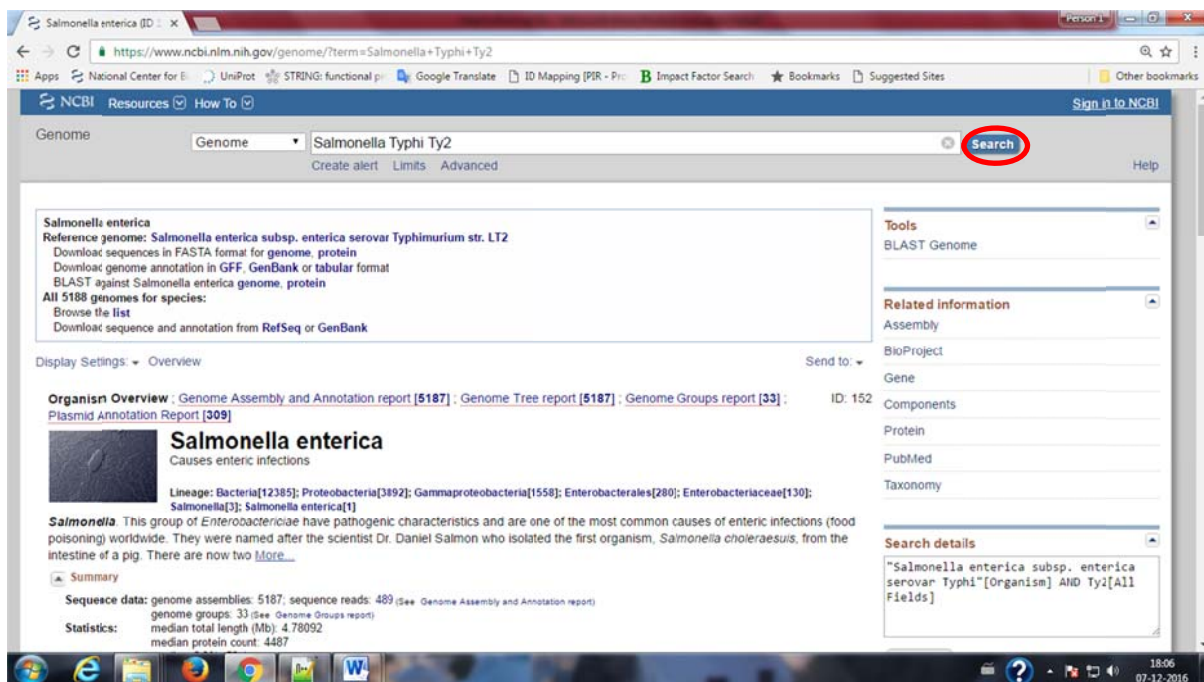# Instruction for parsing protein coding table from coding sequence (CDS)

As we have mentioned in HelpForPredictingsRNAsByBarmanetal.pdf (Step2) that in some cases direct protein coding table is not available in NCBI, so we need to parse all CDS from NCBI in such cases.

**Ex :** If you search *Salmonella* Typhi (*S.* Typhi) Ty2 in NCBI genome section, you will find protein coding table of *S.* Typhi CT18. Therefore we have to find all CDS of *S.* Typhi Ty2. You can find all CDS of *S.* Typhi Ty2 from NCBI nucleotide section. Therefore you have to copy the all CDS related information and paste in notepad or notepad++ or editplus and save in .txt format.

Finally we have developed ParsingCDsFromCompleteGenome.r code to parse all the forward coding start and stop positions of *S.* Typhi Ty2. Then you have to create similar protein coding table mention in **Step4** of HelpForPredictingsRNAsByBarmanetal.pdf by using start and stop positions of all CDS. The other field including Protein name, GeneID, Protein product and Length can be filled by "NA" since they are not playing any role to finding intergenic regions of *S.* Typhi Ty2.

Please see the below screenshots wisely.

