



---

# PCP: An R-based application for visualization data from protein correlation profiling

Prepared for Noah Dephoure, Weill Cornell Medical College  
by the Applied Bioinformatics Core (ABC), Weill Cornell Medical College

revised June 13, 2017

## Contents

<b>1 Installation and updating</b>	<b>2</b>
<b>2 Start</b>	<b>2</b>
<b>3 Using the app - overview</b>	<b>2</b>
<b>4 Data import</b>	<b>2</b>
4.1 Metadata . . . . .	3
4.1.1 Editing the metadata already stored in the data base . . . . .	4
<b>5 Visualizing protein profiles</b>	<b>4</b>
5.1 Interacting with the protein profiles . . . . .	4
5.2 Values for display . . . . .	5
5.3 Customizing the colors . . . . .	5
5.4 Summary plots . . . . .	6
5.5 Individual proteins . . . . .	6
5.6 Spike-in quantification . . . . .	8
<b>6 Tables</b>	<b>9</b>
<b>7 Correlations</b>	<b>9</b>
7.1 Reducing the numbers of displayed proteins . . . . .	10
7.2 Defining proteins of interest for protein profile plots . . . . .	10

## 1 Installation and updating

The routine for installing and updating the package is the same (with the exception of the `devtools` installation, which must be done only once).

1. Open RStudio.
2. In the **Console** section of RStudio, type:

```

1 # install helper package (must be
   done only once!)
2 > install.packages("devtools")
3
4 # download the source of the
   package and install it
5 # note that the version number
   will change depending on the
   update
6 # you can check which version is
   appropriate at http://chagall.
   med.cornell.edu/deplab-pcp/
7 > devtools::install_url(url =
   http://chagall.med.cornell.edu
   /deplab-pcp/DepLab_0.1.1.9000.
   tar.gz")
8
9 # attach the package to your
   current R environment
10 > library(DepLab)
11
12 # start the interactive app
13 > runPCP()

```

## 2 Start

If you have installed the most recent version of the package and you just want to continue using the app, all you need to do is:

1. Open RStudio.
2. In the **Console** section of RStudio, type:

```

1 > library(DepLab)
2 > runPCP()

```

This will start the application in a new window.

## 3 Using the app - overview

There are three main parts to the app:

<b>Data Import</b>	specify the file that contains the results of a protein correlation profiling experiment and add it to the data base
<b>Visualization</b>	generate plots for all values that MaxQuant reported per eluted fraction
<b>Tables</b>	explore the values underlying a specific plot in tabular format

While the Data import panel will always be shown at the start of the app, you can skip right to the visualization or tables part without uploading a new file. You will, however, need to make sure you have selected the correct data base.

## 4 Data import

The app allows for the import of a table of proteins that were identified during a protein correlation profiling experiment. Every file that is uploaded will be stored in a relational `sqlite` database. This allows for storing (and retrieval) of multiple data files in a consistent manner including customized metadata entries, such as information about the specific experiment, the date of data acquisition etc. (see Section ??). This also means that the database will grow over time as more and more files are added.

! The data base should be regularly backed up.

Once you start the app, you will be shown the path to the database that you used during your last session. You will also see the experiment IDs of the data sets that were previously stored. You can change the path to the data base via the **Select database** button, but we strongly recommend to keep all the data in *one* database rather than scattering information from different experiments into different databases.

### SELECT EXISTING DATABASE

Currently selected database: <code>/Users/frd2007/Desktop/proteomics.db</code>
The currently selected database contains the following data sets <code>sc_test_100</code> <code>sc_test_300</code>
<b>Change database</b>
The path to store custom complexes is set to: <code>/Users/frd2007/Desktop/complexes_custom.txt</code>
<b>Change file for custom complexes</b>

If you're fine with these selections and you do not wish to upload new data, simply con

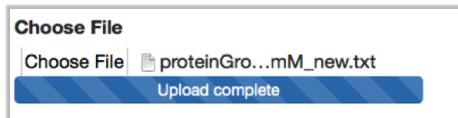
In addition, you can check and change the path to the text file where the custom-defined complexes are stored.

