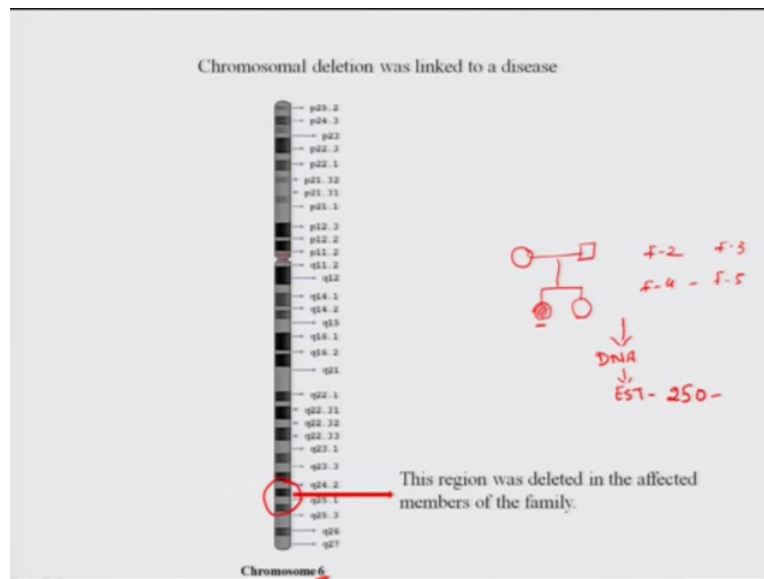**Functional Genomics**
**Professor S Ganesh**
**Department of Biological Sciences & Bioengineering**
**Indian Institute of Technology Kanpur**
**Lecture No 13**
**Tutorial Part 1**

Hello every one and welcome to the tutorial session of functional genomics myself Anupama and as already introduced I am one of the teaching assistance in this course that is functional genomics today we will be learning about whatever information can be arrived at or fetched using few hundred base pairs of DNA sequence.

(Refer Slide Time: 0:50)



So now let us look into that so this is a part of chromosome which was deleted in a disease and linked to that disease. This was identified by a team of scientists so few of you must be wondering that how did they arrived at the region of the chromosomes. So that has been covered already in one of the NPTEL course that is Human Molecular Genetics and also we are rerunning that course but since you are here very briefly I will like to tell you that how do we arrive at such sequences, so what is done that there are a families which come for counseling to the geneticists and they have members who are affected suppose this female was affected and she has a another sister and their parents were unaffected.

Similarly there was a another family, family 2, family 3 with family 4 so when you once you have a certain amount of family for which you can screen, screen for a cause of the disease. So what is done that DNA is extracted from all the family members which are present and then they are EST markers. Now what are EST markers these are small expressed sequence tag, which are unique for particular region of the chromosome.

These are unique for the regions so what do the scientists do they screen for these EST sequence, so suppose approximately around 250 EST sequence can cover almost all the autosomes of the chromosomes and from there on they link, they try to see that what EST region or what EST is getting segregated with the disease by that I mean that which one are always present in the individuals which are affected.

