



HAL
open science

DECOMICS, a shiny application for unsupervised cell type deconvolution and biological interpretation of bulk omic data

Slim Karkar, Ashwini Sharma, Carl Herrmann, Yuna Blum, Magali Richard

► **To cite this version:**

Slim Karkar, Ashwini Sharma, Carl Herrmann, Yuna Blum, Magali Richard. DECOMICS, a shiny application for unsupervised cell type deconvolution and biological interpretation of bulk omic data. 2024. hal-04660629

HAL Id: hal-04660629

<https://hal.science/hal-04660629>

Preprint submitted on 24 Jul 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 DECOMICS, a shiny application for 2 unsupervised cell type deconvolution 3 and biological interpretation of bulk 4 omic data

5
6
7 Slim Karkar^{1,*}, Ashwini Sharma², Carl Herrmann², Yuna Blum^{3,*}⊕, Magali Richard^{4,*}⊕

8
9 ¹IBGC, UMR 5095, University of Bordeaux, CNRS, Bordeaux Bioinformatic Center,
10 Bordeaux, France, ²Health Data Science Unit, Medical

11 ² Faculty Heidelberg and BioQuant, Heidelberg, Germany, ³ IGDR (Institut de Genetique et
12 Developpement de Rennes), UMR 6290, ERL

13 ³ U1305, Equipe Labellisée Ligue Nationale contre le Cancer, Univ Rennes, CNRS, INSERM,
14 35000, Rennes, France

15 ⁴ TIMC, UMR 5525, Univ. Grenoble Alpes, CNRS, F-38700, Grenoble, France

16 *Corresponding authors : slim.karkar@u-bordeaux.fr ; [magali.richard@univ-grenoble-
alpes.fr](mailto:magali.richard@univ-grenoble-
17 alpes.fr) ; yuna.blum@univ-rennes1.fr

18 ⊕Equal co-last contributors

19 keywords: R shiny, deconvolution, unsupervised approaches, transcriptomic data,
20 methylome, tissue heterogeneity

21 Abstract

22

23 Unsupervised deconvolution algorithms are often used to estimate cell composition from
24 bulk tissue samples. However, applying cell type deconvolution and interpreting the
25 results remains a challenge, even more without prior training in bioinformatics. We
26 propose here a tool for estimating and identifying cell type composition from bulk
27 transcriptomes or methylomes. DECOMICS is a shiny-web application dedicated to
28 unsupervised deconvolution approaches of bulk omic data. It provides (i) a variety of
29 existing algorithms to perform deconvolution on the gene expression or methylation-level
30 matrix, (ii) an enrichment analysis module to aid biological interpretation of the
31 deconvolved components, based on enrichment analysis, and (iii) some visualisation
32 tools. Input data can be downloaded in csv format and pre-processed in the web
33 application (normalisation, transformation and feature selection). The results of the
34 deconvolution, enrichment and visualisation processes can be downloaded.

35 **Availability and implementation:** DECOMICS is an R-shiny web application that can be
36 launched (i) directly from a local R session using the R package available here
37 <https://gitlab.in2p3.fr/Magali.Richard/decomics> (either by installing it locally, or via a
38 virtual machine and a Docker image that we provide); or (ii) in the Biosphere - IFB Clouds
39 Federation for Life Science (French Institute of Bioinformatics), a multi-cloud environment
40 scalable for high performance computing: [https://biosphere.france-](https://biosphere.france-bioinformatique.fr/catalogue/appliance/193/)
41 [bioinformatique.fr/catalogue/appliance/193/](https://biosphere.france-bioinformatique.fr/catalogue/appliance/193/), which requires creating an account and
42 logging onto the platform

43 Introduction

44 Identification of the cell composition contributing to bulk molecular signals is a major
45 challenge in molecular analysis in various applications such as cancer¹. The development
46 of in silico deconvolution methods has made it possible to revisit existing bulk omic data
47 from large patient cohorts with regard to intra-sample heterogeneity, and thus to compare
48 sample cell composition with available clinical annotations such as treatment response.
49 Both supervised² and unsupervised³⁻⁷ deconvolution methods have been proposed in the
50 literature. Supervised methods estimate the component proportions using known cell-type
51 reference matrices, whereas unsupervised methods estimate both the reference profiles
52 and the component proportions, without prior knowledge except for the number of
53 components to be considered. Supervised approaches are therefore limited by the quality
54 of the reference signatures, while unsupervised approaches present difficulties in
55 interpreting the inferred components and estimating the number of components to be

56 considered.

57 An intriguing advantage of unsupervised methods is that, unlike supervised methods, they
58 can identify new cell populations or populations that would not have been taken into
59 account a priori. Unsupervised deconvolution approaches have recently been used to
60 identify radiogenomic signatures to predict prognosis of colorectal cancer⁸ , to identify
61 cellular compartment in unknown tumoral samples⁹, or to infer clinical outcomes in
62 melanoma patients¹⁰ .

