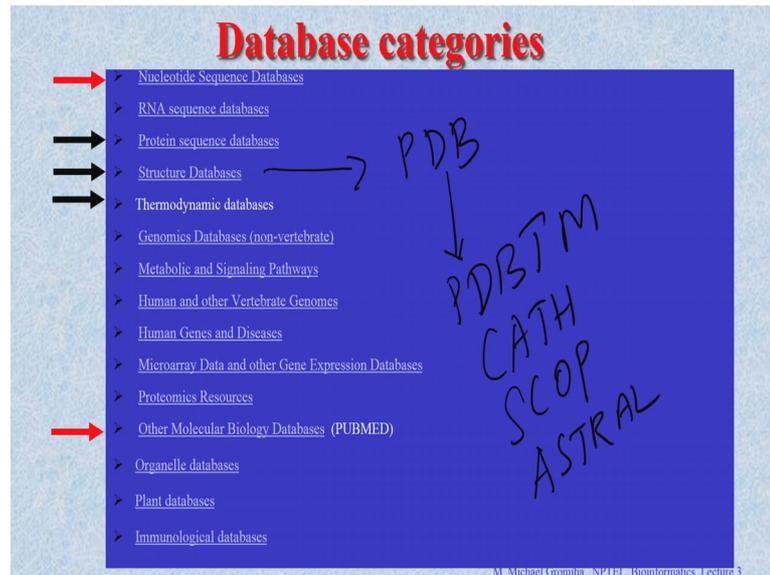**Bioinformatics**
**Prof. M. Michael Gromiha**
**Department of Biotechnology**
**Indian Institute of Technology, Madras**
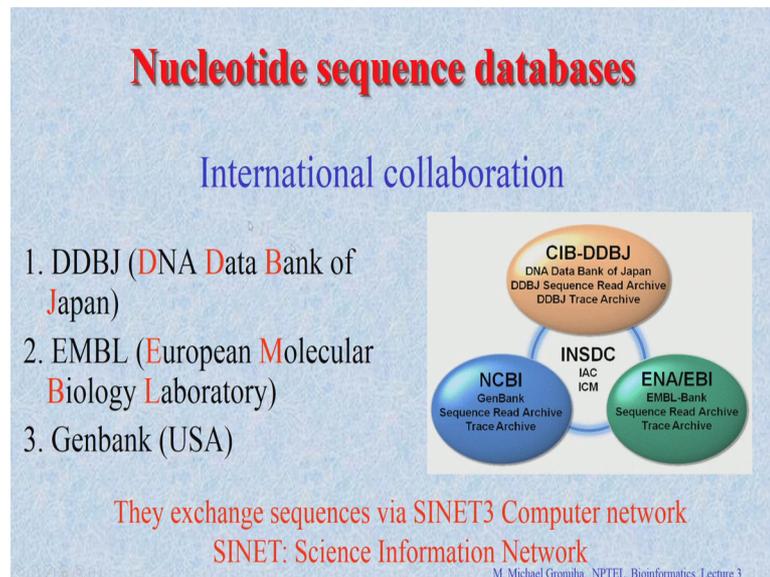
**Lecture - 3b**
**Databases Categories**

(Refer Slide Time: 00:17)



So, this lecture while cover few databases mainly the DNA database and protein databases right. So, in the last class we discussed mainly about the DNA. So, I will move on to little bit about the DNA databases. So, it is very important to get the sequences of a nucleic acids right.
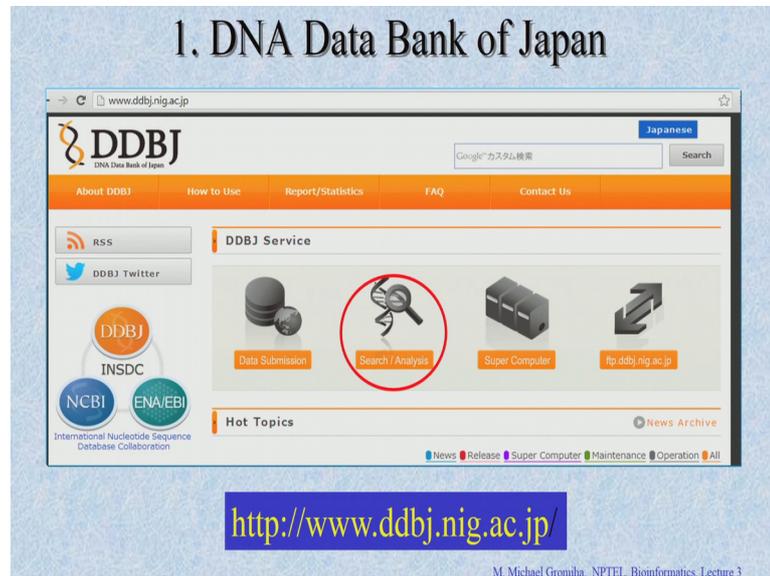
(Refer Slide Time: 00:41)



