**Bioinformatics**
**Prof. M. Michaele Gromiha**
**Department of Biotechnology**
**Indian Institute of Technology, Madras**

**Lecture – 05b**
**Protein sequence databases II**

Now, we will explain more about the contents of UniProt. So, as I discussed earlier it contains a higher notation of the protein sequences. So, it has more information regarding a particular protein. So, what are the major aspects, what are the major contents of Uniprot database? First, they give about the proteins right, the names, origin and attributes and so on and go with the ontologies and then the next the sequence information, right.

(Refer Slide Time: 00:42)



Then once we give this information then they have to support, this information by other means. So, the data they give in the primary data or supported by the bibliographic references because theories are important to assess the reliability of the data as well as the correctness of the data.

So, all the information they provide in the database right or supported with bibliographic references, where they obtained the original data, original information then they give the cross-reference with a lot of our databases I will show some of the important databases and they give the complete information regarding this particular entries.

So, if you want to search the data for any of these specific proteins right then I will show with one example.

(Refer Slide Time: 01:26)



So, hemoglobin B chain what is the importance of hemoglobin B chain? Yeah, transfer protein, oxygen carrier protein right this is (Refer Time: 01:33) hemoglobin B chain. So, if you click the hemoglobin B chain and click on search typically here. So, this will show you all entries the related with the hemoglobin B chain. So, we saw here. So, these are various accession entries right these are the accession number, this entry name and we have the protein names right we have the subunit beta or subunit alpha, subunit gamma and this is a gene name right and we know the organism sometimes from the human, sometime from this mycobacterium tuberculosis.

So, these are the length of the protein, its 147 amino acid residues in this particular hemoglobin subunit beta.

(Refer Slide Time: 02:06)



So, next they give names right the synonymous names what are the alternate names for these particular protein right. Then give a gene name, organism and so on. Then here you get the protein attributes is 147 amino acid its sequence status is complete because they got all the sequences right and its known at the protein level.

(Refer Slide Time: 02:28)



At the protein level, they have the evidence right protein evidence at this protein level. Then you go with general annotation, they try to incorporate try to include most of the data available in the literature.

So, what is the function of this particular protein? They transport oxygen. They transport oxygen from lungs to the tissues. So, this is a major transport. So, then how many subunits or how many small subunits in this particular protein? Totally 4 right because it has 2 alpha chains and 2 beta chains right. We can see the alpha 1, alpha 2, beta 1, beta 2. So, it is the tetramer right. This is heterotetramer right because two different chains this is why hetero, there are 4 chains totally so is a tetramer right.

So, where is it found in tissues? It is in the red blood cells. So, whether any post-translational modifications in this particular protein. So, they mentioned was the glycation and different types of post-translational modifications right at this site.

Mainly the acetylation on Lys 60, Lys 83 and Lys 145 right and they give the proper references also. So, where they have the post-translational modifications; then where this is involved in diseases right they give a lot of information regarding the disease. So, you can see the cause different diseases based on the mutations right. So, they listed up several diseases or we can look all the details right from the UniProt database mainly they give the sickle cell anemia. By the mutation of 6[th] residue (Refer Time: 03:47) Glutamic acid to valine. So, they mention all the details about the diseases from this particular protein.

(Refer Slide Time: 03:53)



So, now they give the ontologies what are the various biological process we see oxygen transports. So, any diversity of this coding region, they have polymorphism right they

have the changes upon mutations or either any other ligands or any other small molecules they bind to this particular protein. So, what are the different small molecules binds to this protein? Heme, iron, metals, and Pyruvate right these are the different ligands they bind with this protein. So, what are the functions and what are the post-translational modifications occurred in this particular protein. So, various types of PTMs: acetylation, glycation, glycoprotein, phosphoprotein and s nitrosylation the various types of TMs, then we go with the ontology. So, there are different types of gene ontology. They are widely used in the literature, one is a biological process cellular component and molecular function. So, they give the different specific subclasses in this different gene ontological process.

