

The oposSOM Package

Henry Löffler-Wirth, Martin Kalcher

October 13, 2015

High-throughput technologies such as whole genome transcriptional profiling revolutionized molecular biology and provide an incredible amount of data. On the other hand, these techniques pose elementary methodological challenges simply by the huge and ever increasing amount of data produced: researchers need adequate tools to extract the information content of the data in an effective and intelligent way. This includes algorithmic tasks such as data compression and filtering, feature selection, linkage with the functional context, and proper visualization. Especially, the latter task is very important because an intuitive visualization of massive data clearly promotes quality control, the discovery of their intrinsic structure, functional data mining and finally the generation of hypotheses. We aim at adapting a holistic view on the gene activation patterns as seen by expression studies rather than to consider single genes or single pathways. This view requires methods which support an integrative and reductionist approach to disentangle the complex gene-phenotype interactions related to cancer genesis and progression. With this motivation we implemented an analysis pipeline based on data processing by a Self-Organizing Map (SOM) (Wirth et al., 2011)(Wirth et al., 2012a). This approach simultaneously searches for features which are differentially expressed and correlated in their profiles in the set of samples studied. We include functional information about such co-expressed genes to extract distinct functional modules inherent in the data and attribute them to particular types of cellular and biological processes

such as inflammation, cell division, etc. This modular view facilitates the understanding of the gene expression patterns characterizing different cancer subtypes on the molecular level. Importantly, SOMs preserve the information richness of the original data allowing the detailed study of the samples after SOM clustering. A central role in our analysis is played by the so-called expression portraits which serve as intuitive and easy-to-interpret fingerprints of the transcriptional activity of the samples. Their analysis provides a holistic view on the expression patterns activated in a particular sample. Importantly, they also allow identification and interpretation of outlier samples and, thus, improve data quality (Hopp et al., 2013a)(Hopp et al., 2013b).

1 Example data: transcriptome of healthy human tissue samples

The data was downloaded from Gene Expression Omnibus repository (GEO accession no. GSE7307). About 20,000 genes in more than 650 tissue samples were measured using the Affymetrix HGU133-Plus2 microarray. A subset of 12 selected tissues from different categories is used here as example data set for the oposSOM-package.

2 Setting up the environment

In order to set the analysis parameters and to create the enclosing environment it is obligatory to use **opossom.new**. If any parameter is not explicitly defined, default values will be used (see also Parameters section):

```
> library(oposSOM)
> env <- opossom.new(list(dataset.name="Tissues",
+                          dim.1stLvlSom=20))
```

The oposSOM package requires input of the expression data. Usually the raw microarray intensity data is preprocessed using appropriate calibration and summarization algorithms (e.g. MAS5, VSN or RMA), and transformed into logarithmic scale prior to utilizing them in the pipeline.

The package then accepts two formats: Firstly a simple two-dimensional numerical matrix, where the columns and rows represent the samples and genes, respectively:

```
> data(opossom.tissues)
> str(opossom.tissues, vec.len=3)

num [1:20957, 1:12] 0.299 2.492 2.293 2.041 ...
- attr(*, "dimnames")=List of 2
 ..$ : chr [1:20957] "ENSG00000115415" "ENSG00000252095" "ENSG00000111640" ...
 ..$ : chr [1:12] "liver" "kidney cortex" "thyroid gland" ...

> env$indata <- opossom.tissues
```

Secondly the input data can also be given as *Biobase::ExpressionSet* object:

```
> data(opossum.tissues)
> library(Biobase)
> opossum.tissues.eset = ExpressionSet(assayData=opossum.tissues)
> opossum.tissues.eset

ExpressionSet (storageMode: lockedEnvironment)
assayData: 20957 features, 12 samples
  element names: exprs
protocolData: none
phenoData: none
featureData: none
experimentData: use 'experimentData(object)'
Annotation:

> env$indata <- opossum.tissues.eset
```

Each sample may be assigned to a distinct group and a respective color to improve data visualization and result presentations. If not defined by the user, the samples will be collected within one group and colored using a standard scheme.

```
> env$group.labels <- c(rep("Homeostasis", 2),
+                       "Endocrine",
+                       "Digestion",
+                       "Exocrine",
+                       "Epithelium",
+                       "Reproduction",
+                       "Muscle",
+                       rep("Immune System", 2),
+                       rep("Nervous System", 2) )

> env$group.colors <- c(rep("gold", 2),
+                       "red2",
+                       "brown",
+                       "purple",
+                       "cyan",
+                       "pink",
+                       "green2",
+                       rep("blue2", 2),
+                       rep("gray", 2) )
```