So, at the same time there are different types of databases have been developed for example. So, one is a DDBJ we say DNA data bank Japan right. So, this is in Japan and the second is EMBL, EMBL is the European molecular biology laboratory. So, they also developed this collecting nucleated data and GenBank USA they also sort out to collect.

Is the similar time because it is several people they think in a similar way right. We cannot say this is very unique that only we will do. Because science is open and research is open. So, it is very competitive right. If you have some information immidiately several people think in a same way, this is the reason we need to be very fast and we very accurate and we have to provide reliable data. So, they started collecting the data and then they also developed some tools to analyze the data right then what happened in this case? Each databases they have different types of data right. So, in this case there is some discrepancies. So, users have to access all the databases, that biggest sometimes you get the data from DDBJ may not be available EMBL.

Sometimes you get the data from EMBL may not be available in GenBank. So, for getting a data set, they have to check all the databases. Second aspect is the search options and the display options and the format may not be same if it is different then again we have to rework right. So, then (Refer Time: 02:09). So, they join together using these science information network. So, here is the DDBJ right. So, NCBI and we have the EMBL. So, joint together. So, they share the data. So, now, the data what we get from
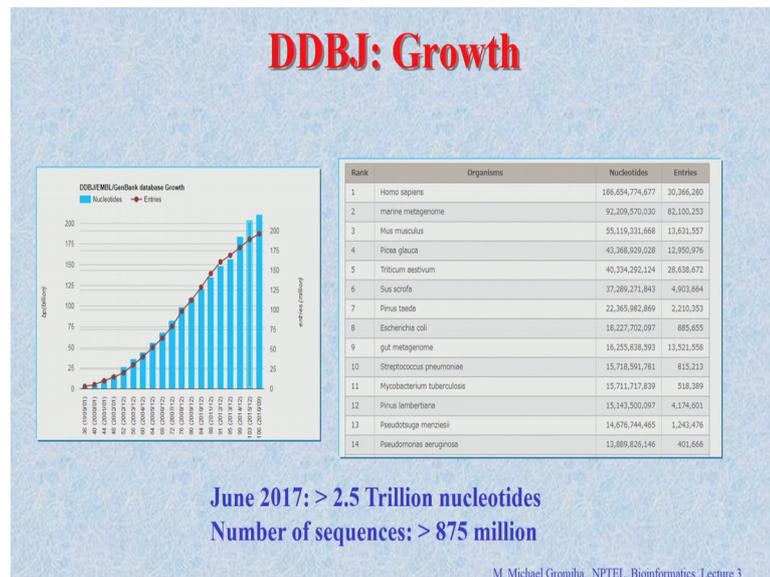
DDBJ you get from EMBL or from the GenBank. So, wherever you use or access the databank database you will get all information right.

(Refer Slide Time: 02:34)



So, now will just demonstrative bit about these databases, DDBJ this is in Japan in (Refer Time: 02:42) Japan a beautiful city in Japan you can have this a (Refer Time: 02:45) and all fine. So, here they have a different options, you can submit the data and you can search and analyze the data and all they have the facilities and also you can FTP other information right. So, this is the website DDBJ dot nig that national institute of genetics this is the institute there maintaining the database and ac dot JP. So, you can access this website and you can get the data of the DNA sequences fine.

