# Instructions for predicting bacterial small non-coding RNAs (sRNAs) from complete genome

**Step1:** Please download complete genome of your bacteria of interest (.fasta extension) from NCBI nucleotide section.

> **Ex:** We have downloaded complete genome of *Escherichia coli* (*E. coli*) K-12 from https://www.ncbi.nlm.nih.gov/nuccore/556503834?report=fasta. Please see the below screenshots.

**Step2:** Please download protein coding table of your bacteria of interest from NCBI genome section. In some cases direct protein coding table is not available in NCBI. Therefore you have to parse the coding table from all coding sequence (CDS) for your bacteria of interest. Please see the HelpForParsingCDS.pdf for parsing all CDS from NCBI.

> **Ex:** We have downloaded protein coding table of *E. coli* K-12 from https://www.ncbi.nlm.nih.gov/genome/proteins/167?genome_assembly_id=161521&gi=556503834. Please see the below screenshots.

**Step3:** Please copy all the tab delimitated .txt file information and paste it in Microsoft excel file (.xlsx). Consider only Protein name, GeneID, Protein product, Length, Start, Stop and Strand columns.

**Ex:** Please see the below screenshots.





**Step4:** Please consider only forward strand since the proposed method has been developed on forward strand specific. Filter the "Strand" column by "+" sign only.

**Ex:** Please see the below screenshots.

**Step5:** We have developed all the codes in R programing language. Therefore we request you to please download latest version of R and install it in your local desktop. We have

