



ELSEVIER

ScienceDirect

SoftwareX | (■■■■) ■■■-■■■

SoftwareX

[www.elsevier.com/locate/softx](http://www.elsevier.com/locate/softx)

# NCBI Mass Sequence Downloader – Large dataset downloading made easy

F. Pina-Martins<sup>a,b,\*</sup>, O.S. Paulo<sup>a</sup>

<sup>a</sup> *Computational Biology and Population Genomics Group, Centre for Ecology, Evolution and Environmental Changes, Departamento de Biologia Animal, Faculdade de Ciências, Universidade de Lisboa, Campo Grande, 1749-016 Lisboa, Portugal*

<sup>b</sup> *Departamento de Biologia e CESAM, Univ. de Aveiro, Portugal*

Received 15 October 2015; received in revised form 6 April 2016; accepted 28 April 2016

## Abstract

Sequence databases, such as NCBI, are a very important resource in many areas of science. Downloading small amounts of sequences to local storage can easily be performed using any recent web browser, but downloading tens of thousands of sequences is not as simple.

*NCBI Mass Sequence Downloader* is an open source program aimed at simplifying obtaining large amounts of sequence data from NCBI databases to local storage. It is written in python (can be run under both python 2 and python 3), and uses PyQt5 for the GUI. The program can be run in either graphical or command line mode.

Source code is licensed under the GPLv3, and is supported on Linux, Windows and Mac OSX. Available at <https://github.com/ElsevierSoftwareX/SOFTX-D-15-00072.git>, [https://github.com/StuntsPT/NCBI\\_Mass\\_Downloader](https://github.com/StuntsPT/NCBI_Mass_Downloader)

© 2016 The Author(s). Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

**Keywords:** Bioinformatics; Genomics; Molecular evolution; BLAST

## Code metadata

Current Code version	3.1
Permanent link to code / repository used of this code version	<a href="https://github.com/ElsevierSoftwareX/SOFTX-D-15-00072.git">https://github.com/ElsevierSoftwareX/SOFTX-D-15-00072.git</a>
Legal Code License	GPLv3; Biopython license
Code Versioning system used	git
Software Code Language used	Python; PyQt5
Compilation requirements, Operating environments & dependencies	Python and (optionally for the GUI) PyQt5
If available Link to developer documentation / manual	<a href="https://github.com/ElsevierSoftwareX/SOFTX-D-15-00072.git">https://github.com/ElsevierSoftwareX/SOFTX-D-15-00072.git</a>
Support email for questions	<a href="mailto:f.pinamartins@gmail.com">f.pinamartins@gmail.com</a> / <a href="mailto:frmartins@ciencias.ulisboa.pt">frmartins@ciencias.ulisboa.pt</a>

## Software Metadata

Current software version	3.1
Permanent link to executables of this version	<a href="https://github.com/ElsevierSoftwareX/SOFTX-D-15-00072.git">https://github.com/ElsevierSoftwareX/SOFTX-D-15-00072.git</a>
Legal Software License	GPLv3; Biopython license
Computing platform / Operating System	GNU/Linux, Microsoft Windows, Mac OSX, any other where python and (optionally for the GUI) PyQt5 are available.
Installation requirements & dependencies	Python and (optionally for the GUI) PyQt5
If available Link to user manual — if formally published include a reference to the publication in the reference list	<a href="https://github.com/ElsevierSoftwareX/SOFTX-D-15-00072.git">https://github.com/ElsevierSoftwareX/SOFTX-D-15-00072.git</a>
Support email for questions	<a href="mailto:f.pinamartins@gmail.com">f.pinamartins@gmail.com</a> / <a href="mailto:frmartins@ciencias.ulisboa.pt">frmartins@ciencias.ulisboa.pt</a>

\* Corresponding author at: Computational Biology and Population Genomics Group, Centre for Ecology, Evolution and Environmental Changes, Departamento de Biologia Animal, Faculdade de Ciências, Universidade de Lisboa, Campo Grande, 1749-016 Lisboa, Portugal.

E-mail address: [f.pinamartins@gmail.com](mailto:f.pinamartins@gmail.com) (F. Pina-Martins).

## 1. Introduction

National Center for Biotechnology Information (NCBI) sequence databases are nowadays a resource of unquestionable

<http://dx.doi.org/10.1016/j.softx.2016.04.007>

2352-7110/© 2016 The Author(s). Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

Please cite this article in press as: Pina-Martins F, Paulo OS. *NCBI Mass Sequence Downloader – Large dataset downloading made easy*. SoftwareX (2016), <http://dx.doi.org/10.1016/j.softx.2016.04.007>

importance for researchers in many areas of science [1]. The current count of sequences available in this database as of 15 December 2015 ascends to over  $18.9 \times 10^7$  sequences and  $20.3 \times 10^{10}$  base pairs (<ftp://ftp.ncbi.nih.gov/genbank/gbrel.txt>), representing roughly 2.5 Tb of compressed data and keeps growing. Due to advances in sequencing technology, the amount of sequence data required by investigators has increased by orders of magnitude in the last few years. This has naturally led to increased use of the NCBI databases by investigators for retrieving sequence data.