63 Applying unsupervised deconvolution approaches to patient cohorts is now possible
64 thanks to recent advances in high-throughput sequencing that have generated an
65 enormous amount of transcriptomic data as well as numerous methylome data. This is
66 particularly true in the field of oncology, where more and more biological samples are
67 being sequenced to help with patient stratification and prognosis. However, analyzing this
68 type of data requires professional coding skills, which are rarely available to clinicians.

69 In order to appeal large data analysis to a wider audience, user-friendly alternatives have
70 been devised, including the development of the R-shiny package, which enables
71 interactive web applications to be built using the R statistical and data mining software¹¹. If
72 web applications exist to apply supervised algorithms on clinical datasets¹², exploration of
73 unsupervised algorithms have been so far limited to bioinformaticians with computing
74 skills and do not provide guidance for biological interpretation of their outputs.

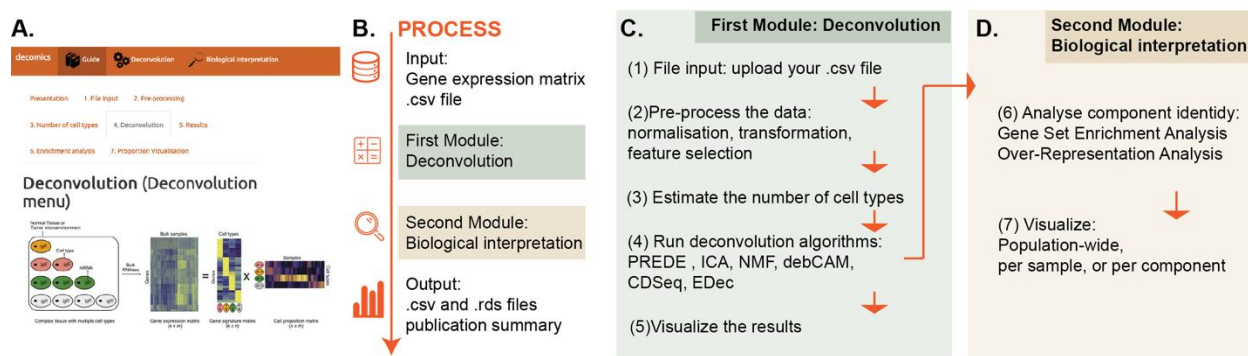
75 Here we propose DECOMICS, a user-friendly shiny interactive web application designed to
76 perform unsupervised deconvolution on transcriptomic and methylome (DNAm) data. Six

77 different unsupervised methods are implemented in DECOMICS, including the most
 78 commonly used (ICA and NMF), and more recent algorithms (CDSeq⁴, debCAM³ and
 79 PREDE⁵, EDec⁷). Our tool also provides guidance during the process and helps with the
 80 biological interpretation of the results, which should be of significant interest for both
 81 bioinformaticians and clinicians.

82 Software description

83 DECOMICS is a Shiny application available as an R package that can be built from source
 84 (GitLab access) or used online through the biosphere cloud of the IFB (Institut Français de
 85 Bioinformatique). DECOMICS workflow is illustrated in Figure 1. It includes a Guide
 86 section, that can be used as a material and methods, and two main modules: the
 87 deconvolution module and a biological interpretation module (Figure 1A-B).

88



89 Figure 1: An overview of DECOMICS workflow to perform deconvolution of omics data. A.
 90 Screenshot of the application. B. Summary of the DECOMICS process. C. Workflow of the
 91 module 'Deconvolution'. D. Workflow of the module 'Biological interpretation'.

92

93 Deconvolution module

94 The deconvolution module (Figure 1C) is used to load the data, carry out the pre-
95 processing, estimate the number of components, run the deconvolution and visualize the
96 results. **File input (1)** requires a .csv file containing the omic data (either gene expression
97 or DNAm) with samples in columns and features in row. Gene expression can be in the
98 form of raw counts, processed counts or processed gene expression in the case of
99 microarray-based technologies. We offer basic pre-processing features for gene
100 expression data in step 2. DNAm data should be provided in the form of β -values. Basic
101 **Pre-Processing (2)** of the gene expression data can be achieved within the application,
102 including normalization adapted to RNA-seq data (Read-per-million or DESeq2) and
103 transformation (log2 or pseudoLog). Pre-processing also includes the option to select a
104 subset of features. Specifically, one can choose the top 1,000 or 5,000 gene expressions,
105 or the top 10,000 or 20,000 β -values, based on the highest coefficients of variation. The
106 **Number of components (3)** to infer (corresponding to the deconvoluted components)
107 have to be estimated by the user. In general, the optimal number of components can be
108 identified through various methods, including Cattell's rule applied to Principal
109 Component Analysis (PCA) eigenvalues¹³, the Minimum Description Length (MDL)³,
110 bootstrapping techniques¹⁴, and cross-validation methods²¹. However, in a prior
111 benchmark analysis, we found that different methods yielded comparable results⁵. For the
112 sake of clarity, we have chosen to present a single method in the application: a guidance
113 plot based on PCA eigenvalues. Then **Deconvolution (4)** is run by one of the six
114 unsupervised algorithms provided in the application. Depending on the type of omic data

115 provided, a subset of algorithms is available: (i) Independent Component Analysis, Non-
116 negative Matrix Factorization, CDSeq, debCAM and PREDE for gene expression data, and
117 (ii) ICA, NMF, debCAM and EDec for DNAm data. Finally, deconvolution **Results (5)** can be
118 visualized by (i) a ‘cell-type’ signature heatmap displaying the 5 top markers of each
119 component, and (ii) a ‘cell-type’ proportion heatmap of each component.