For example in this female and then the other person of the other family which is affected if that EST marker is present or not. So from that they narrowed down to the region of the chromosome and in this disease the case was the chromosom6 and this region, so from here what using this EST and a small probe they designed and then they screened then they screened a cDNA libraries for that.
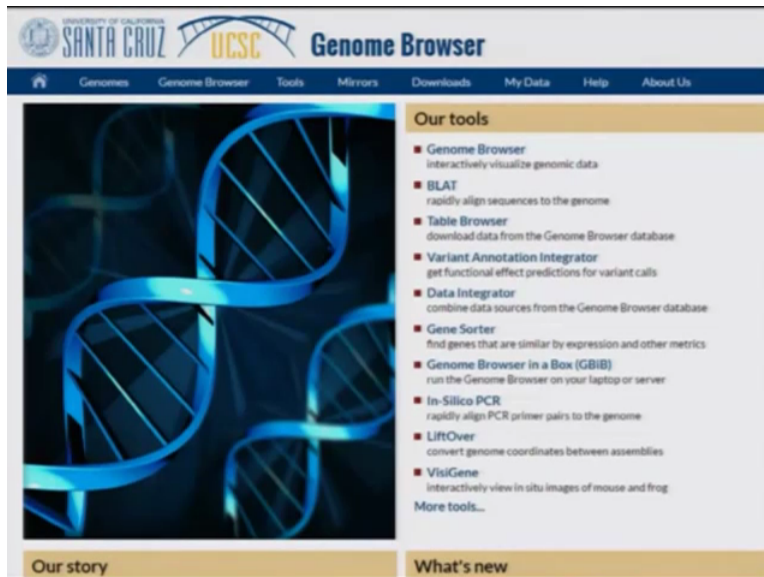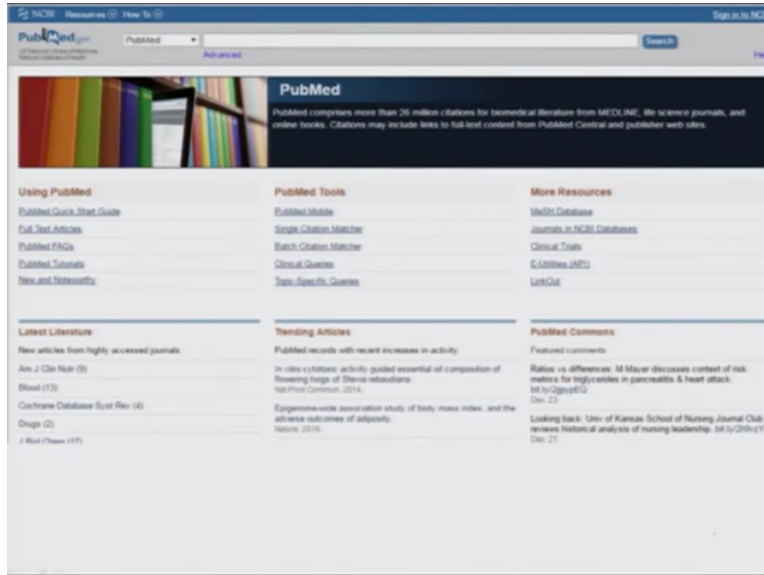
(Refer Slide Time: 3:30)



They try to align them as they screen for many libraries sequence them as sequencing has already being taught and there are various ways by which you can sequence and then you arrive at a

sequence, in this session we will start with the sequence which has been arrived and they have few hundred base pair of the sequence, this is the 5 prime end and this is the 3 prime end of the sequence, so using the current search tools whatever information we can take is we can arrive at or which can be fetched using this few hundred base pair of the sequence that is what we will be learning today.
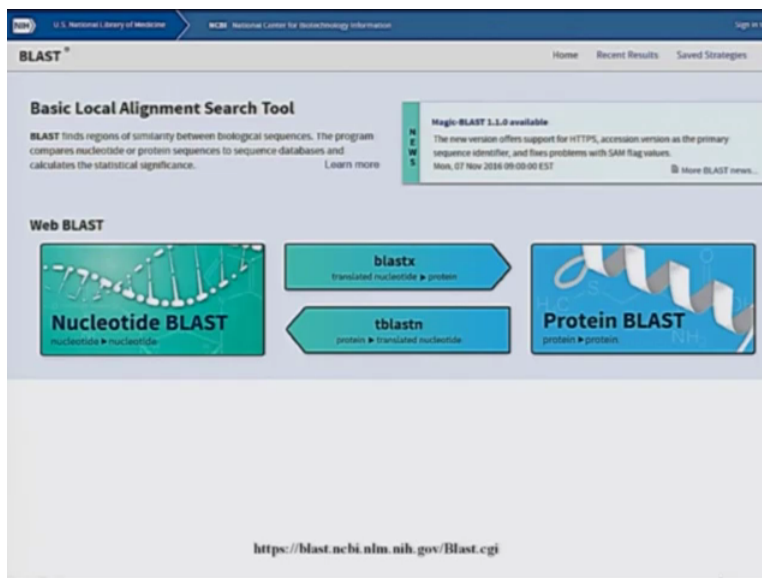
(Refer Slide Time: 4:14)

So using there are few widely used search engines one of those is NCBI another one is UCSC genome browser another one is Ensembl, so today we will be learning about more about NCBI and then we will look online on one of the another search engine that is Ensembl. So since we had a few hundred base pair of sequence and we did not knew from which gene it is or whether this is a coding region or not coding region or this is just a intron of a gene or this a exon of a gene all those information we do not have.

