


# Configure Data for Metagenomics Classification

UGENE integrates a set of tools for taxonomy classification of microorganisms using whole-genome shotgun sequencing data. The tools are Kraken, CLARK, DIAMOND, etc. See, for example, "[Parallel NGS Reads Classification](#)" sample workflow that allows one to classify input FASTQ files with these tools working in parallel and then join their output with another tool WEVOTE.


 The tools are available on 64-bit macOS or Linux operating systems only. Also, as the tools are quite resource-consuming, it is recommended to have at least 16 Gb of RAM available.

To use these tools one should provide appropriate taxonomy data and reference data, specific for a tool. Some reference databases are provided for each of the tool. One can use these data or build a custom database (see, for example, workflow element "[Build Kraken Database](#)").

It is recommended to use the UGENE [Online Installer](#) package to install and automatically configure the data. However, if the Internet is not available on the target computer, or it is required to use another UGENE package for some other reason, follow the instructions below on how to download and configure the data.

## Download data for metagenomics classification

Use links in the "Data for NGS metagenomics classification" section on the web page <http://ugene.net/download-all.html> to download the data.

 Make sure to have enough disk space on the target computer.

See the list of the available downloads in the table below.

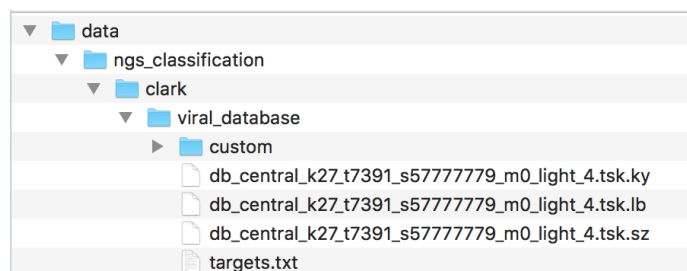
Data	Archive size	Unpacked data size	Description	Data source
NCBI taxonomy classification	2.5 Gb	31 Gb	This includes a set of <a href="#">taxonomy data</a> files from NCBI. <u>These data should be present for any type the NGS classification analysis.</u>	Original data were downloaded from the NCBI FTP ( <a href="ftp://ftp.ncbi.nlm.nih.gov/pub/taxonomy/">ftp://ftp.ncbi.nlm.nih.gov/pub/taxonomy/</a> ).
NCBI RefSeq bacterial genomes	130 Gb	132 Gb	The data can be used to build a database for CLARK-I (light version of CLARK), CLARK, or Kraken.  As UGENE integrates <a href="#">modified version</a> of CLARK/CLARK-I, it is possible to provide *.gz archives as input for building the database. In particular, "CLARK-I DB: RefSeq bacterial+viral genomes" (see below) was generated using the archived data.  Also, keep in mind that changing of some parameters of the " <a href="#">Classify Sequences with CLARK</a> " element may cause re-building of the reference database. The reference data should be present in this case!  For building a Kraken database usage of *.gz archives is not supported, it is required to unpack each *.gz file, so even more disk space will be required.	Original data were downloaded from the NCBI FTP ( <a href="ftp://ftp.ncbi.nlm.nih.gov/refseq/release/bacteria/bacteria*.genomic.fna.gz">ftp://ftp.ncbi.nlm.nih.gov/refseq/release/bacteria/bacteria*.genomic.fna.gz</a> )
NCBI RefSeq viral genomes	77 Mb	77 Mb	Similarly to "NCBI RefSeq bacterial genomes", although the size of the data is rather small.  The reference data are included into "CLARK-I DB: RefSeq bacterial+viral genomes" and "CLARK-I DB: RefSeq viral genomes" databases.	Original data were downloaded from the NCBI FTP ( <a href="ftp://ftp.ncbi.nlm.nih.gov/refseq/release/viral/viral*.genomic.fna.gz">ftp://ftp.ncbi.nlm.nih.gov/refseq/release/viral/viral*.genomic.fna.gz</a> ).
NCBI RefSeq GRCh38 human genome	837 Mb	838 Mb	Similarly to "NCBI RefSeq bacterial genomes".  The data are not included into any database, but provided in case one would like to use them when building a custom database.	Original data were downloaded from the NCBI FTP ( <a href="ftp://ftp.ncbi.nih.gov/genomes/H_sapiens/CHR_*/hs_ref_GRC*chr*.fa.gz">ftp://ftp.ncbi.nih.gov/genomes/H_sapiens/CHR_*/hs_ref_GRC*chr*.fa.gz</a> ).
Kraken DB: MiniKraken 4Gb database	2.5 Gb	4.3 Gb	A sample reference database provided in UGENE for Kraken.  It is a pre-built 4 GB database constructed from complete bacterial, archaeal, and viral genomes in RefSeq (as of Oct. 18, 2017). This can be used by users without the computational resources needed to build a Kraken database. However this contains only 2.7% of kmers from the original database.	Original data were downloaded using a link on the Kraken web site ( <a href="https://ccb.jhu.edu/software/kraken/">https://ccb.jhu.edu/software/kraken/</a> ).
CLARK-I DB: RefSeq bacterial+viral genomes	7.4 Gb	11 Gb	One of the reference databases provided in UGENE for CLARK-I.  The database was build using archived RefSeq bacterial and viral genomes.	See above.