(Refer Slide Time: 03:14)



So, this is a growth get this is the started long term ago. So, it is in 1999. So, these the in January this is the release number 36, before that we started to its initiated database and it is a current status. So, this is 2016 106 release. So, this they we have the 2.5 trillion nucleotides as of June 2017 this let us update and about it; so it 875 million sequences. So, if a plenty of sequences for the analysis and we look into the organisms as we expected. So, we have the top most human and followed the mus musculus and then you can see call a other organisms. So, you can this is the top most organisms it the DDBJ in DDBJ fine. This is nucleotides this is the number of entries available for each organism right.

So now how to get the data; so there are various way is to get the data from DDBJ, the first two condition is the very simple search there is a quick search.

(Refer Slide Time: 04:15)



If you want to get the data for any of these for your DNA or RNA; so here you put the homo sapiens mrna for this protein right. So, if you search right then you will get the link right. When you click on these accession number right, these are the 36833.

(Refer Slide Time: 04:33)



Then you get the complete data. So, now, you see; what are the contents, which are available in DDBJ. So, first you give the name right. So, I have gave some of the data. So, I can give the homo sapiens mrna three and 15 base pairs right and the source, source is homo sapiens and here they have given accession number specifically for this a DDBJ.

So, these accession number and (Refer Time: 04:56) keywords right. So, in this case you can search the data base on any specific keywords.

Right and the others who are though and who sequences data and the reference; you give the complete reference and you give the PUBMED index, they give the link to the PUBMED; PUBMED this literature database ok.

(Refer Slide Time: 05:18)



That I will explain little bit later then and they give the data of number of as T C and G. So, here is the number of A, number of C, number of G number of T and this is the complete sequence right 315 10 15 base points right. So, in these sequence right. So, what is the AG; AG at contents, what is that AT contents?

Student: (Refer Time: 05:38).

At plus 76 divided by 300 and.

Student: 15.

15 right (Refer Time: 05:51) calculate, this is number of as this is number of ts and you can calculate the a t conduct right. Same thing they give the translation protein sequence also right the fine.

(Refer Slide Time: 05:59)



So, that one search; so they give lot of options right some of them are familiar, some of them are art familiar. So, in this case they give the ontology of each terms. Now if you click on this one you will get the details of right the specific term. So, here I put the sequence length 400 to 500 and orgs organism human. You can use any of these search options to get your sequence.

(Refer Slide Time: 06:25)



Now, if you give these are a list of sequences, there are so many error data for extremely that two twenty 26506 entries between 400 and 500 from human right.

(Refer Slide Time: 06:38)



So, then if we click any of this access number we will get the full data right. We discussed earlier. So, they give the definition, they give the keywords and the general name right.

(Refer Slide Time: 06:47)



So, this is the translator protein sequence right, this is the niy where with the last class we discussed about the coding right is the one codes for a protein, I mean amino acid right. So, this is a protein sequence this is DNA sequence, I will get a protein sequence fine.

Now, the DDBJ they have the option to upload your data if we have any sequence right you can also upload the data. So, either you give the single sequence or the multiple sequences.

So, they receive the data and they first validate, and if are all aspects if the data is clean then they will a include your data in the database right that is fine similar way. Now, we discussed about DDBJ, similar data are maintain in the GenBank and EMBL right. So, what is GenBank, how many of you use GenBank? Fine right. So, what is a GenBank right?

(Refer Slide Time: 07:47)



It is a NIH genetic sequence database right, they provide the collection of the information regarding the DNA sequences right. So, they give all the information this is the website. So, consider GenBank will get the information regarding the sequence available in GenBank. So, it is the developed by NIH.

So, in CBI a there maintaining this database right on the all the d-th they publish in your database issues anywhere because anywhere a nucleic acids research is the general as they discussed earlier. So, they publish the details about the databases and you can get all the information regarding that fine.

(Refer Slide Time: 08:20)



So, this is the contents of GenBank. So, like we discussed in DDBJ. So, GenBank also a similar type of contents right; maybe we can say look into this website and get more details just I will go through quickly. So, they has the definition because this is. So, this is saccharomyces cerevisiae right. So, this is accession number right here you give the organism and if the others names right the reference and the features.