120 **Biological interpretation module**

121 The biological interpretation module (Figure 1D) performs an enrichment analysis on the
122 components estimated by the chosen deconvolution algorithm and displays the
123 deconvoluted proportion matrix. First, we propose an **Enrichment analysis (6)** step, to
124 help with the biological interpretation of the components identified by unsupervised
125 approaches. It proposes to perform Enrichment analysis (Gene Set Enrichment analysis¹⁶
126 or Over-representation analysis¹⁷) using various biological databases: ‘CellMatch’¹⁸, a
127 reference database derived from various resources and other reference ones (GO, GTEx,
128 KEGG, Reactome, Tissue Cell Types, Cancer Cell Types and Cancer Cell lines).

129 Second, a **Proportion Visualisation (7)** section offers the possibility to visualize the full
130 component distribution for a single sample, or the distribution of a single component
131 throughout the total cohort.

132 Material and methods

133 Deconvolution Algorithms

134 Unsupervised deconvolution problem applied to omic bulk data consists in solving an
135 equation of form $X = T * A$ where T and A are jointly inferred from X . This is achieved by
136 estimating the mixture of K cell-types, present in different proportion in each sample (cell-
137 type proportion matrix A). Therefore, X can be described as a combination of cell-type
138 specific molecular profiles (cell-type specific gene expression matrix T). The specifics of
139 each existing unsupervised deconvolution algorithm and the reasons for choosing to
140 include them or not in the DECOMICS application are presented in Supplementary Table 1.
141 In DECOMICS, we provide six different algorithms to run deconvolution:

142 - ICA: Independent Component Analysis is a blind source separation algorithm that
143 decomposes signal into statistically independent components. In DECOMICS ICA
144 deconvolution is run using fastICA (*R CRAN*) and the Deconica¹⁹ R packages. By default, 30
145 significant gene markers are selected to get component scores, using the
146 "weighted.mean" summary metric.

147 - NMF²⁰: In non-negative matrix factorization approach, the molecular profile matrix X is
148 factorized into two matrices A and T , with the property that all three matrices have no
149 negative elements. DECOMICS uses the *R CRAN* NMF package with method="snmf/r". The
150 estimated A matrix is constrained to sum the proportion to 1, and T is computed as

151 $T = A^{-1}X$ using the *ginv* inverse function from MASS R package. Finally, all negative values
152 for T are set to 0.

153 - CDSeq⁴ aims at simultaneously estimating A and T matrices using a probabilistic model
154 based on latent Dirichlet allocation (LDA). DECOMICS uses the R implementation of the
155 CDseq method CDSeqR with the following parameters: beta = 0.5, alpha = 5,
156 mcmc_iterations = 300. Reduction factor is computed to avoid expression values $> 10^5$,
157 block numbers and gene block size are computed *s. t.*, a block do not exceed 10^3 genes.

158 - debCAM³ stands for deconvolution by Convex Analysis of Mixtures. This method uses a
159 geometric approach to identify a solution to the NMF problem in the simplex space. Thus,
160 the proposed solution for A is always a proportion matrix. In DECOMICS, the function CAM
161 is called from the debCAM R packages using the following empirical parameters:
162 cluster.num is computed to be 5 times greater than the number of expected components,
163 and dim.rdc set to divide the number of input genes by a tenth.

164 - PREDE⁵ is a method that offers the possibility to conduct partial reference-based
165 deconvolution method solved via an iterative Quadratic Programming procedure. In
166 DECOMICS, PREDE function is used from PREDE R package, with the following
167 parameters: W1 = NULL (which corresponds to a complete deconvolution approach), type
168 = "GE", iters = 100 and rssDiffStrop = 1e-5.

169 - EDec-step¹⁷ estimates both average component methylation profiles and component
170 proportions using an iterative constrained matrix factorization algorithm. This algorithm
171 identifies cell type-specific methylation profiles and constituent cell type proportions by

172 minimizing the Euclidean distance between the reconstituted and original mixed
173 methylation matrices. In DECOMICS, we use the `EDec::run_edec_stage_1` function with
174 the parameters `max_its = 2000` and `rss_diff_stop = 1e-10`.