To upload a file, the current workflow is as follows:

1. Select the “Data Import” tab.



2. Click on “Choose File” and navigate to the MaxQuant output file. Currently, the app accepts the `proteinGroups*.txt` output of MaxQuant [4, 6]. For each protein group, only the founder protein will be extracted. There should be a bar indicating that the upload was successful.



3. Specify a *unique experimental ID* which will help you identify the experiment later on.
4. Select the *organism* from which the samples of that particular experiment originated from.
5. If you know the UniProt ID of the protein you have used as a *standard* for that particular experiment, indicate that in the field below the experimental ID. If you do not specify a UniProt ID, the default is to retrieve the values related to trypsin (from pig) in that data set.
6. Fill out all the relevant metadata connected to that particular experiment (for details, see Section 4.1).
7. Click “Save”. Since this will load all the entries from the MaxQuant output into the database, this may take a while. Once the updating of the database has succeeded, you will see a corresponding message.



**!** It is not possible to upload a file with a previously used experimental ID. Metadata, however, can be updated and corrected after upload (see Section 4.1.1).

## 4.1 Metadata

As of DepLab v0.1.1, you will be required to provide some basic metadata, i.e., additional information about the experiment for which you are uploading the data. The mandatory entries include the name of the experimenter, the genotype, harvest date, lysis method etc.

The screenshot shows a form titled "Sample origin and preparation". At the top right is a "Save" button. A modal window titled "Oops" appears, stating "Experimental ID, Maxquant data, and sample origin metadata required before saving." An arrow points from this modal down to the main form area. The main form contains several dropdown menus and input fields. On the left, under "mandatory metadata", are fields for Name of experimenter (Paola), Genotype (WT), Harvest date (2016-12-13), Lysis method (Dounce), Cell type (MCF10A), Buffer composition (Tris\_HCl), and Digestion enzyme (Trypsin). Below these are "Notes" and "Prefractionation method" (with an "Include?" checkbox). To the right, there are sections for "Mass spectrometry method" and "Data processing", each with its own "Include?" checkbox. A bracket on the right groups these optional sections, labeled "additional, non-mandatory metadata".

Additional types of metadata can be added:

- **Prefractionation method** captures information such as column ID, amount of loaded protein, sample volume, time per and number of fractions.
- **Mass spectrometry method** expects information about the Instrument ID, the run date (default will be the current day) and the length of the method.
- **Data processing** lets you keep track of basic software settings that were used for the processing of the spectra, such as the search and filtering algorithms.

If you need more pre-defined entries – e.g., for the name of the experimenter – get in touch with us ([frd2007@med.cornell.edu](mailto:frd2007@med.cornell.edu)).

You can check the metadata entered for the experiments within the selected database in the Database section (DB Viewer).

The screenshot shows the DB Viewer interface. At the top, there are two tabs: "DB Viewer" (which is selected) and "DB Editor". Below the tabs, a dropdown menu titled "Select expt ID to view" contains the option "mcf10a\_D109N\_01". The main content area is titled "Sample origin and preparation" and lists the following experimental details:

- Experimental ID: mcf10a\_D109N\_01
- Organism: human
- Name of experimenter: Paola Cavaliere
- Genotype: D109N
- Cell type: MCF10A
- Harvest date: 2016-11-18
- Buffer composition: TrisHCl 50mM pH 7.5 KCl 150mM
- Lysis method: Dounce
- Digestion enzyme: Trypsin
- Notes:

#### 4.1.1 Editing the metadata already stored in the data base

To edit the metadata related to experimental data within the selected database, go to the DB Editor in the Database section. Just select the experiment ID of the experiment for which you would like to change the metadata. Then add the respective changes and hit "Save changes".

The screenshot shows the DB Editor interface. At the top, there are two tabs: "DB Viewer" and "DB Editor" (which is selected). Below the tabs, a dropdown menu titled "Select expt ID to edit" contains the option "mcf10a\_D109N\_01". There are two buttons at the bottom: "Edit" (which is highlighted with a black circle) and "Delete - are you sure?".