downloaded latest version of R (3.3.2) from https://cran.r-project.org/bin/windows/base/. You can also find R for Linux and (Mac) OS X from https://cran.r-project.org/.
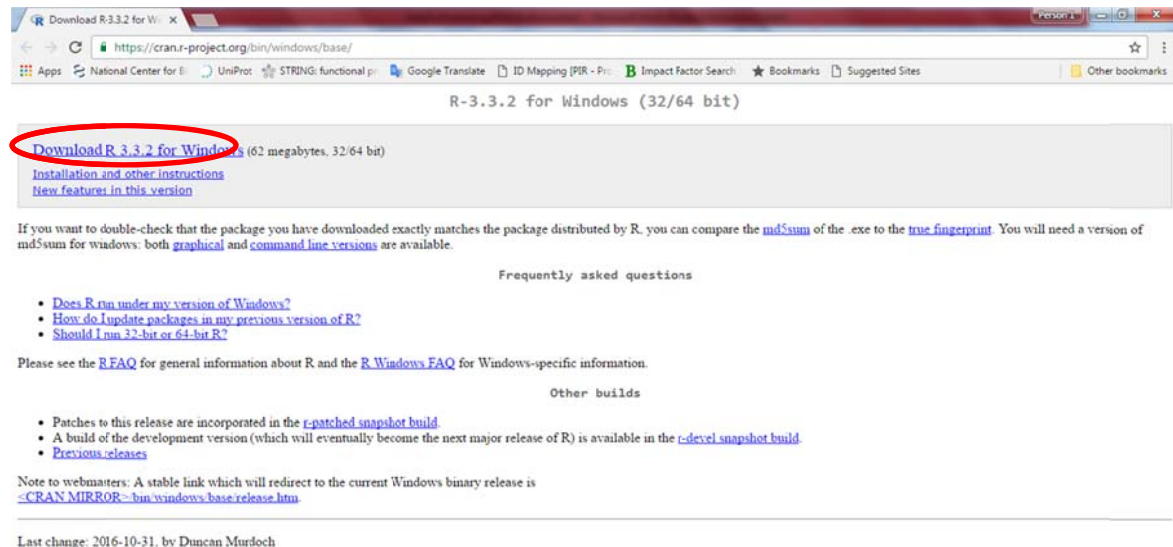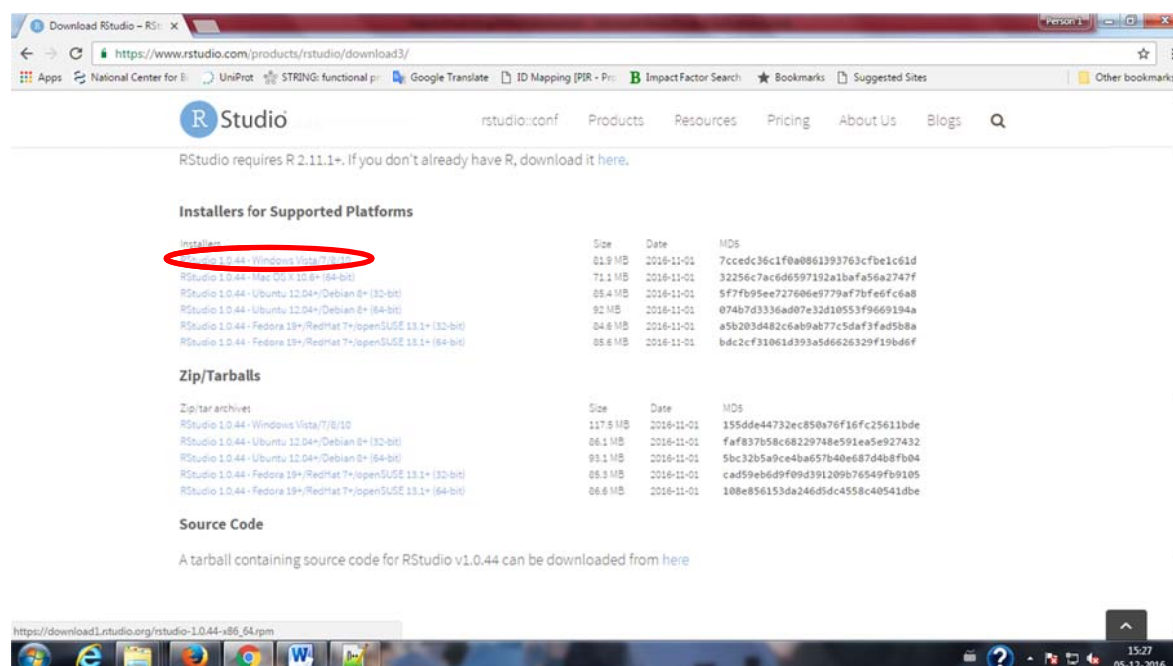
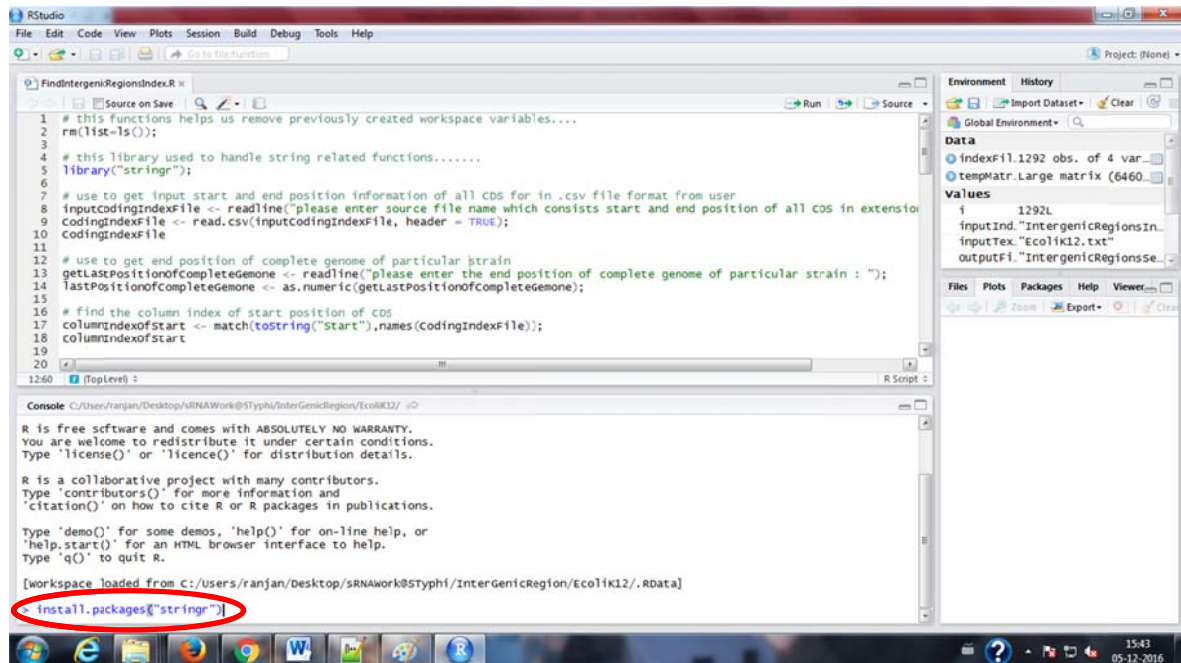**Ex:** Please see the below screenshot.



**Step6:** We always prefer to write code in integrated development environment (IDE). RStudio is a free and open-source IDE for R. You can download RStudio from https://www.rstudio.com/products/rstudio/download3/.