(Refer Slide Time: 5:05)

So one such way is to use blast and here is the page of blast that is that is Basic Local Alignment Search Tool. So what does blast does it finds the region of similarity between biological sequences, this programme it compares the nucleotides or protein sequences to sequence data base and calculate the statistical significance, so they are various blast programmes like one is nuclear type blast, the other one is protein blast. Now you have a stretch of a sequence that can be converted into protein and then again it can be blast to the protein data base that is called blastx.

Similarly the protein can be translated to nucleotide and then it can be blast to the other nucleotide data base and that is called tblastn, so as we have seen that we had a DNA sequence so we will go for nucleotide blast and try to look that what region it is aligning to but before we go into what region it is aligning to very briefly we will look into what does sequence alignment means and what are the various strategies used to align the sequence.

(Refer Slide Time: 6:15)



So there are two key terms which you will encounter more often that is global versus local alignment so very briefly what is global and what is local alignment. Let us look into these two sequences now one you have to align these two sequences one is the one this sequence and the other one is this. Now there are two ways by two ways by which you can align this sequence one is that you just go one by one and try and align for example here if I align this one with this so this would be TACTCA and so on.

And you will try and calculate this core which is based on the various algorithms that what is this core of this alignment, now another way is to that you leave you do not go one by one but you will try and look for small stretch which are aligning the best for example here it would be if you align here with this TACTCACGG and so on you will see that the local stretches are aligning more, so this is called global alignment the first method and the other method is called the local alignment.

Similarly if you have to align these two sequences which are almost of same lengths then it is best to align them globally. So when there is a small stretch of sequence its best to align them local and if there are equal length of equal length of DNA sequences then that can be aligned globally. So on these alignments using different algorithms scores are calculated, few of those are like from where you are starting, what is the gap between the first not aligned and the one where the alignment starts there is a gap penalty then in between there are gap penalty where there are mismatches or you skin few bases and then you align to another stretch.

For example suppose you this is one of this sequence this is quite frank 3 prime and then you have another sequence so it aligns here but then there are some mismatches so you skip that and then again you align. So there would be a certain penalty for this and then that penalty would increase at these regions so on basis of the score is calculated and that is also represented in a blast data.

(Refer Slide Time: 9:13)

So now let us look into how do we go about it, so what happens like this is the blast home page where you have different ways as I have already described blastn, blastp for protein blastx where nucleotide is converted to the protein and then blast to the whole protein data base.

Similarly the nucleotide is translated and then it is blasted to or aligned to to the sequences in the protein data base and so on, so here you can enter your sequence which has that sequence can be entered in a different formats which has accretion numbers which are given by NCBI and that is at the unique suppose you have a gene for which you want to blast, you just need to put that accession number and the tool would take it or you can paste faster sequence.

Now faster sequence is a format of writing the nucleotide or the DNA sequences then on the basis of your of your requirement you can make the... You can use a particular data base suppose you just want to use human the genomics or the transcript or mouse genomics transcript or you can choose from other organism or exclude them add them whatever is your criteria you can do that and then you can just click to blast.

(Refer Slide Time: 10:28)



So in our case this was the query sequence which I have pasted here and then you click on blast just this is the default setting. I have not changed anything and you click on blast.

(Refer Slide Time: 10:45)



So what do you get the results sometimes the software takes time and doing that and this is how the page it comes. This is the job title and this is your query sequence and this is the time you submitted and how much time duration has passed on.

(Refer Slide Time: 11:07)



Then you have a result place this is the blast result page so here are the this is the query the description you had put in nucleic acid which was 421 base pair long and this is the data base through which it has done the search and that is non-redundant data base. Non redundant data

base then it has the description as the nucleotide collections then this is blastn and then you get a distribution of the top 40 blast on the 40 subject sequence, so these are the 40 sequence subjects, subjects sequences against which these are the best scored ones and the score has been here shown in the format of colour. So red shows the highest one and the black one black shows the lowest one and this is your query sequence which was 421 base pair long.

