

Published in final edited form as:

*Curr Opin Immunol.* 2013 October ; 25(5): . doi:10.1016/j.coi.2013.09.015.

# Computational Deconvolution: Extracting Cell Type-Specific Information from Heterogeneous Samples

Shai S. Shen-Orr<sup>a,b,c,\*</sup> and Renaud Gaujoux<sup>c</sup>

<sup>a</sup>Rappaport Institute of Medical Research, Technion-Israel Institute of Technology, Haifa, 31096, Israel

<sup>b</sup>Dept. of Immunology, Faculty of Medicine, Technion-Israel Institute of Technology, Haifa, 31096, Israel

<sup>c</sup>Faculty of Biology, Technion-Israel Institute of Technology, Haifa, 31096, Israel

## Abstract

The quanta unit of the immune system is the cell; yet analyzed samples are often heterogeneous with respect to cell subsets which can mislead result interpretation. Experimentally, researchers face a difficult choice whether to profile heterogeneous samples with the ensuing confounding effects, or a priori focus on a few cell subsets of interest, potentially limiting new discoveries. An attractive alternative solution is to extract cell subset-specific information directly from heterogeneous samples via computational deconvolution techniques, thereby capturing both cell-centered and whole system level context. Such approaches are capable of unraveling novel biology, undetectable otherwise. Here we review the present state of available deconvolution techniques, their advantages and limitations, with a focus on blood expression data and immunological studies in general.

## Keywords

deconvolution; gene expression; systems biology; cell-type; cell-specific

## 1. Introduction

The cellular composition of many biological samples is heterogeneous and varying, that is multiple cell-type subsets are present in each sample, and different samples show high variance between one another in relative cell subset proportions (from hereon, heterogeneous sample). Moreover, many physiological and pathological processes involve cell motility (e.g. infiltration) and differentiation, ultimately resulting in marked shifts in sample cell subset composition (Figure 1A). An example of importance is peripheral blood, which is composed of many different immune cell subsets, comprising the fundamental units of a complex system, whose state and interaction activity reect the type of biological processes taking place whether in health or disease. The ability to measure and interpret phenotypic changes between specific conditions at the cell subset level is therefore critical to obtain a detailed understanding of the role of each cell subset within the immune system.

© 2013 Elsevier Ltd. All rights reserved.

\*Corresponding author, shenorr@technion.ac.il (Shai S. Shen-Orr).

**Publisher's Disclaimer:** This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

While the problem of sample heterogeneity has long been acknowledged [1, 2, 3], researchers have struggled between the choice of focusing on a single cell subset or ignoring the problem and assaying heterogeneous samples. Whereas cell-isolation entails a loss of a systems perspective (i.e. biologically meaningful changes happening in multiple cell subsets and between them), ignoring sample heterogeneity entails misleading and difficult to interpret results. This has been an especially sore point in genomic scale studies (predominantly whole genome gene expression), where it is easy to lose sight of the natural cellular-context of the data amongst thousands of measured features and draw incorrect conclusions.

An emerging solution to this dilemma, comes in the form of computational deconvolution methodologies, capable of extracting cell type-specific information directly from data generated from heterogeneous samples. Research addressing this issue started with the pioneering work of Venet et al. [4], but received relatively little attention. More recently, the increased application of genomic tools to human samples, who exhibit high sample heterogeneity, has spurred further developments. In particular, the availability of data from large efforts that profiled the gene expression of multiple known cell subsets and highlighted distinct differences between them [5, 6] both motivated such computational algorithm development and served as useful input for them.

This article reviews this active area of research and draws attention on the advantages and limitations of the different proposed approaches. For the sake of simplicity, we focus on peripheral blood, which is the primary source of samples in human immunology studies, and in particular on gene expression data and group differences analyses. However research has concomitantly focused on other tissue and data types. Eventually, such techniques have the potential to provide a valuable cell-centered view of the immune system, bringing new insights into inter and intra cellular dynamics, and cell subset states in health and disease.