(Refer Slide Time: 08:45)



So, here you give the various divisions, they give the primary sequences and rodent sequences, mammalian sequences, they have the various conditions, they various

classifications they given the separately. So, if you are interested any given specific sequences for example, viral sequences right they give the database for the data for the viral sequences, they have various classifications right.

(Refer Slide Time: 09:07)



So, then again they give the data about the literature and what the different features and the coding sequence. So, what is the coding sequence?

Student: (Refer Time: 09:17) translated protein.

Translated into protein so that nuclei types that corresponding the sequence of amino acid in protein neither they how they translate. So, this give from which two which. So, we will put n dot dot means from n to m for example, 687 to 3158 see this is the coding region. Then if you partial one they put the less than symbol right we have partial at the 5 dash n then if it is greater than symbol they put partial at the 3 (Refer Time: 09:44). So, they give the information right we also give the complementary stand information right and the translation.

(Refer Slide Time: 09:51)



So, here this is the search option available. So, you put the keyword here right if you click on this go, and then you will get the list.

(Refer Slide Time: 10:00)



So, they give you different types of information, not only the nucleic acid sequences they give the data about the how many literature in the buffered, right if order the how many nucleic rate sequences, I mean GST sequences right to protein sequence and so on. So, if you click on the nuclei type. So, there are 19.

(Refer Slide Time: 10:19)



So, now if you get all the 19 correct these are the 19 sequences you get from GenBank.

(Refer Slide Time: 10:28)



Then if you click on this first one right then it will show the details right fine. So, now, get the data from the GenBank right we discussed about the DDBJ, I am discussed about the GenBank. So, another one is EMBL.

EMBL is maintain from European micro biology laboratory right. So, one is in US one is in Europe and another one is Japan right because they have they very good first commuting facility under facility long time ago this is way first the developed there

right. Using most of the databases called PDB Uniprot and all these databases right they have the (Refer Time: 11:01) between all these developed countries, because internet came there first and they are very first computing facilities and all. So, they started there because familiar trying to use there right because India we got very late right because when I visited Japan in 1997, they ask with you to do with the all everything with the email and all.

But we can only facts at the time right, but currently it is synchronized right whatever we get in the developed countries, we immediately in India also right. This is the reason why almost of the early developed databases right ever thing is from US or in the Europe or in Japan fine; so EMBL. So, they also start we have the database for the nuclei acid sequence database right.
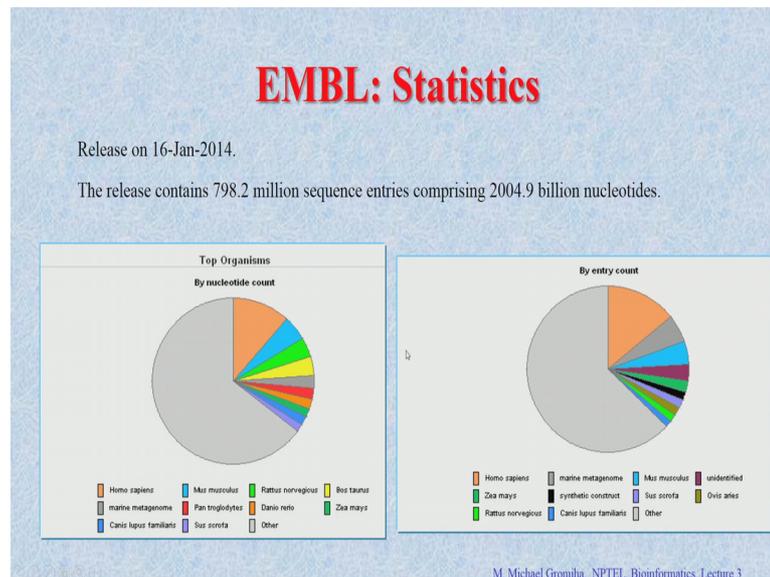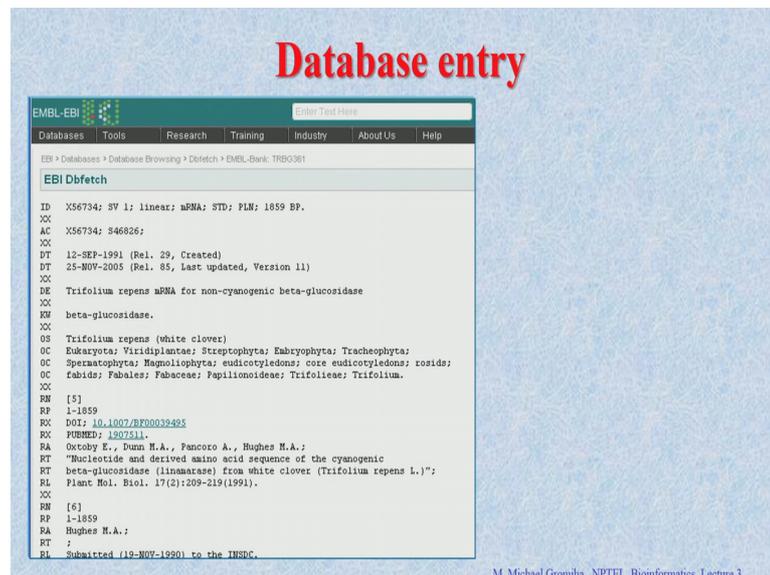
(Refer Slide Time: 11:43)



