# Documentation for 'Analysis Tools' developed at Rishi Biotech, Mumbai

-Nayana Ramachandran

**Introduction**

At Rishi Biotech, Mumbai, we have developed five sequence analysis tools, namely, FormatSeq, ORFreader, TranslORF, InfoProt and RestrDigest. These different tools are integrated and placed on a common platform headed under 'Analysis Tools.'

This documentation will enable a user to utilise the tools to the fullest extent thereby taking advantage of this free edition of Analysis Tools.
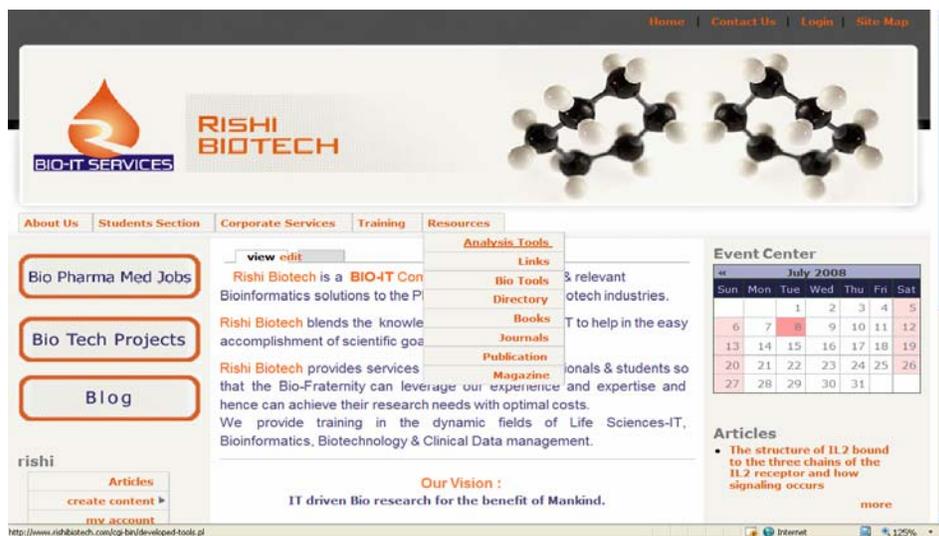
It is advised that proper formats for submitting the nucleotide or protein sequence be followed until further version of 'Analysis Tools' are developed which will take all the standard input formats.

The organisation of this document is based on the five tools and their working.

## A] Analysis Tools

We, at Rishi Biotech, have developed and designed five tools for sequence analysis. The tools can be freely accessed through the home page for 'Analysis Tools' at http://www.rishibiotech.com/cgi-bin/developed-tools.pl

Alternatively, any user can visit the Rishi Biotech home page at http://www.rishibiotech.com   and go to the 'Resources' section and access the link titled 'Analysis Tools.'



On clicking the link, you will be redirected to an integrated tools page titled 'Analysis Tools.'

There is brief description about the various tools given on this page. Scroll to the bottom of the page. You will find a text-area like this:



It is here that you have to input your nucleotide or protein sequence for all the tools that you want to use. After inputting the sequence you have to choose a tool and click on 'Submit' button.

You will be redirected to an intermediary page where you have a chance to change your input sequence, if required. The tool that you have selected will appear in bold red font. If you want to change the tool, simply navigate backwards through the 'back' button on your browser and you will see that your sequence remains intact. Now, you can change your tool, if you wish.
Confirm your actions by clicking on 'Submit.'