## 2. The confounding effects of sample cell-type subset heterogeneity

Most genes are expressed to varying degrees across multiple cell subsets in a tissue or organism, implying that the measured abundance of any such transcript is confounded by the composition of the sample from which it is measured. More precisely, we may break down the total measured abundance of a gene in a sample into three Abundance Components: (1) that due to the characteristic condition of a sample (e.g. disease type etc.), (2) that due to the individual variation, genotype-specific or technical measurement variation, and (3) that due to the average characteristic abundance of a gene as a function of the underlying cell subsets in a sample and their relative proportions (Figure 1B). The extent of contribution of each Abundance Component to total gene expression measured differs between cell subsets, conditions and individuals on a gene by gene basis. In general, however, we note that whereas Abundance Components (1) and (2) describe important biological differences in gene expression (excluding individual technical variation), Abundance Component (3) reflects gene expression variations that are more succinctly captured and interpreted by reporting differences in proportions of cell subsets.

As the composition of relative cell-type subset proportions varies between samples, analysis of sample differences by any methodology oblivious to the underlying cell-type heterogeneity will suffer from several strong drawbacks:

**Misleading information:** if a cell subset's proportion difference is correlated with the phenotype of interest, classical analysis methods (i.e. cell subset agnostic) are prone to produce more false positive differentially expressed genes by detecting differences in total abundance stemming from this shared overall correlation

(Figure 1C). Ideally, differences due to biological condition could be decoupled from those resulting from an artifact of sample heterogeneity (see next section).

**Signal dilution:** differences in genes expressed in cell subsets present at low proportion in a sample may be masked by the signal coming from the same gene expressed in a prevalent cell subset, or alternatively by a high sample-wise variance of the proportions of cell subsets in which the gene is expressed. A prime example of this problem is observed when analyzing gene expression samples derived from whole blood versus PBMCs of the same subjects, where tissue type is a greater determinant of similarity than subject condition [7, 2] (primarily due to the presence of neutrophils, which constitute anywhere from 50–70% of whole blood samples).

**Result interpretation:** without accounting for varying cell subset proportions, it is difficult to attribute the observed effect to any given cell subset (Figure 1D). Often a follow up experiment is required to understand from which cell subset a detected signal originated. Eventually, this makes it harder to build a detailed picture of the immune system and precisely delineate its main driving components – in the considered condition.

Given the above raised difficulties, it may be surprising that sample collection and analysis are still very commonly performed on highly heterogeneous tissues (e.g. peripheral blood in humans, spleen in the case of mice), despite the availability of experimental techniques to isolate and profile fairly homogeneous cell-population, and most recently to profile them at the single cell level [8]. From a research perspective standpoint, effectively choosing whether to profile a heterogeneous tissue or to isolate specific cell subsets should depend on whether one wishes to get a systems perspective by profiling all accessible cell subsets involved in a certain biological process, or conduct a focused study in which cell subsets of interest are a priori chosen. This is especially a problem in genomic level studies which often strive for a systems perspective but where the interpretation of data without cellular context is most difficult. In practice however, based on our own experiences, the decision as to how and what to profile is dictated by limitations in sample material, financial considerations, concerns with respect to biases introduced during isolation procedures [2], or an under-appreciation to the inherent problems in using heterogeneous samples.

### 3. Extracting cell type-specific information from heterogeneous tissue

An attractive approach for gaining insight on cell-subset specific information, is to estimate the proportion and/or gene expression profile of different cell subsets directly from the heterogeneous samples via computational methodologies, thereby preserving the whole-systems perspective, as well as obtaining an unbiased cell-based view (Figure 2).

Present computational methodologies for extracting cell type-specific information from heterogeneous sample data may be divided into five main method classes (Figure 3, Table 1 for a detailed list of methodologies, and Supplementary Material for in depth review of requirements and limitations of each class), which we define based on the type of input they require, and the resolution of the output they offer: Quantifying the presence of different cell subsets in heterogeneous samples may be performed either at low-resolution, by (A) Cell-type enrichment [9, 10, 11, 12, 13, 14], or at high-resolution, by (B) Cell proportion estimation methodologies [15, 16, 17, 18, 19]. Both method classes rely on reference signature expression profiles of the different cell subsets or known cell subsetspecific marker genes [20, 12, 21, 22, 23, 24], and respectively provide detection of the presence/implication of a cell subset in an heterogeneous sample or a disease condition signature (Figure 3A), or actual numeric estimates of the relative proportions of each cell subset in the