Here the main resources sources are the DNA and RNA sequences right. So, they also accept direct submissions, they are also (Refer Time: 11:50) data right. So, mainly from the genome projects as well as the patent applications, this is the website for the EMBL right you can the EBI dot ac dot UK and you got the EMBL. So, this is the website.

(Refer Slide Time: 12:06)



So, now if you got to the EMBL this is the statistics, but similar to DDBJ, now correctly see the mainly this is the organisms right Homo sapiens is the most right and the forward by the other organisms, if are also this is entry count and when nuclei count. So, you can see the similar level of the statistics fine ok.

(Refer Slide Time: 12:25)



Now, go to the database entry, here also if you see the similar to the GenBank or familiar DDBJ, we see the names and the PUBMED entries and we have the coding regions this is the proteins sequence and we have the DNA sequence.

(Refer Slide Time: 12:37)



This has 1859 base pairs and you can see the condense of ACG and T right fine. So, EMBL also have a search options.

(Refer Slide Time: 12:51)



You can search to the homo sapiens gapdm RNA right glyceraldehyd site. Now, go there. So, they have a various options right.

(Refer Slide Time: 13:06)



If you give all the results these are the various results for the EMBL. So, now, we have a very format you check, whether the EMBL format if you click in the EMBL format right.

(Refer Slide Time: 13:15)



So, it gives you the EMBL format, this is the contents of this particular database.

(Refer Slide Time: 13:20)



Then we have the protein sequence and the DNA sequence fine.

(Refer Slide Time: 13:24)



So, if you use the DDBJ right r is the GenBank or is the EMBL now you can get the data for any of the all nucleated sequences. Then you can use the sequences for further analysis. So, you suggested to look into the databases and use the options available in this databases, and then try to get the sequences from different organisms and see whether anything is different or the same once again calculate different properties right.

Last class we discussed about various properties from the nucleotides right what various properties we discussed in the last class?

Student: Flexibility.

Flexibility.

Student: Stacking energy.

Stacking energy; so it depending upon the dinucleotide or trinucleotides right. So, we classified the nucleotide sequence in the overlapping segments, either dinucleotides or trinucleotides depending upon the availability of data, then we can calculate the average values. This will tell you how for this particular sequence is table compared with the other sequence or lot. So, there is database called dinucleotide proper database. Last class we discussed only few properties right mainly the flexibility or rigidity or the base stacking energy or hydrated bond. So, this database serious more than 140 a features right this available is this sets diprodb.

(Refer Slide Time: 14:43)



So, these are the various information available this database here they have a twists stacking energy, rice, bent and so on right.

