

# Phantassus: web-application for visual and interactive gene expression analysis

Maksim Kleverov<sup>1</sup>    Daria Zenkova<sup>1</sup>    Vladislav Kamenev<sup>1</sup>    Margarita Sablina<sup>1</sup>  
Maxim N. Artyomov<sup>1,2</sup>    Alexey A. Sergushichev<sup>1,2,\*</sup>

**1** ITMO University, Computer Tehcnologies Laboratory, Saint Petersburg, 197101, Russia,

**2** Washington University in St. Louis School of Medicine, Department of Pathology and Immunology, St. Louis, MO 63130, USA,

\* Correspondence: [alsergbox@gmail.com](mailto:alsergbox@gmail.com)

## Abstract

Transcriptomic profiling became a standard approach to quantify a cell state, which led to accumulation of huge amount of public gene expression datasets. However, both reuse of these datasets or analysis of newly generated ones requires a significant technical expertise. Here we present Phantassus – a user-friendly web-application for interactive gene expression analysis which provide a streamlined access to more than 84000 public gene expression datasets, as well as allows analysis of user-uploaded datasets. Phantassus integrates an intuitive and highly interactive JavaScript-based heatmap interface with an ability to run sophisticated R-based analysis methods. Overall Phantassus allows to go all the way from loading, normalizing and filtering data to doing differential gene expression and downstream analysis. Phantassus can be accessed on-line at <https://ctlab.itmo.ru/phantassus> or <https://artyomovlab.wustl.edu/phantassus> or can be installed locally from Bioconductor (<https://bioconductor.org/packages/phantassus>). Phantassus source code is available at <https://github.com/ctlab/phantassus> under MIT licence.

## 23 1 Introduction

24 Transcriptomic profiling is an ubiquitous method for whole-genome level profiling of biological  
25 samples [Stark et al., 2019]. Moreover, deposition of these data into one of the public repositories  
26 became a standard in the field, which led to accumulation of huge amount of publicly available  
27 data. The most significant example is NCBI Gene Expression Omnibus (GEO) project [Barrett  
28 et al., 2012], which stores information of more than 180000 studies.

29 Sharing of transcriptomic data opens up possibilities for reusing them: instead of carrying out a  
30 costly experiment, a publicly available dataset can be used, thus decreasing the cost and accelerating  
31 the research [Byrd et al., 2020]. However, the standard approach for gene expression analysis  
32 requires a significant technical expertise. In particular, many analysis methods are implemented  
33 in R as a part of Bioconductor project ecosystem [Gentleman et al., 2004], and thus one has to  
34 have programming skills in R to use them. On the other hand, domain knowledge is beneficial to  
35 improve quality control of the data, which is especially important when working with the publicly  
36 available data, as well as generation of biological hypotheses [Wang et al., 2016].

37 A number of applications have been developed with the aim to simplify analysis of transcriptomic  
38 datasets (see Supplementary File 1 for details). In particular, web-based applications remove the  
39 burden of set-up and configuration from the end users, thus lowering the entry threshold. Shiny  
40 framework [Chang et al., 2022] revolutionized the field as it became easy to create a web interface  
41 for R based pipelines, which led to a significant growth of web-applications for gene expression  
42 analysis [Ge et al., 2018, Iacoangeli et al., 2022, Mahi et al., 2019, Nelson et al., 2016]. However,  
43 such applications generally have limited interactivity due to mainly server-side computations. Shiny-  
44 independent applications can be more interactive, but they suffer from lack of native R support  
45 and require reimplementations of existing methods from scratch [Gould, 2016, Alonso et al., 2015].

46 Here we present Phantasus: a web-application for gene expression analysis that integrates highly  
47 interactive client-side JavaScript heatmap interface with an R-based backend. Phantasus allows  
48 to carry out all major steps of gene expression analysis pipeline: data loading, annotation, nor-  
49 malization, clustering, differential gene expression and pathway analysis. Notably, Phantasus  
50 provides a streamlined access to more than 84000 microarray and RNA-seq datasets from Gene

51 Expression Omnibus database, simplifying their reanalysis. Phantanus can be accessed on-line  
52 at <https://ctlab.itmo.ru/phantanus> or at <https://artyomovlab.wustl.edu/phantanus> or can be in-  
53 stalled locally from Bioconductor. Phantanus is open source and its code is available at <https://github.com/ctlab/phantanus>  
54 under MIT licence.

## 55 2 Results

### 56 2.1 Phantanus web-application

57 We developed a web application called Phantanus for interactive gene expression analysis. Phanta-  
58 nus integrates JavaScript-rich heatmap based user interface originated from Morpheus [Gould, 2016]  
59 with an R back-end via OpenCPU framework [Ooms, 2014]. Heatmap graphical interface provides  
60 an intuitive way to manipulate the data and metadata: directly in a web-browser the user can  
61 create or modify annotations, edit color schemes, filter rows and columns, and so on. On the other  
62 hand the R back-end provides a way to easily run a multitude of computational analysis methods  
63 available as R packages. All together this architecture (Figure 1) provides a smooth experience  
64 for doing all common analysis steps: loading datasets, normalization, exploration, visualization,  
65 differential expression and gene set enrichment analyses.