# Confirm your request

Check your input sequence and modify if you wish to.

```
ORIGIN
        1 ccctgcactt gggagccggt agcactccta tcactgcttc tcaacccgtg agctaccagc
       61 tgtgtcatga gctgcagaca gttctcctcg tcctacttga gccgcagcgg cgggggtggc
      121 gggggcggcc tgggcagcgg gggcagcata aggtcttcct acagccgctt cagctcctca
      181 gggggcggtg gaggagggggg ccgattcagc tcttctagtg gctatggtgg gggaagctct
      241 cgtgtctgtg ggagggagg cggtggcagt tttggctaca gctacggcgg aggatctggg
      301 ggtggtttta gtgccagtag tttaggcggt ggctttgggg gtggttccag aggtttttggt
      361 ggtgcttctg gaggaggcta tagtagttct gggggtttg gaggtggctt tggtggtggt
      421 tctggaggtg gctttggtgg tggctatggg agtgggttg gggggtttgg gggctttgga
      481 ggtggtgctg gaggaggtga tggtggtatt ctgactgcta atgagaagag caccatgcag
      541 gaactcaatt ctcggctggc ctcttacttg gataaggtgc aggctctaga ggaggccaac
      601 aacgacctgg agaataagat ccaggattgg tacgacaaga agggacctgc tgctatccag
      661 aagaactact cccttatta taacactatt gatgatctca aggaccagat tgtggacctg
      721 acagtgggca acaacaaaac tctcctggac attgacaaca ctcgcatgac actggatgac
      781 ttcaggataa agtttgagat ggagcaaaac ctgcgggcaag gagtggatgc tgacatcaat
```

The chosen Analysis Tool is: **'FormatSeq'**

Submit

---

Home      Analysis Tools      Information / Contact

# FormatSeq

This is a nucleotide sequence formatting tool.

This tool formats your output, i.e. removes all the spaces, gaps and unknown bases from your input nucleotide sequence. (Note: Any 'U's in the sequence will be converted to 'T').

This tool renders your sequence fit to apply to many tools found over the World Wide Web. However, the condition is that you have to submit an unformatted sequence, in either of these formats:

1. GenBank: full-text starting with 'LOCUS' or sequence information starting with '1' or 'ORIGIN.'
2. EMBL: full-text starting with 'ID' or sequence information starting with 'SQ' or the sequence itself.

The output displays various information like: Formatted sequence, base constituents, their percentages, AT: GC and GC:AT ratio, which can be analysed for comparing between two sequences.

# ORFreader

ORFreader scans the input nucleotide sequence for the Open Reading Frames in all the six frames. The ORF's start with a start codon ATG (or AUG in RNA) and end with one of the three stop codons: TAA, TGA or TAG.

The tool takes nucleotide input as mentioned for FormatSeq as well as an additional FASTA input. The information is displayed in two levels:

1. Complete ORF information which segregates the ORF's frame-wise for all the six frames.
2. Partial ORF information which segregates the ORF with respect to its length, in six levels, of or above 50, 100, 150, 200, 250 and 300 base pairs.

This enables the user to view selectively the ORF's and use them for further analysis.



The positions and the length of ORF are sufficient for the user to analyze the sequence, especially genomic sequences.

In further versions of ORFreader, we propose to include a graphical interface to view the ORF's.

# TranslORF

This tool is an extension of ORFreader and translates the given sequences wherever there is an ORF.

The input format is same as above, namely, GenBank, EMBL or FASTA nucleotide sequence.

TranslORF considers only those ORF's which are equal to or more than 50 base pairs in length as only these can form stable peptides or proteins. The output is organized frame-wise with all the lengths of ORF equal to or above 50 bp. The details of output are: frame information, ORF start, end and length information, the translated product and the peptide length.

```
--- For frame -3: ---

ORF Start: 149
ORF End: 93
ORF Length: 57 bp
Translated product:   *KLRLPPPPPPPRPLPPLM
Peptide length: 18


ORF Start: 536
ORF End: 444
ORF Length: 93 bp
Translated product:   *PLPNPPNPPKPPPAPPPSPPIRVALSFLVM
Peptide length: 30


ORF Start: 872
ORF End: 678
ORF Length: 195 bp
Translated product:   *LVISSRLSWITSRVTPLLLVRRSMSLVRMVSSSKLIFNSISCFRRCPTSASMLPRRCTSSLRVM
Peptide length: 64


ORF Start: 1043
ORF End: 957
ORF Length: 87 bp
Translated product:   *SVPCFLPSTLTSIFTAGPLSRVLVRLSM
Peptide length: 28


ORF Start: 2087
ORF End: 1995
ORF Length: 93 bp
Translated product:   *LRQLSRRGLGLQADGSPCVGSSAQKNSGSM
Peptide length: 30


--- For frame -2: ---

ORF Start: 1824
```

In future versions of TranslORF we plan to include gene organization information into introns and exons.

