

GEOsearch User Interface Manual

Zhicheng Ji, Hongkai Ji

July 10, 2017

1 GEOsearch Online GUI

GEOsearch Online GUI can be launched directly online: <https://zhiji.shinyapps.io/GEOsearch/>. The only needed software is the web browser and no additional software is required to be installed.

2 Installation of GEOsearch Software

GEOsearch software can be installed via Github. Users should first have R installed on their computer before installing GEOsearch. R can be downloaded here: <http://www.r-project.org/>. Users can then install the latest version of GEOsearch package via Github by running the following commands in R (the lastest version of R is recommended):

```
if (!require("devtools"))
  install.packages("devtools")
devtools::install_github("zji90/GEOsearch")
```

After that, one can launch the user interface by running the following commands in R:

```
library(GEOsearch)
GEOsearchui()
```

3 GEOsearch GUI Overview

The GEOsearch GUI consists of three main parts: the main menu on the top of the GUI where users can switch between different functions (Figure 1A), the sidebar panel on the left side where users can specify options (Figure 1B), and the main panel where the results will appear (Figure 1D).

GEOsearch has four main functions: searching GSE or GSM terms, key word enrichment analysis, obtaining details of all samples in GSE series, and generating a shell file to batch download samples in GSE series.

4 Performing Search

The first and most important step is to perform the search (Figure 1). GEOsearch provides two modes: traditional search and extended search, and GSM samples can also be searched in addition to GSE series (Figure 1B). Note that searching GSM samples could take a long time in some cases. To start the search one can simply enter the search term (“Oct4” in this example) in the text box “Enter Search Term” and click “Start Search” button (Figure 1C). By default GEOsearch will perform traditional search and will return exactly the same records as searching the same term in GEO. One can also switch to “Extended Search” using the radio button above the text box (Figure 1B). For extended search GEOsearch will split the term being searched into separate words and check whether each word is a gene name. For gene names GEOsearch will find all its gene aliases (e.g. Pou5f1, Oct3, Oct3/4, etc.) and perform additional searches by replacing the original gene name with each gene alias. The search results are then integrated and displayed in a compact table (Figure 1D). By default GEOsearch shows four columns that contain the most important information. One can use the “Column visibility” tab (Figure 1D) to display additional columns.

After an initial search, one can further narrow down the search results by performing a second-round search. For example, one can type in “pluripotency” in the text box above the “Title” column (Figure 1E) and type “array” in the text box above the “Type” column to only keep the search results that have the word “pluripotency” in title and the word “array” in type. This function allows one to conveniently focus on the GSE records of most interest without having to perform complicated programmatic search in GEO. One can also select the results that contain the keyword in any of the field by using the search box (Figure 1F).

One can change the order of the rows and columns using the computer mouse. To change the order of the columns, click and hold the title of a column and hover to a new place. To change the order of the rows, click and hold the ID of a row and hover to a new place. To sort the results according to a column, click on the title of a column and the results will be sorted in ascending or descending order. One can select the GSE records by clicking on the rows. The details of each samples for the selected GSE can be found in “Sample Details”.

One can download the contents of the current page as CSV, Excel or PDF (Figure 1D). One can also download all search results using the download button (Figure 1G). Click “Download Selected Part” if only the selected rows are to be downloaded.

5 Key Word Enrichment

GEOsearch enables one to effectively summarize the biological contexts (cell types, tissues and diseases) of the search results. Lists of cell types, tissues and diseases commonly used in biological experiments were precompiled from ATCC (www.atcc.org). After one selects which lists to use (Figure 2A), GEOsearch displays the number of GEO records whose titles or descriptions contain each key word (Figure 2B). GEOsearch will also perform a simple enrichment analysis to test whether each key word is significantly enriched in the search results

This screenshot shows the 'Perform Search' interface. At the top, there are tabs for 'GEOsearch', 'Perform Search' (which is active), 'Key Word Enrichment', 'Sample Details', and 'Batch Download'. Below the tabs is a search bar labeled 'Search GEO' with options for 'Traditional Search' (selected) and 'Extended Search'. A checkbox for 'Search GSM samples in addition to GSE series (may take a long time)' is unchecked. The search term 'Oct4' is entered in the search field. Below the search bar is a table of search results. The table has columns for 'Series', 'Organism', 'Title', and 'Type'. The results are as follows:

Series	Organism	Title	Type
A	A	All	All
1 GSE99631	Homo sapiens	The functional enhancer repertoire of human embryonic stem cells	Genome binding/occupancy profiling by high throughput sequencing; Other
2 GSE99630	Homo sapiens	The functional enhancer repertoire of human embryonic stem cells [STARR-RNA-seq]	Other
3 GSE99629	Homo sapiens	The functional enhancer repertoire of human embryonic stem cells [plasmid DNA-seq]	Other
4 GSE99628	Homo sapiens	The functional enhancer repertoire of human embryonic stem cells [isolated plasmid DNA-seq]	Other
5 GSE99627	Homo sapiens	The functional enhancer repertoire of human embryonic stem cells [ChIP-seq]	Genome binding/occupancy profiling by high throughput sequencing
6 GSE84009	Mus musculus	Vitamin C Induces Specific Demethylation of H3K9me2 in Mouse Embryonic Stem Cells via Kdm3a/b	Genome binding/occupancy profiling by high throughput sequencing

Figure 1: Screenshot demonstrating: Perform Search

This screenshot shows the 'Key Word Enrichment' interface. At the top, there are tabs for 'GEOsearch', 'Perform Search', 'Key Word Enrichment' (which is active), 'Sample Details', and 'Batch Download'. Below the tabs is a section for 'Key Word Enrichment Analysis' with a checkbox for 'Include Following Information' which is checked for 'Disease', 'Cell Type', and 'Tissue'. The search term 'Oct4' is entered in the search field. Below the search bar is a table of enriched terms. The table has columns for 'Term', 'Frequency', 'Logfoldchange', 'FDR', and 'All'. The results are as follows:

Term	Frequency	Logfoldchange	FDR	All
1 pluripotent	260	2.99449193096773	4.7110446863438e-247	
2 embryonic	291	1.49935383344215	8.32618983657471e-112	
3 stem cell	197	1.92372703725857	3.74600350963499e-101	
4 embryonic stem cell	104	2.64688161555774	1.84886916587981e-79	
5 pluripotent embryonic stem cell	20	6.763140297552	1.11997459642389e-38	

Figure 2: Screenshot demonstrating: Key Word Enrichment

compared to all samples in GEO (Figure 2C). In this example the top enriched key words include “embryonic stem cell” and “pluripotent embryonic stem cell”, which are well-known contexts related to Oct4. Thus the key word enrichment analysis can provide a quick overview of related biological contexts. One can select one or multiple biological contexts in the table and GEOsearch will return all GSE records whose titles or descriptions contain the selected contexts (Figure 2D). Again, one can select GSE records and explore the sample details in “Sample Details”.

6 Sample Details

For GSE records selected in “Perform Search” or “Key Word Enrichment”, GEOsearch will return the details of all samples related to these GSE records. In this way one can conveniently compare samples from different GSE records, which could be quite tedious in GEO. Users can choose selected GSE from “Perform Search” or “Key Word Enrichment” or enter new GSE ID in the text box (Figure 3A). The GSE ID can be entered in the text box (Figure 3B) if users have chosen “Enter new GSE”. Details of all samples in the given GSE will be shown in Figure 3C.

The screenshot shows the 'Sample Details' tab selected in the top navigation bar. On the left, there's a sidebar with options to choose from selected GSE in step 1 or 2, or enter new GSE. Below that is a text input field for 'Enter GSE name' containing 'GSE53531;GSE53532'. A red letter 'A' is placed over the 'Choose from selected GSE in step 1' option. A red letter 'B' is placed over the 'Enter GSE name' input field. On the right, a main panel displays a table of sample details. The table has columns for Experiment, Sample, Title, Type, Source, and Description. The data shows five entries, with the first two being GSE53531 and the last three being GSE53532. A red letter 'C' is placed over the table header.

Experiment	Sample	Title	Type	Source	Description
All	All	All	All	All	All
1 GSE53531	GSM1295590	siCTL_ChIP-Seq	ChIP-Seq	mammary gland, adenocarcinoma	
2 GSE53531	GSM1295591	siERα_ChIP-Seq	ChIP-Seq	mammary gland, adenocarcinoma	
3 GSE53531	GSM1295592	IgG	ChIP-Seq	mammary gland, adenocarcinoma	
4 GSE53532	GSM1295593	siCTL_RNA-Seq	RNA-Seq	mammary gland, adenocarcinoma	siCTL; polyA RNA
5 GSE53532	GSM1295594	siERα_RNA-Seq	RNA-Seq	mammary gland, adenocarcinoma	siERα; polyA RNA

Figure 3: Screenshot demonstrating: Sample Details

The screenshot shows the 'Batch Download' tab selected in the top navigation bar. On the left, there's a sidebar with instructions for running the function and a text input field for 'Enter GSE or GSM name' containing 'GSE53531;GSE53532'. A red letter 'A' is placed over the 'Enter GSE or GSM name' input field. A red letter 'B' is placed over the large text area containing the shell script. On the right, a large text area displays the generated shell script, which is a series of wget commands for downloading SRA files for the specified GSEs.

```

mkdir /home/data/GEO/GSM1295590/
wget -O /home/data/GEO//GSM1295590/data.sra ftp://ftp-trace.ncbi.nlm.nih.gov/sra/sra-instant/reads/ByExp/sra/SRX/
mkdir /home/data/GEO//GSM1295591/
wget -O /home/data/GEO//GSM1295591/data.sra ftp://ftp-trace.ncbi.nlm.nih.gov/sra/sra-instant/reads/ByExp/sra/SRX/
mkdir /home/data/GEO//GSM1295592/
wget -O /home/data/GEO//GSM1295592/data.sra ftp://ftp-trace.ncbi.nlm.nih.gov/sra/sra-instant/reads/ByExp/sra/SRX/
mkdir /home/data/GEO//GSM1295593/
wget -O /home/data/GEO//GSM1295593/data.sra ftp://ftp-trace.ncbi.nlm.nih.gov/sra/sra-instant/reads/ByExp/sra/SRX/
wget -O /home/data/GEO//GSM1295594/
wget -O /home/data/GEO//GSM1295594/data.sra ftp://ftp-trace.ncbi.nlm.nih.gov/sra/sra-instant/reads/ByExp/sra/SRX/
mkdir /home/data/GEO//GSM1295594/

```

Figure 4: Screenshot demonstrating: Batch Download

7 Batch Download

GEOsearch provides a function to generate a shell file for downloading all the raw files corresponding to the GSE from GEO. Enter GSE or GSM name for which the raw data will be downloaded in Figure 4A. Download path can also be designated. The shell file generated by GEOsearch to download the SRA raw data for the given GSE or GSM will be shown in Figure 4B. The generated shell file can be downloaded in Figure 4C. Users can run the shell file on their own computer to download all the files.

8 Contact

To report bugs and provide suggestions for the GEOsearch GUI as well as the GEOsearch package, please contact the maintainer Zhicheng Ji (zji4@jhu.edu).