data (Figure 3B). Estimating measured abundance accounting for sample heterogeneity may be performed either at low-resolution, by (C) Sample heterogeneity correction [25, 26, 27, 28, 29], or at high-resolution, by (D) cell type-specific deconvolution methodologies [25, 30, 26]. Both methodologies rely on the availability of the proportions of one or more cell subset present in samples, and respectively provide a correction of the measured sample data for biases introduced due to cell subset sample heterogeneity (Figure 3C), or estimates of cell subset-specific gene expression profiles for each cell subset – for which proportion information was provided, which can be subsequently used to infer cell type-specific differential expression between groups of samples (Figure 3D). We note that the input cell-subset proportions in these two method classes, may come either from actual measurements or computationally estimated. In fact, a fifth category (E) consists of *complete deconvolution* methods which estimate both proportions and cell type-specific expression profiles, often using a combination of deconvolution methods (B and D), and require some limited prior knowledge on proportions [31] or expression profiles [30, 32, 19, 33, 34, 35, 36].

#### 4. Gain in biological insights

Computational deconvolution methods aim at providing a cell-centered view of heterogeneous molecular data, by decoupling the effect of proportion from cell type-specific phenotype. In particular, they have the potential to mine high-throughput data in a way that even upcoming laboratory techniques may not yet or ever handle, e.g., due to limitations of cell surface markers for cell-sorting or to the mere unavailability of biological material for past studies. Notably, they have already proved to be able to provide new insights in complex diseases such as autoimmune disease or cancer. Abbas et al. [15] deconvolve whole blood samples from Systemic Lupus Erythematosus (SLE) patients, identifying specific changes in leukocyte proportions and activation (NK, T and monocytes, in particular), as well as correlation of proportions with treatment type and other clinical measures. Deconvolution based cell-type specific differential expression of acute rejection versus stable patients identified a previously unsuspected role for monocytes in both kidney [25] and cardiac transplant [39], respectively undetectable or only mildly detectable from whole blood. Liu et al. [27] showed that DNA methylation is a potential mediator of genetic risk in Rheumatoid Arthritis (RA), and highlighted the importance of correcting for sample heterogeneity in blood DNA methylation data. Quon et al. [34] showed that significant improvements in outcome prediction of lung and prostate cancer can be achieved when building a classifier on computationally purified tumor data, as opposed to data from bulk biopsies. Given the large and increasing number of cell-types known, and the desire to capture their difference in behavior, we foresee deconvolution methodologies, which offer increased resolution and interpretability at little or no extra costs, being increasingly utilized such that they become part of main stream analysis pipelines in human profiling studies.

#### 5. Limitations of computational approaches

Despite the various successful application of computational deconvolution methodologies, we believe several open issues remain to be investigated before they become widely adopted. First, a better understanding of the accuracy lower bound for estimates of cell subsets proportions or differential gene expression detection must be developed. This general accuracy is difficult to assess because of the many factors to consider (proportion dependencies, individual variation, clinical condition, cell-cell interaction), a large scale evaluation on simulated and public data could provide much information with respect to their power. Second, and particularly relevant in the case of blood, is the development of algorithms capable of performing “*deep deconvolution*”, i.e. accurately estimating from a whole blood or PBMC sample, the proportions and expression patterns of a greater number of cell subsets, going further down into the hematopoietic tree (T-regs, naive, memory, and

effector cell subsets). One challenge in achieving this is the minimum required sample size that increases together with the number of considered cell subsets, a restriction that may be lifted using regularization methodologies. Finally more research on how to efficiently perform single-sample deconvolution, where cell type-specific profiles are estimated for each individual sample rather than groups, as already proposed in cancer tissues [34], would enable moving towards more personalized tests. Of note, one of us [RG], has recently published an R package, *CellMix* [40], which complies together many of the published computational gene expression deconvolution methodologies, and facilitates future algorithm development and benchmarking.

## 6. Conclusion

Cell subset heterogeneity is inherent to most primary biological samples, which may confound downstream data analysis if it is not taken into account and strongly restricts result interpretability. From a systems immunology perspective to health and disease, it is critical to be able to assess each cell subset's state and interactions, over a range of condition and molecular environment. In this respect, computational deconvolution methodologies showed to be powerful tools capable of providing novel high-resolution system-wide insights.