(Refer Slide Time: 12:05)



Now in this you can always get your mouse there and you can find that what that sequence is what is the alignment, what is the accession number and what and what is the E values score which I will be explaining in a short duration.

(Refer Slide Time: 12:24)



So this is our description page where we have seen that our query sequence has aligned to these many sequences and there were many more however through the print screen I have got only this one this much. But when you will be doing it online you can scroll down and see that there were many more sequences from various species to which your queries sequence has aligned so what could be the maximum score the total score and how much the query has been covered is given here.

Now what is E value so E value defines the number of hits which you can expect by change, so that here it is zero which means that all the base pairs of your query were aligning hundred percent and it was not a chance event and here you have accession numbers by which you can access this particular transcript which is here it is like EPM2A which is associated with I suppose lafora disease and the protein is laforan and it is a transcript variant X6 mRNA and you can access these by this accession ID.

So once you have that this accession ID then you can also look into how your query sequence was aligning to the subject sequence or to this sequence which the search engine has searched for. So here you see that your this is the query and this is the subject sequence. Your query first nucleotide is aligning with with the 2393 or 2393 base pair of the subject sequence and from there on it is aligning to till 2813 base pair.

Now from here you can access the sequence what do you find that this the first hit which the search engine has given, it is a protein EPM2A or laforan EPM2A or Laforan and it is associated with a progressive myoclonus type of epilepsy. The below you have all the information about this entry in whatever organisms it is found it is found in the (())(14:49) cod data then there are commons about it. What genome annotation data it has been provided for example from NCBI the annotation status and then as you will scroll down you will get lots of information about that.

(Refer Slide Time: 15:23)



If you click on the graphics of this you can see the alignment in the form of a graph here this is your query sequence and this is EPM2A and how they are aligning. You can always explore going to these arrows or you can use the tools and tracks and whatever you can explore all this online.

(Refer Slide Time: 15:41)



Similarly then there you can go for the FASTA sequence and you will find the sequence of the gene to which your query has aligned at the most, from here as you can see on the right hand side there are various ways to analyse these sequences for example you can run a blast in which this now this sequence of the subject which you have arrived at that is EPM2A gene or lafora can be used to used to run a blast against all the nucleotide sequence.

Similarly you can use primers to amplify this particular gene you can highlight the sequence speeches what are the characteristics sequence speeches and you can also find something in this sequence suppose you are looking for a particular stretch of a sequence you can use this tool to identify, similarly there are articles about EPM2A gene as you are trying to find out the gene for the linked disease so you have come so far so you can always look into what are the recent articles related to with gene, what are the pathways in which this genes is involved like for example glycogen , synthesis or glucose metabolism.

You can always from here link to the phenotype of the disease is this are theses pathways defective in those patience in which this chromosomal region is deleted. Do you see something abrupt with glycogen metabolism in those patients? So with this you can link the information in the tool.

(Refer Slide Time: 17:20)





When you run blast using that query now again this is just taken the accession number. As I have already described that anything you can use as a query sequence in the NCBI blast and then you can again using the default as I have done here you can run blast again the search engine takes time and thus the page appears on your screen and then this is the way the result is again.

(Refer Slide Time: 17:54)



So this is the similar result as I have described so here these are top 109 blast hits on the hundred subject sequences and this is how it is aligning mostly are red almost all are red and these are the stretch of sequences now these are the stretch of sequences which were not aligning and the last part was aligning.

