

Rich Analysis of Variable gene Expressions in Numerous tissues



RAVEN

Documentation and user manual

Author:	Julia Sophia Gerke, MSc
Supervisor:	Thomas Grünewald, MD, PhD
Current version:	v1.0.0_RAVEN
Published:	October 2017

Contents

Contents	1
1 Installation	2
1.1 Requirements	2
1.2 Installation	3
1.3 Start RAVEN	8
1.4 Login	9
2 Upload Data	10
2.1 Upload	10
2.2 Processing Time	11
2.3 Available Input Files	12
2.4 Creating a new Input File	13
3 General Navigation	16
3.1 Application	16
3.2 Home - Working Directory	18
3.3 Table	19
4 Methods	21
4.1 Extracting Subsets	22
4.2 Direct Gene Expression Comparison	23
4.3 Cancer Specific Gene (CSG) Score	26
4.4 Peptide Matching Pipeline (PMP)	28
5 Save Data & Results	31
5.1 Images & Plots	31
5.2 Results & Data	31
6 About & Contact	34
6.1 Citation	34
6.2 Authors	34
6.3 Questions & Error Reporting	35
6.4 FAQ	36
6.5 License	39

1 Installation

1.1 Requirements

Raven runs on Windows , Linux  and Mac .

For Windows and Linux a **64-bit version** is obligatory. A random access memory (**RAM**) of at least **16 GB** is necessary. However, for big data files as provided from us, 32 GB are recommended for the easy and fast handling. You are not sure if your computer fulfills the requirements? Have a look at Table 1.1 and Figure 1.1 to find out more about your operating system (OS) and RAM or ask your system administrator for help.

Raven is a Java application and requires **Java 8** (see more in the next Section 1.2).

When using a laptop, your screen should have a resolution of at least **1440x900p** to display the application correctly. (Desktop computers normally have a higher resolution).

For the usage of the Peptide Matching Pipeline as described in 4.4 a **stable internet connection** is essential. Make sure you have the permission to let RAVEN send and receive queries to and from internet servers. This is sometimes restricted in some institutions with high data security. In this case talk to your IT-administrator.

Tip: The requirements really scare you or the Installation Chapter overburdens you? Ask your IT-administrator for help to check the requirements mentioned above and to install or update Java on your computer. Afterwards you can continue with Section 1.3 to start RAVEN.










You have OS and use and look at
	 → 'About This Mac'	Figure 1.1a
	Ctrl + Alt + T → write: <code>uname -a</code> → write: <code>less /proc/meminfo</code>	Figure 1.1b
 10 	 + X → 'system' check out the Microsoft help ¹ for older versions	Figure 1.1c

Table 1.1: Keyboard shortcuts depending on your operating system (OS) to open the information panel about your system. Your RAM is shown there. For Windows OS also the version (64-bit or 32-bit) can be seen.

1.2 Installation

Check your Java version!

Before you run RAVEN the first time make sure Java is installed on your computer and you have the right version of it (> 1.8.0_91 → Java 8). If you have Windows make sure you have installed the 64-bit version. To check this out use one of the following approaches:

- You can display your Java version via **commandline**. Open your terminal ( , ) or command prompt () with the keyboard shortcuts shown in Table 1.2 and type the command below:

```
java -version
```

The shown version should be 1.8.0_91 or higher. Additionally, it should include '64-Bit Server VM' both as shown in Figure 1.2. Otherwise you probably have an old version or the wrong system. If you get an error message for unknown command, java is not installed on your computer at all. In both cases update your version or install the right java version as described in the next Section 1.2.

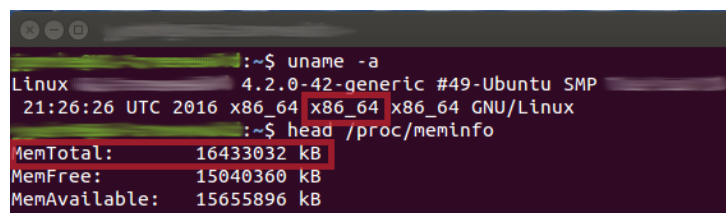
- If you are not familiar with the command line, you can follow the instructions from the [java website description](#)² on your **browser** to find your installed version of Java.

¹<https://support.microsoft.com/en-us/help/827218/how-to-determine-whether-a-computer-is-running-a-32-bit-version-or-64-bit-version-of-the-windows-operating-system>

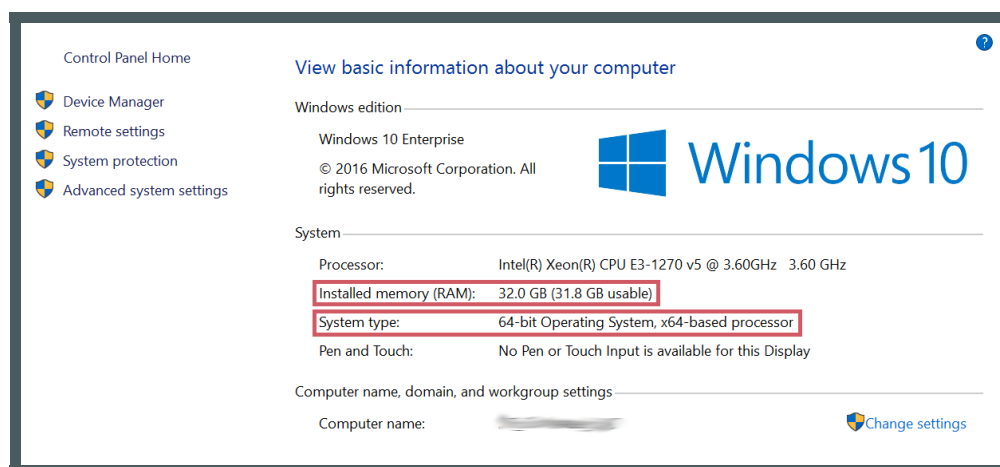
²https://java.com/en/download/help/version_manual.xml



(a) MacOS 🍏



(b) Linux 🐧



(c) Windows 🪟

Figure 1.1: Where to find your RAM and OS version (both framed red) for different operating systems (🍏, 🐧, 🪟). These information panels can be opened as described in Table 1.1.



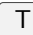









You have OS and use ...
	 + 
	 +  + 
 10	 +  → 'command prompt'
 7	 menu → search for 'command prompt'

Table 1.2: Keyboard shortcuts depending on your operating system (OS) to open the terminal (MacOS, Linux) or command prompt (Windows).

```
$ java -version
java version "1.8.0_91"
Java(TM) SE Runtime Environment (build 1.8.0_91-b14)
Java HotSpot(TM) 64-Bit Server VM (build 25.91-b14, mixed mode)
```

Figure 1.2: Output for 'java -version' on Windows command prompt with red framed java version and its OS dependencies. This output can vary subject to your OS.

Download Java 8

If you do not have Java 8 already installed on your computer you can download the Java Runtime Environment (JRE) on the [website](#)³ of its developer ORACLE. A detailed view of the website is shown in Figure 1.3, though make sure you accept the license agreement in order to start the download (red framed). Choose the corresponding executable file highlighted in orange in Figure 1.3. After saving the executable application, start the installation process with double click and follow the shown instructions until the installation process is finished. Done!



Tip: It is also possible to download Java from its own website. However, this mostly refers to the Java 32-bit version. Be careful when you are not familiar with it. Rather download java from ORACLE as described above!

³<http://www.oracle.com/technetwork/java/javase/downloads/jre8-downloads-2133155.html>

Java SE Runtime Environment 8u144



You must accept the [Oracle Binary Code License Agreement for Java SE](#) to download this software.

☒ Accept License Agreement
 ☐ Decline License Agreement

Product / File Description	File Size	Download
Linux x86	59.13 MB	jre-8u144-linux-i586.rpm
Linux x86	75.01 MB	jre-8u144-linux-i586.tar.gz
Linux x64	56.48 MB	jre-8u144-linux-x64.rpm
Linux x64	72.41 MB	jre-8u144-linux-x64.tar.gz
Mac OS X	63.94 MB	jre-8u144-macosx-x64.dmg
Mac OS X	55.56 MB	jre-8u144-macosx-x64.tar.gz
Solaris SPARC 64-bit	52.12 MB	jre-8u144-solaris-sparcv9.tar.gz
Solaris x64	49.95 MB	jre-8u144-solaris-x64.tar.gz
Windows x86 Online	0.7 MB	jre-8u144-windows-i586-iftw.exe
Windows x86 Offline	54.57 MB	jre-8u144-windows-i586.exe
Windows x86	60.2 MB	jre-8u144-windows-i586.tar.gz
Windows x64 Offline	62.34 MB	jre-8u144-windows-x64.exe
Windows x64	63.99 MB	jre-8u144-windows-x64.tar.gz

Figure 1.3: JRE download panel on the ORACLE website, choose the right one (orange) to your operating system ,  or .