<b>CLARK-I DB: RefSeq viral genomes</b>	16 Mb	72 Mb	One of the reference databases provided in UGENE for CLARK-I.  The database was build using archived RefSeq viral genomes.	See above.
<b>DIAMOND DB: UniRef50</b>	5.2 Gb	13 Gb	One of the reference databases provided in UGENE for DIAMOND.  Note that unlike Kraken and CLARK, DIAMOND requires protein reference sequences as input.	Original data were downloaded from the Uniprot FTP ( <a href="http://ftp.uniprot.org/pub/databases/uniprot/current_release/uniref/uniref50/uniref50.fasta.gz">http://ftp.uniprot.org/pub/databases/uniprot/current_release/uniref/uniref50/uniref50.fasta.gz</a> ). Then a DIAMOND database was built.
<b>DIAMOND DB: UniRef90</b>	13 Gb	34 Gb	One of the reference databases provided in UGENE for DIAMOND.	Original data were downloaded from the Uniprot FTP ( <a href="http://ftp.uniprot.org/pub/databases/uniprot/current_release/uniref/uniref90/uniref90.fasta.gz">http://ftp.uniprot.org/pub/databases/uniprot/current_release/uniref/uniref90/uniref90.fasta.gz</a> ). Then a DIAMOND database was built.
<b>MetaPhlAn 2 embedded DB: mpa_v20_m200</b>	1 Gb	1.2 Gb	This database is provided with the MetaPhlAn2 tool. It was built on ~1M unique clade-specific marker genes identified from ~17,000 reference genomes (~13,500 bacterial and archaeal, ~3,500 viral, and ~110 eukaryotic).	The database was downloaded here: <a href="https://bitbucket.org/biobakery/metaphlan2/downloads/">https://bitbucket.org/biobakery/metaphlan2/downloads/</a> .
<b>Total:</b>	<b>162 Gb</b>	<b>227 Gb</b>		

## Configure data

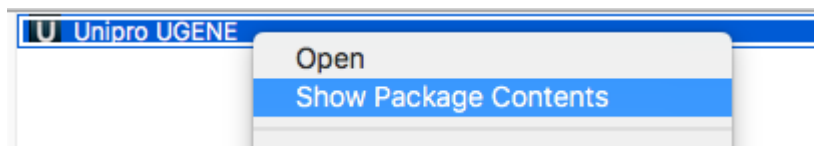
Data described above are stored as 7zip archives. After a file download, unpack it using an appropriate file archiver (for example, [Keka](#) on macOS).

The unpacked data are stored in a folders structure with the root folder called "data". For example, for "NCBI RefSeq viral genomes" the archive is called "ngs\_classification.clark.viral\_database.7z" and the unpacked data look as follows:



It is required to move these data to the UGENE data folder, following the hierarchical data structure:

- On Linux it is "data" folder, located in the UGENE installation folder.
- On macOS the "data" folder is located inside the "Unipro UGENE.app" bundle. Right-click on the bundle, select "Show Package Contents", select "Contents -> MacOS -> data" folder.



Thus, all required data will be placed to the "ngs\_classification" sub-folder of the UGENE data folder.



Kraken, CLARK, DIAMOND and WEVOTE are integrated as [external tools](#). So, also make sure the tools executables are set in the UGENE Application Settings.