Alternatively, the *group.labels* and *group.colors* can also be defined within the phenotype information of the ExpressionSet:

```
> group.info <- data.frame(
+   group.labels = c(rep("Homeostasis", 2),
+   "Endocrine",
+   "Digestion",
+   "Exocrine",
+   "Epithelium",
+   "Reproduction",
+   "Muscle",
+   rep("Immune System", 2),
+   rep("Nervous System", 2) ),
+   group.colors = c(rep("gold", 2),
+   "red2",
+   "brown",
+   "purple",
+   "cyan",
+   "pink",
+   "green2",
+   rep("blue2", 2),
+   rep("gray", 2) ),
+   row.names=colnames(opossum.tissues))

> opossum.tissues.eset = ExpressionSet(assayData=opossum.tissues,
+   phenoData=AnnotatedDataFrame(group.info) )
> opossum.tissues.eset

ExpressionSet (storageMode: lockedEnvironment)
assayData: 20957 features, 12 samples
  element names: exprs
protocolData: none
phenoData
  sampleNames: liver kidney cortex ... cerebral cortex (12 total)
  varLabels: group.labels group.colors
  varMetadata: labelDescription
featureData: none
experimentData: use 'experimentData(object)'
Annotation:

> env$indata <- opossum.tissues.eset
```

Finally the pipeline will run through all analysis modules without further input. Periodical status messages are given to inform about running and accomplished tasks. Please note that the tissue sample will take approx. 30min to finish, depending on the users' hardware:

```
> opossom.run(env)
```

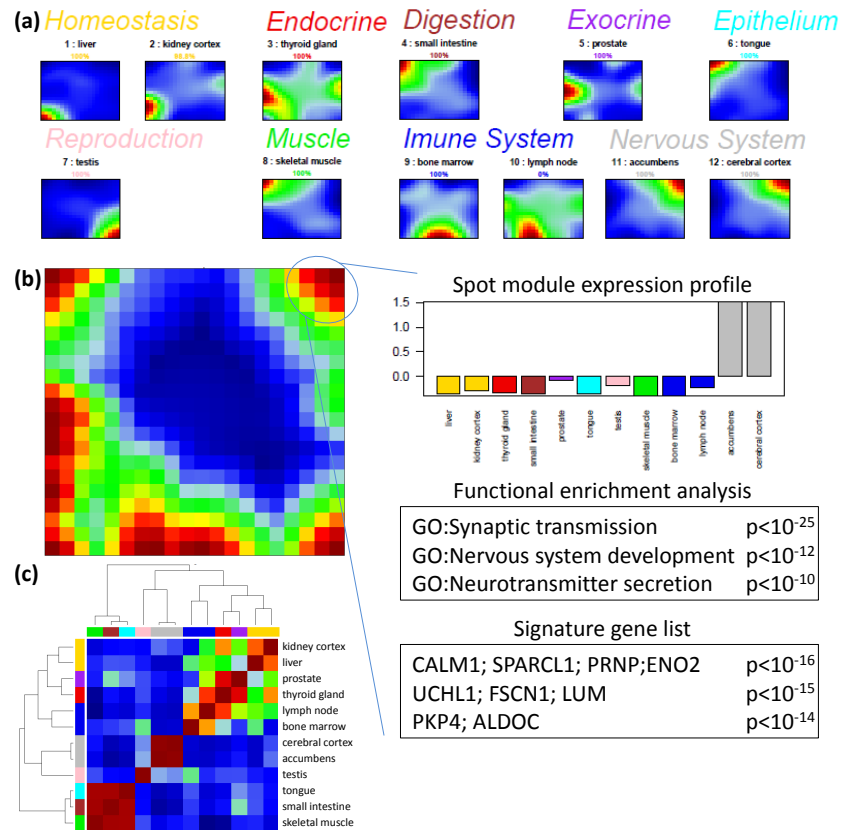


Figure 1: Few selected results provided by the oposSOM package: (a) Expression landscape portraits represent fingerprints of transcriptional activity. The *group.labels* and *group.colors* parameters are used to arrange and represent the samples throughout all analyses. (b) Functional expression modules are identified in the expression landscapes and described using appropriate summary portraits (left part), and expression profiles, enrichment analyses and differential gene lists (right part). (c) Sample similarity structure is analysed using different algorithms and distance metrics. Here a clustered pairwise correlation matrix is shown.