(Refer Slide Time: 18:14)



So once you have the sequence or the all the sequence which are aligning to that gene that was EPM2A, EPM2A gene and these are the genes which are now aligning to that gene. You can
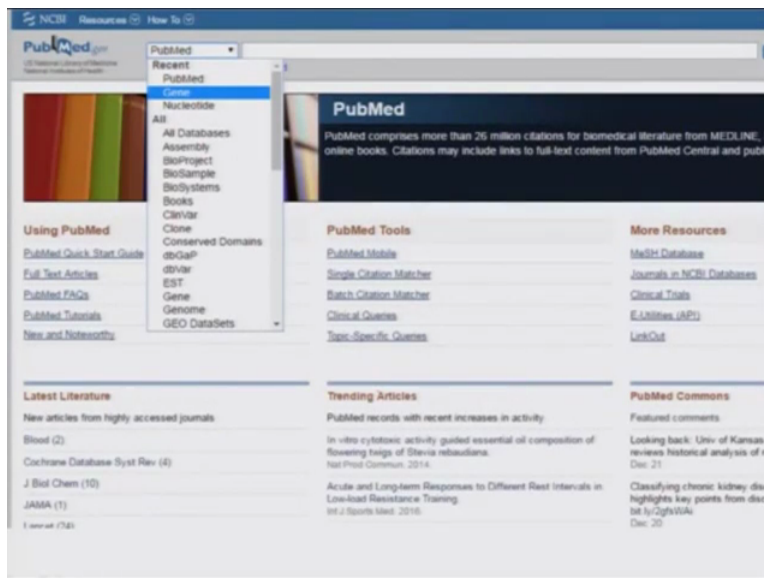
select all clicking on all and then you can go for distance tree results so from here you can get the phylogenetic tree of the gene of your interest.

(Refer Slide Time: 18:46)



So this tree shows how many how the gene has evolved and how close to one or the other organisms genes is how it is. So for example primates are more closer to rodents and so on, so once you have the sequence.

(Refer Slide Time: 19:03)

So you can also access to the protein of the sequence you can go to PubMed and you can go on gene.
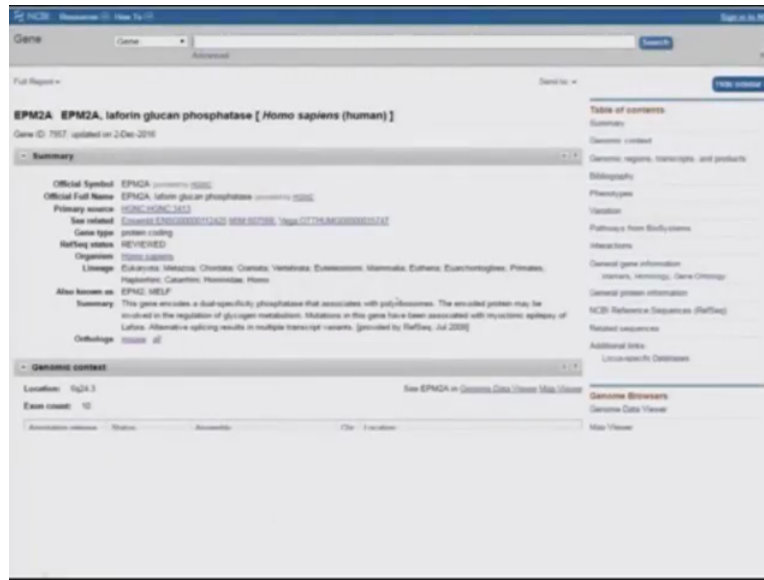
(Refer Slide Time: 19:16)
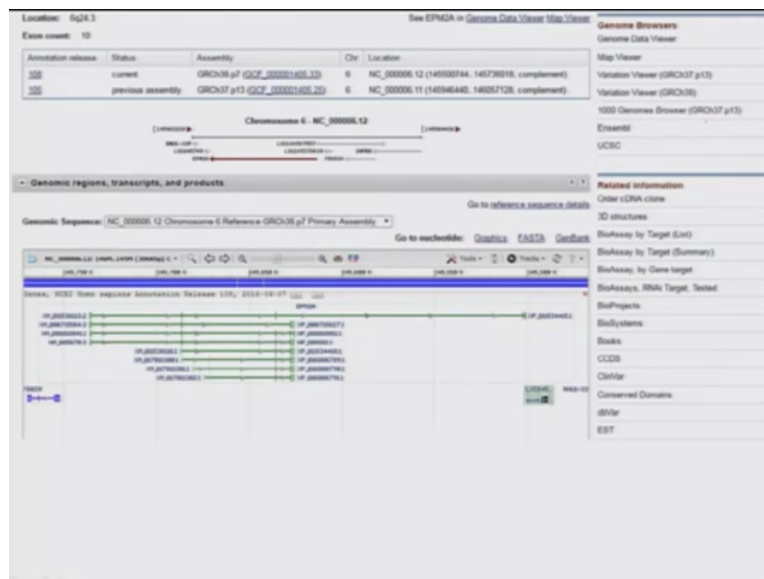