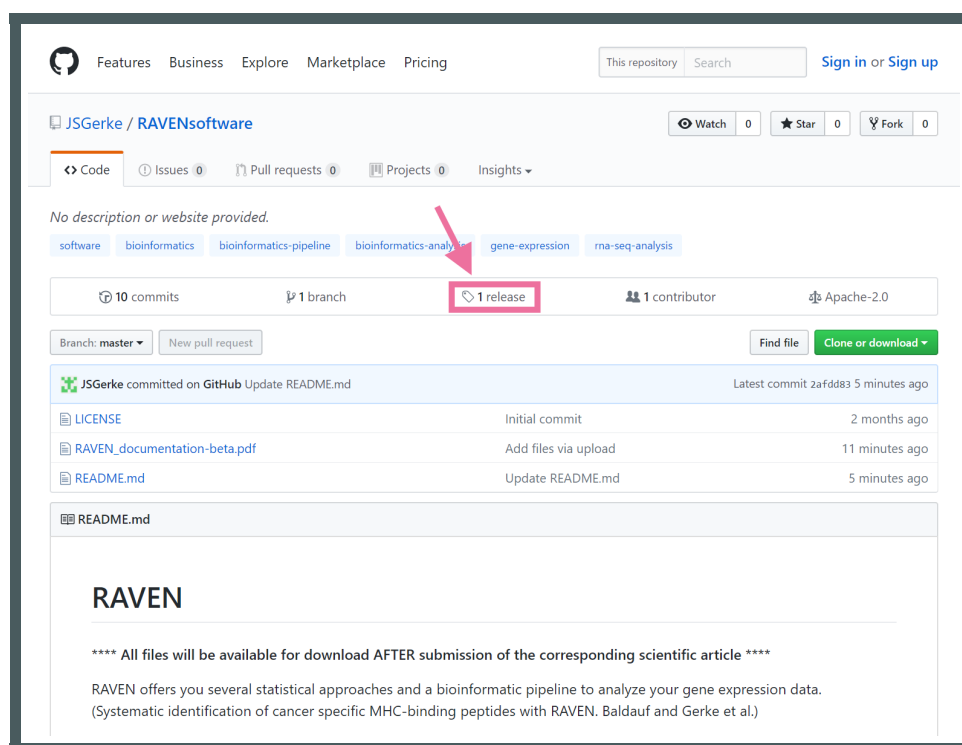
Download RAVEN software

Download the current release version of RAVENsoftware (vX.X.X_RAVEN) from [GitHub](https://www.github.com/JSGerke/RAVENsoftware)  (<https://www.github.com/JSGerke/RAVENsoftware>). Also its manual and already prepared gene expression datasets are available there on Github [GitHub](https://www.github.com/JSGerke/RAVENsoftware) ⁴(Figure 1.4a).

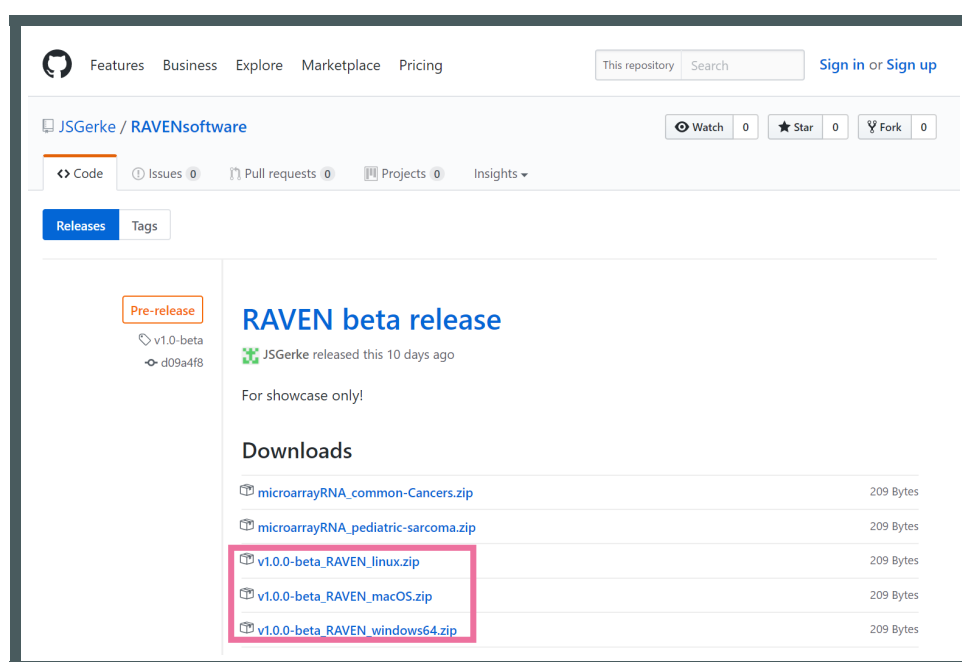
To install RAVEN, download the vX.X.X_RAVEN_ddd.zip corresponding to your operating system (, , ) as shown in red in Figure 1.4b.

Unpack your downloaded .zip file via **right click** > **extract all** and choose a (new) directory to save the files. The current release contains two files, the user manual and the executable file. Please keep **both** files from your unpacked .zip in the **same directory** to always have a documentation of RAVEN at your fingertips via the softwares 'help' button.

⁴<https://github.com/JSGerke/RAVENsoftware>



(a) RAVEN download main page on GitHub. Beside a short overview, you can download RAVENs documentation here. To download the application continue via release marked red. (<https://github.com/JSGerke/RAVENsoftware>).



(b) release: Download current RAVEN version vX.X.X_RAVEN_ddd.zip for your operation system (Apple, Windows, Linux). This figure displays the beta version as an overview for you. Please use the current version which download site is constructed the same way. (<https://github.com/JSGerke/RAVENsoftware/releases>).

Figure 1.4: Screen shots of the platform GitHub where you can download RAVEN and gene expression data.

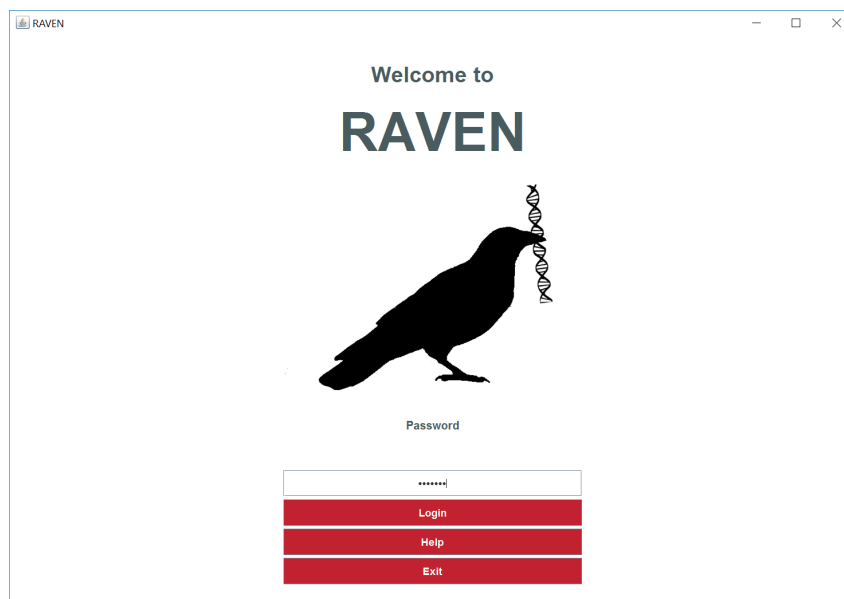





Figure 1.5: RAVEN start window. Enter the password to proceed (Section 1.4).

1.3 Start RAVEN

You can easily start RAVEN by double click on the executable file RAVEN_vX.X.X (.exe → , .app → , .bin → ) which you unpacked from your downloaded .zip file. The Login Window as shown in Figure 1.5 will open immediately to welcome you to RAVEN.

Tip: RAVEN quick start! You are annoyed of the exhausting navigation through all your directories or even have troubles remembering where you saved RAVEN? Create a shortcut to your desktop!

Copy RAVEN_vX.X.X.exe file » go to your desktop » right click on desktop » create shortcut

Now you can easily start the application with a double click on your newly created shortcut. Be, careful: do not copy + paste the file! By doing so the help manual cannot be found by RAVEN anymore.

If you prefer to run RAVEN from command line, navigate into the directory where you saved RAVEN and start the java application as described below. Replace X with your downloaded version.

```
cd yourDirectoryPath\  
  
.\RAVEN_vX.X.X.exe
```

1.4 Login

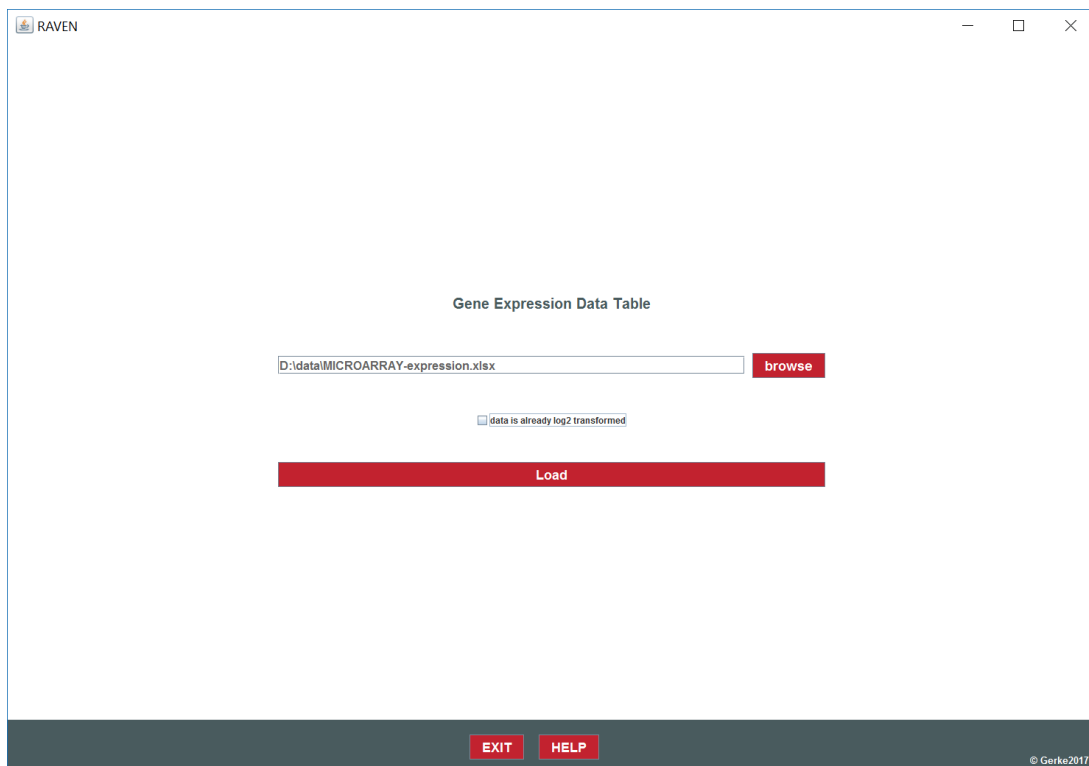
RAVEN can only be accessed with a password (shown in Figure 1.5), which is free for academic use only. Please contact us (thomas.gruenewald@med.uni-muenchen.de) to receive the login password to start exploring RAVEN.