3 Browsing the results

The pipeline will store the results in a defined folder structure. These results comprise a variety of PDF documents with plots and images addressing the input data, supplementary descriptions of the SOM generated, the metadata obtained by the SOM algorithm, the sample similarity structures and also functional annotations. The PDF reports are accompanied by detailed CSV spreadsheets which render the complete information richness accessible.

Figure 1 shows few selected outputs generated by the pipeline. The expression landscape portraits (Figure 1a) represent fingerprints of transcriptional activity. They are used to identify functional expression modules, which are further visualized and evaluated (Figure 1b). Sample similarity structure is analysed using different algorithms and distance metrics, for example by clustering the pairwise sample correlation matrix (Figure 1c).

Detailed description of the respective algorithms and visualizations would exceed the scope of this outline. We therefore refer to our publications aiming at methodical issues and application of the pipeline (Wirth et al., 2011)(Wirth et al., 2012b)(Wirth et al., 2012a)(Wirth, 2012)(Steiner et al., 2012)(Binder et al., 2012)(Hopp et al., 2013a)(Hopp et al., 2013b).

HTML files are generated to provide straightforward access to this great amount of analysis results (see Figure 2). They guide the user in terms of giving the most prominent links at a glance and leading from one analysis module to another. The **Summary.html** is the starting point of this browsing and can be found in the results folder created by the oposSOM pipeline.

(a) General Information

Number of samples: 12
 Number of groups: 9
 Number of genes: 20825
 Dimension: 1st level SOM: 20 x 50
 Analysis finished: 12 Jun 13 16:26:20 2014 CEST

Results

- Raw Data (PDF)

1st Level SOM Analysis

These reports comprise the SOM portraits in standard and alternative color scales, as well as supporting maps and profiles which provide supplementary information about the 1st level SOM.

- 1st Level SOM Expression Profiles (PDF)
- Alternative Color Scales: absolute, raw, loglogFC (PDF)
- Rank Profiles: FC, WAO, shrinkage-1 (PDF)
- Supporting Maps (PDF)
- Entropy Profiles (PDF)
- Topology Profiles (PDF)

Sample Summaries

Summary page for the individual samples.

- Sample Reports (HTML)

Geneset Enrichment Analysis

Functional analyses using predefined gene sets. The results are visualized in terms of heatmaps, profile plots and population maps.

- Functional Analysis (HTML)

2nd Level Analysis

Sample similarity analyses based on different metrics applied, using the metadata as input.

- 2nd Level SOM (PDF)
- Similarity Based Methods: Neighbor Joining, Hierarchical Clustering (PDF)
- Correlation Based Methods: Spanning Tree, Networks, Maps (PDF)
- Component Based Methods: SASCA, SASCA (PDF)

3rd Level Analysis

Different criteria of spot module definition such as overexpression or mutual correlations between the metagenes where applied. The reports comprise integrated portraits, functional analyses.

- Spot Reports (HTML)

Group Analyses

Analyses based on group wise aggregated data, including portraits, clustering and functional analyses.

- Group Analysis Reports (HTML)

(b) Group Overview

Groups	Number of Samples
metastasis	2
Endocrine	1
Digestion	1
Esocrine	1
Epithelium	1
Reproduction	1
Muscle	1
Immune System	2
Nervous System	2

Sample Summary Sheets

For each sample a report sheet is created which summarizes the most relevant information using the global and local perspective. The global summary shows the ranked list of differentially expressed genes for the whole sample, the ranked list of over- and underexpressed gene sets after GSG-overexpression analysis and the respective p-value distributions. The local summary sheets present the analogous information for each single spot detected. The gene and gene set lists are provided as tables.