From here you will write the gene which you have identified that is EPM2A and then you can select whatever organisms organisms you want for example you want the gene which is there sequenced in the humans or in the mouse or in the rat here we our query was about a sequence which we have identified from a part of a chromosome which was deleted in the family of of the affected people.

(Refer Slide Time: 19:45)



So we will go for humans and then you can have all the information about this gene in the humans from what is the lineage what is the summary what does this gene does what are the orthologs present.

(Refer Slide Time: 20:02)



Similarly from the on the right hand side you have all the other information then you can scroll down and look into what are the other cDNA clones, 3D structures and everything. What are the EST associated with this gene and so on, now here you see on the screen that this is showing
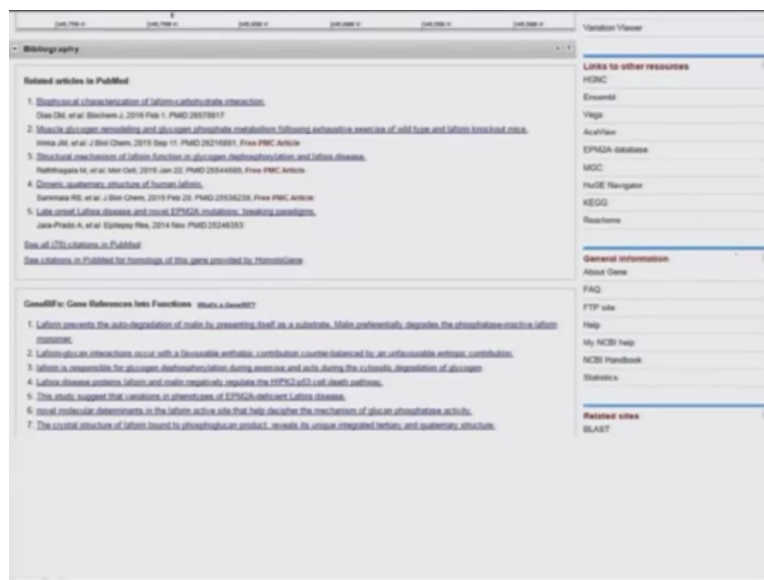
what are the transcripts present so there are I, 2 3, 4, 5, 6, 7 and 8, 8 transcripts for this gene that is EPM2A.
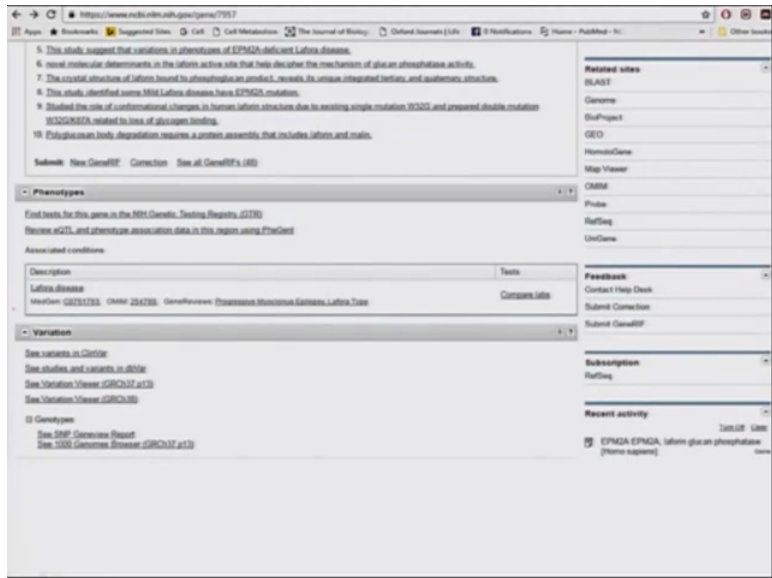
(Refer Slide Time: 20:38)