Upload Data

2.1 Upload

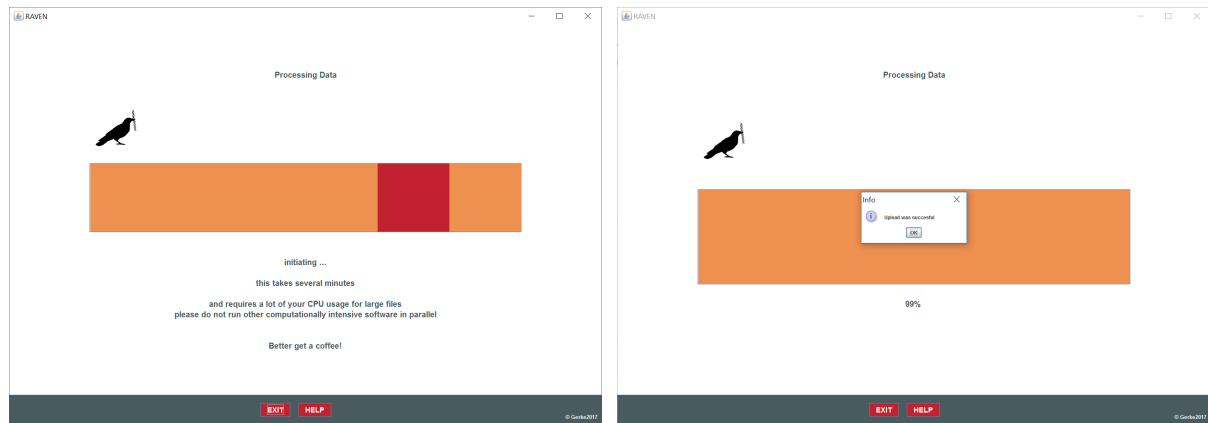
Enter or browse your input data file as shown in Figure 2.1. For later calculations it is important to know if your data is already log2 transformed or still in its natural scale. Activate the checkbox accordingly before pushing the 'Load' button. The input data table can be a **text file** (.txt) or in **Excel format** (.xlsx). For more details and requirements regarding your input data have a look at Section 2.3 and 2.4.

After the upload you will receive a short notification if your upload was successful (Figure 2.2b) or not. A failed upload mostly derives from a formatting error of your data. An error message will notify you about the flaw to help you to correct it. For this, also check out Section 2.3 and 2.4.



The screenshot shows a web application window titled "RAVEN". The main content area is titled "Gene Expression Data Table". It features a text input field containing the file path "D:\data\MICROARRAY-expression.xlsx", followed by a red "browse" button. Below this is a checkbox labeled "data is already log2 transformed", which is currently unchecked. A large red "Load" button is positioned below the checkbox. At the bottom of the window, there is a dark grey footer bar containing two red buttons labeled "EXIT" and "HELP", and a small copyright notice "© Gerke2017" on the right side.

Figure 2.1: Gene Expression data upload panel for log2 transformed or natural scaled files in text or Excel format which will directly appear after your login.



(a) View on upload progress ⌘. Watch RAVEN working or get a ☕ break!

(b) Info message for successful upload.

Figure 2.2: RAVEN processing your input data.

2.2 Processing Time

Depending on the size of your input file and the available RAM, the processing time can vary from a few seconds to several minutes (Figure 2.2a). In worst case your file is too large and cannot be processed at all. Check out Table 2.1 to estimate how many samples you can process dependent on your RAM. Make sure your file size does not diverge a lot from those in the table. Otherwise, consider the listed suggestions below.

How to reduce processing time

Your file is big and cannot be processed or you want to speed up the loading process? Try the following suggestions:

- use text format instead of Excel format
- reduce your decimals to 2-3 digits
- close other parallel running computational intensive software
- change to another computer with higher performance (RAM)


# sample size [file size]	RAM			
	4 GB	8 GB	16 GB	32 GB
500 [62MB]	✓	✓	✓	✓
1000 [124MB]	✗	✓	✓	✓
1500 [185MB]	✗	✓	✓	✓
2000 [247MB]	✗	✗	✓	✓
2500 [360MB]	✗	✗	✗	✓
3000 [420MB]	✗	✗	✗	✓
2400 [500MB]	✗	✗	✗	✓

Table 2.1: Approximate number of samples (including their file size) that can be processed by RAVEN dependent on the available RAM of your computer, assuming that for every sample around 20000 genes were measured. The files are text files and all measurements (natural scale) are rounded to 3 decimals. Though, the last table row refers to natural scale raw data in Excel format.

✓ file can be processed; ✗ file cannot be processed

2.3 Available Input Files

Together with RAVEN, we also published two datasets each consisting of several cancer and normal tissue samples from Affymetrix microarray data downloaded from Gene Expression Omnibus ([GEO](#)). While for both files the normal tissue samples are the same, they differ in their tumor types.

All gene expression files are already converted to the right format to be accepted as input for RAVEN. They are available in natural scale and rounded to 3 decimals. You can download all .txt files zipped on [GitHub](#)  where you also downloaded the application before as described in Section 1.2. A complete list of the covered cancer types is included in the zipped files. All samples underwent normalization and a quality check described by *Baldauf and Gerke et al.*

Normal Tissues

The normal tissues in our files cover 71 different types of in total 929 samples including the most important normal tissues.

Most Common Cancer

This file contains samples from 15 of the most common cancers in human of in total 1462 patients. The tumors include breast, bladder, cervical, colon, endometric, esophagus, kidney, liver, lung, oral, ovarian, pancreatic, prostate, stomach and thyroid cancer.

Pediatric Sarcoma

This file focuses mainly on sarcomas and pediatric cancers. Beside the normal tissues, 1749 samples of 50 tumor entities are assembled in this dataset, such as Ewing sarcoma.

2.4 Creating a new Input File

It is also possible to upload your own data to RAVEN. Make sure your data file matches the specifications and requirements listed below, otherwise your data file upload will fail. An excerpt of an (affymetrix) microarray example file is shown in Figure 2.3.

	A	B	C	D	E	F
1	Study	GSE1111	GSE1111	GSE1111	GSE1111	GSE1111
2	Tissue	NORMAL_ADRENAL	NORMAL_ADRENAL	NORMAL_ADRENAL	BREAST_CANCER	BREAST_CANCER
3	Probesets	GSM27.CEL	GSM28.CEL	GSM29.CEL	GSM30.CEL	GSM31.CEL
4	100009676_at	33.332562	28.081642	44.485829	33.091241	26.995848
5	10000_at	61.889346	72.09631	50.17429	46.002119	53.288426
6	10001_at	220.894017	132.348928	65.006652	92.446847	116.872781
7	10002_at	23.894864	20.511832	36.115859	30.654433	24.532586

Figure 2.3: Example excerpt of an input file for RAVEN. Please note that the first three cells (A1-3) are required to be exactly as shown here!

Input file specification

- the **first line** (or header) starts with the term 'Study', followed by each samples study ID (e.g. GSE123). As this is never processed in any method you can also store some sample information in these cells. Do not leave them empty!
- the **second line** starts with the term 'Tissue' followed by the normal tissue or cancer of each sample. Normal tissues **always** have to begin with 'NORMAL_' ! (e.g. NORMAL_BREAST, BREAST_CANCER)
- the **third line** starts with the term 'Probesets' followed by a unique id for each sample. (e.g. GSM3333)
- the forth and **consecutively lines** contain the measured gene expressions (Affymetrix) for all samples beginning with the gene ID of the corresponding gene in the first column.
- the **gene ID** must be an entrez gene id (GeneID) followed by the suffix '_at', as often in Affymetrix data. (e.g. 1000_at). Otherwise, check out the description for ID translation below this list. If your Affymetrix data sets has probesets starting with 'AFFX-', delete them from your data table!

- all columns (except first one) must be **sorted alphabetically** by their samples' tissue (second line).
- within cells of the first 3 rows no whitespace is allowed. Use _ to join several words to one per cell (COLON_CANCER ✓; ~~COLON CANCER~~ ✗)
[in general: always avoid whitespace!]
- if your data includes **gender specific normal tissues** make sure they match the following **annotation**. Otherwise some methods which take gender specific normal tissues into account will not work properly.
→ NORMAL_UTERUS_ENDOMETRIUM, NORMAL_UTERUS_MYOMETRIUM,
NORMAL_CERVIX, NORMAL_OVARY, NORMAL_VAGINA,
NORMAL_PROSTATE, NORMAL_TESTES, NORMAL_PENIS
- for **decimal mark** the application supports both point (e.g. US, UK) and comma (e.g. D, EU), but no thousands separator for digit grouping (1000.1 ✓; ~~1,000.1~~ ✗)
- **no empty cells** within the data matrix/table. Use NA or NaN instead.
- if your file is a text file use **tab separated columns**
- if your file is an Excel file make sure your data table is on the **first sheet**. All other sheets will be ignored.
- if your file is an Excel file it has to **end with .xlsx**. Excel files created with older versions than Excel2007 have a different ending and are not accepted.

Tip: RAVEN is optimized to analyze microarray gene expression data. However, you can also run the software with your own RNAseq data, as long as you prepare your data accordingly. Make a quality assessment on your data and preprocess it correctly. Afterwards, normalize your reads and adapt the input file specifications necessary to upload data to RAVEN. Though, you probably have to adjust the IDs accordingly to match the specifications (see below)

ID translation

RAVEN only accepts the entrez gene ID (GeneID) for genes. If you have another ID, you can convert it to GeneID with the ID converter (second tab!) of biodb.jp¹, [DAVID](https://david.ncifcrf.gov/)² or [bioDBnet](https://biodbnet-abcc.ncifcrf.gov/db/db2db.php)³.

It is also important that your gene IDs end with the suffix '_at'. If the suffix is missing you can add it via Excel or on command line.