In this review, we focused our discussion on gene expression studies, where the majority of research on the effects of sample cell subset heterogeneity on analysis and the benefit of addressing it has been done to date. Expressed transcripts are of course not the only molecular species whose interpretation of measured abundance is affected by sample cell subset heterogeneity. For example, computational techniques have already been applied to DNA methylation data, which revealed new biology otherwise masked behind sample heterogeneity [28, 27]. These are not the exception, and we would rather propose that the opposite is the case, namely the majority of molecular species assayed vary in their total abundance from one cell subset to another, and perhaps even for secreted molecular species. Thus their measurement and analysis would be affected by sample composition and benefit from computational techniques aimed at providing a cell-centered view of the system. We envision that different molecular species and measurement modalities would require the tailoring of their own particular methodology, but expect that any new methodology development will likely fall into one of five meta-methods classes we defined here for obtaining cell type-specific information from cell subset heterogeneous gene expression samples.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

This work was supported by US National Institutes of Health (NIH) (U19 AI057229). SSO is a Taub Fellow. RG is supported by the Lady Davis Fellowship.

## References

1. Davey HM, Kell DB. Flow cytometry and cell sorting of heterogeneous microbial populations: the importance of single-cell analyses. *Microbiological reviews*. 1996; 60(4):641–696. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=239459&tool=pmcentrez&rendertype=abstract>. [PubMed: 8987359]
2. Whitney A, Diehn M, Popper S, Alizadeh A, Boldrick J, Relman D, et al. Individuality and variation in gene expression patterns in human blood. *Proceedings of the National Academy of Sciences of the United States of America*. 2003; 100(4):1896. URL: <http://www.pnas.org/content/>



100/4/1896.short. [PubMed: 12578971] \*\* The authors provide a comprehensive analysis of factors of variation in whole blood and PBMC gene expression, notably showing that a significant fraction of the observed variation reflects differences in cell type proportions.