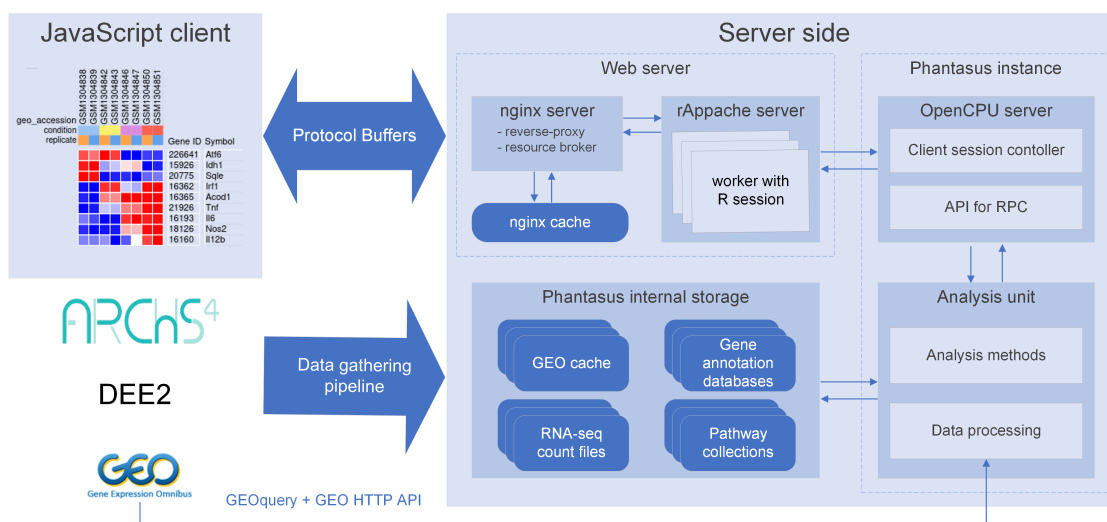


Figure 1: Overview of Phantanus architecture. The front-end interface is a JavaScript application, that requests the web-server to load the data and perform resource-consuming tasks. The core element of the back-end is the OpenCPU-based server which triggers execution of R-based analysis methods. Protocol Buffers are used for efficient client-server dataset synchronization.

66 Several options to load the gene expression data into Phantasus are available. First, datasets  
67 from Gene Expression Omnibus [Barrett et al., 2012] can be loaded by their identifier. Phantasus  
68 supports microarray datasets, which are loaded directly from GEO, as well as RNA-sequencing  
69 datasets, for which counts data from third-party databases are used (see section Section 2.3 for  
70 details). Second, datasets can be loaded from a gene expression table file in GCT, TSV and XLSX  
71 formats. Finally, a set of curated datasets are available directly from the home page.

72 A number of methods can be used to prepare, normalize and explore the gene expression table.  
73 In particular, it is possible to aggregate microarray probe-level data to gene levels, transform and  
74 filter the data, do a principal component analysis (PCA), do a k-means or hierarchical clustering,  
75 etc. These tools allows to do a thorough quality control of the dataset and remove the outliers if  
76 they are present.

77 When the dataset is properly filtered and normalized, differential expression analysis using limma  
78 [Ritchie et al., 2015] or DESeq2 pipelines [Love et al., 2014] can be carried out. These results  
79 can then be used with other web-services for downstream analysis, with shortcuts for pathway  
80 analysis with Enrichr [Kuleshov et al., 2016] and metabolic network analysis with Shiny GAM  
81 [Sergushichev et al., 2016]. Additionally, gene set enrichment analysis can be done directly in  
82 Phantasus as implemented in fgsea package [Korotkevich et al., 2021].

83 All of the plots produced by Phantasus during the data exploration and analysis can be exported  
84 as vector images in SVG format. This includes heatmaps, PCA plots, gene profiles, enrichment  
85 plots etc. The obtained images can be used for publications as is or adjusted in a vector graphics  
86 editor.

87 Another option for presenting final or intermediate results is a session link sharing. When a link is  
88 generated, a snapshot with the current dataset and its representation: annotations, color scheme,  
89 sample dendrograms, etc, is saved on the server. The link can be shared with other users, and,  
90 when opened, restores the session.