#### Basic sample info

You can also delete all information related to an entire experiment by clicking the Delete button.

## 5 Visualizing protein profiles

The screenshot shows the Visualization tab interface. At the top, there are three tabs: "Data Import" (selected), "Visualization", and "Tables".

There are three kinds of plots that can currently be generated to visualize the fraction-wise protein profiles:

1. **Summary plot:** Depicting the sums of values for all protein groups per fraction.

2. **Individual proteins:** Tracking the measured values across all fraction for individual proteins and complexes.

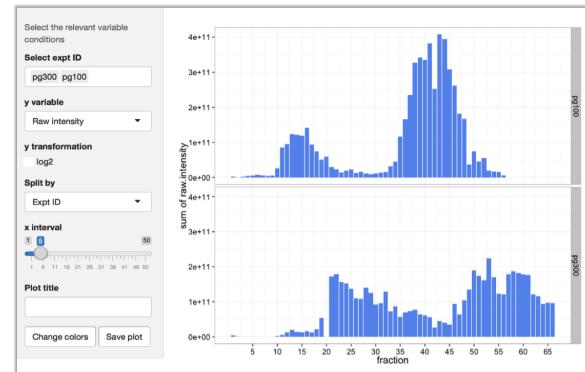
3. **Spike-in:** Tracking the values for the protein spike-in controls across the fractions.

The general steps to obtain a plot are very similar for all three types:

1. Select the *type of plot* that you would like to generate.
2. Select *one or more experiments*. The field "Select expt ID" will contain all experimental IDs stored in the database. You can choose by scrolling through the list or typing the specific ID.

The screenshot shows the plot selection interface. At the top, there are three tabs: "Summary plots" (selected), "Individual proteins" (highlighted with a blue border), and "Spike-in quantification". A callout points to the tabs with the text "Select the type of plot". Below the tabs, there is a section titled "Select the relevant variable conditions" and a "Select expt ID" dropdown. The dropdown contains the options "pg\_300mM", "pg300", and "pg100". A callout points to the dropdown with the text "Select the experiment".

Once you select the experiment ID, a plot will immediately be generated using default settings that currently include to use the raw intensity for the *y* variable and to make a separate plot for each selected experiment ID.



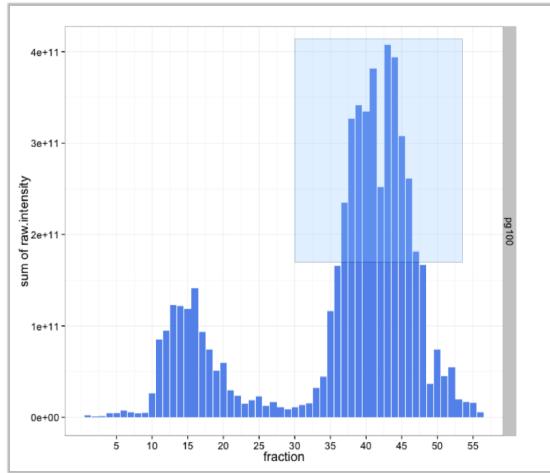
All images can be saved as pdf or png files ("Save plot").

#### 5.1 Interacting with the protein profiles

The plots can be customized in various ways. Every change will immediately refresh the image.

The *size of the plot* is connected to the size of the app's display. If you reduce the size of the app's window, the plot will be adjusted accordingly.

To *zoom* into certain fractions, simply use the mouse to draw a rectangle around the region of interest. Then double-click into the light blue box.



To zoom out to the original scale, double click anywhere within the image.

To *retrieve basic information* from UniProt for each depicted molecule, click on the respective link underneath each plot (scroll down if you don't see it immediately). Since the download from UniProt takes some time, this list is not automatically updated if you change the selection of proteins. You will have to refresh that list manually. If you click on the individual UniProt IDs within the result table, the corresponding web site will be opened in your Internet browser.