Here these are the various dinucleotides, we have the property names. So, you can analyze various features; that means we do not know which feature is important for the

flexibility, which features important for the binding affinity right. So, you can use any of the features and try to understand why this is important what properties are important for the binding affinity or if the any mutation is casting the diseases, you have the mutation data here right then you can relates with this decision information so on fine, so this database available in the website. So, these are download able. So, you can download the data right for all the features and you can use this features to understand the different DNA sequences.

Fine; so till away discuss about the database; different types of databases, the collection of databases right and only for the DNA or the where databases we develop for the DNA.

Student: DDBJ

DDBJ, EMBL, GenBank, and this, then nuclei proper database. Now, we will discuss about literature database right what is literature database?

Student: (Refer Time: 15:56).

And discuss data about the published articles, that the earlier days right.

(Refer Slide Time: 16:06)



## Literature database

Chemical abstracts

Physical abstracts

Science citation index

PUBMED

Mainly for biological sciences and are freely available

Useful for getting the references of any work and related papers

M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 3

If you (Refer Time: 00:00) because they started to have this some abstracts like chemical abstracts the physical abstracts right. So, in this case this is very big volume very small letters, very difficult worried that. Even I have not sure the publish a number base in the

computer readable form. Publish in the chemical abstract and the physical abstracts they are the very famous abstracts. Then we started the science citation index right they started to rank the different journals as the different papers, and see how many citations each article or each general published articles in generals are cited by other others.

Then PUBMED is the widely used database for the live senses right they include mainly a live sense papers, not the physics or chemistry right from several lading generals right with the reference to the all the papers published in the literature ok.

(Refer Slide Time: 16:50)



So, I will explain about the proper because its widely used database and we get almost all the information, there going to the biology and as well as for the medicine. So, if you want to get the information regarding DDBJ nucleic acid sequence database. So, you want to get the article published about this a database. If you search with the keyword right, you will show you the all the data right. From this you can get this is the collection, because this is also includes the DDBJ this is listed here and you can say somewhere DDBJ progress report and various other information right.

They listed other articles also, because they also linked with this DDBJ, this is the real why they get the other articles right.

(Refer Slide Time: 17:28)



So, we go DDBJ progress report right. So, then we can get the data. So, when you display the data there are various options; whether they need the format this summary or the abstract or any other formats and how many items are page? 5 or 10 or 20 and whether you can sort whether you need the first other or the general name or the recently added and so on right. So, when you apply this right, then we can get the information; this is what we need. So, they aspect DDBJ nucleic acid sequence database, here we have DDBJ, we will get the data. For if you want to get the full text (Refer Time: 18:03) full text then you get the complete paper, because this is the title.

If you click here you will get the abstract and if you click here you will get the full text right, then we are the related citations what is the meaning of related citations?

Student: (Refer Time: 18:17).

The papers which are similar to DDBJ. So, we can get a GenBank, you can get the EMBL also other databases right. So, can get the different entries publish articles, which related to this DDBJ right. So, this is a website for the PUBMED and you can access this website to get the literature database right.

(Refer Slide Time: 18:36)



So, now you get the click on the DDBJ progress report, you get the abstract and here you can see get this full paper. Either you get this Oxford University process this is the publish in your site or you can get the proper central right what is open access? So, they put is open access right.

What is the meaning of open access?

Student: You can use the free earlier.

You can use the free earlier right because to publish and article publish a general issue. So, it will was start of cast, who will bear the cast two options; one option is the subscribers they have to bear the cast, in this case (Refer Time: 19:13) have to pay anything. We have to publish the article when this accepted with the review or in the other procedures publishers to publish, how to earn the money because they we have to subscribe the generalist. When you subscribe you have to pay, though those who want to read this articles they have to pay right. These go they maintain till few years ago. In the open access means the others will play right in this case the read us do not go to pay because others will pay the charges and the tell the publishers to make it offer. So, anybody can access this article.

So, this is call though open access right there are two options, some journals they can only open access, some journals they can only sub by subscription some journals they

have the option given to the others. Others can choose whether you want to make it open or not. Where is open you have to pay buts available to everyone, where is not open its not available to everyone. So, only restricted users can read the article right. So, depends on the others they can you make it as open access or not fine, this is open access means you can use it right.