## 91 **2.2 Stand-alone phantasm distribution**

92 Aside from using the two official mirrors: <https://ctlab.itmo.ru/phantasm> and [https://](https://artyomovlab.wustl.edu/phantasm)  
93 [artyomovlab.wustl.edu/phantasm](https://artyomovlab.wustl.edu/phantasm), there is a possibility to set up phantasm locally. Phantasm can  
94 be installed as an R package from Bioconductor (<https://bioconductor.org/packages/phantasm>)  
95 or loaded as a Docker image (<https://hub.docker.com/r/asergushichev/phantasm>). In both cases  
96 almost all of the Phantasm functions will be available from the start.

97 Some of Phantasm features require additional server-side set up. Extended support of GEO  
98 datasets requires preprocessed expression matrices and platform annotations. Identifier mapping  
99 requires organism annotation databases. Pathways enrichment requires pathway databases. For  
100 the initial set up, all these files can be downloaded from [https://ctlab.itmo.ru/files/software/](https://ctlab.itmo.ru/files/software/phantasm/minimal-cache)  
101 [phantasm/minimal-cache](https://ctlab.itmo.ru/files/software/phantasm/minimal-cache).

102 Important feature of a stand-alone version of Phantasm is an ability to share manually curated  
103 datasets. Similar to Phantasm session link sharing, one can generate a named session consisting  
104 of a dataset and its visual representation. Link to this named session (e.g. [https://ctlab.itmo.](https://ctlab.itmo.ru/phantasm/?preloaded=GSE53986.Ctrl.vs.LPS)  
105 [ru/phantasm/?preloaded=GSE53986.Ctrl.vs.LPS](https://ctlab.itmo.ru/phantasm/?preloaded=GSE53986.Ctrl.vs.LPS)) can then be shared for the other users to view.  
106 Such predictable display of the data can be particularly useful in a publication context.

## 107 **2.3 Available datasets**

108 Phantasm provides a streamlined access to more than 84379 GEO datasets. For these datasets the  
109 expression values and gene identifiers (Entrez, ENSEMBL or Gene Symbol) are readily available  
110 (Figure 2). Moreover, these datasets are used to populate initial Phantasm cache, and thus they  
111 have low loading times.

112 From these 84379 datasets 49666 are microarrays based on 2767 platforms. For 1347 platforms  
113 GEO datasets have a machine-readable annotations in the *annot.gz* format with Entrez gene and  
114 Gene symbol columns, which correspond to 39689 datasets. The remaining 9977 datasets are  
115 obtained from platforms that do not have a GEO-provided annotation. For these 1420 platforms we  
116 have automatically marked up user-provided annotations in *SOFT* format to extract gene identifiers  
117 and convert the annotations into *annot.gz* format.

118 RNA-seq subset of the datasets with a streamlined access consists of 34713 datasets. As GEO does  
 119 not store expression values for RNA-seq datasets, we rely on other databases for the expression  
 120 data. The first-priority database for RNA-seq gene counts is ARCHS4 (Human, Mouse and Zoo  
 121 versions), which covers 25833 datasets. The other source is DEE2 database (human, mouse and  
 122 other available organisms), which covers an additional 8880 datasets. DEE2 database contains  
 123 transcript-level quantification, so it has been preprocessed to sum read counts into gene-level tables.

Source	Full support	Limited support	Total
NCBI GEO	84379	39044	123423
GEO Microarray	49666	25470	75136
GEO curated annotation	39689	676	40365
Parsed user-provided annotation	9977	677	10654
Other user-provided annotation	0	24117	24117
GEO RNA-seq	34713	13574	48287
ARCHS4	25833	9977	35810
<i>Homo sapiens</i>	12081	5206	17287
<i>Mus musculus</i>	12476	4473	16949
Other organisms	1276	298	1574
DEE2 (not in ARCHS4)	8880	3597	12477
<i>Homo sapiens</i>	2547	1879	4426
<i>Mus musculus</i>	2287	845	3132
Other organisms	4046	873	4919

Figure 2: Dataset availability in Phantasus. For fully supported datasets gene expression data is accompanied by gene annotations in a standardized format. Limited support datasets have either incomplete gene expression matrix or gene annotations.

## 124 3 Implementation

### 125 3.1 Web application architecture

126 Phantasus is a web-application that combines interactive graphical user interface with an access to  
 127 a variety of R-based analysis methods (Figure 1). The front-end is JavaScript-based, and is derived  
 128 from Morpheus web-application for matrix visualization and analysis [Gould, 2016]. The back-end  
 129 is written in R, with an OpenCPU server [Ooms, 2014] translating HTTP-queries from the client  
 130 into R procedure calls.

131 The JavaScript client is responsible for the matrix visualization, as well as certain analysis methods.

132 In particular, steps like subsetting the dataset, working with annotations, basic matrix modification  
133 (e.g., log-transformation, scaling, etc) have client-side implementation. Furthermore, the client  
134 supports additional visualization methods such as row profile plots, volcano plot, and others.