¹<http://biodb.jp/>

²<https://david.ncifcrf.gov/tools.jsp>

³<https://biodbnet-abcc.ncifcrf.gov/db/db2db.php>

Excel: Open your file in Excel (also possible with a text file but choose the right separator!). Add a new sheet and rename it to 'suffix'. Copy your first column to the new sheet. There, start with cell B4 and enter the first formula shown below in the black box. Drag this cell (by clicking on the small square in the bottom right cell corner) over the whole column B. The formula will adapt accordingly. Switch back to your original data sheet and override cell A4 with the formula as shown below (second one). Also drag this cell over the whole column A. Save your modifications.

```
=A4&"_at"
```

```
=suffix!B4
```

cmd line: (for Windows only with cygwin installed) Open the terminal and navigate into the directory of your data file (as already described in Section 1.2) and write the following (only working for text files!).

```
cd /yourDirectory/

head -n3 inputData.txt > newInputData.txt

awk -v OFS='\t' 'NR>3 { $1=$1"_at"; print }'
inputData.txt >> newInputData.txt
```


General Navigation

3

3.1 Application

RAVEN offers you several statistical analysis methods and bioinformatic pipelines to analyse your gene expression data.

Methods implemented in RAVEN:

- **Extracting Subsets** (see 4.1)
- **Direct Gene Expression Comparison** (see 4.2)
- **Cancer Specific Gene (CSG) Score** (see 4.3)
- **Peptide Matching Pipeline (PMP)** (see 4.4)

Setting Parameters

In general, all methods are based on the same operational procedure. See the stepwise procedure illustrated with black circled numbers in Figure 3.1.

- ❶ selecting one of the tabs on the left side for the method you want to run (🍏 → tabs are shown at the top)
- ❷ specify if you want to perform this method on a gene or tissue
- ❸ enter your favorite gene or tissue
- ❹ specify the genes or tissues you want to compare/test against your favourite against (gene vs tissues or tissue vs gene only!)
- ❺ select some additional settings
- ❻ click on 'view' button to run the analysis method
- ❼ the result will show up here after RAVEN is done with computing
- ❽ click on 'save' button to save your results or data used for the analysis

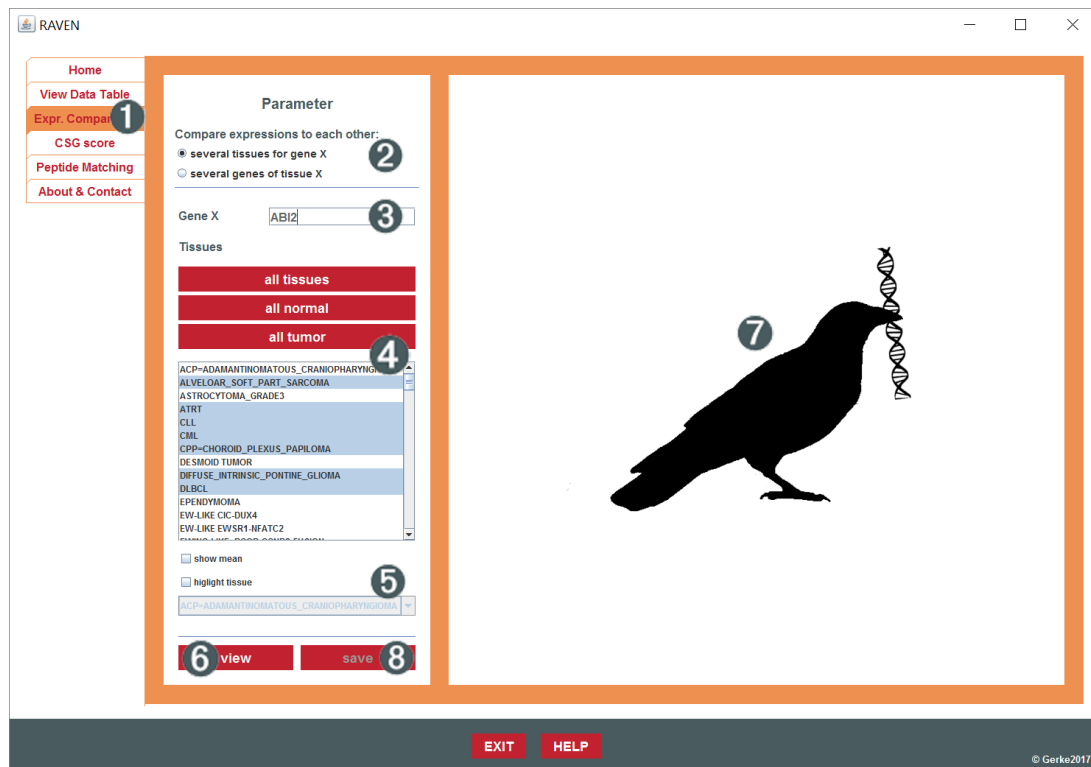


Figure 3.1: Step wise procedure to run an analysis with RAVEN shown on an Example.

While you can add only one gene to the one line text fields in the parameter section, you can add several genes in the larger text areas to apply a method with same settings on multiple genes at once. Thereby, you can use the gene symbol as well as the gene ID (entrez) as identifier. However, for the peptide matching pipeline you can only use gene symbol. Gene symbols are not case sensitive.

Tissue types cannot be entered, but only selected from the available ones of your input file.

Additional Information

Some parameters or options, which might not be intuitive have a small information symbol **i** next to it. A information panel will pop up when you move your mouse pointer onto the symbol (Figure 3.2).

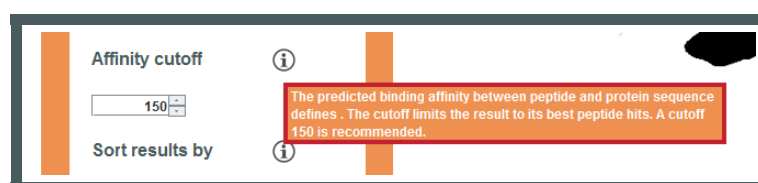


Figure 3.2: Inbuilt information with pop up text.

Help

If you need help while using RAVEN and do not have this manual at your fingertips you can open this documentation via the red HELP button in the menubar at the bottom of the application.

3.2 Home - Working Directory

Setting a working directory is optional but saves you a lot of time and effort saving multiple files. Click on the 'Home' tab and browse or type (non existing paths are colored red) your favorite directory where you want to store your data later. To confirm your directory click the 'set directory' button (see Figure 3.3). Then the file manager will automatically open this directory for you when saving data. Only existing directories are allowed. Otherwise, the button will be disabled with its writing fading to grey.

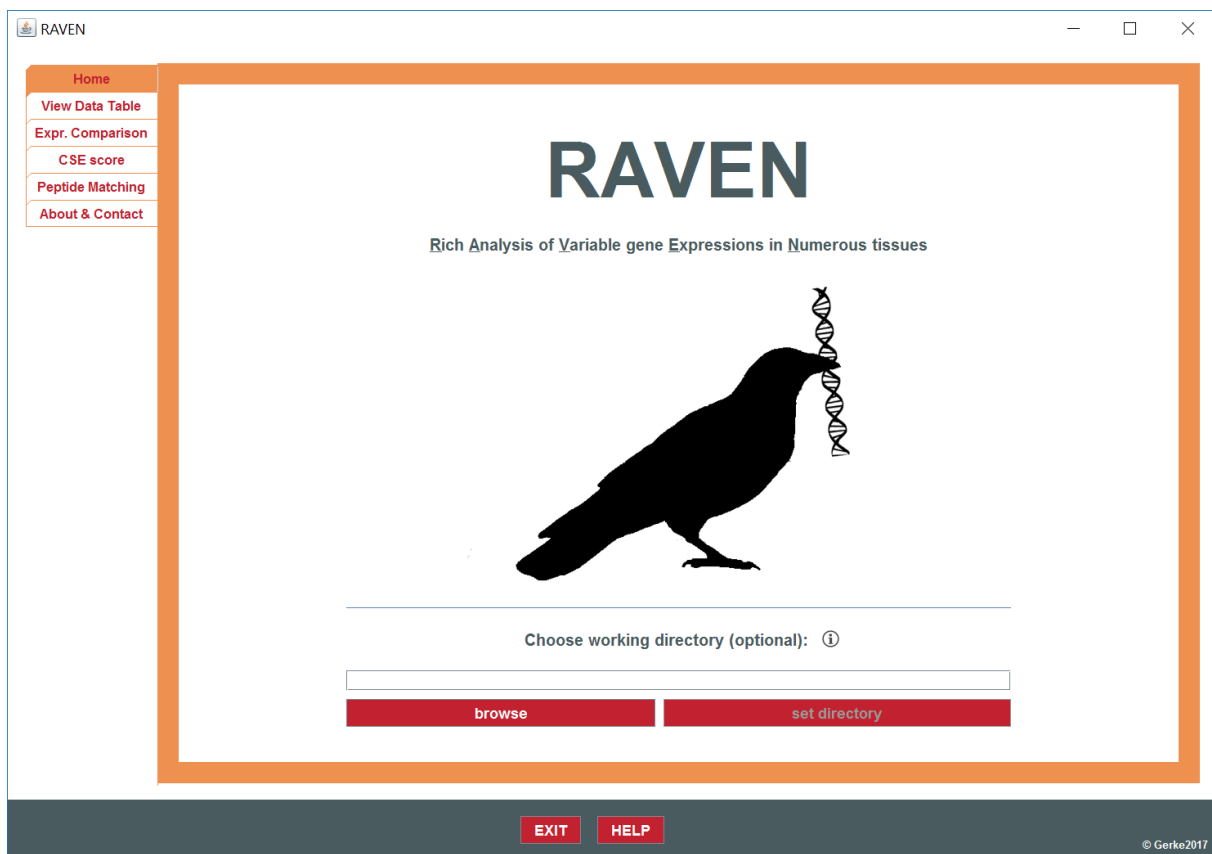


Figure 3.3: Raven Home: set the working directory via the 'set directory' button (disabled here because of missing/unvalid directory in the textfield)

3.3 Table

The results are always shown or accompanied by a table.

Navigation and Selection

Mostly a table is sorted by the result's p-value or most important score. The current sorting column is marked with a small adumbrated arrow, pinpointing to its sorting direction. To resort the table by values of another column, click on the corresponding column header. A second click on the same column changes from increasing order to decreasing or other way around. It is also possible to change the order of the columns. Drag and drop the columns by their header into your desired order. When saving table data the reorder is taken.

You can select several rows to accentuate them from the remaining ones or to save a subset of data. For an easy navigation you can use the keyboard as described in table 3.1. On the search bar (Figure 3.4) below the table you can also see the number of selected rows.