Refresh UniProt info...						
gene	orf	sgd_id	protein	organism	status	size
ABP1	YCR088W	S000000684	P15891	ABP1_YEAST	Reviewed;	592
ACT1	YFL039C	S000001855	P60010	ACT_YEAST	Reviewed;	375

↑ link to UniProt entry

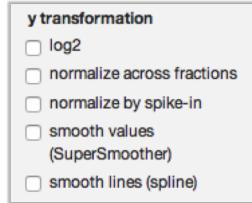
## 5.2 Values for display

The following measurements are available for all three kinds of plots ("y variable"):

- Raw intensity
- LFQ intensity
- MS count
- Peptides count
- Unique peptides only
- Razor and unique peptides
- Sequence coverage

These are all the numerical values that MaxQuant supplies per fraction.

For the *individual protein plots*, these numerical values can be modified to improve the effectiveness of the visualization.

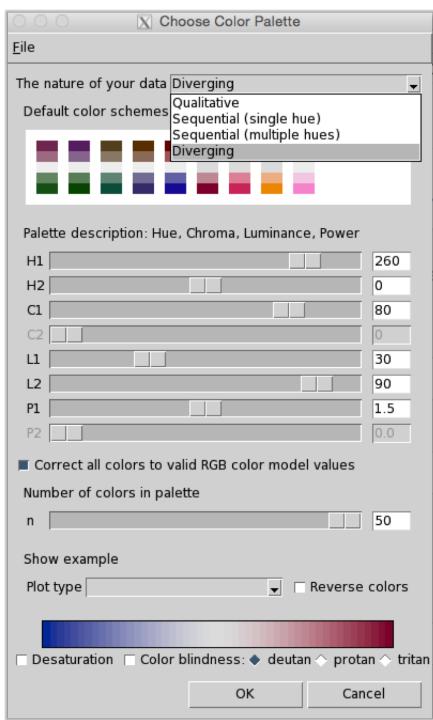


- **Normalize across fractions:** Each measured value,  $m_f$ , is normalized so that proteins with equal measurements in each fraction,  $f$ , have a normalized value of 1 per fraction:  $m_{norm} = \frac{m_f \times F}{\sum_{f=1}^F m}$ .
- **Spike-in control:** Each measured value,  $m_f$ , is normalized to the corresponding value,  $s_f$ , of the spiked-in protein:  $m_{norm} = \frac{m_f}{s_f}$ .
- **Smooth values (SuperSmoother):** This will apply Friedman's SuperSmoother, which is a sophisticated moving average method [5].
- **Smooth lines (spline):** This applies a function that interpolates the numerical values using polynomial functions [3].

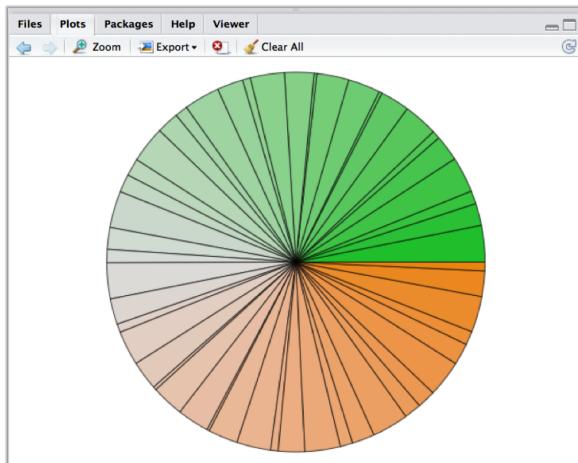
## 5.3 Customizing the colors

You can change the color scheme via the button "Change colors". Once you click on it, a separate window will appear which will allow you to change the type of color scheme (qualitative, sequential, diverging), choose the colors, manipulate hue<sup>1</sup>, chroma<sup>2</sup>, luminance<sup>3</sup> and power, and specify the numbers of colors that will be used.

1. Hue: The degrees on the RGB color wheel where Red = 0, Green = 120, Blue = 240.  
 2. Chroma: The saturation of the hue indicated with values 0 (black) to 255.  
 3. Luminance: Perceived brightness of a color; depends on the hue as well as the saturation.