*NCBI Mass Sequence Downloader* provides a user friendly interface and automated error checking for downloading large sets of sequence data.

## 2. Problems and background

Although downloading sequences from NCBI can be done in a simple fashion using any standards compliant web browser via the *Entrez* [2] web portal, this method does not scale well, and downloading large amounts of sequences (in the order of the tenths of thousand) from these databases can cause problems when performed this way (<https://www.biostars.org/p/43970/>). Furthermore, manually performing this type of tasks is time consuming and error prone, which may hamper the reproducibility of scientific work.

For retrieving large sets of data, NCBI provides the *E-utilities* API [2], although it can be difficult to use by investigators without an IT related background, despite it's through documentation. Frameworks exist, written in various popular languages, such as *Python*, *Perl* or *Ruby*, that provide some level of abstraction for using this API, such as *Biopython* [3], *BioPerl* [4], or *BioRuby* [5], respectively. However, these too, require some degree of programming knowledge to use, rather than providing end-user packages ready to use for a specific purpose. This leaves investigators without a simple, ready-made solution. Although this is not much of a problem for someone with a bioinformatics background, it poses a serious issue for someone with a molecular biology background, who may frequently require this kind of data, but lack the programming skills to use one of the mentioned frameworks or the API.

By using NCBI's API, our program intends to solve the problem of retrieving large datasets, in a user friendly, automated, and reproducible way. The tool is therefore aimed at molecular biologists that do not have an IT related background, but need to download large datasets from the NCBI databases.

## 3. Software framework

### 3.1. Software architecture

*NCBI Mass Sequence Downloader* is written in python (<http://www.python.org>) and can be run under both python 2 and python 3. The command line interface (CLI) version of the program can be run on any OS that has python available. The GUI version further requires PyQt5 (<http://www.riverbankcomputing.com/software/pyqt/intro>) available, which means all major currently used operating systems such as

GNU/Linux, MS Windows and Mac OSX are supported. The program uses a slightly altered (changed the *import* statements) module from *Biopython* [3] to *Entrez* [2], which is included with the software. This avoids the need to have *Biopython* installed, which, despite being a popular library in the bioinformatics community, is not usually so for molecular biologists. The consequence of this convenience for the user is a higher maintenance requirement since it makes it necessary to keep up with the upstream *Entrez* module. However, this *Biopython* module has not had many recent changes, and merging them into *NCBI Mass Sequence Downloader* has so far been trivial.

The program consists of essentially three modules — a back-end, a front-end and the *Entrez* module. If the program is run without arguments, the GUI version is launched (Fig. 1), but if the program is run with arguments, the command line version will be run instead. This makes the program quite flexible to use in different environments.

The program's source code is available on github (<https://github.com/ElsevierSoftwareX/SOFTX-D-15-00072.git>), along with binaries for GNU/Linux, MS Windows and Apple OSX.

### 3.2. Software functionalities and limitations

*NCBI Mass Sequence Downloader* is made to solve a single task — downloading sets of sequences from the NCBI databases. For this, the user should provide an email address for eventual contact from NCBI (which is sent only to NCBI), the database to be queried, the search query, and a path to the file for the downloaded sequences. Download progress is indicated in both user interfaces.

Currently, the program is limited to downloading sequences in the FASTA format and to NCBI databases, but data from several databases can be retrieved: *nucleotide*, *nucore*, *nucgss*, *protein*, *genome* and *popset*.

### 3.3. Internal routines and error handling

Once the program is requested to start the download, it queries the selected NCBI database for the inputted search term. It will then store the returned sequence IDs in memory, and begin downloading the respective records in batches of 3000 sequences. Every batch is temporarily stored in memory, and once 3000 sequences are downloaded, they are immediately stored in the output file, flushed from memory, and then, the next batch is processed.

After all the records are downloaded, the output file is parsed and its sequences' IDs matched to the originally retrieved sequence IDs. If any sequences are missing, a new pass is made to retrieve them. This process is repeated as often as necessary until all the requested sequences are stored in the output file. Stopping the program at any time will not affect any sequences already stored in the output file.

If a pre-existing FASTA file is selected as the output file, instead of overwriting it, the file is parsed, and the sequences' IDs are retrieved and compared to those returned by NCBI for the requested query. Any already present sequences are not downloaded again, and any missing sequences are appended to the

Download English Version:

<https://daneshyari.com/en/article/4978427>

Download Persian Version:

<https://daneshyari.com/article/4978427>

[Daneshyari.com](https://daneshyari.com)