135 The analysis methods that require external data or algorithms are implented in the form of the  
136 **phantasus** R package to be carried out on the server side. The operations include differential gene  
137 expression analysis, principal component analysis, pathway analysis, and others. Commonly, these  
138 methods rely on functions which are already available in the existing R packages, for such methods  
139 only wrapper R functions are implemented.

140 OpenCPU server is a core component of the Phantasus back-end. The server provides an HTTP  
141 API for calling computational methods implemented in R. For each call OpenCPU creates a new R  
142 environment with the required data, in which the method is then executed. OpenCPU can manage  
143 these R environments both in a standard single-user R session, and, with the help of rApache, in  
144 a multi-user manner inside an Apache web-server.

145 The transfer of large objects between the server and the client exploits a binary Protobuf proto-  
146 col. The Phantasus back-end uses protolite R package [Ooms, 2021] for object serialization and  
147 deserialization. The front-end relies on protobuf.js module [Coe, 2020].

148 For the further performance improvement Nginx server is used to wrap OpenCPU server. Nginx  
149 server caches the results of the OpenCPU method calls. If the same method with the same data is  
150 called again the cached result can be returned without any additional computations. Furthermore  
151 Nginx is used to serve static content and to manage permissions.

## 152 **3.2 Data sources and data gathering**

153 The main data source for Phantasus is NCBI GEO database [Barrett et al., 2012]. All of the GEO  
154 datasets are identified by a GSEnnnnn accession number (with a subset of the datasets having an  
155 additional GDSnnnnn identifier). However, depending on the type of the dataset the processing  
156 procedure is different.

157 The majority of gene expression datasets in GEO database can be divided into two groups: microar-  
158 ray data and RNA-seq data. While the experiment metadata is available for all of the datasets, the

159 expression matrices are provided only for the microarray datasets. Phantasus relies on GEOQuery  
160 package [Davis and Meltzer, 2007] to load the experiment metadata (for all of the datasets) and  
161 expression matrices (for microarray datasets) from GEO.

162 When a GEO RNA-seq dataset is requested by the user, Phantasus refers to precomputed gene  
163 counts databases available in the internal storage. In particular, data from ARCHS4 [Lachmann  
164 et al., 2017] and DEE2 [Ziemann et al., 2019] projects are used. Both of these projects contain gene  
165 counts and metadata for RNA-seq samples related to different model organisms including but not  
166 limited to mouse and human. For any requested RNA-seq dataset the gene counts are loaded from  
167 a single database, whichever covers the highest number of samples.

168 Next, Phantasus stores gene annotation databases which are used to map genes between different  
169 identifier types. These databases are stored in sqlite format compatible with *AnnotationDbi* R  
170 package [Pagès et al., 2022]. Currently only human and mouse databases are available, which are  
171 based on *org.Hs.eg.db* and *org.Mm.eg.db* R packages respectively.

172 Pathway databases are stored to be used for gene set enrichment analysis. Currently gene set  
173 collections include GO biological processes database [Ashburner et al., 2000], Reactome database  
174 [Gillespie et al., 2021] and MSigDB Hallmark database [Liberzon et al., 2011] for human and mouse.

175 Finally, for a faster access, Phantasus dataset cache is automatically populated by a large com-  
176 pendium of datasets. The automatic pipeline is Snakemake-based and consists of four steps. First  
177 of all the pipeline converts DEE2 files to ARCHS4-like HDF5 files. During this procedure tran-  
178 script expression provided by DEE2 is summed up to gene level. Second step checks for which  
179 microarray platforms GEO contains a curated machine-readable annotation in *annot.gz* format.  
180 Third step tries to generate the machine-readable annotation for the rest of microarray platforms  
181 from the annotations available in the *SOFT* format. Currently, this step produces an additional  
182 1300 *annot.gz* files. The last step goes over all of the microarray datasets with a machine-readable  
183 annotation and over all of the RNA-seq datasets with the counts available in ARCHS4 or DEE2.  
184 For each such dataset the cached entry with all of the data and metadata is created and stored.

185 A snapshot of Phantasus internal storage is available at [https://ctlab.itmo.ru/files/software/  
186 phantasus/minimal-cache](https://ctlab.itmo.ru/files/software/phantasus/minimal-cache). It contains preprocessed count files, automatically marked-up anno-



187 tations, gene and pathways databases. This snapshot can be used for a local Phantasus set  
188 up.

## 189 4 Case study

### 190 4.1 Basic usage

191 To illustrate the basic usage of Phantasus we will consider dataset GSE53986 [Noubade et al., 2014]  
192 from GEO database. This dataset consists of 16 samples of bone marrow derived macrophages, un-  
193 treated and treated with three stimuli: LPS, IFN $\gamma$  and combined LPS+IFN $\gamma$ . The gene expression  
194 was measured with Affymetric Mouse Genome 430 2.0 Array array. Here we give an overview of  
195 the steps, the full walk-through for the analysis is available in Supplementary File 2.