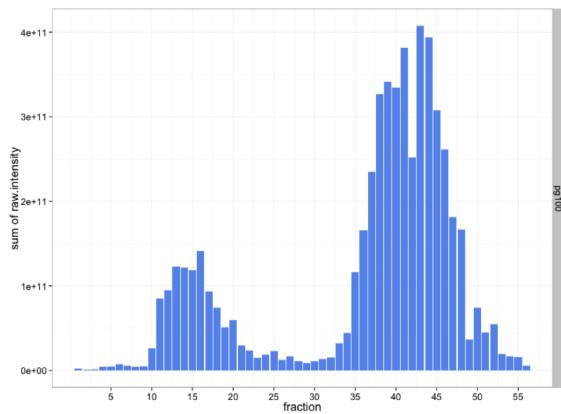


The bar at the bottom will reflect your choices and should give you a good idea of how your plot will be affected. In addition, you can choose to see an example plot, e.g. a pie chart. Note that this example will not be shown within the app, but in RStudio's Plot environment!



## 5.4 Summary plots

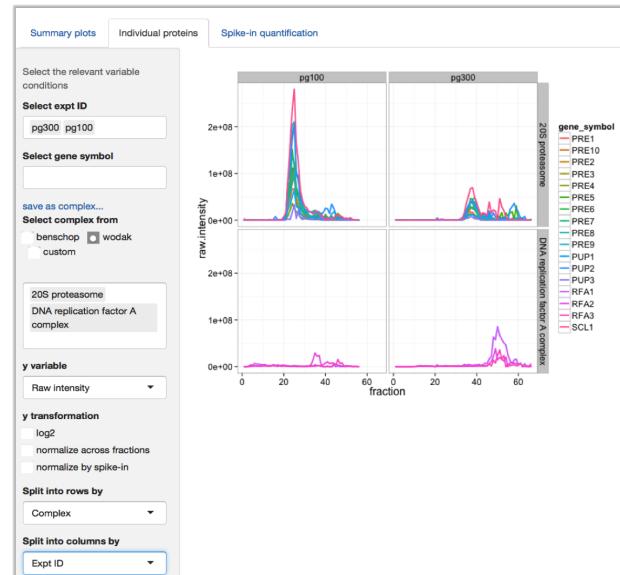
Summary plots display the sum of all measurements per fraction. The default setting will make one plot per selected experiment ID. If you would like to visually distinguish the different fractions, choose a different color scheme.



Since the summary plots will have to retrieve and compute *all* values per experiment, the generation of the image may take a couple of seconds.

## 5.5 Individual proteins

The individual protein plots are useful to track the measurements for single protein across the different fractions.



As for the summary plots, start by *selecting the experiment ID(s)*. Once you have selected the first experiment ID, the list of possible additional experiment IDs will be reduced based on the organism of your first selection. This way, you will always only be able to compare samples from the same organism. This also means that you should wait for the successful loading of the first data set which will be indicated by a message stating the organism of your choice.

Select expt ID  
sc\_500  
Data from: yeast

To display protein profiles of interest, select them by their gene name or their UniProt ID, e.g. "GAPDH".

P04406 (GAPDH)

**Customizing the plots** There are numerous options to customize the display to allow for useful comparisons, e.g., across experiments and between the members of distinct protein complexes.

Typically, it will make sense to *assign the color based on the gene symbol*. The legend will always be sorted alphabetically, and so will the table of basic information that you can retrieve from UniProt.

To compare multiple complexes across different experiments, we recommend to *split by rows* using the experiment IDs and to *split by columns* using the complex' IDs.

You can also choose to differentiate different gene symbols by the point shape or line type, but note that this makes only sense for very limited numbers of variables (up to 5 different genes, for example).

Point shape  
dot  
plus  
cross  
diamond  
Assign to gene symbol  
Assign to ORF  
Assign to SGDID  
Assign to expt ID  
Assign to complex

Split into rows by  
Expt ID  
None  
Gene symbol  
ORF  
SGDID  
Expt ID  
Complex

If you have replicates for the same experimental condition, you can specify both, the name of the condition and numeric replicate labels. For this, you will have to check the box **Specify replicates** underneath the box for the selection of the experiment IDs. You can then manually define the name of the experimental condition, which can then be used to, for example, split the plots.