RStudio

File   Edit   Code   View   Plots   Session   Build   Debug   Tools   Help

ParsingCDsFromCompleteGenome.r

```
1   # this functions helps us remove previously created workspace variables....
2   rm(list=ls());
3
4   # this library used to handle string related functions.......
5   library("stringr");
6
7   # use to get complete list of all CDs file in .txt file from user
8   inputTextfile <- readline("please enter the complete list of all CDs file in .txt file format: ");
9
10  #Read text data with all lines
11  readTextFile <- paste0(readLines(inputTextfile));
12  readTextFile
13
14  # This section will find all the line numbers with keyword "CDS" and return the particular line for CDS
15  allCDs <- grep("CDS",readTextFile,value = TRUE);
```

Console   C:/Users/ranjan/Desktop/sRNAWork@STyphi/UserInstructionsForPredictingsRNAs/SourceCodesForParsing/

```
R version 3.1.0 (2014-04-10) -- "Spring Dance"
Copyright (c) 2014 The R Foundation for Statistical Computing
Platform: i386-w64-mingw32/i386 (32-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> source('C:/Users/ranjan/Desktop/sRNAWork@STyphi/UserInstructionsForPredictingsRNAs/SourceCodesForParsing/ParsingCDsFromCompleteGenome.r')
please enter the complete list of all CDs file in .txt file format: STty2CDS.txt
please enter the output file name for CDS start and end position of coding part in .csv file format: StartAndStopPositionsofAllCDs.csv
```

Environment    History

Values

Forwa.chr  [1:2249...
Forwa.chr  [1:2249...
Forwa.chr  [1:2249...
allCD.chr  [1:4324...
input."STty2CDS.t...
readT.Large chara...

---

Excel spreadsheet: StartAndStopPositionsOfAllCDs

| | A | B |
|---|---|---|
| 1 | 190 | 255 |
| 2 | 337 | 2799 |
| 3 | 2801 | 3730 |
| 4 | 3734 | 5020 |
| 5 | 7665 | 8618 |
| 6 | 8729 | 9319 |
| 7 | 11594 | 13510 |
| 8 | 13596 | 14735 |
| 9 | 15020 | 15967 |
| 10 | 16094 | 16438 |
| 11 | 17873 | 19972 |
| 12 | 20004 | 21155 |
| 13 | 21198 | 23060 |
| 14 | 23341 | 24045 |
| 15 | 24473 | 25015 |
| 16 | 25116 | 25802 |
| 17 | 28429 | 29436 |
| 18 | 29437 | 29982 |
| 19 | 29998 | 30516 |
| 20 | 30509 | 31213 |
| 21 | 31278 | 32123 |
| 22 | 32162 | 32449 |
| 23 | 33368 | 34363 |
| 24 | 35339 | 37057 |
| 25 | 40128 | 41621 |

| | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| 1 | Protein name | GeneID | Protein product | Length | Start | Stop | Strand |
| 2 | NA | NA | NA | NA | 190 | 255 | + |
| 3 | NA | NA | NA | NA | 337 | 2799 | + |
| 4 | NA | NA | NA | NA | 2801 | 3730 | + |
| 5 | NA | NA | NA | NA | 3734 | 5020 | + |
| 6 | NA | NA | NA | NA | 7665 | 8618 | + |
| 7 | NA | NA | NA | NA | 8729 | 9319 | + |
| 8 | NA | NA | NA | NA | 11594 | 13510 | + |
| 9 | NA | NA | NA | NA | 13596 | 14735 | + |
| 10 | NA | NA | NA | NA | 15020 | 15967 | + |
| 11 | NA | NA | NA | NA | 16094 | 16438 | + |
| 12 | NA | NA | NA | NA | 17873 | 19972 | + |
| 13 | NA | NA | NA | NA | 20004 | 21155 | + |
| 14 | NA | NA | NA | NA | 21198 | 23060 | + |
| 15 | NA | NA | NA | NA | 23341 | 24045 | + |
| 16 | NA | NA | NA | NA | 24473 | 25015 | + |
| 17 | NA | NA | NA | NA | 25116 | 25802 | + |
| 18 | NA | NA | NA | NA | 28429 | 29436 | + |
| 19 | NA | NA | NA | NA | 29437 | 29982 | + |
| 20 | NA | NA | NA | NA | 29998 | 30516 | + |
| 21 | NA | NA | NA | NA | 30509 | 31213 | + |
| 22 | NA | NA | NA | NA | 31278 | 32123 | + |
| 23 | NA | NA | NA | NA | 32162 | 32449 | + |
| 24 | NA | NA | NA | NA | 33368 | 34363 | + |
| 25 | NA | NA | NA | NA | 35339 | 37057 | + |

ForwardCodingPartOfSTty2

Rest of the steps are same as mention in HelpForPredictingsRNAsByBarmanetal.pdf.