196 As the first step of the analysis the dataset can be loaded and normalized. The dataset is loaded  
197 straightforwardly by the *GSE53986* identifier. Because this is a microarray dataset, internally the  
198 gene expression values are obtained from GEO. In this particular case the expression values have  
199 not been normalized, but it can be done in Phantasus. From the available normalization options  
200 we select log2 scaling and quantile normalization. Further, we can aggregate microarray probe-  
201 level expression values into gene-level expression. We chose *Maximum Median Probe* method which  
202 retains only a single probe per gene, the one that has highest median expression value. Finally, we  
203 can filter out lowly expressed genes, for example, by keeping only the top 12000 expressed genes.

204 After the normalization step we can apply a number of exploratory techniques. In particular we  
205 can do a Principal Component Analysis, k-means gene clustering and hierarchical clustering of the  
206 samples. From these analysis we can discover, that there is an overall good concordance between  
207 the replicates of the same treatment, with an exception of the first replicate in each group. We can  
208 conclude that these samples are outliers and remove them before the downstream analysis.

209 Finally, we can do a comparison between the sample groups, for example by comparing untreated  
210 and LPS-treated samples. As the data has been normalized, we can apply *limma* for differential  
211 gene expression analysis. The result appears as additional gene annotation columns: P-values,  
212 log-fold-changes and other statistics. Next we can use differential expression results for a pathway  
213 enrichment analysis: for example, we can use R-based gene set enrichment analysis via *fgsea* package

214 or we can use external tools, such as Enrichr.

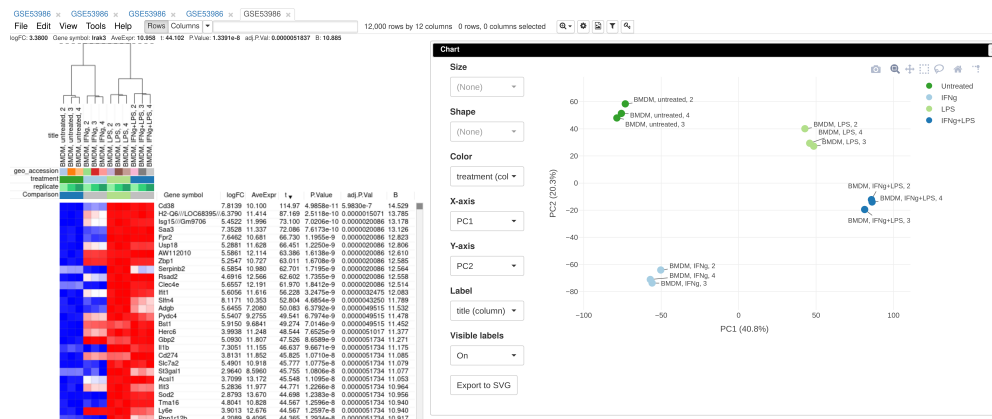


Figure 3: Example of analysed dataset GSE53986 with normalized gene expression values, filtered outliers, hierarchically clustered columns, and rows annotated with differential expression analysis between untreated and LPS-treated macrophages

## 215 4.2 Data reanalysis

216 To highlight Phantasia ability to reanalyze publicly available data in a context of a biological study  
 217 let us consider a study by Mowel and colleagues [Mowel et al., 2017]. The study considers a genomic  
 218 locus *Rroid* linked by the authors to homeostasis and function of group 1 innate lymphoid cells  
 219 (ILC1). The authors hypothesized that *Rroid* locus controls ILC1s by promoting the expression  
 220 of *Id2* gene, a known regulator of ILCs. To confirm this hypothesis the authors generated an *Id2*-  
 221 dependent gene signature based on an existing transcriptomic data [Shih et al., 2016] and showed  
 222 its deregulation in *Rroid* deficient cells.

223 The described above computational analysis linking *Rroid* and *Id2* can be replicated in Phantasia  
 224 in a straightforward way (see Supplementary File 3 for the detailed walk-through).

225 First, we can open GEO dataset GSE76466 [Shih et al., 2016], containing gene expression data for  
 226 *Id2* deficient NK cells. Notably, GSE76466 is an RNA sequencing dataset, without  
 227 gene expression values stored directly in GEO database, however Phantasia loads the dataset  
 228 leveraging precomputed expression values from ARCHS4 project [Lachmann et al., 2017]. Then we  
 229 can run differential gene expression analysis with DESeq2, comparing *Id2*-deficient and wild-type  
 230 NK cells. *Id2*-dependent gene signature can be obtained by sorting the genes by *stat* column.