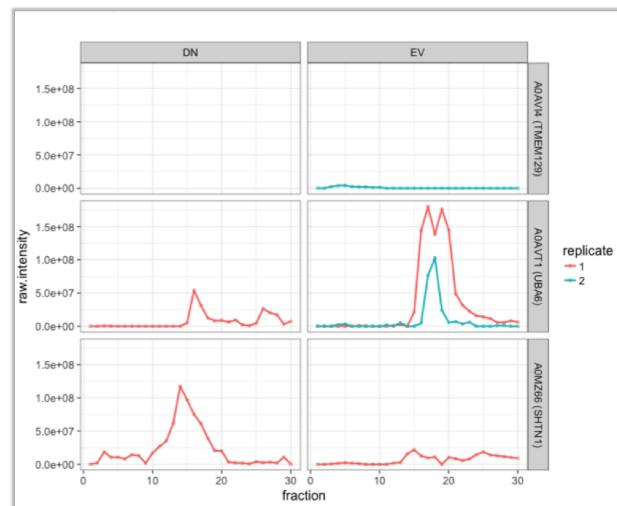
Reset inputs  
Select expt ID  
mcf10a\_EV\_01  
mcf10a\_EV\_02  
mcf10a\_D109N\_01  
Data from: human  
 Specify replicates  
Condition\* Rep.\*  
mcf10a\_EV\_01  
EV 1  
mcf10a\_EV\_02  
EV 2  
mcf10a\_D109N\_01  
DN 1

define names for the conditions and choose integers to distinguish different replicates

A useful combination of settings may be:

- Split into rows by: Gene symbol
- Split into columns by: Condition
- Color by: Replicate

This would result in the following figure:



**Normalization** The values displayed on the y-axis can be *log*-transformed or normalized using two approaches: either across fractions or by using a spike-in control. When normalizing based on the values for spike-in controls, you will be able to choose the specific protein if more than one standard was indicated in the metadata section for any experiment.

The screenshot shows the 'Select expt ID' interface. At the top, there's a search bar with 'pg100 pg300'. Below it, a message says 'You have selected yeast'. A 'Select gene symbol' section contains 'ATP20 CYC8 ERP1'. Under 'y variable', 'Raw intensity' is selected. In 'y transformation', 'normalize by spike-in' is checked. On the left, 'experiment IDs' are listed: 'pg100' (selected), 'P00761', 'pg300' (selected), and 'P00761'. A note on the right says: 'once you select "spike-in" normalization, the corresponding available proteins are indicated'.

For details on the normalization of the y variable, see Section 5.2.

**Complexes** Instead of specifying individual gene symbols that should be displayed, you can also select a set of proteins that have been reported to be part of common complexes. We currently supply three lists of complexes: two for yeast based on publications by Benschop et al. [2], and Wodak et al. [7], and one for human complexes, CORUM Core [1].

To see the profiles for proteins of a complex defined by either one, e.g., select “Wodak”, then use the empty field below to search for the complex of interest, e.g. “20S Proteasome”. You can select multiple protein complexes. To add additional proteins, use the field “Select gene symbol”.

The screenshot shows the 'Select complex from' interface. A radio button 'wodak' is selected. Below it, a list of complexes is shown: '1,3-beta-glucan synthase complex (Fks1p/Rho1p)', '1,3-beta-glucan synthase complex (Gsc2p/Smk1p)', '19/22S regulator', '20S proteasome', and '6-phosphofructokinase complex'. A note on the right says: 'list of complexes within the selected list'.

Many complexes may encompass proteins for which your experiment does not provide any values. This will cause NA entries.

If you want to limit the complexes that are available in the auto-select box to those for which many proteins were recovered in your experiment(s), you can specify the minimum number of proteins per complex that should be present in *every* experiment.

