

## Data retrieval and download

### 10.1 Downloading a set of gene results and associated data.

There are several ways to obtain gene results and associated data. You can download this information via:

- (1) gene results section by clicking on Download Genes function;
- (2) Sequence Retrieval tool in the Tools Section on the right side of the Home page;
- (3) Downloads menu at the top and Data files submenu.

For this exercise you can start with any result list you generated this morning, or pick a strategy from the list below:

<http://microsporidiadb.org/micro/im.do?s=0ae5b5bd79a40271>

Download this list of genes with the following associated data: Genomic Location, Product Description, Transcript Length and Predicted GO Function. Hint: click on the Download ## Genes link.

The screenshot shows the Microsporidiadb web interface. At the top, there are navigation tabs: "My Strategies: New Opened (1) All (1) Basket Public Strategies (3) Help". Below this, a search strategy is defined: "Strategy: Genes with BAMH1 site up- and downstream but not within". The strategy is visualized as a flowchart with four steps: Step 1 (DNA Motif 24728), Step 2 (Organism 45129 Genes), Step 3 (DNA Motif 24728 Segments), and Step 4 (Organism 45129 Genes). An "Add Step" button is visible. Below the flowchart is an "Expanded View of Step Organism" showing a diagram of a gene with a DNA motif and an organism. At the bottom, a table of 245 genes is displayed. The table has columns for "All Results", "Ortholog Groups", and various taxonomic groups: "Eumetazoa", "Eucnidia", "E. acido", "E. cucullif. (nr Genes: 37)", "E. hellem (nr Genes: 21)", "E. imbricata", "E. romaleae", "E. dieneusi", "N. ganisi (nr Genes: 2)", "N. sp. 1 (nr Genes: 6)", "N. bombycis", and "N. ceranae". The "All Results" column shows 245 genes. A red circle highlights the "Download 245 Genes" link in the top right corner of the table area.

All Results	Ortholog Groups	E. acido	E. cucullif. (nr Genes: 37)	E. hellem (nr Genes: 21)	E. imbricata	E. romaleae	E. dieneusi	N. ganisi (nr Genes: 2)	N. sp. 1 (nr Genes: 6)	N. bombycis	N. ceranae					
245	116	0	35	32	34	18	15	23	21	12	3	1	3	3	10	0

Hint: select the Tab delimited type of report to download and then click on the boxes to customize your report. The gene ID is automatically downloaded and so is not an option in the popup.

- **Tab delimited (Excel): choose from columns** – create a file with one row per gene and unlimited columns per gene. Any data that is available as a column on the result page can be downloaded with this option.
  - **Tab-delimited** text, also known as **tab-separated** values (TSV), is a format that can be created or viewed by most spreadsheet programs and text editors. The TSV format follows these rules: Each entry in the **file** takes up a single line. The first line in the **file** is the header line, which labels each field.
- **FASTA (sequence retrieval, configurable)** – create a multi-fasta file of your sequences. Each sequence begins with a single-line description, which contains greater-than (“>”) symbol, followed by lines of sequence data. You have the option to configure the start and end points of the sequence
- **GFF3: Gene models and optional sequences** –a simple **tab delimited** format for describing genomic features in a 9-column text file. GFF stands for *Generic Feature Format*. GFF3 allows multi-level grouping and multi-level descriptive attributes.
- **Text: choose from columns and/or tables** – create a text file of data associated with your sequences. This option allows you to download data that has multiple associations per gene, such as multiple GO terms assigned to one gene. The file structure is NOT one row per gene.
- **XML: choose from columns and/or tables** – create an xml file of columns or tables of data for your sequences. **XML** (extensive mark-up language) stores data in plain text format and is both human- and machine-readable format. This is also a software- and hardware- independent way of storing, transporting, and sharing data. XML is commonly used for the interchange of data over the Internet.

- **json** (JavaScript Object Notation): choose from columns and/or tables – Create a json formatted file that can be used as an alternative to xml. It is a way to store information in an organized, easy-to-access manner. In a nutshell, it gives us a human-readable collection of data that we can access in a really logical manner.

## 10.2 Download the sequences of genes in a list of results.

What if you are interested in examining the 5' flanking sequences of these genes? How can you easily get these sequences for subsequent analysis? What kind of sequences can you retrieve? Protein? Genomic? Coding?

Hint: use same list of results as in 10.1. Choose Download ### Genes again but this time choose **FASTA (sequence retrieval, configurable)**.

Now, retrieve the 500 nucleotides upstream of the start site of your genes.

**Download 245 Genes from the search:**  
*Combine Gene results*

Please select a format from the dropdown list to create the download report.

\*\*Note: IDs will automatically be included in the report and the report will be sorted by ID.

**This reporter will retrieve the sequences of the genes in your result.**

Choose the type of sequence:  genomic  protein  CDS  transcript

Choose the region of the sequence(s):

begin at   nucleotides

end at   nucleotides

Download Type:  Save to File  Show in Browser

\*\* Note: If UTRs have not been annotated for a gene, then choosing "transcription start" may have the same effect as choosing "translation start".