231 Second, RNA-sequencing dataset GSE101459 [Mowel et al., 2017], generated by Mowel and col-

232 leagues for Rroid-deficient NK cells, can also be opened in Phantasus. There we can do differential  
233 gene expression analysis with DESeq2 and remove lowly expressed genes. Finally we can enter  
234 the generated Id2-dependent gene signature into Phantasus gene search field and use GSEA plot  
235 tool to obtain an enrichment plot, similar to one presented by Mowel and colleagues, confirming a  
236 potential regulation via Id2.

## 237 **5 Conclusion**

238 Phantasus is a tool for visual and interactive gene expression analysis that allows in an easy and  
239 streamlined manner to go from loading, normalizing and filtering data to differential gene expression  
240 and downstream analysis. Additionally, due to its tight integration with R environment, Phantasus  
241 can be extended with other analysis methods, in particular the ones available at Bioconductor.  
242 Phantasus can be both used on-line at <https://ctlab.itmo.ru/phantasus> or be installed locally from  
243 Bioconductor.

## 244 **6 Acknowledgements**

245 The project was supported by Ministry of Science and Higher Education of the Russian Federation  
246 (Priority 2030 Federal Academic Leadership Program).

## 247 **References**

- 248 R. Stark, M. Grzelak, and J. Hadfield. RNA sequencing: the teenage years. *Nat Rev Genet*, 20  
249 (11):631–656, 11 2019.
- 250 Tanya Barrett, Stephen E. Wilhite, Pierre Ledoux, Carlos Evangelista, Irene F. Kim, Maxim Toma-  
251 shevsky, Kimberly A. Marshall, Katherine H. Phillippy, Patti M. Sherman, Michelle Holko, An-  
252 drey Yefanov, Hyeseung Lee, Naigong Zhang, Cynthia L. Robertson, Nadezhda Serova, Sean  
253 Davis, and Alexandra Soboleva. NCBI GEO: archive for functional genomics data sets—update.  
254 *Nucleic Acids Research*, 41(D1):D991–D995, 11 2012. ISSN 0305-1048. doi: 10.1093/nar/gks1193.  
255 URL <https://doi.org/10.1093/nar/gks1193>.

- 256 J. B. Byrd, A. C. Greene, D. V. Prasad, X. Jiang, and C. S. Greene. Responsible, practical genomic  
257 data sharing that accelerates research. *Nat Rev Genet*, 21(10):615–629, 10 2020.
- 258 Robert C Gentleman, Vincent J Carey, Douglas M Bates, Ben Bolstad, Marcel Dettling, Sandrine  
259 Dudoit, Byron Ellis, Laurent Gautier, Yongchao Ge, Jeff Gentry, Kurt Hornik, Torsten Hothorn,  
260 Wolfgang Huber, Stefano Iacus, Rafael Irizarry, Friedrich Leisch, Cheng Li, Martin Maechler,  
261 Anthony J Rossini, Gunther Sawitzki, Colin Smith, Gordon Smyth, Luke Tierney, Jean YH  
262 Yang, and Jianhua Zhang. Bioconductor: open software development for computational biology  
263 and bioinformatics. *Genome Biology*, 5(10):R80, 2004. doi: 10.1186/gb-2004-5-10-r80. URL  
264 <https://doi.org/10.1186/gb-2004-5-10-r80>.
- 265 Zichen Wang, Caroline D Monteiro, Kathleen M Jagodnik, Nicolas F Fernandez, Gregory W Gun-  
266 dersen, Andrew D Rouillard, Sherry L Jenkins, Axel S Feldmann, Kevin S Hu, Michael G Mc-  
267 Dermott, et al. Extraction and analysis of signatures from the gene expression omnibus by the  
268 crowd. *Nature communications*, 7(1):1–11, 2016.
- 269 Winston Chang, Joe Cheng, JJ Allaire, Carson Sievert, Barret Schloerke, Yihui Xie, Jeff Allen,  
270 Jonathan McPherson, Alan Dipert, and Barbara Borges. *shiny: Web Application Framework for*  
271 *R*, 2022. URL <https://CRAN.R-project.org/package=shiny>. R package version 1.7.3.
- 272 Steven Xijin Ge, Eun Wo Son, and Runan Yao. idep: an integrated web application for differential  
273 expression and pathway analysis of rna-seq data. *BMC bioinformatics*, 19(1):1–24, 2018.
- 274 Alfredo Iacoangeli, Richard Dobson, and Ammar Al-Chalabi. Geoexplorer: a webserver for gene  
275 expression analysis and visualisation. *Nucleic Acids Research*, April 2022. ISSN 0305-1048.
- 276 Naim Al Mahi, Mehdi Fazel Najafabadi, Marcin Pilarczyk, Michal Kouril, and Mario Medvedovic.  
277 Grein: An interactive web platform for re-analyzing geo rna-seq data. *Scientific reports*, 9(1):  
278 1–9, 2019.
- 279 Jonathan W Nelson, Jiri Sklenar, Anthony P Barnes, and Jessica Minnier. The START App: a web-  
280 based RNAseq analysis and visualization resource. *Bioinformatics*, 33(3):447–449, 10 2016. ISSN  
281 1367-4803. doi: 10.1093/bioinformatics/btw624. URL [https://doi.org/10.1093/bioinformatics/](https://doi.org/10.1093/bioinformatics/btw624)  
282 [btw624](https://doi.org/10.1093/bioinformatics/btw624).

