# Using BLAST To Compare the PDC gene sequence across different Saccharomyces cerevisiae strains



This is the standard view of the BLASTn. It is divided into three sections, namely: "Enter Query Sequence," "Choose Search Set", and" Program Selection."
In the section "Enter Query Sequence," we have to Enter the input query i.e.- the sequence for which we want to perform the BLAST. Here we can either paste the sequence directly or give the accession number of that sequence and provide a subrange for that sequence.
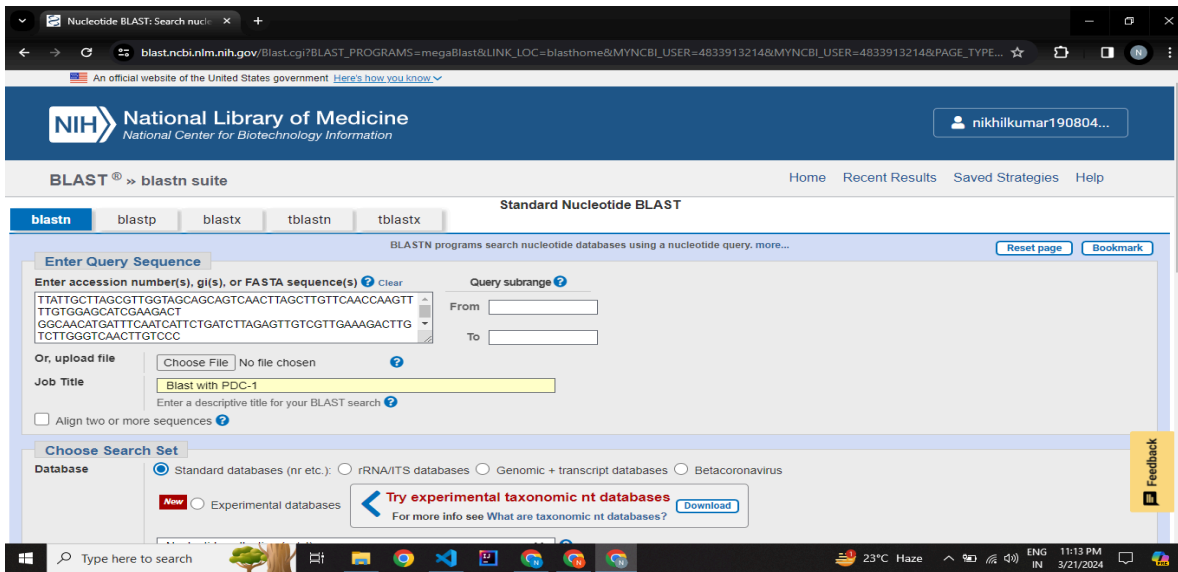


Now, in the section "Choose Search Set" we can select the database we want to use for BLAST. By default, it is selected as "Nucleotide" and standard database (non-redundant). Also here, we have the option to add an organism or to exclude an organism from the database. And we also have the option to filter our database using an Entrez Query.
And in the last section, "Program Selection" we have some optimization options to enhance our searches.
And lastly we can change some algorithm parameters according to our needs.

# BLAST With PDC-1



Now, here i am performing BLAST with input as the PDC-1 gene. I pasted the sequence of the PDC-1 gene, which I downloaded using Python, and gave a suitable Job Title to this.
The accession number for PDC-1 is NC_001144.5 and the sequence starts from 234081 to 232390. It is in reverse complement form.



Here I selected the database as Standard database, which is a by default option, and selected the Nucleotide collection (nr/nt) which again was by default. After that i added the Organism "Saccharomyces cerevisiae," and I selected the top among them with taxid as 4932.
After that, in section "Program Selection" I kept the options by default, and in algorithm parameters I kept everything by default only.
And finally, i clicked on the button for BLAST.
It took some time to complete this request, and after some time, the result page was shown.

# Result of BLAST



This is the result page of the BLAST. On the left side, it provides the details related to the search like Job Title, a unique ID for every search, the BLAST program that was selected (blastn.blastp,tblastx etc.), the database which was selected, and the query length which was 1692 bases long for the PDC-1.
On the right side, it gives the option to filter out the results based on some parameters like percent identity and E value or by using an organism name.



Now, here we can clearly see the BLAST of PDC-1 against various Saccharomyces cerevisiae strains. Along with this we also have their alignments with PDC-1 and with percent identity, E value, scores, etc.
There were more than 600 strains displayed on keeping the number of records to a maximum allowed value in NCBI BLAST. Hence, this shows the BLAST of Saccharomyces cerevisiae strains with the PDC-1 gene.
A similar approach can be used to perform the BLAST with every PDC gene (I performed this in my code).

Now for further SNP analysis, I selected some of the Saccharomyces cerevisiae strains and wrote a Python code to find the SNPs between the PDC gene and strain sequence.

I selected these two strains: Saccharomyces cerevisiae YJM996 and Saccharomyces cerevisiae YJM1199. The script I wrote prints the position of SNP change along with the base being changed.