175 Gene Set Enrichment Analysis

176 In order to biologically characterize each of the unsupervised components identified,
177 biological enrichment analyses are performed. For each component, the first step consists
178 in ranking the genes according to their coordinates on the component, in order to identify
179 the most contributing genes of the component. For methylation data, CpG coordinates are
180 aggregated at the gene level, taking the maximum value observed for CpGs of the same
181 gene. This approach is a way of considering a gene as strongly contributing to the
182 component if it has at least one strongly contributing CpG. In a second step, an
183 enrichment analysis is performed either based on Gene Set Enrichment Analysis (GSEA)¹⁶
184 or Over-Representation Analysis (ORA)¹⁷. If the coordinates contain sufficient non-
185 duplicate values (threshold set at 30% by default) to enable reliable ordering of the values,
186 a GSEA analysis is performed using the `fgsea` R package²¹; otherwise, an ORA analysis is
187 performed, taking as gene selection the top 20% of the component's most contributing
188 genes and as gene universe all the genes available in the user's dataset.

189 Various biological databases can be queried. DECOMICS includes the CellMatch database
190 restricted to the human species and cell types with at least 3 marker genes per cell type.
191 After filtering, it provides marker genes for 120 different cell types across 103 normal

192 tissues and 26 tumoral tissues. DECOMICS also includes the latest versions of the
193 following biological databases provided on the Enrichr²² tool website:

194 - *Cancer Cell Line Encyclopedia* (967 terms)

195 - *CellMarker Augmented 2021* (1097 terms)

196 - *GO Biological Process 2023* (5407 terms)

197 - *GO Cellular Component 2023* (474 terms)

198 - *GO Molecular Function 2023* (1147 terms)

199 - *GTEX Tissues V8 2023* (511 terms)

200 - *KEGG 2021 Human* (320 terms)

201 - *MSigDB Oncogenic Signatures* (189 terms)

202 - *NCI 60 Cancer Cell Lines* (93 terms)

203 - *Reactome 2022* (1818 terms)

204 Previously to the enrichment analysis, components obtained from ICA-based method are
205 subjected to a reorientation, as proposed in the deconica R package based on the
206 hypothesis that the highest absolute values of a component weight should be positive

207 Availability

208 Installation instructions can be found on DECOMICS gitlab webpage:
209 <https://gitlab.in2p3.fr/Magali.Richard/decomics>. There are three installation options: (i)
210 full local installation, (ii) running locally a virtual machine (VM), or (iii) using the Biosphere-
211 IFB cloud. Local installation requires several packages to be loaded. To help users, we
212 propose a conda recipe on the DECOMICS gitlab webpage. To offer the possibility to run
213 DECOMICS on a local VM, we built a docker container. The user simply needs to install
214 docker on his machine and to launch the image we provide. Finally, DECOMICS is
215 deployed on the Biosphere portal (searchable through the RAINBio catalogue). To use the
216 clouds of IFB-Biosphere, users need to create an account and get membership of an active
217 group (more information can be found here : <https://ifb-elixirfr.github.io/biosphere/signin>).
218 Then users can deploy and connect to VM using the web interface (tutorial here: [https://ifb-](https://ifb-elixirfr.github.io/biosphere/vm_connect)
219 [elixirfr.github.io/biosphere/vm_connect](https://ifb-elixirfr.github.io/biosphere/vm_connect)).

220 Application and results

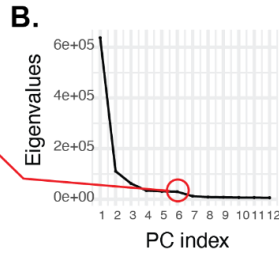
221 We have provided two use cases to illustrate the DECOMICS pipeline. The first use case is
222 based on gene expression data (GSE64385²³) dataset is available for download in .csv
223 format from the DECOMICS application. This dataset consists of a mixture of six cell types:
224 HTC116, neutrophils, natural killer cells, monocytes, B cells, and T cells. As shown in
225 Figure 2, the unsupervised components identified by DECOMICS are highly consistent with
226 the constituent cell types. In this example, the deconvolution algorithm employed is

227 debCAM and the functional enrichment analyses was performed using the CellMatch
228 database. The second use case utilizes DNAm profiles from reconstituted mixtures of six
229 purified immune cells derived from human blood samples (GSE77797²⁴). This example
230 further showcases the efficacy of the DECOMICS pipeline for DNAm unsupervised
231 deconvolution and biological interpretation of the unsupervised components. Detailed
232 information on this use case is provided in Supplementary Figure 1, and the corresponding
233 .csv dataset is available for download from the DECOMICS application.

234 In the DECOMICS pipeline, the input data consist of simple count tables in .csv format,
235 which can be uploaded directly to the application. Example input files demonstrating the
236 input format are available for download from the application. The pre-processed data
237 (.csv), deconvolution results (.rds), top 100 contributing genes of each component (.csv),
238 estimated typical gene expressions for components (.csv), and proportion estimates (.csv)
239 can also be downloaded from the application. Additionally, enrichment analyses can be
240 downloaded in the form of a .csv table containing the enrichment scores and *P*-values for
241 the queried database.

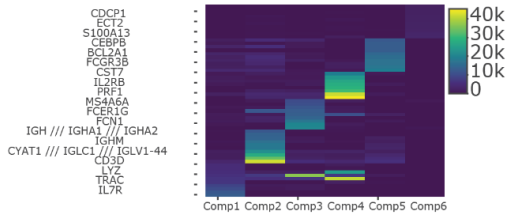
A.