3. De Ridder D, Van Der Linden CE, Schonewille T, Dik WA, Reinders MJT, Van Dongen JJM, et al. Purity for clarity: the need for purification of tumor cells in DNA microarray studies. *Leukemia official journal of the Leukemia Society of America Leukemia Research Fund UK*. 2005; 19(4): 618–627. URL: <http://www.ncbi.nlm.nih.gov/pubmed/15744349>.
4. Venet D, Pecasse F, Maenhaut C, Bersini H. Separation of samples into their constituents using gene expression data. *Bioinformatics*. 2001; 17(suppl 1):S279. URL: [http://bioinformatics.oxfordjournals.org/content/17/suppl\\_1/S279.short](http://bioinformatics.oxfordjournals.org/content/17/suppl_1/S279.short). [PubMed: 11473019]
5. Novershtern N, Subramanian A, Lawton LN, Mak RH, Haining WN, McConkey ME, et al. Densely interconnected transcriptional circuits control cell states in human hematopoiesis. *Cell*. 2011; 144(2):296–309. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3049864&tool=pmcentrez&rendertype=abstract>. [PubMed: 21241896] \*\* D-MAP: the authors generated a dataset of 39 human cell types across the whole hematopoietic tree, showing, in particular, that gene expression programs get reused across lineages and cell subsets are generally more similar to one another in deep branches of the hematopoietic tree.
6. Shay T, Kang J. Immunological Genome Project and systems immunology. *Trends Immunol*. 2013 \*\* The ImmGen project generated cell-specific profiles of 249 immune cells in mice, and provide analytical tools to browse and mine the data for a variety of cell-specific information.
7. Cobb JP, Mindrinos MN, Miller-Graziano C, Calvano SE, Baker HV, Xiao W, et al. Application of genome-wide expression analysis to human health and disease. *Proceedings of the National Academy of Sciences of the United States of America*. 2005; 102(13):4801–4806. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=555033&tool=pmcentrez&rendertype=abstract>. [PubMed: 15781863] \* Comparison of gene expression from whole blood and sorted cells (here T cells and Monocytes) suggests that differences in leukocyte proportions constitutes a major source of variation of the apparent global gene expression.
8. Ramsköld D, Luo S, Wang YC, Li R, Deng Q, Faridani OR, et al. Full-length mRNA-Seq from single-cell levels of RNA and individual circulating tumor cells. *Nature Biotechnology*. 2012; 30(8):777–782. URL: <http://www.nature.com/doi/10.1038/nbt.2282>.
9. Barbie DA, Tamayo P, Boehm JS, Kim SY, Moody SE, Dunn IF, et al. Systematic RNA interference reveals that oncogenic KRAS-driven cancers require TBK1. *Nature*. 2009; 462(7269): 108–112. URL: <http://www.ncbi.nlm.nih.gov/pubmed/19847166>. [PubMed: 19847166]
10. Verhaak RGW, Hoadley KA, Purdom E, Wang V, Qi Y, Wilkerson MD, et al. Integrated genomic analysis identifies clinically relevant subtypes of glioblastoma characterized by abnormalities in PDGFRA, IDH1, EGFR, and NF1. *Cancer Cell*. 2010; 17(1):98–110. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2818769&tool=pmcentrez&rendertype=abstract>. [PubMed: 20129251]
11. Bolen CR, Uduman M, Kleinstein SH. Cell Subset Prediction for Blood Genomic Studies. *BMC Bioinformatics*. 2011; 12(1):258. URL: <http://www.biomedcentral.com/1471-2105/12/258>. [PubMed: 21702940]
12. Shoemaker JE, Lopes TJ, Ghosh S, Matsuoka Y, Kawaoka Y, Kitano H. CTen: a web-based platform for identifying enriched cell types from heterogeneous microarray data. *BMC Genomics*. 2012; 13(1):460. URL: <http://www.biomedcentral.com/1471-2164/13/460>. [PubMed: 22953731]
13. Chikina MD, Huttenhower C, Murphy CT, Troyanskaya OG. Global prediction of tissue-specific gene expression and context-dependent gene networks in *Caenorhabditis elegans*. *PLoS computational biology*. 2009; 5(6):e1000417. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2692103&tool=pmcentrez&rendertype=abstract>. [PubMed: 19543383]
14. Nakaya H, Wrammert J, Lee E. Systems biology of seasonal influenza vaccination in humans. *Nature*. 2011; 12(8):786–795. URL: <http://www.ncbi.nlm.nih.gov/pmc/articles/pmc3140559/>.
15. Abbas AR, Wolslegel K, Seshasayee D, Modrusan Z, Clark HF. Deconvolution of blood microarray data identifies cellular activation patterns in systemic lupus erythematosus. *PloS one*. 2009; 4(7):e6098. URL: <http://www.ncbi.nlm.nih.gov/pubmed/19568420>. [PubMed: 19568420] \*\* The authors built an set of signature gene expression profiles for 17 immune resting/ activated

cell subsets, optimized for their ability to accurately estimate cell proportions in whole blood gene expression microarray data.