(Refer Slide Time: 04:53)



Then you go with this interactions, they have the interaction in other proteins. So, this the database for the protein interactions that intact.

So, they give the interaction with the other proteins, how they are interacting with the different other proteins you had different sites. What are the metal binding sites, where different binding sites they give the information regarding the binding site of a particular protein?

Now the variants, there are various variants in this particular protein right. So, they give the wild-type residue, this is the residue which are real in the original residue in the protein.

This is the residue which is replaced. So, this is replaced, for example, valine is replaced to alanine, then they will happen to what is the major disease. So, here give you all the mutations 'from' this is 'to' and which type of diseases, and they give the references. Then they can give all the information essentially if you see UniProt is the unique resource which contains all the information regarding any particular protein right.

Now, you go with the next step, now still now they mention about a general information right what are the functions and what are the gene ontology, what are the binding sites what are the PTM sites all the information they give. Now they go with the sequence level right the major aspect of this UniProt. So, UniProt contains not only the sequence, it contains the data for other different functions and the structure of a particular protein.

(Refer Slide Time: 06:22)



So, here this is the sequence. How many residues is in this protein? 41, 42, 43, 44, 45, 46, 47 they have 147 residues right. So, it gives the complete sequence right and here they give the secondary structure because as you know this protein right this is the alpha-helical protein this is predominantly with the alpha helices, I will discuss the secondary structure in the later classes.

So, if you see the blue once they are mainly helices right and the green one strands, but here it does not have any strand and we have some turns here and there we can see some turns right. This is the secondary structure information regarding their particular protein. So, you can see this sequence in two different formats, either you can see this is the UniProt format and also you can see the FASTA format what is FASTA format? FASTA format is a format which starts with a greater than symbol, and here then we can see this is the command line and here the sequence will start.

So, this is a kind of format which you adopt in bioinformatics generally for the treating the sequences as well as for our large-scale analysis. So, to separate two sequences they use this specific format. Then you can also see if you see a sequence any other sequences in any other organisms they have similar to this sequence right. So, there are various tools if you see here there are various tools available, one is BLAST I will discuss the details in later classes right this will help you to see the proteins, which are related with your original protein sequence.

So, you have your own sequence, if you click on this blast this will give you what are the other sequences related with your own sequence. So, now, we have a reference because this is an important part, because how they obtain the information. Because they cannot get directly from again is from outside right that is why we have to use any reliable resources. So, only one major resource reliable resource, so that is the published articles, where shall we get this information.

Student: PUBMED.

PUBMED database right now we discussed in the previous class previous classes right. So, PUBMED provides the resource for all the published articles mainly in biology right and medicine. So, we have this different papers we can get the information from PUBMED and they integrate all the data in the database and wherever they collect the information from the literature. So, they give the references.

(Refer Slide Time: 18:42)



Right these are the references. To provide all the information its very time-consuming. So, manual curation is a very hard this is a reason why the manually curated sequences are less compared with the computer translated sequences. So, where they can collect the information by keyword searching and they put the collect the data from the resources and put it as it is right there we do not we have to work on that, but what the reliability. So, we need to do the manual curation.

(Refer Slide Time: 09:07)



So, now give the sequence databases right we can see the; what EMBL Genbank and DDBJ which sequence database? Nucleic acid sequence database, where they give the complete sequence databases and the translated sequences they provide the information.

(Refer Slide Time: 09:21)



Then this is the 3D structures. So, this is a method how they obtained a data right and the resolution and the positions and so on.

More details I will discuss in later classes.

(Refer Slide Time: 09:32)



Right.

(Refer Slide Time: 09:34)



So, these are the protein-protein interaction databases, MINT, IntAct, and STRING. So, it will give you the particular proteins and how this interact with other proteins. Then they give other protein databases and the post translational modification database and so on right.

(Refer Slide Time: 09:49)



Then also they provide genome annotations, organism-specific database as well as the phylogenomic databases. This is the one part; the first part is a general information, the second part is a sequence and the structure and the functional information plus the links and now they give the different one type annotations.