Source : GSE64385
 Input.Type : RNA - Gene expression
 Normalisation : None
 Transformation : linear
 Feature.Selection : cv1000
 Number.of.Component : 6
 Method : debCAM



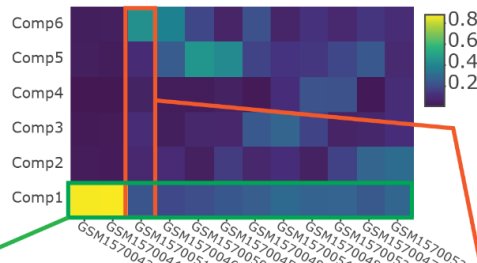
C.

Component Signature : top markers

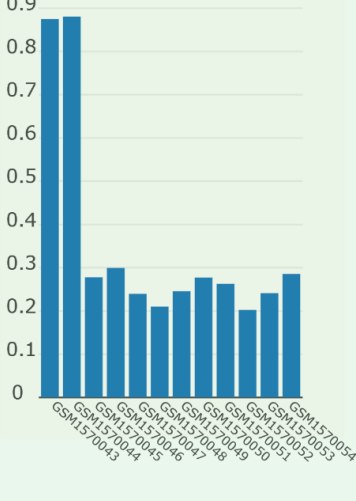


D.

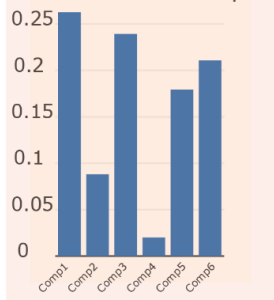
Component proportions for each Samples



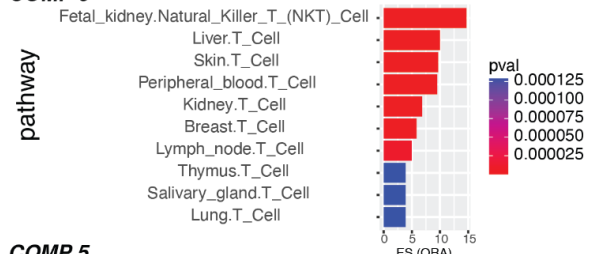
E. Component distribution in the cohort



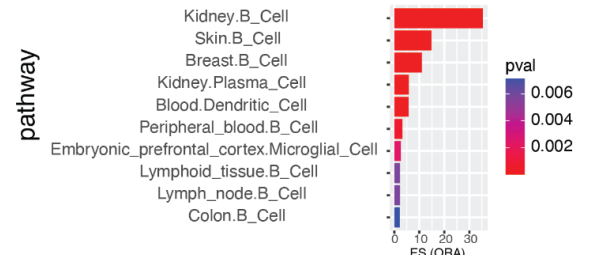
F. Components proportion in the selected sample