16. Gong, T.; Szustakowski, JD. DeconRNASeq: A Statistical Framework for Deconvolution of Heterogeneous Tissue Samples Based on mRNA-Seq data; Bioinformatics (Oxford, England). 2013. p. 1-2. URL: <http://www.ncbi.nlm.nih.gov/pubmed/23428642>
17. Song S, Nones K, Miller D, Harliwong I, Kassahn KS, Pinese M, et al. qpure: A tool to estimate tumor cellularity from genome-wide single-nucleotide polymorphism profiles. PloS one. 2012; 7(9):e45835. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3457972&tool=pmcentrez&rendertype=abstract>. [PubMed: 23049875]
18. Carter SL, Cibulskis K, Helman E, McKenna A, Shen H, Zack T, et al. Absolute quantification of somatic DNA alterations in human cancer. Nature Biotechnology. 2012; 30(5):413–421. URL: <http://www.nature.com/doi/10.1038/nbt.2203>.
19. Zhong Y, Wan YW, Pang K, Chow LM, Liu Z. Digital sorting of complex tissues for cell type-specific gene expression profiles. BMC Bioinformatics. 2013; 14(1):89. URL: <http://www.biomedcentral.com/1471-2105/14/89>. [PubMed: 23497278]
20. Liberzon A, Subramanian A, Pinchback R, Thorvaldsdóttir H, Tamayo P, Mesirov JP. Molecular signatures database (MSigDB) 3.0. Bioinformatics. 2011; 27(12):1739–1740. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3106198&tool=pmcentrez&rendertype=abstract>. [PubMed: 21546393]
21. Liu W, Yuan K, Ye D. Reducing microarray data via nonnegative matrix factorization for visualization and clustering analysis. Journal of biomedical informatics. 2008; 41(4):602–606. URL: <http://www.ncbi.nlm.nih.gov/pubmed/18234564>. [PubMed: 18234564]
22. Yang X, Ye Y, Wang G, Huang H, Yu D, Liang S. VeryGene: linking tissue-specific genes to diseases, drugs, and beyond for knowledge discovery. Physiological genomics. 2011; 43(8):457–460. URL: <http://physiolgenomics.physiology.org/cgi/content/abstract/43/8/457>. [PubMed: 21245417]
23. Xiao SJ, Zhang C, Zou Q, Ji ZL. TiSGeD: a database for tissue-specific genes. Bioinformatics (Oxford, England). 2010; 26(9):1273–1275. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2859128&tool=pmcentrez&rendertype=abstract>.
24. Kogenaru S, Val C, Hotz-Wagenblatt A, Glatting KH. TissueDistributionDBs: a repository of organism-specific tissue-distribution profiles. Theoretical Chemistry Accounts. 2009; 125(3-6): 651–658. URL: <http://www.springerlink.com/index/10.1007/s00214-009-0670-5>.
25. Shen-Orr SS, Tibshirani R, Khatri P, Bodian DL, Staedtler F, Perry NM, et al. Cell type-specific gene expression differences in complex tissues. Nature methods. 2010; 7(4):287–289. URL: <http://www.ncbi.nlm.nih.gov/pubmed/20208531>. [PubMed: 20208531] \*\* Computational method to estimate cell-specific differential expression between two groups of samples (using SAM). A number of other authors used it to successfully identified novel cell-driven disease signatures.
26. Qiao W, Quon G, Csaszar E, Yu M, Morris Q, Zandstra PW. PERT: a method for expression deconvolution of human blood samples from varied microenvironmental and developmental conditions. PLoS computational biology. 2012; 8(12):e1002838. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3527275&tool=pmcentrez&rendertype=abstract>. [PubMed: 23284283]
27. Liu Y, Aryee MJ, Padyukov L, Fallin MD, Hesselberg E, Runarsson A, et al. Epigenome-wide association data implicate DNA methylation as an intermediary of genetic risk in rheumatoid arthritis. Nature biotechnology. 2013; 31(2):142–147. URL: <http://www.ncbi.nlm.nih.gov/pubmed/23334450>. \*\* The authors show that DNA methylation correction for heterogeneity significantly reduces the association signals of CpG methylation with disease (here rheumatoid arthritis (RA)), and highlight the importance of taking into account sample heterogeneity in the analysis of these data.
28. Houseman EA, Accomando WP, Koestler DC, Christensen BC, Marsit CJ, Nelson HH, et al. DNA methylation arrays as surrogate measures of cell mixture distribution. BMC bioinformatics. 2012; 13:86. URL: <http://www.biomedcentral.com/1471-2105/13/86/http://www.biomedcentral.com/1471-2105/13/86http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3532182&tool=pmcentrez&rendertype=abstract>. [PubMed: 22568884] \* The authors propose a deconvolution