Sample Name	Group	Summary Sheet	Global Gene List	Local Gene List	Gene Set List
liver	Homocestas	PDF	CSV	CSV 1 CSV 2 CSV 3	CSV
Kidney cortex	Homocestas	PDF	CSV	CSV 1 CSV 2 CSV 3 CSV 4	CSV
Thyroid gland	Endocrine	PDF	CSV	CSV 1 CSV 2 CSV 3	CSV
small intestine	Digestion	PDF	CSV	CSV 1 CSV 2	CSV
prostate	Esocrine	PDF	CSV	CSV 1 CSV 2 CSV 3	CSV
testis	Reproduction	PDF	CSV	CSV 1 CSV 2 CSV 3	CSV
axillary muscle	Muscle	PDF	CSV	CSV 1 CSV 2	CSV
bone marrow	Immune System	PDF	CSV	CSV 1 CSV 2 CSV 3	CSV
lymph node	Immune System	PDF	CSV	CSV 1 CSV 2	CSV
adipocytes	Nervous System	PDF	CSV	CSV 1 CSV 2 CSV 3	CSV
cerebral cortex	Nervous System	PDF	CSV	CSV 1 CSV 2 CSV 3	CSV

(c) Spot Module Report Sheets

Reports contain the spot module expression profiles and assignments of the spots to samples and to groups.

- Overexpression Spot Report (PDF)
- Underexpression Spot Report (PDF)
- K-Means Cluster Report (PDF)
- Group Overexpression Report (PDF)

Spot Module Network Analysis

Networks of spot association are visualized as graphs. WTC, correlation networks and correlation spanning trees, are given for individual spots and spot patterns.

- Overexpression Networks (PDF)
- Underexpression Networks (PDF)
- K-Means Cluster Networks (PDF)
- Group Overexpression Networks (PDF)

Chromosomal Enrichment

For each spot, enrichment of chromosomal positions (chromosome bands) is visualized as overview heatmaps and individual chromosome plots.

- Overexpression Chromosomal Enrichment (PDF)

(d) Gene Sets

Enrichment profiles of individual predefined gene sets are shown as bar plots across all samples. Additionally the top FC-overexpression profiles of the leading metagenes are shown. Further, members of each gene set are given as population maps and tables.

Category BP

GeneSet name	Category	Profile	Population Map	Members
10 ribon' posttranslational protein folding	BP	PDF	PDF	CSV
2-regulatory metabolic process	BP	PDF	PDF	CSV
3-protein-protein-@-phosphorylate metabolic process	BP	PDF	PDF	CSV
5-UTR-mediated mRNA stabilization	BP	PDF	PDF	CSV
T-methylguanine mRNA capping	BP	PDF	PDF	CSV
acrosome assembly	BP	PDF	PDF	CSV
acrosome reaction	BP	PDF	PDF	CSV
actin cytoskeleton organization	BP	PDF	PDF	CSV
actin cytoskeleton reorganization	BP	PDF	PDF	CSV
actin filament based movement	BP	PDF	PDF	CSV
actin filament bundle assembly	BP	PDF	PDF	CSV

Figure 2: HTML files allow browsing all results provided by the oposSOM package: (a) The central *Summary.html* serves as starting point and contains general information and results, as well as links to other HTML files such as (b) the sample summary page, (c) the spot module summary page and (d) the functional analyses page.

4 Parameter settings

All parameters are optional and will be set to default values if missing. However we recommend to adapt the following parameters according to the respective analysis:

- *dataset.name* (character): name of the dataset. Used to name results folder and environment image (default: 'Unnamed').
- *dim.1stLvlSom* (integer): dimension of primary SOM (default: 20). Given as a single value defining the size of the square SOM grid.
- *feature.centralization* (boolean): enables or disables centralization of the features (default: TRUE).
- *sample.quantile.normalization* (boolean): enables quantile normalization of the samples (default: TRUE).

Database parameters are required to enable gene annotations and functional analyses (details are given below):

- *database.dataset* (character): type of ensemble dataset queried using biomaRt interface (default: "auto"). Use "auto" to detect database parameters automatically.
- *database.id.type* (character): type of rowname identifier in biomaRt database (default: ""). Obsolete if *database.dataset*="auto".

The parameters below are secondary and may be left unattended by the user:

- *dim.2ndLvlSom* (integer): dimension of the second level SOM (default: 20). Given as a single value defining the size of the square SOM grid.
- *training.extension* (numerical, >0): factor extending the number of iterations in SOM training (default: 1).
- *rotate.SOM.portraits* (integer {0,1,2,3}): number of rotations of the primary SOM in counter-clockwise fashion (default: 0). This solely influences the orientation of the portraits.
- *flip.SOM.portraits* (boolean): mirroring the primary SOM along the bottom-left to top-right diagonal (default: FALSE). This solely influences the orientation of the portraits.
- *geneset.analysis* (boolean): enables or disables geneset analysis (default: TRUE).