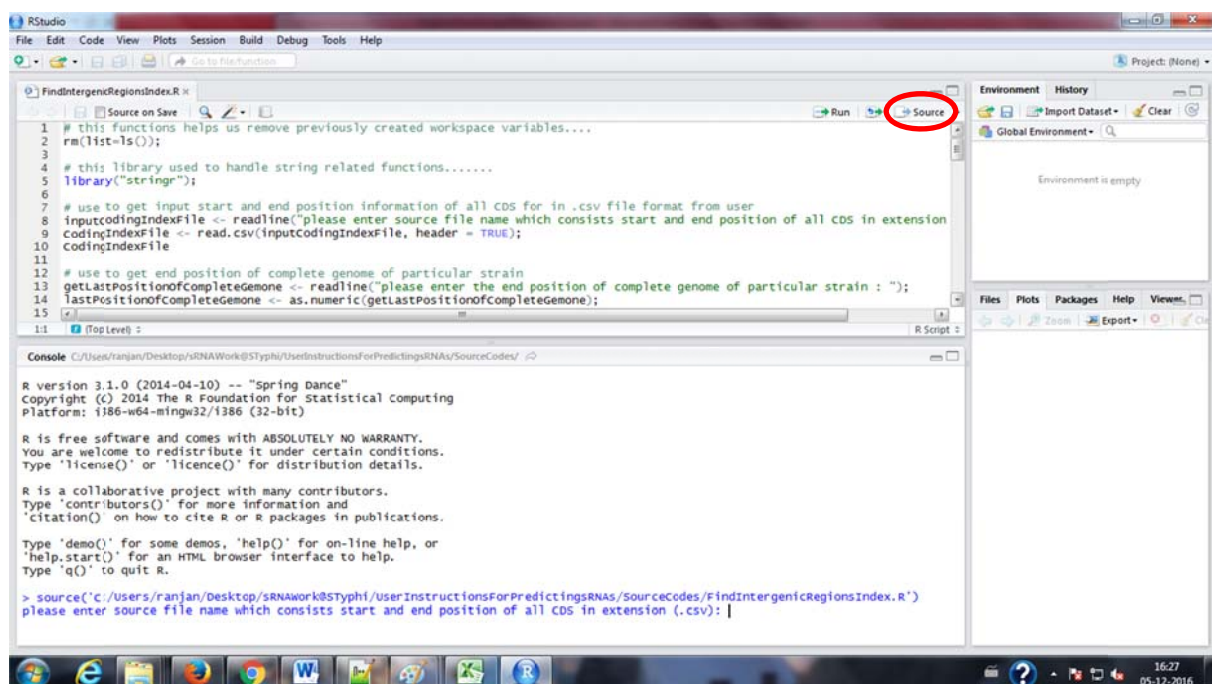**Ex:** Please see the below screenshot.

**Step7:** We have used R package 'stringr' for handling string operation. Therefore we request you to please install R package 'stringr' from console of Rstudio by using install.packages("stringr") comment.

  **Ex:** Please see the below screenshot.



**Step8:** We have developed FindIntergenicRegionsIndex.r code to find all the intergenic regions of a particular strain. This code will accept input as protein coding table (created in **Step4**) and end position of complete genome of particular strain. This code will produce all the intergenic index of particular strain.

  **Ex:** Please see the below screenshots.

**Step9:** We have considered only those intergenic regions having length greater than 50nts. You can filter intergenic regions by length greater than 50.

    **Ex:** Please see the below screenshots.

**Step10:** Next we have developed FindIntergenicRegionSequencesFromCompleteGenome.r code to find corresponding sequence of intergenic regions. This code will accept input as complete genome sequence (created in **Step1**) and index file of start and stop position of intergenic regions (created in **Step9**). This code will produce sequence corresponding to the intergenic region. Please note that this code will take few minutes since it will operate on complete genome.

**Ex:** Please see the below screenshots.

**Step11:** Following this we have developed ExtractWindowFromIntergenicRegions.r code to extract window from intergenic region. This code will accept input as intergenic region sequence file (created in **Step10**), window size and step size of the window. This code will produce single and multiple windows based on intergenic region length, window size and step size.

 **Ex:** Please see the below screenshots.

**Step12:** Afterwards we have developed IntergenicRegionsSequenceFrequencyCalculator.r code to calculate mono (4), di (16) and tri (64) nucleotide composition features for each window of intergenic region. This code will accept input as an intergenic region sequence file (created in **Step10**) and produce sequence composition (84) for all windows. Please note that this code will take few hours for calculating composition features for each window of each intergenic region. The execution time is directly proportional to the number of intergenic region and length of the intergenic region.

**Ex:** Please see the below screenshots.

**Step13:** Next we have developed GenerateSVMDataForIntergenicWindow.r code to generate SVM format data for each window of each intergenic region. This code will take input as intergenic region sequence file (created in **Step10**) and decision (1 and 0 indicates Yes and No respectively) of generating mono, di and tri nucleotide features. This code will produce data in SVM format. We found that trinucleotide composition features performed well, therefore we have generated only tri nucleotide features in this step.

    **Ex:** Please see the below screenshots.

**Step14:** Please download our proposed best model (ModelTriNucleotideCompositionOneistoTwo7_7_4_3) and svm_classify.exe files from SourceCodes. You can also download SVM^light ( svm_learn.exe and svm_classify.exe)from http://svmlight.joachims.org/ (under Source Code and Binaries).

**Step15:** Finally we have developed PredictIntergenicRegionWindows.r code to predict each window of each intergenic region using proposed best model. This code will take only input as an intergenic region sequence file (created in **Step10**). This code will produce prediction score of each window of each intergenic region.

**Ex:** Please see the below screenshots.