- 283 Joshua Gould. Morpheus: Versatile matrix visualization and analysis software.  
284 <https://software.broadinstitute.org/morpheus/index.html>, 2016. URL <https://software.broadinstitute.org/morpheus/index.html>.  
285
- 286 Roberto Alonso, Francisco Salavert, Francisco Garcia-Garcia, Jose Carbonell-Caballero, Marta  
287 Bleda, Luz Garcia-Alonso, Alba Sanchis-Juan, Daniel Perez-Gil, Pablo Marin-Garcia, Ruben  
288 Sanchez, et al. Babelomics 5.0: functional interpretation for new generations of genomic data.  
289 *Nucleic acids research*, 43(W1):W117–W121, 2015.
- 290 Jeroen Ooms. The opencpu system: Towards a universal interface for scientific computing through  
291 separation of concerns, 2014.
- 292 Matthew E Ritchie, Belinda Phipson, Di Wu, Yifang Hu, Charity W Law, Wei Shi, and Gordon K  
293 Smyth. limma powers differential expression analyses for RNA-sequencing and microarray studies.  
294 *Nucleic Acids Research*, 43(7):e47, 2015.
- 295 Michael I. Love, Wolfgang Huber, and Simon Anders. Moderated estimation of fold change and  
296 dispersion for rna-seq data with deseq2. *Genome Biology*, 15, 2014.
- 297 M. V. Kuleshov, M. R. Jones, A. D. Rouillard, N. F. Fernandez, Q. Duan, Z. Wang, S. Koplev, S. L.  
298 Jenkins, K. M. Jagodnik, A. Lachmann, M. G. McDermott, C. D. Monteiro, G. W. Gundersen,  
299 and A. Ma’ayan. Enrichr: a comprehensive gene set enrichment analysis web server 2016 update.  
300 *Nucleic Acids Res.*, 44(W1):W90–97, 07 2016.
- 301 A. A. Sergushichev, A. A. Loboda, A. K. Jha, E. E. Vincent, E. M. Driggers, R. G. Jones, E. J.  
302 Pearce, and M. N. Artyomov. GAM: a web-service for integrated transcriptional and metabolic  
303 network analysis. *Nucleic Acids Res.*, 44(W1):194–200, 07 2016.
- 304 Gennady Korotkevich, Vladimir Sukhov, Nikolay Budin, Boris Shpak, Maxim N. Artyomov, and  
305 Alexey Sergushichev. Fast gene set enrichment analysis. *bioRxiv*, 2021. doi: 10.1101/060012.  
306 URL <https://www.biorxiv.org/content/early/2021/02/01/060012>.
- 307 Jeroen Ooms. *protolite: Highly Optimized Protocol Buffer Serializers*, 2021. URL <https://CRAN.R-project.org/package=protolite>. R package version 2.1.1.
- 309 Benjamin Coe. *protobuf.js is a pure JavaScript implementation with TypeScript support for node.js*

310 *and the browser*, 2020. URL <http://protobufjs.github.io/protobuf.js>. JS npm package version  
311 6.11.3.

312 Sean Davis and Paul Meltzer. Geoquery: a bridge between the gene expression omnibus (geo) and  
313 bioconductor. *Bioinformatics*, 14:1846–1847, 2007.

314 Alexander Lachmann, Denis Torre, Alexandra B. Keenan, Kathleen M. Jagodnik, Hyojin J. Lee,  
315 Moshe C. Silverstein, Lily Wang, and Avi Ma’ayan. Massive mining of publicly available RNA-seq  
316 data from human and mouse. <https://www.biorxiv.org/content/early/2017/09/15/189092>, 2017.  
317 URL <https://www.biorxiv.org/content/early/2017/09/15/189092>.

318 Mark Ziemann, Antony Kaspi, and Assam El-Osta. Digital expression explorer 2: a repository of  
319 uniformly processed RNA sequencing data. *GigaScience*, 8(4), 04 2019. ISSN 2047-217X. doi:  
320 10.1093/gigascience/giz022. URL <https://doi.org/10.1093/gigascience/giz022>. giz022.

321 Hervé Pagès, Marc Carlson, Seth Falcon, and Nianhua Li. *AnnotationDbi: Manipulation of*  
322 *SQLite-based annotations in Bioconductor*, 2022. URL [https://bioconductor.org/packages/](https://bioconductor.org/packages/AnnotationDbi)  
323 *AnnotationDbi*. R package version 1.58.0.