Select...	Windows / Linux	Mac
all rows	Ctrl + A	⌘ + A
consecutive rows	↑ + mouse last row to select	↑ + mouse last row to select
scattered rows	Ctrl + mouse unselected rows	⌘ + mouse unselected rows
deselect a certain row	Ctrl + mouse selected single row	⌘ + mouse selected single row

Table 3.1: Keyboard combinations to select or deselect your desired table rows dependent on your operating system.

Search

You are looking for a certain gene or tissue? Search for it via the search bar below the table. All rows matching the search word will be filtered in the table. Separate several search terms by comma. If part of the table, you can search for categories like entrez gene ID, genesymbol, tissue, tumor, peptide, allele and protein. Be careful when using the 'select all' button after a gene search. As this is a simple text search you will also filter other genes containing your search terms (e.g. S0X1 → S0X1, QS0X1, S0X12 ...).

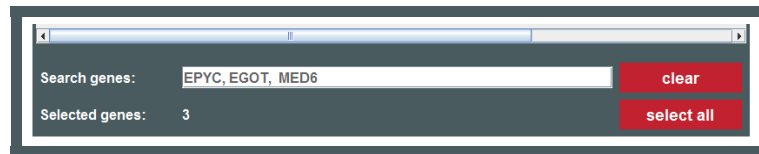


Figure 3.4: Inbuilt information with pop up text.

Gene transfer between methods

You identified multiple genes in your current method with which you want to continue in another RAVEN method? To save you much time and writing

Do Select rows of current table right click on this selection copy genes navigate to another method
right click on text area enter genes

4 Methods

RAVEN offers you several statistical approaches and a bioinformatic pipeline to analyze your gene expression data. In this chapter you learn more about use and of those methods. By combining them all (see workflow in Figure 4.1) you can receive a comprehensive gene expression analysis. But you can also run each method independently from the others.

- **Extracting Subsets** (see 4.1)
- **Direct Gene Expression Comparison** (see 4.2)
- **Cancer Specific Gene (CSG) Score** (see 4.3)
- **Peptide Matching Pipeline (PMP)** (see 4.4)

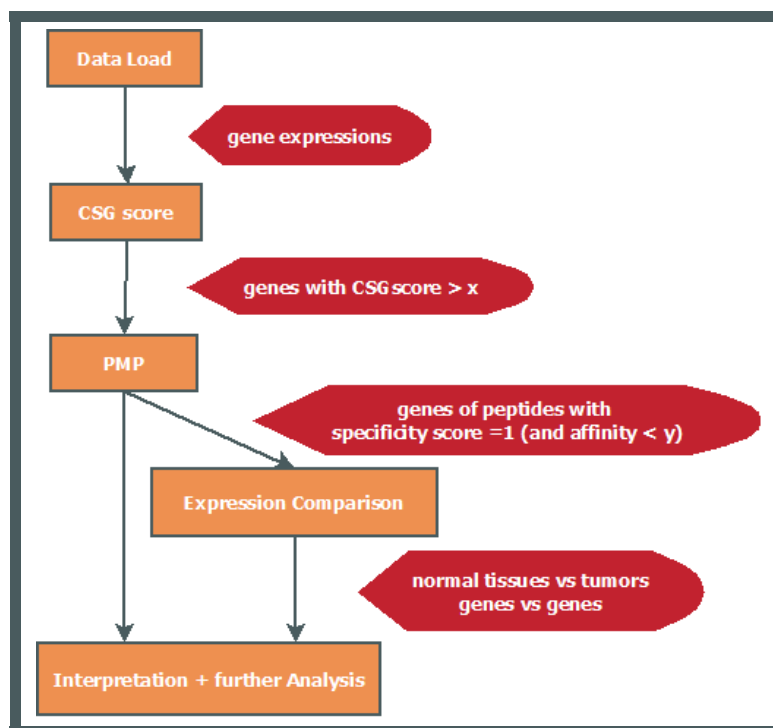


Figure 4.1: Example workflow of data analysis to fully exploit RAVEN's possibilities. Every method (orange) can be performed independently.

[Please note: All images and examples in this manual were randomly chosen without any scientific reason. The aim of those is to illustrate and explain the handling and support understanding of RAVEN and its included approaches only.]

4.1 Extracting Subsets

Via the 'View Data Table' tab you can display your uploaded data. If your data contains samples from different tissues or tumors you can display a subset of your whole data set. Select all tissue types you want to extract and view the corresponding table. You can also focus on specific genes by searching and selecting them. Save your table via the 'save' button and choose how your data table should look like.

Beside a classical table, you can also save your subset in RAVEN style which enables you to upload your extracted subset into RAVEN saving you the time to manually annotate your table to the standards of RAVEN input files (described in Section 2.4). See in Section 5.2 for more information on saving your data and results.

The screenshot shows the RAVEN web application interface. On the left, a sidebar contains navigation links: Home, View Data Table (active), Expr. Comparison, CSG score, Peptide Matching, and About & Contact. The main area is titled 'Parameter' and 'View and export tissue dependent subsets of your data.' It features a 'Tissues' section with three red buttons: 'all tissues', 'all normal', and 'all tumor'. Below this is a scrollable list of cancer types, including ACP-ADAMANTINOMATOUS_CRANIOPHARYNGIOMA, ALVEOLAR_SOFT_PART_SARCOMA, ASTROCYTOMA_GRADE3, ATRT, CLL, CML, CPP-CHOROID_PLEXUS_PAPILOMA, DESMOID TUMOR, DIFFUSE_INTRINSIC_PONTINE_GLIOMA, DLBCL, EPENDYMOMA, EW-LIKE_CIC-DUX4, EW-LIKE_EWSR1-NTF2C2, EWING-LIKE_BCOR-CCNB3-FUSION, EWING_SARCOMA, FOLLICULAR_LYMPHOMA, and GERMINOMA. At the bottom of this list are 'view' and 'save' buttons. The right section displays a table with columns: Gene ID, Genesymbol, CLL.GSM977137_SDG_00021.CEL, and CLL.GSM977131. The table lists various genes and their expression values across different samples. At the bottom of the table is a search bar with the text 'Search genes:' and a 'clear' button, and a 'Selected genes: 0' label with a 'select all' button. The footer of the application includes 'EXIT' and 'HELP' buttons, and a copyright notice '© Gerke2017'.

Gene ID	Genesymbol	CLL.GSM977137_SDG_00021.CEL	CLL.GSM977131
100009676_at	ZBTB11-AS1	35.589	
10000_at	AKT3	42.602	
10001_at	MED6	196.411	
10002_at	NR2E3	26.506	
10003_at	NAALAD2	9.303	
100048912_at	CDKN2B-AS1	21.252	
10004_at	NAALADL1	177.134	
10005_at	ACOT8	249.31	
10006_at	ABI1	947.341	
10007_at	GNPDA1	135.398	
10008_at	KCNE3	31.098	
100093630_at	SNHG8	1544.441	
100093698_at	GS1-600G8.3	13.211	
10009_at	ZBTB33	337.593	
1000_at	CDH2	42.378	
100101467_at	ZSCAN30	88.954	
100101938_at	ANKRD26P3	48.383	
10010_at	TANK	836.763	
100113407_at	TMEM170B	17.353	
10011_at	SRA1	348.093	
100124700_at	HOTAIR	10.292	
100125288_at	ZGLP1	128.709	
100126784_at	LOC100126784	26.563	
100126791_at	EGOT	21.361	
100126793_at	GHRLOS	42.554	
100127888_at	SLC04A1-AS1	48.899	
100127983_at	C8orf88	15.149	
100128025_at	WVTR1-AS1	19.069	
100128071_at	FAM229A	246.302	
100128098_at	ST8SIA6-AS1	19.875	
100128124_at	HGC6.3	48.432	
100128126_at	STAU2-AS1	27.641	

Figure 4.2: Overview of RAVEN method 'View Data Table' to extract a st of randomly selected cancer types.

4.2 Direct Gene Expression Comparison

The direct gene expression comparison approach enables the direct visualization of gene expression data via box plots. You can decide to analyze expression differences **between multiple genes** in the same tissue or cancer as well as differences in expression values of the same gene in **several tissue or/and cancer types** as shown in Figure 4.4.

The boxes of the plot are sorted descending by their median. You can spice up your box plot by marking the mean or highlighting your favorite gene, tissue or cancer type in red.

Parameter Setting

What you have to do:

1. **choose the category type** of your boxplot: several genes in cancer X or rather gene X in various cancers/normal tissues
2. **note your X** of which you want to compare your categories
3. **select your categories** which should be compared (genes or cancers/tissues)
4. **approve** if you also want to display the **mean** in your boxplot; only this will include the **extreme outliers** as well
5. **approve** if and which category you want to **highlight** in red
6. **Submit your query** Via the 'view' button.
7. **Save** your result table or plot data to import them in another program and recreate your plot via the 'save' button.

Tip: Especially for box plots with a high number of bars it is recommendable to accentuate the most important bar in the graph.

Results

The results include a plot as well as a table summarizing each box of it in numbers. A box represents a category either a gene or tissue/cancer, depending on your chosen option. A detailed description on the box and its calculation is shown in Figure 4.3.