Then on again you can look for OMIM which is Online Mendelian Inheritance of Man that is this is another search tool by you can search for EPM2A and you can see that to what all disease it is associated with.
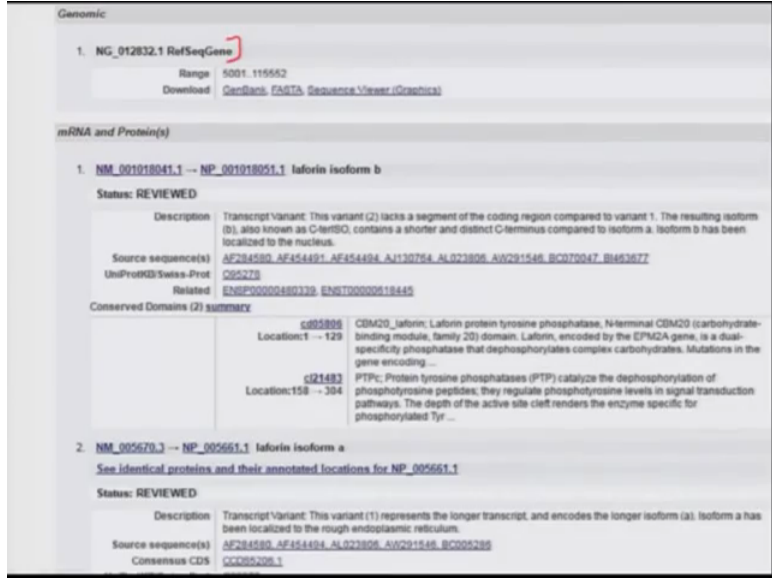
(Refer Slide Time: 21:10)

Then again you have all the related articles and then gene reference and there is a huge amount of information already available on the internet which can help you arrive if all kind of information you are looking at but you need to know what you are what do you what exactly you want to search.

(Refer Slide Time: 21:21)



Similarly when you will scroll down you wills see that there are the reference sequences of the gene you can go to the chromosomal region from which this gene was coded and there you can have all the information about introns, exons what are the sequences as the human chromosome

or the human genome has been sequenced aligned and deposited on internet for our ease. Similarly there is a protein sequence and mRNA sequence NM generally is for the mRNA sequence and NP is mainly for proteins but sometimes you will also see accession IDs with xp and all, these are for the proteins.

(Refer Slide Time: 22:17)



So from here you can go to the protein sequence which this gene is coding for and then you have all the information about the protein data base that is you can have the conserve domain known you can highlight some sequence feature you can use the FASTA sequence and you can run a blast as in against all the protein data base and so on.

(Refer Slide Time: 22:38)



So with this we have the FASTA sequence and we can run blast this protein is aligning to what all species of the protein sequences it is aligning to the protein sequences of what all species.

(Refer Slide Time: 22:47)

Similarly you can identify conserve domains. So if you click here you will get to know what are the conserved domains or motives in this proteins sequence for example in this protein which was being coded by EPM2A it has CBM20, it is carbohydrate binding domain. So as the function of EPM2A was in glycogen metabolism and in glycogen synthesis path way so it will goes with the the functions which it does that it would bind to the starch or the glycogen, similarly there is another region that is PTPC super family. What is PTPC super family? It is a protein tyrosine phosphatase which catalyses the de-phosphorylation of phospho tyrosine peptide.

So this protein which was being encoded by EPM2A it has two domains one is CBM20 which has ability to bind two glycogen like molecules and similarly it has an another domain that is PTPC which is for the which is a phosphatase domain by that we mean that it can de-phosphorylate the phosphate group from the peptides.

(Refer Slide Time: 24:08)



Similarly you can highlight the sequence features of this peptide sequence.

(Refer Slide Time: 24:14)



And this is the sequence feature which is highlighted which shows that it is the iso form B is encoded by transcript variant 2 and so on.