- *geneset.analysis.exact* (boolean): enables or disables p-value and fdr calculation in geneset analysis (default: TRUE). Obsolete if *geneset.analysis=F*.
- *spot.threshold.samples* (numerical, between 0 and 1): expression threshold for the spot regions in single sample portraits (default: 0.65).
- *spot.threshold.modules* (numerical, between 0 and 1): spot detection in summary maps, expression threshold (default: 0.95).
- *spot.coresize.modules* (integer, >0): spot detection in summary maps, minimum spot size (default: 3).
- *spot.threshold.groupmap* (numerical, between 0 and 1): spot detection in group-specific summary maps, expression threshold (default: 0.75).
- *spot.coresize.groupmap* (integer, >0): spot detection in group-specific summary maps, minimum spot size (default: 5).
- *pairwise.comparison.list* (list of group lists): group list for pairwise analyses (default: empty list). Each element is a list of two character vectors containing the sample names to be analysed in pairwise comparison. The sample names must be contained in the column names of the input data matrix. For example, the following setting will compare the homeostasis (liver, kidney) to the nervous system samples (accumbens, cortex), and also tongue to the nervous system:

```

> env$preferences$pairwise.comparison.list <-
+   list(list(c("liver","kidney cortex"),
+             c("accumbens","cerebral cortex")),
+        list(c("tongue"),
+             c("accumbens","cerebral cortex")))

```

5 Biomart database settings

Two parameters are required to access gene annotations and functional information via biomaRt interface:

database.dataset defines the Ensembl data set to be queried, e.g. "hsapiens_gene_ensembl", "mmusculus_gene_ensembl" or "rnorvegicus_gene_ensembl". A complete list of possible entries can be obtained by

```
> library(biomaRt)
> mart<-useMart("ensembl")
> listDatasets(mart)
```

The default setting "auto" will cause oposSOM to test frequently used settings of *database.dataset* and *database.id.type*. If this automatic download of annotation data fails, a warning will be given and manual definition of the parameters will be necessary to enable functional analyses.

database.id.type provides information about the identifier type constituted by the rownames of the expression matrix, e.g. "ensembl_gene_id", "refseq_mrna" or "affy_hg_u133_plus_2". A complete list of possible entries can be obtained by

```
> library(biomaRt)
> mart<-useMart(biomart="ensembl", dataset="hsapiens_gene_ensembl")
> listFilters(mart)
```

6 New functionalities introduced with oposSOM 1.0 on Bioconductor

The oposSOM-package release on Bioconductor is highly superior to the version released on CRAN in 2011:

- Structure of the source code was thoroughly revised to meet the requirements of Bioconductor.
- Organization and presentation of the results output was improved, accompanied with an extended HTML interface to access all results.
- A package vignette was introduced.
- New analysis modules were implemented:
 - Metagene entropy and portrait topology analyses
 - Neighbor-joining clustering of the samples
 - Correlation Network analysis of the samples
 - GSZ-profiles for the individual gene sets
 - Overview heatmaps summarizing enrichment of a large number of gene sets
 - Cancer hallmark enrichment analyses
 - Enrichment analyses for genes sets relating to chromosomal positions
 - Spot report sheets and spot correlation (wTO) networks
 - Expression portraits, differential expression analyses and functional characteristics summarized for the groups defined
 - Stability analyses of the groups using correlation silhouette methods
 - Differential expression analyses for pairs of samples or groups of samples, including differential expression portraits and functional characterization
- Primary input data can be given as Bioconductor 'ExpressionSet' object.

7 Citing oposSOM

Please cite (Löffler-Wirth et al., 2015) when using the package.