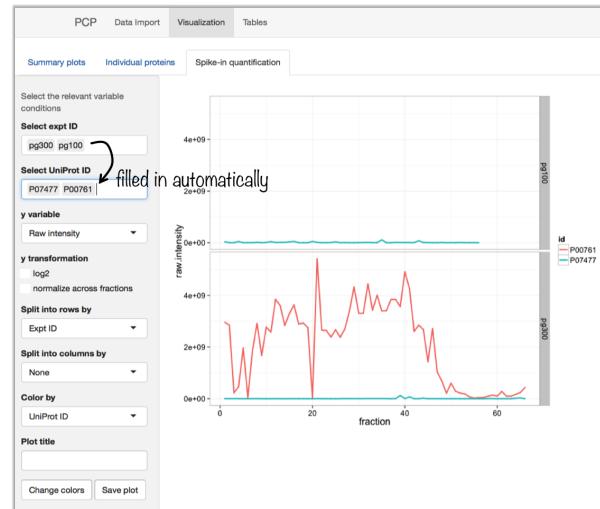
The screenshot shows the 'Available complexes' interface. It has a text input 'Define the min. # of recovered proteins per complex' with the value '2'. A note on the right says: 'the number of recovered proteins will limit the number of complexes that can be selected'.

In addition to the pre-defined complexes, you can *generate your own groups of proteins*. Simply choose the individual proteins that you would like to group together via the “Select gene symbol” line, then hit “*save as complex...*”. You will be asked to specify a name for the protein group. Once you have saved a group, it will be available in all future sessions: simply select the “*custom*” field and use the empty line below to select your custom-made group of interest.

## 5.6 Spike-in quantification

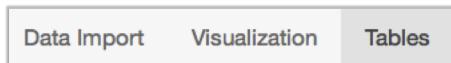
Ideally, each experimental ID should have the corresponding UniProt ID for the spike-in protein that was used. Therefore, all you will have to do here is to indicate the experiment(s) of interest and the profiles for the spike-ins related to the experiment(s) will be shown.

If no UniProt ID was specified during the upload of the file (as described in Section 4), Trypsin from pig (UniProt ID P00761) will be used.



The field “Select UniProt ID” will be automatically populated depending on the experiment ID. You can delete individual entries, but you cannot add spike-ins. Use the “Individual proteins” plots for assessing individual protein profiles.

## 6 Tables



The tables section allows you to see the values that underlie the summary plots and individual protein profiles. Conversely, if the visualization section is empty (e.g., because you haven't selected an experimental ID), the corresponding table section will be empty, too. The following example shows the table for the summary plot.

gene_symbol	fraction	value	measurement	expt_id
AAC3	1	0	raw.intensity	pg100
AAC3	2	0	raw.intensity	pg100
AAC3	3	0	raw.intensity	pg100
AAC3	4	0	raw.intensity	pg100
AAC3	5	0	raw.intensity	pg100
AAC3	6	0	raw.intensity	pg100
AAC3	7	0	raw.intensity	pg100
AAC3	8	0	raw.intensity	pg100
AAC3	9	0	raw.intensity	pg100
AAC3	10	0	raw.intensity	pg100
AAC3	11	0	raw.intensity	pg100

Each row corresponds to *one* fraction per protein per experiment ID, which is different from MaxQuant's table where each row corresponds to *all* fractions per protein and experiment ID.

Select/indicate the relevant variable conditions. Note that the entries for Condition and Replicate will determine the labels printed in the heatmap.

Cond.*	Repl.*	Select Expt. ID*	
WT	1	mcf10a_WT_01	<a href="#">add</a>
WT	2	mcf10a_WT_02	<a href="#">delete</a>
EV	1	mcf10a_EV_01	<a href="#">delete</a>
EV	2	mcf10a_EV_02	<a href="#">delete</a>
D	1	mcf10a_D109N_01	<a href="#">delete</a>