They give the enzyme in pathways, where they have this reaction and gene expression databases right and family domain database and as well as the other resources, for example, the drug bank and other resources. Fine, then they have the other relevant documents like chromosome and other polymorphism disease mutation and so on.

Now, you are familiar with the all the contents of UniProt database right when you have time you will look into the UniProt and take your own protein and if you search and if you read, then you will be comfortable with understanding all the aspects of a particular protein. So, here the hemoglobin B chain, in the first part we give all the functional information regarding the general information and second part with the protein sequence along with the secondary structure or the tertiary structure and the links with other databases, and the last part they give the enzymes and pathways and the interactions and so on. If you look search with the hemoglobin B chain, for example, it will give you lot of data. So, for example, if you see this, many data and from that some of them are redundant. If you see the 2 sequences there are sometimes 90 percent identity, sometimes 80 percent identity right.

(Refer Slide Time: 11:14)



For example, if you have these two sequences, what is the identity? 100 percent right. So, we can see the 100 percent identity. So, in this case, this UniProt provides an option to select these sequences with some level of redundancy for example, if you want to reduce to 90 percent; that means, if two sequences if the similarities identity is less than 90 percent you can get, but more than 90 percent you will take one and discard this one. So, also you can get the redundant reduce redundancy up to 50 percent. They will check the two sequences and see the identity if it is more than 50 percent they will keep one and discard the other.

How to reduce this redundancy I will explain in the subsequent classes right. So, you can get that if you do this here now for example, currently you get 1243 results, when you look into this specific redundancy right now we can get 240 results. So, here we use the identity of 0.5. So, we reduce the data from 1243 to 240. Because to do any analysis it is very important to have the data from different ways, not to use the same several times because you will introduce bias because of the same sequence right. So, in this case, it is important to have a non-redundant data. So, here you can get with the two different cutoffs right.

For example, this is 100 percent, 90 percent, and 50 percent right, if you want to have other redundancies we have to use other programs available in the literature.

(Refer Slide Time: 12:50)

So, these are the entries now you have 240 results. So, these are the different entries right obtained from this particular search.

(Refer Slide Time: 12:57)

Now we can also download the data the various options to download the data right they have the various options. For example, you can see a tab delimited and you can use excel, you can download in FASTA and you can download XML and you can list the accession numbers right as well as you can get the XML or RDF format.

So, various formats you can download the data. So, if you give the FASTA format right. So, you give the FASTA format this is the result. So, as we discussed earlier in the FASTA format its start with the greater than symbol and this is the name and followed by the sequence you can see from here to here this amino acid sequence, sequence number 1 and here again the greater than symbol started this is sequence 1 right and here this symbol started again. So, it is started with sequence 2. So, and this is a sequence.

So, you can get the all the 240 results based on the greater than symbol, you can understand the first sequence starts from' and end with the 'RSF'. The second sequence started from 'MAN' and ended with this 'AMRF'. Then you can separate data and you can use these data for the analysis. Now I have two questions. The first question is obtain the sequences of transcription factors with less than 50 percent sequence identity.

(Refer Slide Time: 14:10)



Transcription factors are DNA binding proteins right they have specific functions. So, to understand the features of this transcription factors, sequences, I would like to collect the sequence from the transcription factors and reduce the data with 50 percent sequence identity. This is question number one.

Question number 2 I am interested to understand the function of a particular protein, for example, this human mitochondrial beta-barrel membrane protein VDAC, two questions. So, to understand the first question, how to obtain this information how to obtain the transcription factors with less than 50 percent sequence identity, first you have to go to

UniProt, right. The first step is to go to UniProt database and search for transcription factors.

(Refer Slide Time: 14:59)



So, go to there and we can see the UniProt. So, we have to search with transcription factors. So, we search this is the second step.