8 Details

This document was written using:

```
> sessionInfo()
```

```
R version 3.2.2 (2015-08-14)
Platform: x86_64-pc-linux-gnu (64-bit)
Running under: Ubuntu 14.04.3 LTS
```

```
locale:
```

```
[1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
[3] LC_TIME=en_US.UTF-8      LC_COLLATE=C
[5] LC_MONETARY=en_US.UTF-8  LC_MESSAGES=en_US.UTF-8
[7] LC_PAPER=en_US.UTF-8     LC_NAME=C
[9] LC_ADDRESS=C             LC_TELEPHONE=C
[11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
```

```
attached base packages:
```

```
[1] parallel stats graphics grDevices utils datasets methods
[8] base
```

```
other attached packages:
```

```
[1] Biobase_2.30.0 BiocGenerics_0.16.0 oposSOM_1.6.0
```

```
loaded via a namespace (and not attached):
```

```
[1] lattice_0.20-33 ape_3.3 IRanges_2.4.0
[4] XML_3.98-1.3 bitops_1.0-6 grid_3.2.2
[7] nlme_3.1-122 DBI_0.3.1 magrittr_1.5
[10] stats4_3.2.2 RSQLite_1.0.0 KernSmooth_2.23-15
[13] som_0.3-5 pixmap_0.4-11 fdrtool_1.2.15
[16] scatterplot3d_0.3-36 S4Vectors_0.8.0 fastICA_1.2-0
[19] tools_3.2.2 igraph_1.0.1 biomaRt_2.26.0
[22] RCurl_1.95-4.7 AnnotationDbi_1.32.0
```

References

Hans Binder, Lydia Hopp, Volkan Cakir, Mario Fasold, Martin von Bergen, and Henry Wirth. Molecular phenotypic portraits - Exploring the *ÄYOMESÄZ* with individual resolution. In Jens Allmer, editor, *Health Informatics and Bioinformatics (HIBIT), 2011 6th International Symposium*, pages 99–107. IEEE Xplore, 2012. ISBN 978-2-4673-4394-4. doi: 10.1109/HIBIT.2011.6450817.

Lydia Hopp, Kathrin Lembcke, Hans Binder, and Henry Wirth. Portraying the Expression Landscapes of B-Cell Lymphoma - Intuitive Detection of Outlier

- Samples and of Molecular Subtypes. *Biology*, 2(4):1411–1437, 2013a. doi: 10.3390/biology2041411.
- Lydia Hopp, Henry Wirth, Mario Fasold, and Hans Binder. Portraying the expression landscapes of cancer subtypes: A glioblastoma multiforme and prostate cancer case study. *Systems Biomedicine*, 1(2):1–23, 2013b. URL <http://www.landesbioscience.com/journals/systemsbiomedicine/toc/volume/1/issue/2/>.
- Henry Löffler-Wirth, Martin Kalcher, and Hans Binder. oposSOM: R-package for high-dimensional portraying of genome-wide expression landscapes on Bioconductor. *Bioinformatics (Oxford, England)*, June 2015. ISSN 1367-4811. doi: 10.1093/bioinformatics/btv342. URL <http://www.ncbi.nlm.nih.gov/pubmed/26063839>.
- Lydia Steiner, Lydia Hopp, Henry Wirth, Jörg Galle, Hans Binder, Sonja Prohaska, and Thimo Rohlf. A global genome segmentation method for exploration of epigenetic patterns. *PLoS one*, 7(10), 2012.
- Henry Wirth. *Analysis of large-scale molecular biological data using self-organizing maps*. Dissertation, Leipzig University, 2012. URL <http://www.qucosa.de/fileadmin/data/qucosa/documents/10129/DissertationHenryWirth.pdf>.
- Henry Wirth, Markus Löffler, Martin von Bergen, and Hans Binder. Expression cartography of human tissues using self organizing maps. *BMC Bioinformatics*, 12(1):306, 2011. ISSN 1471-2105. doi: 10.1186/1471-2105-12-306. URL <http://www.biomedcentral.com/1471-2105/12/306>.
- Henry Wirth, Martin von Bergen, and Hans Binder. Mining SOM expression portraits: feature selection and integrating concepts of molecular function. *BioData Mining*, 5(1):18, October 2012a. ISSN 1756-0381. doi: 10.1186/1756-0381-5-18. URL <http://www.ncbi.nlm.nih.gov/pubmed/23043905>.
- Henry Wirth, Martin von Bergen, Jayaseelan Murugaiyan, Uwe Rösler, Tomasz Stokowy, and Hans Binder. MALDI-typing of infectious algae of the genus *Prototheca* using SOM portraits. *Journal of microbiological methods*, 88(1): 83–97, January 2012b. ISSN 1872-8359. doi: 10.1016/j.mimet.2011.10.013. URL <http://www.ncbi.nlm.nih.gov/pubmed/22062088>.