(Refer Slide Time: 20:18)



So, you can go the full text right now this the PDF format right. So, what is PDF format? Portable dockman format right in this case you can get the full paper you can read in the format, then you can like a print right.

Here we can read it online whether if figure if PDF format, this just look at of print like (Refer Time: 00:00) in general article right you can do that.
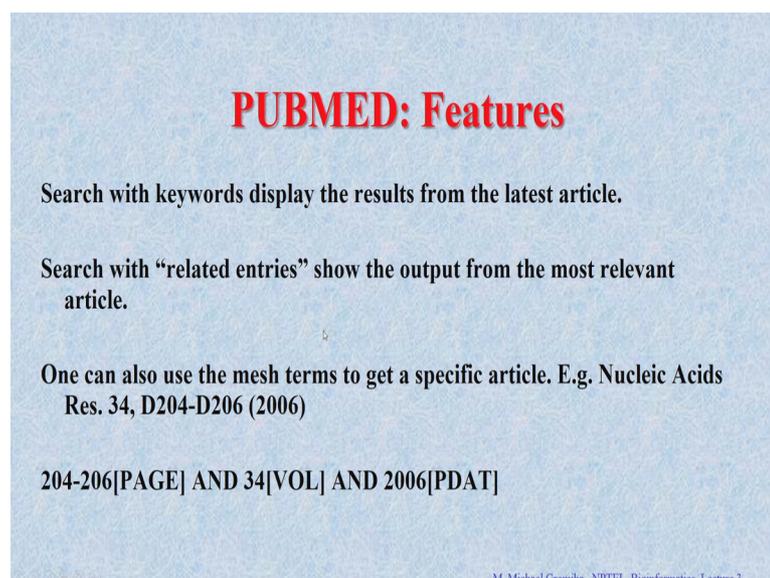
(Refer Slide Time: 20:45)



So, now it is related article. So, if you look into these related click on related articles. So, it show the all the articles, which are relevant to your search. So, for example, this is your search right the man article, if you go to related citations these are the other papers which publish in the literature which related to the article which our query article right you get all the related articles you can get that fine.

(Refer Slide Time: 21:09)



So, what is the various features of the PUBMED. So, you can search with the any keywords, when get data for the very latest article right and search with the related

entries to give the articles which are relevant to the related articles, and also you can some mesh terms to get any of these articles. So, can you they get with the page numbers, you can with the with the general name, you can search with others and so on; so main difference between currently available PUBMED and previous years well.

Previous days, unless we get the article we cannot site, but currently if you type any keyword you get all the articles. So, main problem is currently if you write a man script just you type the keyword, get all the papers, they do not data full paper, some thermal people they do not get the abstract also right. They do not know the others, just they take for topmost 5 15 and site everything.

The earlier days the original others cut proper credit, because they do not get all the papers only the important papers we have to send a request and they will send the refract, then we read, then we understand, then only be site number of citation if this is very less, because all of papers almost state by the others, but currently they said 100s, but they do not read in other papers. That is a difference between the advantage of this PUBMED right very easily they others use the get the information fine.

(Refer Slide Time: 22:42)



So, this is another database called disease database right this is the 3D inside this database which contains the various information regarding the proteins sequence and they there are thermodynamics and (Refer Time: 22:49) and so on they are such added some of the information regarding the diseases. So, how are the mutations, they change

the of the mutations right because any of these diseases. They obtain the data from the other databases and then develop the database for the diseases.

(Refer Slide Time: 23:06)



This is another database for the protein function database, we developed few years ago. So, here it will tell you these important residues which are perform in different functions in membrane proteins.

So, you have a search options and we have the display options and we have go they have get the results. For any mutations you can say whether (Refer Time: 23:23) and. So, what is the function of the protein as well as how many specific mutation, which alters this specific function or not; likewise if you do little literature the very databases based on the protein sequence, protein structure and the thermodynamics, diseases and the literature and so on. So, look into this nucleic acid research website, we will get the information regarding the all the information so that you can use it.

Thanks for your attention.