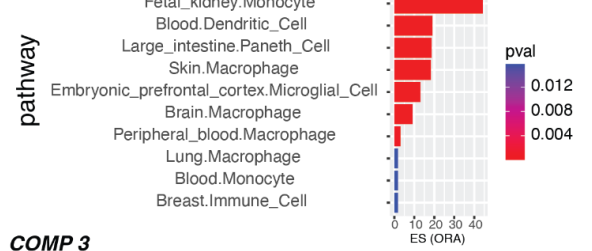
G. COMP 6



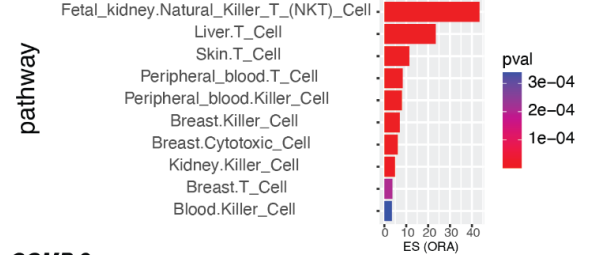
COMP 5



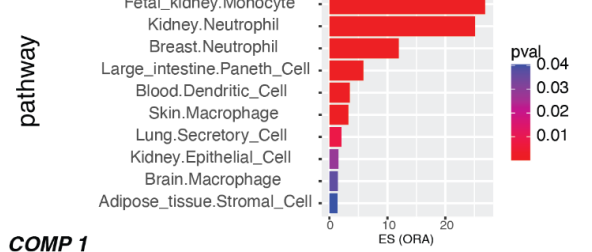
COMP 4



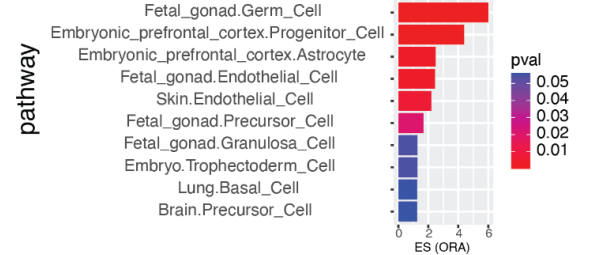
COMP 3



COMP 2



COMP 1



244 Figure 2: Illustration of DECOMICS based deconvolution of gene expression data. **A.**
245 Parameters used in the DECOMICS application are described. **B.** Scree plot illustrating the
246 selection process for the number of components to be deconvolved (module 1, step 3:
247 Number of cell types). **C-D.** Deconvolution results are presented: **(C)** Component
248 signatures are plotted, and **(D)** a heatmap showing component proportions (module 1,
249 step 5: Results). **E-F.** Visualization of component distribution: **(E)** Distribution of a specific
250 component across the cohort, and **(F)** distribution of all components within a given sample
251 (module 2, step 7: Proportion visualization). **G.** Enrichment plot displaying the enrichment
252 analysis for each component (module 2, step 6: Enrichment analysis) ORA indicates if an
253 Overrepresentation Analysis has been performed, GSEA indicates if a Gene Set
254 Enrichment Analysis has been performed for a given sample. Pval corresponds to the
255 adjusted p-values of the Enrichment Score (ES) after correction for multiple testing.

256 Discussion

257 Despite their advantages over supervised methods, the use of unsupervised deconvolution
258 methods is not trivial, as it requires a priori knowledge of the number of cell populations to
259 be considered, as well as a biological interpretation of the estimated components. Here
260 we provide a user-friendly interactive web-application to perform unsupervised
261 deconvolution by assisting the user in the choice of the number of components and
262 biological interpretation of the results.

263 Significant efforts have been made by colleagues to integrate unsupervised deconvolution
264 into analysis pipelines, including determining the number of components to infer and
265 interpreting the biological data^{25,26}. However, executing these pipelines in their entirety
266 necessitates the installation of R packages and the use of command line interfaces, which
267 our web application avoids. Additionally, these pipelines offer a limited selection of
268 deconvolution methods, restricted to those developed by the authors of the integrative
269 pipeline. In this work, we propose an unbiased approach that incorporates several
270 deconvolution methods.

271 DECOMICS stands out for its ease of use, speed, and comprehensive documentation. It is
272 designed to be accessible to users without expertise in R or biostatistics. With its
273 integrated biological interpretation feature, users can seamlessly perform both
274 deconvolution and interpretation within the same platform, eliminating the need for
275 external tools. Moreover, DECOMICS is highly adaptable, allowing for easy updates and
276 redeployment as new reference-free methods or enrichment analysis databases become
277 available. This flexibility ensures that DECOMICS remains at the cutting edge, capable of
278 incorporating the latest technical and methodological advancements.

279 Competing interests

280 No competing interest is declared.

281 Authors contributions and statements

282 S.K., Y.B and M.R. conceived the project. S.K., A.S., C.H., Y.B and M.R. contributed to code
283 development and results analysis. S.K., Y.B. and M.R. wrote the manuscript.

284 Acknowledgements

285 This work is a contribution of the EIT Health program COMETH. It has been partially
286 supported by MIAI @ Grenoble Alpes (ANR-19-P3IA-0003), and by the French Agency for
287 National Research (CauseHet // ANR-22-CE45-0030). Finally, it has also been carried out
288 with financial support from ITMO Cancer of Aviesan within the framework of the 2021-2030
289 Cancer Control Strategy, on funds administered by Inserm (ACACIA project AAP-MIC-
290 2021). We thank E. Amblard, H. Barbot, F. Pittion and L. Lamothe for testing the
291 application.

292

293 Bibliography

- 294 1. Nguyen, H., Nguyen, H., Tran, D., Draghici, S. & Nguyen, T. Fourteen years of cellular
295 deconvolution: methodology, applications, technical evaluation and outstanding
296 challenges. *Nucleic Acids Res.* **52**, 4761–4783 (2024).
- 297 2. Avila Cobos, F., Alquicira-Hernandez, J., Powell, J. E., Mestdagh, P. & De Preter, K.
298 Benchmarking of cell type deconvolution pipelines for transcriptomics data. *Nat.*