What RAVEN tells you:

- **visualization** → boxplot of your complete (selected) data

- **mean** → mean of all values (per category)
- **median** → median of all values (per category)
- **IQR** → inter quartile range is the distance between the Q1 and Q3 covering 50% of the data around the median; corresponds to the actual box (red in Figure 4.3) in a boxplot
- **regular range** → distance between the whiskers covering all non-outliers (per category)
- **regular min + max** → smallest/highest (non-outlier) value of all values (per category) → without considering the outlier
- **samples** → number of samples n of your category
- **number of outliers** → total number of outliers and farouts (extreme outlier)

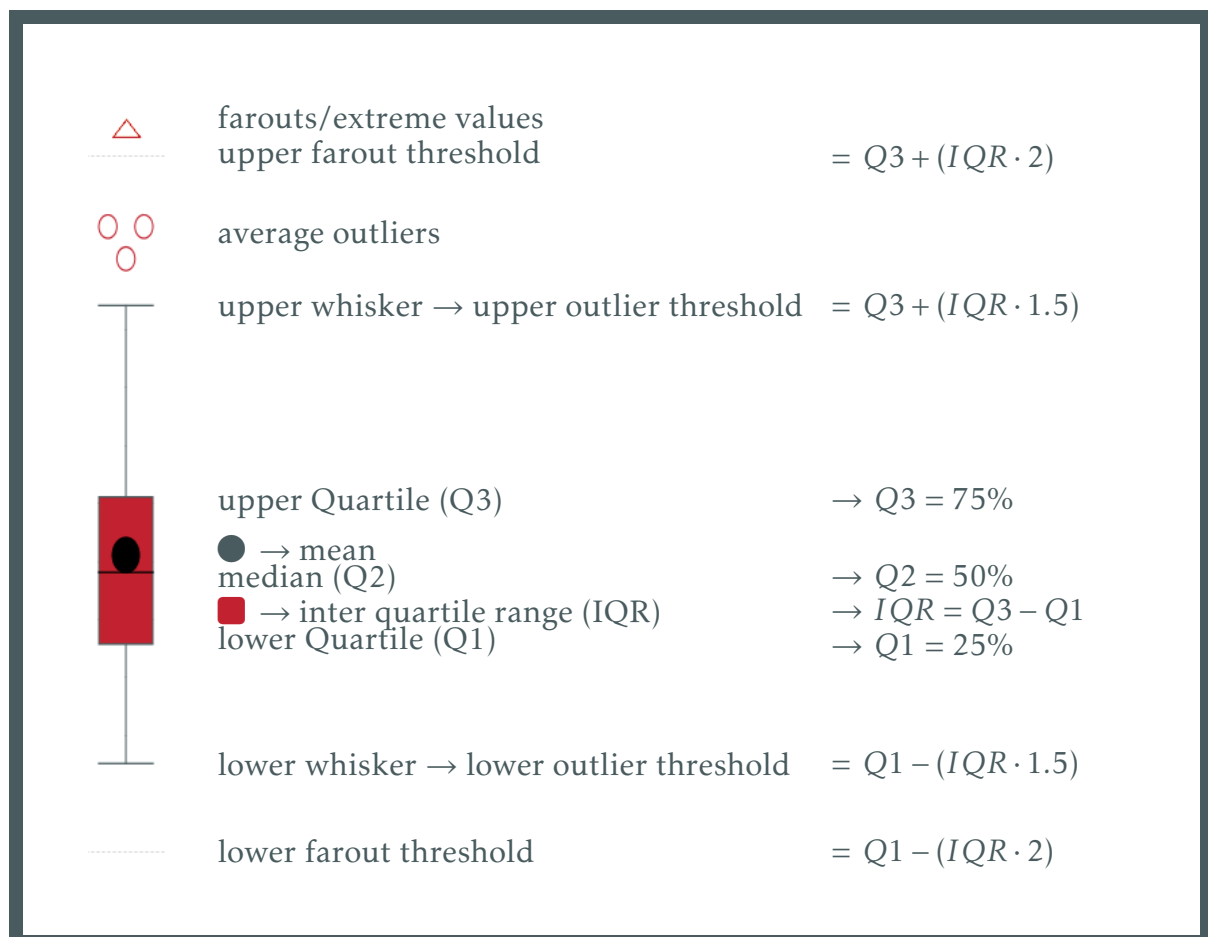
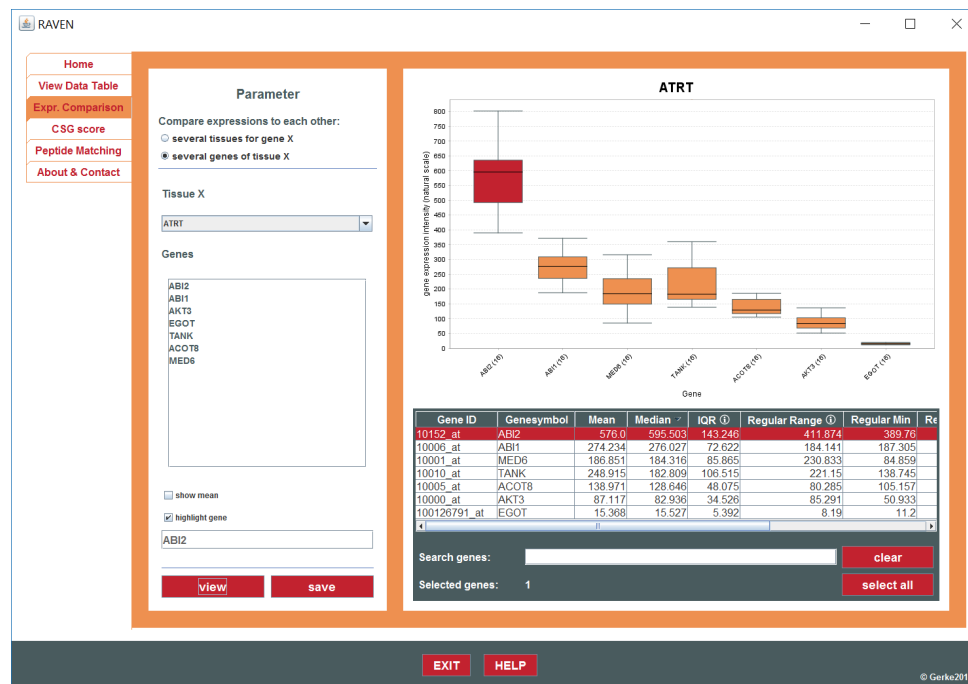
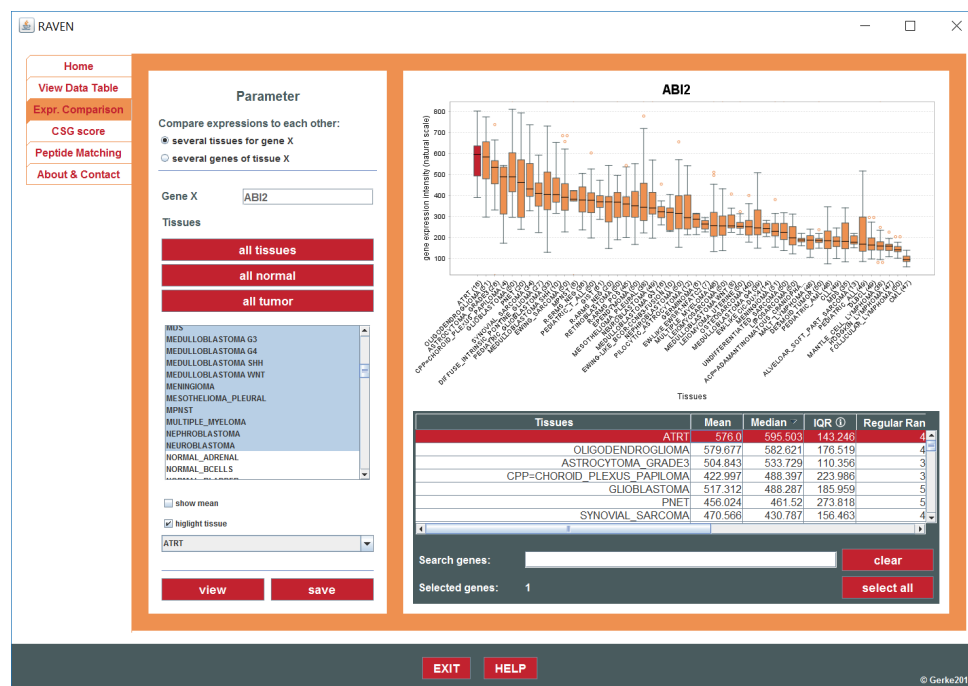


Figure 4.3: Schematic description of the boxplot bar visualized by RAVEN on an artificial example.



(a) Comparing expression of multiple genes with focus on ABI2 (red) in a specific cancer (ATRT).



(b) Comparing expression of the same gene (ABI2) in different cancer types with focus on ATRT (red)

Figure 4.4: RAVEN results of the direct gene expression comparison. Example of analyzing differences in expression values with focus on ABI2 (gene) and ATRT (cancer) compared to other genes, cancer and normal tissues.

4.3 Cancer Specific Gene (CSG) Score

The CSG score indicates if a certain gene is higher or lower expressed in one cancer compared to normal tissues. The CSG scores for CSGs potentially suitable for immunotherapeutic targets in a given cancer entity were usually greater than 2. CSG scores greater than 3 were empirically considered as high and those greater than 4 as very high. Read more about the CSG score, its calculation and meaning, in our paper [Baldauf and Gerke et al].

CSG Score workflow

What RAVEN does for you:

1. **identification** of 'Cancer Specific Genes':
2. calculate outlier expression score (*OS*) of gene X in cancer C
3. calculate penalty score (*PS*) for gene X based on N different types of normal human tissues.
4. CSG score of gene X in cancer C: $CSG(X, C) = OS(X, C) - PS(X, N)$

This enumeration is only a short headword list, to make you curious. A detailed description of the CSG score is written in our paper [Baldauf and Gerke et al.].

CSG Score parameter setting

What you have to do:

1. **choose a type** to calculate CSG scores for gene X in several cancers, or for several genes in cancer X
2. **set gene/tumor X**
3. **select your genes/tumors** for which you want to get your score regarding to X
4. **submit your query** Via the 'view' button
- ... You want to continue with several specific genes of your result in another method?
5. **select your favorite genes** in result table
6. **copy their gene symbols** via right click » copy selected gene symbols
7. **use them for another method** navigate to another method » right click into large text area
» paste genes here » continue with remaining parameters ...

Results

Every CSG score is calculated twice. To avoid a gender specific bias RAVEN calculates a female score considering female specific normal tissues and excluding male specific normal tissues. The male score is accordingly calculated the other way around.

What RAVEN tells you:

- **CSG-score** → the resulting score which can be positive or negative indicates a potential suitability of the CSG as immunotherapeutic target in a specific cancer; RAVEN gives you two gender specific scores
[>2 : potentially; >3 : high; >4 : very high]
- **highest tissue** → normal tissue with the highest penalty used to calculate the final CSG score
- **female specific normal tissues** → cervix, ovary, vagina, uterus endometrium, uterus myometrium
- **male specific normal tissues** → prostate, testes, penis

The male and female CSG score is mostly the same. They only differ if normal tissue used for the calculation is one of the gender specific normal tissues as listed above.