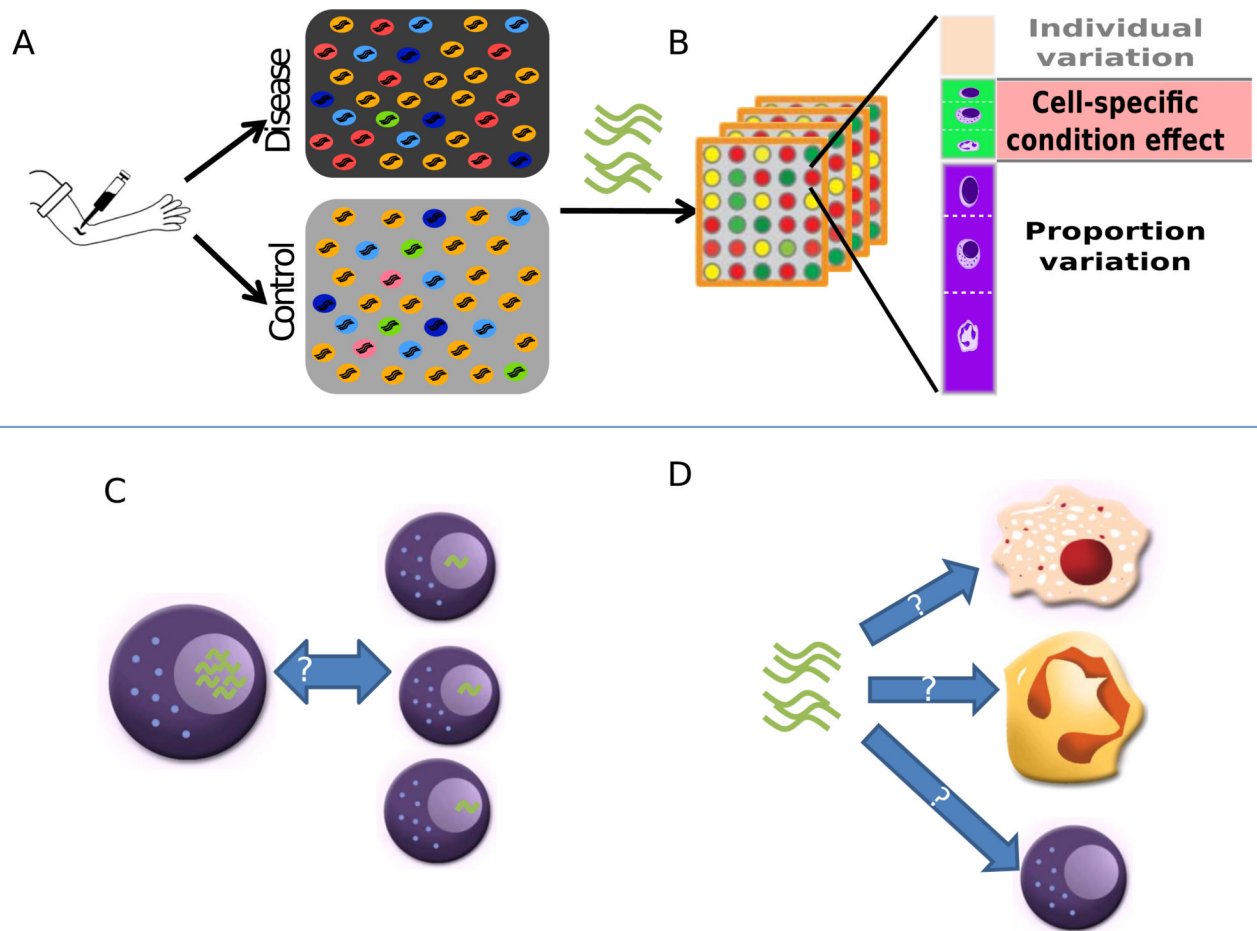
method for DNA methylation data, that accurately estimates the proportions of 8 immune cell subsets from whole blood samples.