Similarly you can run a blast and what do you find that you find that again this blast would take NP as a query sequence and using the default settings of the search engine when we run blast. This page appears showing that it is running the the job which we have assigned to it and then you get the results again in a similar fashion. How well your query sequence is aligning to the subjects sequence what are the top hits what are organisms it is aligning all those information we can get from here.

Now you have all the sequences which has aligned to your query sequence and from there you can select all and then again you can go to see the distance tree results which will bring the phylogenetic tree for this protein along the various phylums of the various texonomic groups along the various texonomic groups. For example frog, toads or turtle or insectivorous rabbits and so on, so you can see that how your protein of interest has evolved through this stages of evolution and how conserved it is.

(Refer Slide Time: 25:50)



You can also select few sequences for example here I selected the proteins sequence from the humans then monkey then mouse and chicken and then you can see how well it is aligning and you can also identify what all regions are what all regions are more conserved.

(Refer Slide Time: 26:09)



So you see that in these four groups which I which I have selected mainly most of it is aligning and only like suppose the one of the humans is not aligning to the other three that is monkey, mouse or and the chicken. The few last base pairs are not aligning that well.

This is how this huge amount of information we just got from few hundred base pair of the DNA sequence so what all we arrived at that we had few hundred base pairs few hundred base pair of the DNA sequence. From there we arrived to a gene then we looked into what are the features of that gene, we can look into that. What is the current literature about that we looked into in what all species it is present that is by using blast. We can also look into the how conserved that is that is by looking into the phylogenetic tree in and what all organisms 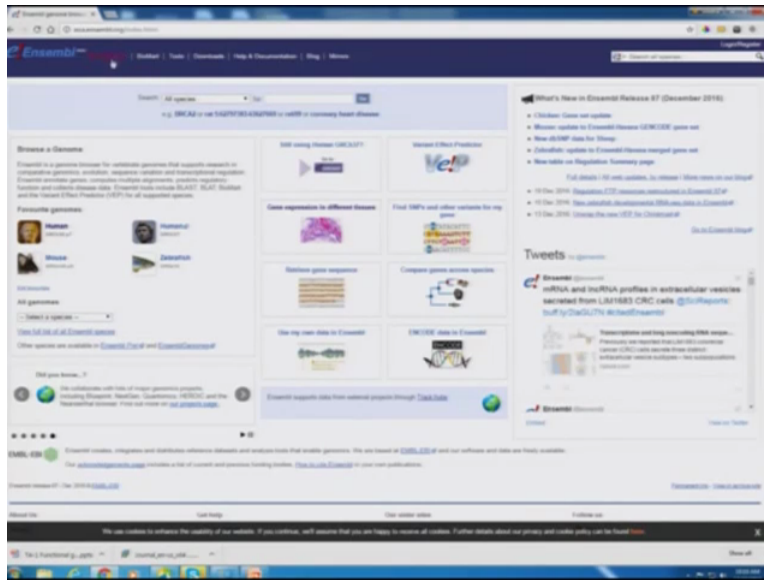it is there. From that gene you can look go to the protein sequence, you can see into what all motives or conserved motives it has how well it is conserved in the other species and so on. What is the current literature about that protein you can look into the structure of that protein and many more.

So this is the power of Bio-informatics and computer science which has been given to the biologists that when you have a small sequence this all you can search it against all the data base present. Now we will look into another search engine that is Ensembl search engine and very quickly we will go we have almost understood what all we can do but let us look into what we can do. So now we will look into another search engine that is Ensembl and we will go online and look into what all information can be fished from that search engine.

(Refer Slide Time: 28:32)



So this is the Ensembl home page and you can do here blast or what are other tools are available similarly you can have your favourite genome selected if you want the gene to be only looked into humans or in zebra-fish or mouse information, similarly you can go for gene expression in different tissues how you gene expresses in different tissues of your organism of your interest, you can compare the gene across the species which we have also done with the help of NCBI using the phylogenetic tree which we created selecting all the sequences.

You can use your own data or you can also upload the data which you have created so depending on your requirement you can search the available search engines for the information and you can explore whatever you wished to. There are various tutorials already available on the YouTube which we can use to do so that is all thank you I hope this session has been somewhat informative for you and which will help you in using these search engines that is all. Thank you.