4.4 Peptide Matching Pipeline (PMP)

The Peptide Matching Pipeline (PMP) combines several bioinformatic webserver to predict T-cell binding MHC-I peptides and calculate their corresponding binding affinity. Furthermore, it also identifies all other proteins with this peptide which indicates the specificity of the peptide. A peptide specifically present in only a single protein can indicate a potential target for immunotherapy. By automatizing this time-consuming and tedious process and therefore minimizing slips, the PMP of RAVEN supports your research on targeted immunotherapy.

By using both the CSG score (see Section 4.3) calculation followed by the PMP a first draft for research on targeted cancer specific immunotherapy can be modeled.

To use the peptide matching pipeline (PMP), a stable internet connection is required. It is also important that your system allows RAVEN to access the internet to send and receive queries. Otherwise the firewall will block your request and the pipeline will break in the first step. In this case talk to your IT-administrator.

Systematic workflow of the peptide matching pipeline (Figure 4.5)

What RAVEN does for you:

1. Map the symbol of your favourite gene to UniProt to receive its uniprot accession ID (AC)
2. Get the protein sequence in FASTA format for your gene from UniProt with the AC
3. Send Query with protein sequence, allele and peptide length to IEDB
4. Receive T-cell binding MHC-I peptides predicted with artificial neural networks (ANN)
5. Search each peptide against the UniProtKB database with **Apaches Lucene**¹ text search provided by Protein Information Resource (PIR) to find all proteins comprising these peptides
6. Determine each peptides' specificity by its number of matching proteins

PMP parameter setting

What you have to do:

1. **Enter at least one query gene.** Here, RAVEN only accepts genesymbols (e.g. MYBL2, PTPN18, ABI2). Otherwise, you will receive an error message.

¹<https://lucene.apache.org/>

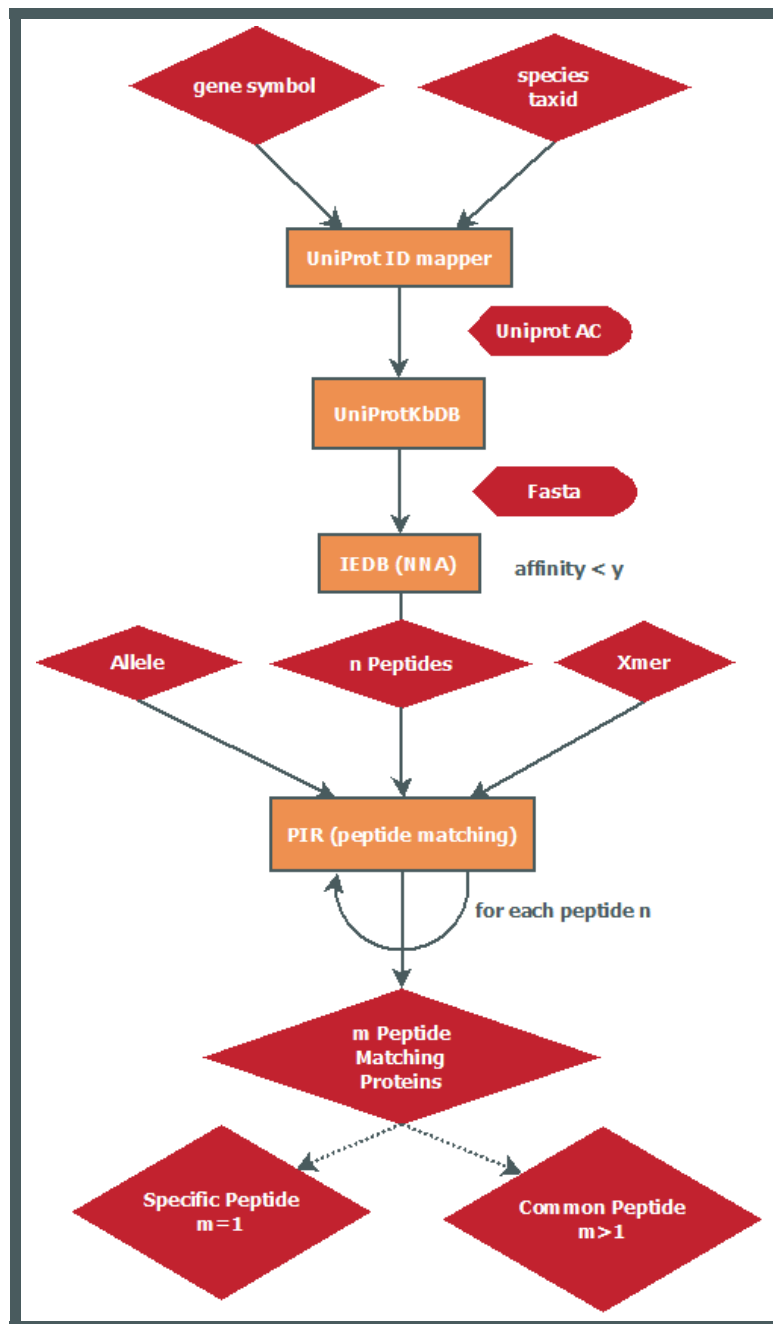


Figure 4.5: Systematic workflow of the peptide matching pipeline.

2. **Select your species.** Choose between human and mouse.
3. **Choose a database.** You can limit the proteins, which match the peptides, to Swiss-Prot entries or use the whole UniprotKB (Swiss-Prot + TrEMBL). SwissProt only contains manually annotated and reviewed proteins. Contrary, TrEMBL assembles automatically annotated and unreviewd proteins. When using both, proteins may be found repeatedly under several different IDs, which mostly refer to multiple unreviewed and obsolete sequence versions of an already validated one. To exclude such cases and avoid false positives, the IDs have to be revised manually. Thus, we recommend to use Swiss-Prot only.

4. **Specify the peptides.** Chose your favourite allele from the most common ones and confine the peptides length.
5. **Set an affinity cutoff** With this cutoff you limit your results to peptides with a stronger affinity than your set cutoff. The smaller the affinity score the stronger the MHC-I binding.
6. **Sort the results** Sorts the peptides in the result table by the chosen column. While the affinity score indicates the strength of the MHC-binding, the specificity score refers to the number of proteins comprising the peptide.
7. **Submit your query.** Via the 'view' button.

Tip: By sending too many query genes at once, RAVEN has to communicate multiple times with several web services. This susceptible procedure can fail after some time due to overloaded servers or connection issues resulting into a breakdown of the PMP and thereby a loss of already received and processed data.
To avoid this, better split your query into several batches of **not more than 100 genes** to reduce such vulnerability.

PMP result table


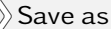
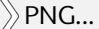
What RAVEN tells you:

- **Peptide** → sequences of predicted T-cell binding MHC-I peptides of your query gene(s) (result from the ANN)
- **Query Gene** → gene encoding the protein which comprises the peptide
- **Uniprot ID** → UniProt accession ID of the protein encoded by the query gene
- **Affinity (nM)** → Binding affinity in nM between the T-cell binding MHC-I and peptides (result from the ANN)
- **Specificity Score** → number of proteins binding to this specific peptide
- **Peptide specific proteins** → list of proteins which bind this specific peptide. For convenience to interpret the results we list the proteins by their symbols.

Save Data & Results


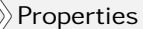
5.1 Images & Plots

You can save plots and images directly by clicking on it. With a right click of your mouse there will open a small menu which offers you different action such as print or zoom.

Do   

Altering Plots

You do not like the plot and want to change it? Alter labeling, font, colors etc. via the plot property menu. The color of points of boxes inside the plot cannot be changed.

Do  

You are still not happy and want to alter more? Save the data on which the plot is drawn and import it in another statistical visualization software such as Prism. Look at the next Section to see how you should save your data to easily import them in other applications.

5.2 Results & Data

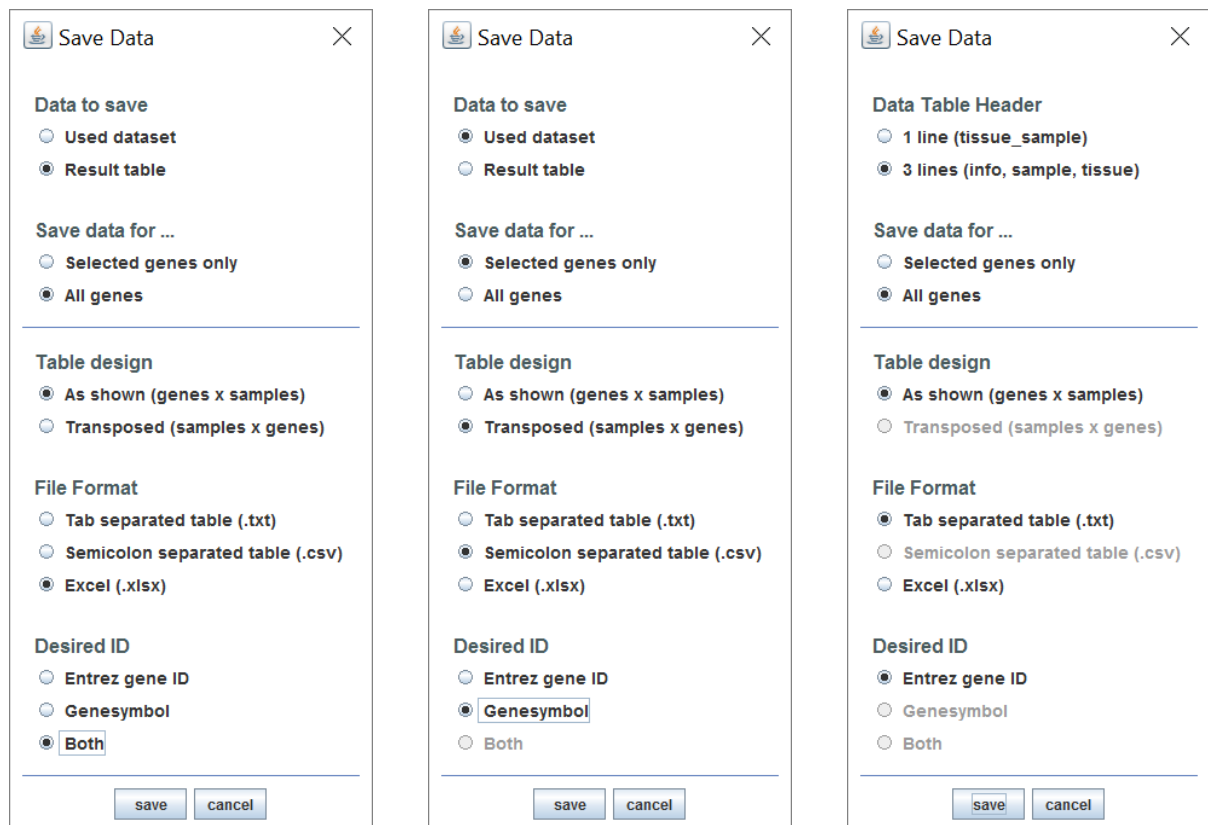
Save your results or the used data by clicking on the 'save' button. This is only possible after you clicked on the 'view' button to run and produce the results. Before this the 'save' button is disabled.

