N159, icn; chongLing, lon inN115, F75]-235; \$0;20; 1149. Envirien in fontion 18 19(A5) (10: 135, 31013) n. 6CV11c.11471ags (1b%) is(5; 11004, 30, 69, 400, PioLeStio. 9L55; in 10a, 11in, Jalesr, inltogeraa)) revel\s; hy, cikier Juwing, falen, DC _J , /maiial); y250/en. Enttlien Restdz. 865. Intc. in, 20011		C Kelet Etace Ciey of
<pre>te, 2013) F. A4, 20.2) S. LAST, A:on Infors, cesperi urifich cesups .11:107, 966-1055, (2159) %and of the W159, icn; chongcing. Ion], imW115, P75]-235; \$0;20;1144. <u>Enxistan in fonition .18:17(45)(10-155,510:13) n, 66(W114,11471ags (15%) is(5,11004,10,69,400, FipleStio.9L55; in loa, 11in, Taleer, inItegeraa)) revel;; ny, ciAier Juring, falen, DC;, /maiia1); y12uion, .Enttlien Restdz.865.intc.in, 20011</u></pre>	71, 109,	Abblicgion
S. LASY, L:on Infors, cesperi urldth cesups 2 11:107, 966-1055, [2159) %and of the W159, icn; chong.ing. Ion], imW115, P75J-235; \$0;20;1149. <u>Enxietza</u> in fonition .18:17(45) (10:155, 510:13) n, 66W114.11471ags (1bW) is15:11004, 10, 63,400, FipLeStio.9L55; in Ioa, 11in, Taleer, inltogeraa)) revel;; ny, clAier Juwing, falen, DC _J , /maiial); yizuion. .Enttlien Restdz. 865. inte. in, 20011	te, 2013) F. A4, 20.2)	
10/073, cesser3 > Point Calls wildth cesups > Feig retestation intess .11:107, 966-1055, [2159) > %and of the WI59, icn; > chongring. > lon], fmWli5, P75J-235; \$0;20;1149. > @nxieiza in fonition - .16:17(45) (10-125, 510:13) > n.6(4114, 11471ags (1b%) > ist5, 11004, 10, 69, 400, FipLeStio. 9L55; . .in lon, 11in, Taleer, inltogerau)) - revel\$; ny, clAler Juwing, falen, DC _J , /maiial); ylzuion, . .Enttlien Restdz. 865. Inte. in, 20011 -	S. LAST, Alon	Prod Colle
aflan 2 aflan 3 aflan 3 aflan aflan aflan aflan aflan	larors, cespera	Pont Calls
.11:107,066:1055,(2159) Send of the NI59.icn; chongring, lon].imWli5.F75]-235;50;20;1149. <u>Sinvivies</u> in f <u>onition</u> .1819(45)(10-155,310113) n.664114.11471ags (15%) is(5:11004,10,69,400,F1pLe5tic.0155; in lon,11in, faleer, inltegerau)) revels; ny,clAler Juwing, falen,0C;,/maiial); yizuion. .Enttlien Restdz.865.inte.in,20011	artor it ceseps	af lan
lon], unWl15, F73)-235, 50720; 1149. <u>Snylwies</u> in F <u>onition</u> -18 17(45) (10-155, 310113) n, 6CWl14. 11471ags (15%) is(5, 11084, 10, 59, 400, FloLe5tio. 9L55; in log, 11in, faleer, inltegerau)) revels; ny, clAlet Juwing, falen, DC,, (maiial)) ylsuion, . Enttlien Resids. 865. Inte. in, 20011	.11:107,066/1055,(2155) 9and of the N159.icn; changcing,	
<u>Gnvieica</u> in f <u>onition</u> -1617(45)(1V-105,3101(3)) on,6CVII4.1147lags (16%) -161(5:11004,30,61,400,FloLeStio.9L55; ,in 100,11fn,faler, inltegeraa)) revels; ny,cikler Juwing,falen,0C ₇ ,/maifal)} yEsuion. - Enttlien Reside.865.intc.in,20011	lon), unW115, 075]-355; \$0;20;1149.	
ris(S,11004,30,59,400,FipLeStio.9L55; ,in lon,11in,Jaleer, inltegeraa)) Trevel%; ny,ciAler Juring,Falen,DC7,/maiTal)? yEsuian, .Enttlien Reside.865.Inte.in,20011	<u>ânvieten</u> in f <u>onition</u> -18 17(45) (10-195, 310113) n. 6(4)15, 11471aas (15%)	
:15(5:11004,30,53,400,PlpLeStio.0155; ,in lon,11in,falerr,inltegeraa)) Trevel%; ny,ciAlet Juring,falen,DC;,/maiIal)} y2suion, ,Enttllen Restdz.065.1ntc.in,20011		
ny, cikler Juring, falen, DC7, /mailal)} y zzuien, , Entklien Reside. 865. Inte. in, 20011	15(5,11004,30,63,400,PipLeStic.QL55; ,in loa,11in,Jalesr,inltogerau)) revelt:	
y Esulen, . Entklien Resids. 865. intc. in, 20011	ny, cifier Juring, falen, DC7, (mailal)	
	ylsuien. .Entillen Reside.865.Inte.in,20011	