So, when you search transcription factors what will happen? We will get a list of sequences which contain 'transcription factors' in any field because if you use the symbol if you symbol search right you can get the 'transcription factors', which are matched with anywhere or you can they search in this you can have the different options, if you want to get the only the title you can give title or any other MESH term (Refer Time: 15:30) you want to use it we can use any MESH(Refer Time: 15:32) terms.

So, in a transcription factors how many data we get? So, more than 10,000 data right now if you use now you will get more than that. So, now, what is the question? So, we need the transcription factors or with less than 50 percent sequence identity. So, what to do with this? You have to choose the 50 percent and if you click we will get the data. So, to download the data. So, there is an option called download here, click here download what will happen? We will ask for the option which format do you want, if you ask for the FASTA format, we will get the data in FASTA format right.

(Refer Slide Time: 16:00)

This is the one and here this is the second one and the third one and so on. So, you get all the data in FASTA format right the first question now it is done, now we can use this for the analysis right. What is the second question? Yeah, I want to see amino acid sequence as well as some of the functions of this specific protein, because this protein is a recently published one. So, it has several functions right, eukaryotic protein, it is a first eukaryotic protein right in beta barrel membrane protein.

(Refer Slide Time: 16:32)

So, if you do this right just for go to the UniProt and if you search with this keyword, mitochondrial beta barrel membrane protein and human and VDAC right then you will get these specific entries and if you click on any of these things right, you will get the data this is a sequence. So, these are your sequences if you want to get the FASTA format you can click into FASTA format and then you will get the sequence in FASTA format. Now again in the UniProt, you can obtain the information on different perspectives; for example, if you are interested on any of the particular protein, you want to do that. So, we can do it and get the information. If you want to get the information regarding post translational modification sites right if you give the codes, if you can get all the post translational modification sites or if you want to collect the sequences of a set of data.

For example DNA binding proteins or the RNA of binding proteins or any information right if you give the see the in the search, you search the correctly you get the data and it is also possible to get the data with any specific sequence redundancy. So, it is a unique resource, it contain lot of information regarding protein sequences right. So, this is the reason why several researchers they are using the UniProt database in their research. So, now, we will recollect again, what are the various aspects we discussed in the class today. Primary structures, what is a primary structure? Specific arrangement of amino acid residues in a protein right like it has a main chain and it has a site chain. So, main chain it is the same right what is the main chain?

(Refer Slide Time: 17:57)

Student: (Refer Time: 17:59).

The amino acid side there is a NH2 right.

Student: C alpha.

C alpha.

Student: C

C N C alpha C it goes and right then you can put the COOH, this is the amino terminal N terminal this called the N terminal or amino terminal. This is C terminal. So, we can see this is the information we will get, but C alpha is connected one side with the.

Student: R group.

R group. Another with the?

Student: H.

H right. So, this is H here this is R1 this is R2. So, now, this is the main chain this is chain is the same only the side chains are different. So, you get R 1, we will get R 2, R 3 and so on this can be same amino acid or the different amino acid because this determine the sequence. So, here if you see this one what is R1 here in this sequence this is M R 2.

Student: Alanine.

Alanine right because M is here right and then this is 1 2. So, this how we get the sequence there is a primary sequence we know the link, we know the sequence arrangement, but we do not need anything else if we do not know the locations right.

So, now what is the primary resource for the protein sequences?

Student: PIR.

PIR and the (Refer Time: 19:20) SWISS-PROT right this is the earlier developed databases right then they merge together to form.

Student: UniProt.

UniProt database UniProt is the universal protein database right now it is widely used in literature. Now what are the different contents of UniProt sequences? Functions and the structural information and the various interactions and the pathways and the assemblies right and they supplemented with the original literature, then what are the major applications of the UniProt? Yeah different functions and the different sequences different distribution of amino acids say different functions right we can use for higher order specific sequence or any correctly sequences right and you can use it for the further applications right.

This UniProt is the unique resource for protein sequence databases. Next class I will discuss about what are the various aspects you can derive from this amino acid sequences that we will discuss in the later classes.

Thanks for the kind attention.