The saving menu shown in Figure 5.1 will pop up for you to define which data should be saved how.

Saving options:

Data to save	choose whether saving the result table or the data producing it.
Save data for...	do you want to save the whole table or only those you selected.
Table design	The first option will save the table as displayed on the screen (or: rows → genes, columns → samples; as in Affymetrix data). Choosing the transposed table design, the shown table will be mirrored by its diagonal (or: rows → samples, columns → genes) as favored by GraphPad Prism.

File Format	RAVEN offers 3 different formats to save your file. When saving your table to a text file (.txt) the columns will be tab separated and can be read by any text editor (e.g. Notepad++) but also in Excel. Saving your table directly as Excel file (.xlsx) is also a possibility. Some programs such as GraphPad Prism favor a comma separated text file (.csv).
Desired ID	When your table corresponds to genes you can choose between entrez gene ID or gene symbol as identifier. The first is easier to process with other programs the later might be more intuitive for your interpretation. For result tables it is convenient to include both IDs.
Table Header	This option is only available for the subset creating method. You can export a classical one line header table, where the column name includes both tissue and sample ID (e.g. tissue_sample). Instead also a header of 3 lines is possible splitting this information into 3 cells per sample (info → tissue → sample •) as required for RAVEN input files.



(a) How to save whole result table.

(b) How to save data to import and use with GraphPad Prism.

(c) How to save data to upload in RAVEN.

Figure 5.1: 'Save Data' windows to choose your saving options for data files and result tables.

Save data to use with Prism

If you want to redo the plots or perform other statistical methods with GraphPad Prism on selected genes you identified with RAVEN methods, you can save the original data in a special data format to easily import it into Prism. Save your file with the options shown in Figure 5.1b.

How to import into Prism:

open Prism » start new project » choose your 'table & graph' type » select 'empty data table' » create
File » Import... » choose your file (.csv)

Save data to use with RAVEN

Saving a data set which can be uploaded in RAVEN can only be done with the subset creating method in 'View Data Table' (see Section 4.1). Create a subset from your current data set and click the 'save' button. The table options differ from those for the other methods. Save your file with the options shown in Figure 5.1c.

About & Contact

6.1 Citation

When using our software, please cite our paper:

Systematic identification of cancer-specific MHC-binding peptides with RAVEN

Baldauf and Gerke *et al.*

DOI: <https://doi.org/10.1101/193276>

6.2 Authors

Development, Implementation & Design

Julia Sophia Gerke, MSc

Algorithm

Michaela Baldauf

Idea & Supervision

Thomas G.P. Gruenewald, MD, PhD

Visit us...

... online on 

www.lmu.de/sarkombiologie



... or in person 

AG Gruenewald

Max-Eder Research Group for Pediatric Sarcoma Biology

Institute of Pathology

Ludwig-Maximilian-University Munich

Thalkirchner Str. 36

80337 Munich, Germany

✉ thomas.gruenewald@med.uni-muenchen.de

☎ ++49-89-2180-73716

6.3 Questions & Error Reporting

If you have problems with RAVEN or unsolved questions that could not be answered with this manual, feel free to write us. For handling or technical issues please have a look at Section 6.4 first and contact our developer (julia.gerke@med.uni-muenchen.de). With any other questions Dr. Gruenewald (thomas.gruenewald@med.uni-muenchen.de) will be happy to help you.

Error Reporting

When reporting technical issues or problems with the handling please send us the following information regarding your problem if possible. This will help us to solve your problem much faster. Thank you in advance.

- Which operating system do you have? How much RAM do you have? See Table 1.1 or Figure 1.1 (only if you have trouble with the installation or data upload process!)
- What did you do?
- What did you want to achieve with it?
- Did you get an error message? Which one?
- If you started the application via command line you might see an error message written on the command line. If this is the case, please copy it into your email.

6.4 FAQ

I have RNAseq data instead of microarray. Can I analyze it with RAVEN anyway?

RAVEN is optimized to analyze microarray gene expression data. However, you can also run the software with your own RNAseq data, as long as you prepare your data accordingly. Make a quality assessment on your data and preprocess it correctly. Afterwards, normalize your reads and adapt the input file specifications necessary to upload data to RAVEN.

Does RAVEN only accept human gene expression data?

RAVEN is optimized to work with gene expression data from human tissues and cancers. You can also upload data from other species but you have to abandon gene symbols. Gene symbols are only available for humans. As the PMP operates independent from your data set you can use gene symbols from both human and mouse.

Why are some notes or pop-ups not in English?

Some implemented notes or pop-ups automatically fit to your computer systems language. Change it to English to have a uniform language in RAVEN.

Why can I see only one outlier circle in my boxplot for an outlier number of three

The circle in the boxplot is an average outlier summarizing several outliers. In your table you can look up the total number of outliers. This number also may include extreme outliers (farouts) which are shown as triangles in your boxplots.

Why did the peptide matching pipeline break early?

The pipeline broke when the red progress bar is not moving anymore.

After sending a query to a web server, RAVEN waits for its response. Sometimes the server needs too long to respond and cancels the job itself after some time. Without the job's result the pipeline cannot continue and terminates early. In this case try it again later or even the next day.

Why is the PMP not making progress anymore?

The time RAVEN needs to process a gene via the PMP can vary a lot. Depending on the amount of predicted peptides per gene the time can vary from several seconds to half an hour and more. In this case the red progress bar is still moving but not increasing

the displayed percentage. By decreasing the affinity cutoff you can reduce the amount of peptides.

However, if your cutoff is less than 150 and the displayed percentage did not increase during the last hour the peptide matching process of the server got stuck. Quit it and try the PMP again. After a restart it should work again.

Why does the program not respond anymore?

If you do not have the minimum requirement to run RAVEN, the program most likely will get stuck during the data upload phase. Decrease the size of your input file or increase your RAM.

Why is the entered gene ID, name or symbol not accepted?

First, check if you used the allowed IDs which are entrez gene ID or gene symbol only. Control your spelling. Otherwise, this can occur due to a less common synonym or alias of your gene symbol. You can verify your gene symbol via [GeneCards](http://www.genecards.org/)¹. If all this does not apply to, the gene might be missing in your data set.

Why is my gene missing from the peptide table in the PMP?

A missing gene is either excluded from the beginning (see below) or is dropped during the pipeline process depending on your parameter. If all detected peptides have a higher affinity than your set affinity cutoff, the corresponding gene will be directly excluded due to its missing results. Increase the affinity cutoff accordingly.

Why did RAVEN exclude some genes from my query list in the peptide matching pipeline (PMP)?

This can have two different reasons:

- **non-coding gene:** The PMP only works for protein-coding genes. If your query is a RNA gene or pseudo gene. The pipeline will not find a protein fasta sequence. Therefore, we excluded them from the query. In this case however, you will get a notification at the end of your processed job, informing you which "genes" were excluded from your query genes in advance.
- **filtering** Missing protein-coding genes were filtered out by your set parameters. With the affinity cutoff you limit the amount of resulting peptides which should be checked for their specificity. If your cutoff is set too low, no peptide with such a strong binding affinity can be found and thus the gene will not show up in your result table.

¹<http://www.genecards.org/>

Why is the button not working?

This can have different reasons:

- **the button is disabled:** If RAVEN disables a button (→ grey writing) you are not allowed to perform this action at the moment. For example, the 'save' button is disabled if there is no data (table) there which can be saved. Run the method first, before saving its results or used data.
- **several methods are running at the same time:** You are running too many jobs in RAVEN at the same time. Wait until RAVEN is done before starting a new job or method. To avoid such an overload run one RAVEN method after another.
- **the previous query was not finished yet:** The current job with your query is still running, RAVEN will not allow you to start a new one until the previous is finished. Wait until RAVEN is finished before sending the next query. You can cancel a running job by closing RAVEN.

Why do I get the following error messages when using the PMP?

- **WebServiceConnectionException:** the PMP could not connect to the server providing the web service. In most cases, this is a temporary problem. If this error still occurs 2 days later, please contact us.
- **TimeoutException:** The web service takes too long to respond. Try it again a few hours later.

Why is the header/row name, containing the sample IDs, missing when saving the data of Expression Comparison

Each expression value belongs to another sample. Values in the same row/column do not belong to the same sample. Contrary, the samples in the same column/row represent a group. So, the sample IDs are not necessary for this group visualization (box plot) and are rather misleading information.

Why did the table not update after changing the parameter?

After changing the parameters, you always have to press the 'view' button to display the updated table or plot.

RAVEN is not responding anymore! What can I do?

Currently RAVEN seems to be too busy to accept new intake from you. Do not click around and confuse RAVEN even more. Give it some time, it probably will be back in a few minutes. After 15 minutes without response, you can abort and restart RAVEN. In this case it is recommendable to use a computer with more power (see Section 1.1).

6.5 License

Copyright 2017 Gerke

Licensed under the Apache License, Version 2.0 (the "License"); you may not use this file except in compliance with the License. You may obtain a copy of the License at

<http://www.apache.org/licenses/LICENSE-2.0>

Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.