Drand Leonis 00000 Keccicat

Compretions



Multiple Sequence Alignment: Basics

Definition and Purpose

Multiple sequence alignment (MSA) is the alignment of three or more sequences to identify conserved regions and evolutionary relationships. MSA is crucial for understanding sequence variability and functional conservation.

Progressive Alignment Method

The progressive alignment method is a common approach for MSA. It starts with pairwise alignments of the most similar sequences and progressively adds more sequences to the alignment based on similarity scores.

Scoring Matrices and Gap Penalties

Scoring matrices (e.g., PAM, BLOSUM) assign scores to different amino acid or nucleotide matches. Gap penalties are used to penalize the introduction of gaps in the alignment, reflecting evolutionary events such as insertions or deletions.



CLUSTALW

SEQUENCE S_EMPIMANTS" AITROMANS.

T COFFEE +

COFFEE

Multiple Sequence Alignment Tools

ClustalW/Omeg a

ClustalW and its successor, Clustal Omega, are widely used tools for MSA. They employ progressive alignment methods and offer various options for parameter tuning.

MUSCLE

MUSCLE (Multiple Sequence Comparison by Log-Expectation) is a fast and accurate MSA tool. It uses an iterative refinement approach to improve the alignment quality.



T-Coffee

T-Coffee (Tree-based Consistency Objective Function for alignment Evaluation) is an MSA tool that combines local and global alignment information to produce accurate alignments.





Visualizing Multiple Sequence Alignments

Consensus Sequences

Consensus sequences represent the most common nucleotide or amino acid at each position in the alignment. They highlight conserved regions and provide insights into the functional importance of specific residues.

Conservation Analysis

Conservation analysis involves quantifying the degree of conservation at each position in the alignment. Highly conserved regions are often critical for protein structure and function.

Color Schemes and Formatting

Color schemes and formatting are used to enhance the visual representation of MSAs. Different colors can represent different amino acids or levels of conservation, making it easier to identify patterns and trends.





Applications of Sequence Alignments

Evolutionary Studies

2

3

Sequence alignments are used to infer evolutionary relationships between different organisms. By comparing sequences, researchers can construct phylogenetic trees and understand the evolutionary history of genes and proteins.

Functional Annotation

Sequence alignments are used to annotate the function of unknown genes and proteins. By comparing sequences to those with known functions, researchers can infer the function of the target sequence.

Protein Structure Prediction

Sequence alignments can aid in predicting protein structures. By identifying homologous proteins with known structures, researchers can model the 3D structure of the target protein.

Challenges in Sequence Analysis

Big Data in Genomics

The increasing volume of genomic data presents significant challenges for sequence analysis. Handling and processing large datasets require substantial computational resources and efficient algorithms.



Computational Complexity

Sequence alignment algorithms can be computationally intensive, especially for large datasets. Optimizing these algorithms and developing faster methods are crucial for efficient sequence analysis.

Low-Complexity Regions

Dealing with low-complexity regions, such as repetitive sequences, poses challenges for sequence alignment. These regions can lead to spurious alignments and require specialized methods for accurate analysis.





Future Directions and Conclusion

Integration of Multi-**Omics Data**

3

Integrating sequence data with other omics data (e.g., transcriptomics, proteomics, metabolomics) provides a more comprehensive understanding of biological systems. This integration requires advanced bioinformatics tools and databases.

Machine Learning

Importance of Databases and Tools

Biological databases and sequence analysis tools are essential for advancing biological research. Continued development and innovation in these areas are critical for unlocking the full potential of genomic and proteomic data.



Machine learning is playing an increasingly important role in sequence analysis. Machine learning models can be used to predict gene function, protein structure, and other biological properties from sequence data.