- 299 *Commun.* **11**, 5650 (2020).
- 300 3. Chen, L. *et al.* debCAM: a bioconductor R package for fully unsupervised
301 deconvolution of complex tissues. *Bioinforma. Oxf. Engl.* **36**, 3927–3929 (2020).
- 302 4. Kang, K., Huang, C., Li, Y., Umbach, D. M. & Li, L. CDSeqR: fast complete
303 deconvolution for gene expression data from bulk tissues. *BMC Bioinformatics* **22**, 1–
304 12 (2021).
- 305 5. Qin, Y. *et al.* Deconvolution of heterogeneous tumor samples using partial reference
306 signals. *PLoS Comput. Biol.* **16**, e1008452 (2020).
- 307 6. Sompairac, N. *et al.* Independent Component Analysis for Unraveling the Complexity of
308 Cancer Omics Datasets. *Int. J. Mol. Sci.* **20**, 4414 (2019).
- 309 7. Onuchic, V. *et al.* Epigenomic Deconvolution of Breast Tumors Reveals Metabolic
310 Coupling between Constituent Cell Types. *Cell Rep.* **17**, 2075–2086 (2016).
- 311 8. Zhong, M.-E. *et al.* CT-based radiogenomic analysis dissects intratumor heterogeneity
312 and predicts prognosis of colorectal cancer: a multi-institutional retrospective study. *J.*
313 *Transl. Med.* **20**, 1–12 (2022).
- 314 9. Peng, X. L., Moffitt, R. A., Torphy, R. J., Volmar, K. E. & Yeh, J. J. De novo compartment
315 deconvolution and weight estimation of tumor samples using DECODER. *Nat.*
316 *Commun.* **10**, 4729 (2019).
- 317 10. Nazarov, P. V. *et al.* Deconvolution of transcriptomes and miRNomes by independent
318 component analysis provides insights into biological processes and clinical outcomes
319 of melanoma patients. *BMC Med. Genomics* **12**, 1–17 (2019).
- 320 11. Jia, L. *et al.* Development of interactive biological web applications with R/Shiny. *Brief.*

- 321 *Bioinform.* **23**, bbab415 (2022).
- 322 12. Li, T. *et al.* TIMER2.0 for analysis of tumor-infiltrating immune cells. *Nucleic Acids Res.*
323 **48**, W509–W514 (2020).
- 324 13. Cattell, R. B. The Scree Test For The Number Of Factors. *Multivar. Behav. Res.* **1**, 245–
325 276 (1966).
- 326 14. Houseman, E. A. *et al.* Reference-free deconvolution of DNA methylation data and
327 mediation by cell composition effects. *BMC Bioinformatics* **17**, 259 (2016).
- 328 15. Lutsik, P. *et al.* MeDeCom: discovery and quantification of latent components of
329 heterogeneous methylomes. *Genome Biol.* **18**, 1–20 (2017).
- 330 16. Mootha, V. K. *et al.* PGC-1 α -responsive genes involved in oxidative phosphorylation are
331 coordinately downregulated in human diabetes. *Nat. Genet.* **34**, 267–273 (2003).
- 332 17. Goeman, J. J. & Bühlmann, P. Analyzing gene expression data in terms of gene sets:
333 methodological issues. *Bioinformatics* **23**, 980–987 (2007).
- 334 18. Shao, X. *et al.* scCATCH: Automatic Annotation on Cell Types of Clusters from Single-
335 Cell RNA Sequencing Data. *iScience* **23**, 100882 (2020).
- 336 19. Czerwinska, U. UrszulaCzerwinska/DeconICA: DeconICA first release. Zenodo
337 <https://doi.org/10.5281/zenodo.1250070> (2018).
- 338 20. Gaujoux, R. & Seoighe, C. A flexible R package for nonnegative matrix factorization.
339 *BMC Bioinformatics* **11**, 367 (2010).
- 340 21. Korotkevich, G. *et al.* Fast gene set enrichment analysis. 060012 Preprint at
341 <https://doi.org/10.1101/060012> (2021).
- 342 22. Chen, E. Y. *et al.* Enrichr: interactive and collaborative HTML5 gene list enrichment

- 343 analysis tool. *BMC Bioinformatics* **14**, 128 (2013).
- 344 23. Becht, E. *et al.* Estimating the population abundance of tissue-infiltrating immune and
345 stromal cell populations using gene expression. *Genome Biol.* **17**, 1–20 (2016).
- 346 24. Koestler, D. C. *et al.* Improving cell mixture deconvolution by identifying optimal DNA
347 methylation libraries (IDOL). *BMC Bioinformatics* **17**, 1–21 (2016).
- 348 25. Li, Z. & Wu, H. TOAST: Improving reference-free cell composition estimation by cross-
349 cell type differential analysis. *Genome Biol.* **20**, (2019).
- 350 26. Scherer, M. *et al.* Reference-free deconvolution, visualization and interpretation of
351 complex DNA methylation data using DecompPipeline, MeDeCom and FactorViz. *Nat.*
352 *Protoc.* **15**, 3240–3263 (2020).

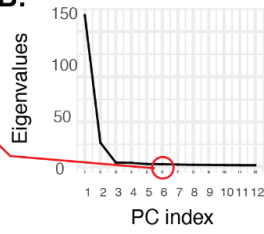
353 Supplementary material

354

A.

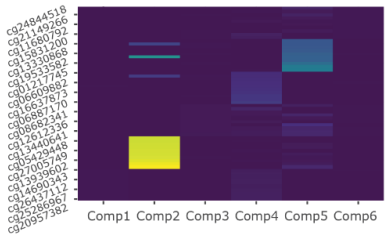
Source: GSE77797
 Input.Type: DNAmethylation 450K
 Normalisation: NA
 Transformation: NA
 Feature.Selection: cv20,000
 Number.of.Component: 6
 Method: debCAM

B.