# InfoProt

This is currently the sole protein sequence information tool at Rishi Biotech. This tool requires protein sequence as input in the following formats:

1. Swiss-Prot sequence format (consisting of sequence, gaps and numbers)
2. FASTA protein sequence format

The output is a text-cum-tabulated format displaying sequence information, number of amino acids, the peptide constituents in terms number and percentages of amino acids, atoms and species.
At the end, there are two formulae assigned to the protein:
   a. Based on atomic information
   b. Based on amino acid information

| Asparagine | Asn | N | 10 | 5.15 |
|---|---|---|---|---|
| Glutamine | Gln | Q | 5 | 2.58 |
| Aspartate | Asp | D | 5 | 2.58 |
| Glutamate | Glu | E | 14 | 7.22 |
| Lysine | Lys | K | 19 | 9.79 |
| Histidine | His | H | 5 | 2.58 |
| Arginine | Arg | R | 12 | 6.19 |
| Grand Total | | | 194 | 100.00 |

****** Amino acid constituents of protein based on species (number and percentage) ******

| Group | Amino acids | Number | Percentage |
|---|---|---|---|
| Negatively charged | $D_i$ | 19 | 9.79 |
| Positively charged | $R_i$ | 31 | 15.98 |
| Total charged | G, A, V, L, I, F, W, C, M, S, T, P, Q | 139 | 71.65 |
| - | G, A V, L, | 61 | 31.44 |
| - | F, Y, | 18 | 9.28 |
| -Sulphur-containing | $C_i$ | 15 | 7.73 |
| -Oxy- | $S_i$ | 25 | 12.89 |
| - | P | 5 | 2.58 |
| -Acidic | $N_i$ | 15 | 7.73 |
| Grand Total | | 194 | 100.00 |

****** Atomic constituents of protein based on type (number and percentage) ******

| Atom | Symbol | Number | Percentage |
|---|---|---|---|
| Carbon | C | 988 | 31.44 |
| Hydrogen | H | 1579 | 50.25 |
| Oxygen | O | 283 | 9.01 |
| Nitrogen | N | 277 | 8.82 |
| Sulphur | S | 15 | 0.48 |
| Grand Total | | 3142 | 100 |

****** Formulae ******

| Based on | Formula for protein |
|---|---|
| Atoms | $C_{988} H_{1579} O_{283} N_{277} S_{15}$ |
| Amino Acid | $F_5 S_{11} T_{14} N_{10} K_{16} E_{14} Y_{10} V_{11} Q_5 M_4 C_7 L_{15} A_{11} W_3 P_5 H_5 D_5 R_{12} I_{12} G_{12}$ |

Done

This is a very useful tool which can be used to predict restriction digestion by restriction endonucleases. If you choose this tool on 'Analysis Tools' page, on the confirmatory page you will have to select the restriction enzyme with which you want to predict the digestion.

The user can select from a range of twenty enzymes available (in future versions, all enzymes will be made available) and proceed with 'Submit.'

The result page will display the fragments created due to restriction digestion and fragment information like the positions, length, point of cleavage, etc.

If there is no restriction digestion, only a single fragment will be displayed to the user.



Thus, the user can carry out in-silico restriction digestion with any given sequence. The user can also navigate back and forth and try different restriction enzymes on the same sequence.

**CONCLUSION**

The analysis tools have been designed using Perl 5.8.8 as the programming language.
They are under constant upgradation and development.
We plan to release the future versions of 'Analysis Tools' by improving on the existing tools and adding new tools to the website.


If any queries or suggestions regarding the tools, please e-mail to:
info@rishibiotech.com