324 Michael Ashburner, Catherine A. Ball, Judith A. Blake, David Botstein, Heather L. Butler,  
325 J. Michael Cherry, Allan Peter Davis, Kara Dolinski, Selina S. Dwight, Janan T. Eppig, Mi-  
326 dori A. Harris, David P. Hill, Laurie Issel-Tarver, Andrew Kasarskis, Suzanna E. Lewis, John C.  
327 Matese, Joel E. Richardson, Martin Ringwald, Gerald M. Rubin, and Gavin Sherlock. Gene  
328 ontology: tool for the unification of biology. *Nature Genetics*, 25:25–29, 2000.

329 Marc Gillespie, Bijay Jassal, Ralf Stephan, Marija Milacic, Karen Rothfels, Andrea Senff-Ribeiro,  
330 Johannes Griss, Cristoffer Sevilla, Lisa Matthews, Chuqiao Gong, Chuan Deng, Thawfeek  
331 Varusai, Eliot Ragueneau, Yusra Haider, Bruce May, Veronica Shamovsky, Joel Weiser, Tim-  
332 othy Brunson, Nasim Sanati, Liam Beckman, Xiang Shao, Antonio Fabregat, Konstantinos  
333 Sidiropoulos, Julieth Murillo, Guilherme Viteri, Justin Cook, Solomon Shorser, Gary Bader,  
334 Emek Demir, Chris Sander, Robin Haw, Guanming Wu, Lincoln Stein, Henning Hermjakob,  
335 and Peter D’Eustachio. The reactome pathway knowledgebase 2022. *Nucleic Acids Re-*  
336 *search*, 50(D1):D687–D692, 11 2021. ISSN 0305-1048. doi: 10.1093/nar/gkab1028. URL  
337 <https://doi.org/10.1093/nar/gkab1028>.

- 338 Arthur Liberzon, Aravind Subramanian, Reid Pinchback, Helga Thorvaldsdóttir, Pablo Tamayo,  
339 and Jill P. Mesirov. Molecular signatures database (MSigDB) 3.0. *Bioinformatics*, 27(12):1739–  
340 1740, 05 2011. ISSN 1367-4803. doi: 10.1093/bioinformatics/btr260. URL [https://doi.org/10.](https://doi.org/10.1093/bioinformatics/btr260)  
341 [1093/bioinformatics/btr260](https://doi.org/10.1093/bioinformatics/btr260).
- 342 R. Noubade, K. Wong, N. Ota, S. Rutz, C. Eidenschenk, P. A. Valdez, J. Ding, I. Peng, A. Sebrell,  
343 P. Caplazi, J. DeVoss, R. H. Soriano, T. Sai, R. Lu, Z. Modrusan, J. Hackney, and W. Ouyang.  
344 NRROS negatively regulates reactive oxygen species during host defence and autoimmunity. *Na-*  
345 *ture*, 509(7499):235–239, May 2014.
- 346 Walter K. Mowel, Sam J. McCright, Jonathan J. Kotzin, Magalie A. Collet, Asli Uyar, Xin  
347 Chen, Alexandra DeLaney, Sean P. Spencer, Anthony T. Virtue, EnJun Yang, Alejandro Vil-  
348 larino, Makoto Kurachi, Margaret C. Dunagin, Gretchen Harms Pritchard, Judith Stein, Cyn-  
349 thia Hughes, Diogo Fonseca-Pereira, Henrique Veiga-Fernandes, Arjun Raj, Taku Kambayashi,  
350 Igor E. Brodsky, John J. O’Shea, E. John Wherry, Loyal A. Goff, John L. Rinn, Adam Williams,  
351 Richard A. Flavell, and Jorge Henao-Mejia. Group 1 innate lymphoid cell lineage identity is  
352 determined by a cis-regulatory element marked by a long non-coding rna. *Immunity*, 47(3):  
353 435–449.e8, 2017. ISSN 1074-7613. doi: <https://doi.org/10.1016/j.immuni.2017.08.012>. URL  
354 <https://www.sciencedirect.com/science/article/pii/S1074761317303709>.
- 355 Han-Yu Shih, Giuseppe Sciumè, Yohei Mikami, Liying Guo, Hong-Wei Sun, Stephen R. Brooks,  
356 Joseph F. Urban, Fred P. Davis, Yuka Kanno, and John J. O’Shea. Developmental acquisition  
357 of regulomes underlies innate lymphoid cell functionality. *Cell*, 165(5):1120–1133, 2016. ISSN  
358 0092-8674. doi: <https://doi.org/10.1016/j.cell.2016.04.029>. URL [https://www.sciencedirect.com/](https://www.sciencedirect.com/science/article/pii/S0092867416304238)  
359 [science/article/pii/S0092867416304238](https://www.sciencedirect.com/science/article/pii/S0092867416304238).