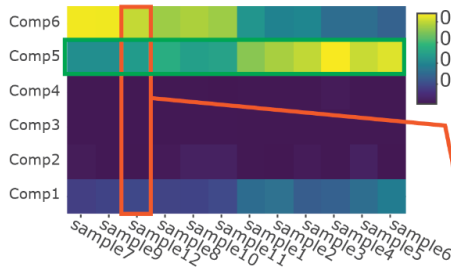


C.

Component Signature : top markers

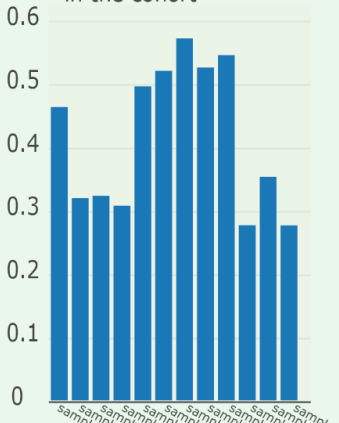


D. Component proportions for each Samples



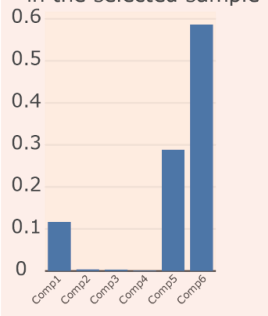
E.

Component distribution in the cohort

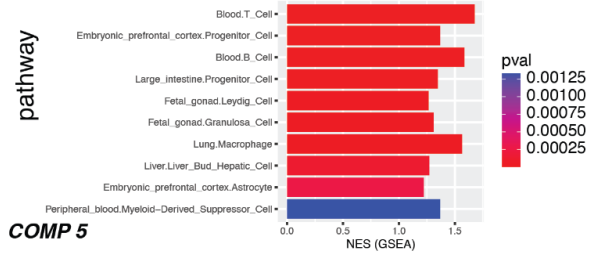


F.

Components proportion in the selected sample



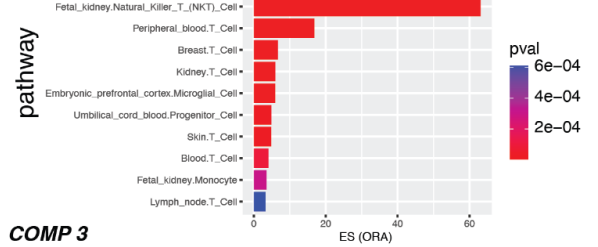
G. COMP 6



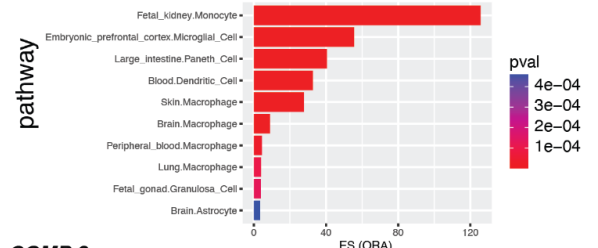
COMP 5



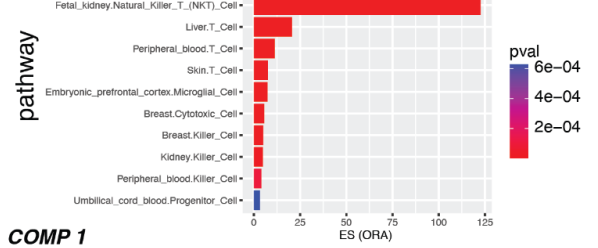
COMP 4



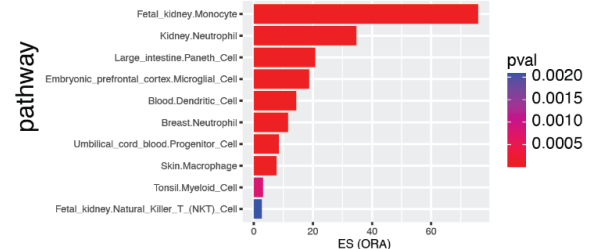
COMP 3



COMP 2



COMP 1



356 Figure S1: Illustration of DECOMICS based deconvolution of DNA methylation data. **A.**
357 Parameters used in the DECOMICS application are described. **B.** Scree plot illustrating the
358 selection process for the number of components to be deconvolved (module 1, step 3:
359 Number of cell types). **C-D.** Deconvolution results are presented: **(C)** Component
360 signatures are plotted, and **(D)** a heatmap showing component proportions (module 1,
361 step 5: Results). **E-F.** Visualization of component distribution: **(E)** Distribution of a specific
362 component across the cohort, and **(F)** distribution of all components within a given sample
363 (module 2, step 7: Proportion visualization). **G.** Enrichment plot displaying the enrichment
364 analysis for each component (module 2, step 6: Enrichment analysis) ORA indicates if an
365 Overrepresentation Analysis has been performed, GSEA indicates if a Gene Set
366 Enrichment Analysis has been performed for a given sample. Pval corresponds to the
367 adjusted p-values of the Enrichment Score (ES) after correction for multiple testing.