**Help**

The diagram illustrates the structure of a gene. It shows a 5' UTR (Untranslated Region) followed by an exon, an intron, another exon, and a 3' UTR. The 3' UTR ends with a stop codon and a polyA tail. Below the gene structure, the CDS (Coding Sequence) is shown as a pink bar, and the protein sequence is shown as a blue bar. The CDS starts at the ATG start codon and ends at the stop codon.

Tools:

**BLAST**  
Identify Sequence Similarities

**Sequence Retrieval**  
Retrieve Specific Sequences using IDs and coordinates

**PubMed and Entrez**  
View the Latest Pubmed and Entrez Results

Now retrieve 5' UTR sequences for the list of genes. Begin with setting a transcription start parameter at 0 and end at a translation site (ATG) and parameter (-1). Setting a translation start (ATG) site parameter to (-1) would eliminate incorporating "A" of the start codon into the 5' UTR sequence.

## 10.3 Use the Sequence Retrieval Tool to download the genomic sequence for your genes.

Note that you can download sequence with the sequence retrieval tool (SRT) accessed from the tools menu on the **home page**:

- Retrieve Sequences By Gene IDs (Gene IDs can be accessed from gene result page and the uploaded to the retrieve sequence tool that can be accessed from the home page).

- Retrieve Sequences By Genomic Sequence IDs (Retrieve Genomic Sequence IDs by accessing Home > Genomic Sequence > Organism > Nosema > Get answer).
- Retrieve Multiple Sequence Alignments by Contig / Genomic Sequence IDs (follow instructions above, then select a contig of interest and follow prompts to Retrieve contig sequence or perform multiple sequence alignment in the Sequence section).
- Retrieve Sequences By Open Reading Frame IDs (hint: you can find ORFs using gene IDs).

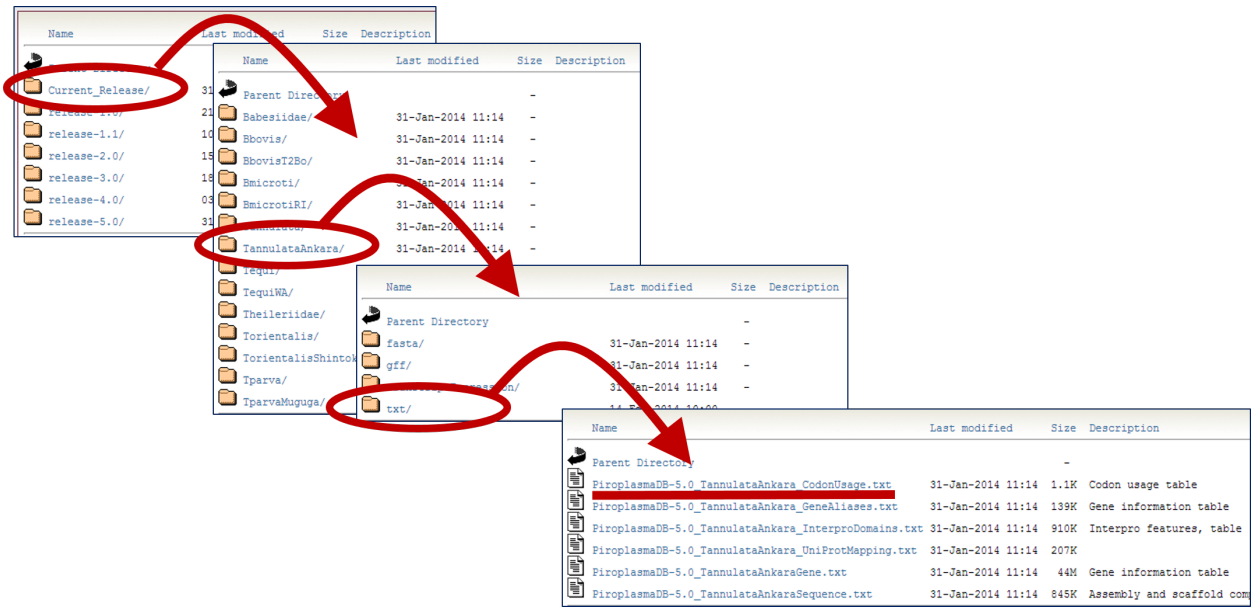
#### 10.4 Downloading large data files such as all coding sequences or all protein sequences for an entire genome.

For this exercise use any EuPathDB site. The example below illustrates a use case in PiroplasmaDB: <http://piroplasmadb.org>

Files are available from the Download section of all EuPathDB sites  
Hint: select “Data Files” under the “Download” menu in the grey tool bar.



Hint: navigate through the subfolders and find the txt files containing codon usage information for *T. annulata* Ankara. Folders without a strain designation contain species level data.



What other data are available for download? Do the directories make sense ... fasta, gff, transcript Expression, txt? Is there any data in the transcript Expression folder for *T. annulata*? Look at the Transcript Expression searches to determine which of the organisms have this data type.