**Measurement**

Smoothed

**Correlation method**

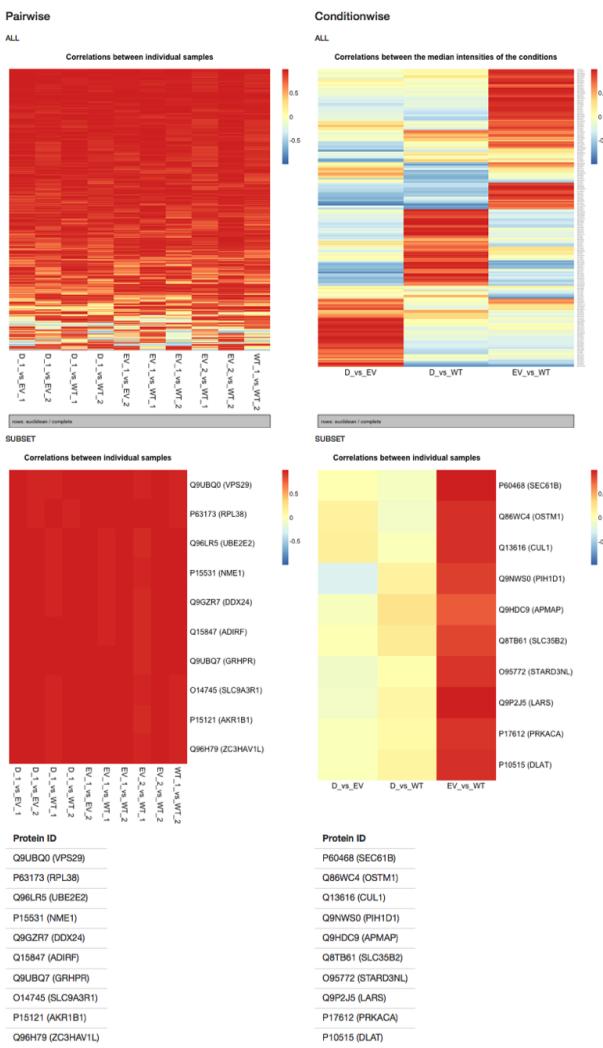
Pearson

[Calculate correlations and plot heatmap](#)



Define Condition (Cond) and replicate numbers (Repl.) that correspond to the experimental ID.

Currently saving has to be done via right-mouse click



## 7.1 Reducing the numbers of displayed proteins

Specify the range for the correlation values; this will limit the number of proteins shown in the heatmap.

### Corr. range for pairwise

comp.:  
1



Corr. range for  
conditionwise comp.:



Select a subset of proteins from the top or bottom of the first set of heatmaps which will be shown in more detail below.

- top
- bottom

Limit the range of correlation values

specify the number of proteins shown in the subsets

## 7.2 Defining proteins of interest for protein profile plots

## References

- [1] CORUM: The comprehensive resource of mammalian protein complexes-2009. *Nucleic Acids Research*, **38**(SUPPL.1), 2009. doi:10.1093/nar/gkp914.
- [2] Benschop JJ, Brabers N, van Leenen D, Bakker LV, van Deutekom HWM, van Berkum NL, Apweiler E, Lijnzaad P, Holstege FCP, and Kemmeren P. A consensus of core protein complex compositions for *Saccharomyces cerevisiae*. *Molecular Cell*, **38**(6):916–928, 2010. doi:10.1016/j.molcel.2010.06.002.
- [3] Blanc C and Schlick C. X-splines: A spline model designed for the end-user. In *Proceedings of the 22nd Annual Conference on Computer Graphics and Interactive Techniques*, pp. 377–386, 1995. doi:<http://doi.acm.org/10.1145/218380.218488>.
- [4] Cox J and Mann M. MaxQuant enables high peptide identification rates, individualized p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nature Biotechnology*, **26**(12):1367–1372, 2008. doi:10.1038/nbt.1511.
- [5] Friedman JH. A Variable Span Scatterplot Smoother. Tech. Rep. October, Stanford University, 1984.
- [6] Kirchner M and Selbach M. In vivo quantitative proteome profiling: Planning and evaluation of SILAC experiments. *Methods in Molecular Biology*, **893**:175–199, 2012. doi:10.1007/978-1-61779-885-6{\\_}13.
- [7] Pu S, Wong J, Turner B, Cho E, and Wodak SJ. Up-to-date catalogues of yeast protein complexes. *Nucleic Acids Research*, **37**(3):825–31, 2009. doi:10.1093/nar/gkn1005.