29. Smyth GK. Linear models and empirical bayes methods for assessing differential expression in microarray experiments. *Statistical applications in genetics and molecular biology*. 2004; 3(1) Article3. URL: <http://www.ncbi.nlm.nih.gov/pubmed/16646809>.
30. Kuhn A, Thu D, Waldvogel HJ, Faull RLM, Luthi-Carter R. Population-specific expression analysis (PSEA) reveals molecular changes in diseased brain. *Nature methods*. 2011; 8(11):945–947. URL: <http://www.ncbi.nlm.nih.gov/pubmed/21983921>. [PubMed: 21983921]
31. a TE, Lehmusvaara S, Ruusuvaari P, Visakorpi T, Shmulevich I, Lahdesmaki H. Probabilistic analysis of gene expression measurements from heterogeneous tissues. *Bioinformatics (Oxford, England)*. 2010; 26(20):2571–2577. URL: <http://www.ncbi.nlm.nih.gov/pubmed/20631160http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2951082&tool=pmcentrez&rendertype=abstract>.
32. Gaujoux R, Seoighe C. Semi-supervised Nonnegative Matrix Factorization for gene expression deconvolution: a case study. *Infection, genetics and evolution : journal of molecular epidemiology and evolutionary genetics in infectious diseases*. 2012; 12(5):913–921. URL: <http://www.ncbi.nlm.nih.gov/pubmed/21930246>.
33. Repsilber D, Kern S, Telaar A, Walzl G, Black GF, Selbig J, et al. Biomarker discovery in heterogeneous tissue samples -taking the in-silico deconvolution approach. *BMC bioinformatics*. 2010; 11:27. URL: <http://www.ncbi.nlm.nih.gov/pubmed/20070912>. [PubMed: 20070912]
34. Quon G, Haider S, Deshwar AG, Cui A, Boutros PC, Morris Q. Computational purification of individual tumor gene expression profiles leads to significant improvements in prognostic prediction. *Genome medicine*. 2013; 5(3):29. URL: <http://www.ncbi.nlm.nih.gov/pubmed/23537167>. [PubMed: 23537167] \*\* Single sample deconvolution: given a set of separate profiles from tumor sample and their surrounding healthy tissue (not necessarily matched), the proposed method estimates, for each individual tumour sample, the fraction of normal tissue as well as the “pure” expression profiles of each tissue type (normal/cancer).
35. Ahn J, Yuan Y, Parmigiani G, Suraokar MB, Diao L, Wistuba II, et al. DeMix: deconvolution for mixed cancer transcriptomes using raw measured data. *Bioinformatics (Oxford, England)*. 2013; 29(15):1865–1871. URL: <http://www.ncbi.nlm.nih.gov/pubmed/23712657http://bioinformatics.oxfordjournals.org/content/early/2013/05/26/bioinformatics.btt301.short>.
36. Zuckerman NS, Noam Y, Goldsmith AJ, Lee PP. A self-directed method for cell-type identification and separation of gene expression microarrays. *PLoS computational biology*. 2013; 9(8):e1003189. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3749952&tool=pmcentrez&rendertype=abstract>. [PubMed: 23990767]
37. a Miller J, Cai C, Langfelder P, Geschwind DH, Kurian SM, Salomon DR, et al. Strategies for aggregating gene expression data: The collapseRows R function. *BMC Bioinformatics*. 2011; 12(1):322. URL: <http://www.biomedcentral.com/1471-2105/12/322>. [PubMed: 21816037]
38. Gong T, Hartmann N, Kohane IS, Brinkmann V, Staedtler F, Letzkus M, et al. Optimal deconvolution of transcriptional profiling data using quadratic programming with application to complex clinical blood samples. *PloS one*. 2011; 6(11):e27156. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3217948&tool=pmcentrez&rendertype=abstract>. [PubMed: 22110609]
39. Shannon CP, Hollander Z, Wilson-McManus J, Balshaw R, Ng RT, McMaster R, et al. White blood cell differentials enrich whole blood expression data in the context of acute cardiac allograft rejection. *Bioinformatics and biology insights*. 2012; 6:49–61. URL: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3329187&tool=pmcentrez&rendertype=abstract>. [PubMed: 22550401]
40. Gaujoux R, Seoighe C. CellMix: A Comprehensive Toolbox for Gene Expression Deconvolution. *Bioinformatics (Oxford, England)*. 2013:1–2. URL: <http://www.ncbi.nlm.nih.gov/pubmed/23825367>. \* R package that provides a standardised interface to many gene expression deconvolution methods, as well as related data like cell-specific marker gene sets or benchmark datasets, where cell proportions or sorted cell profiles are available.



**Highlights**

- Primary biological samples are often heterogeneous tissues, particularly blood samples
- Cell isolation entails a loss of system perspective and is not currently feasible for all cell subsets
- Computational deconvolution methods can provide a cell-centered system perspective
- The use of deconvolution methods requires pre-planning on experimental design
- Such approaches are applicable to many types of molecular species

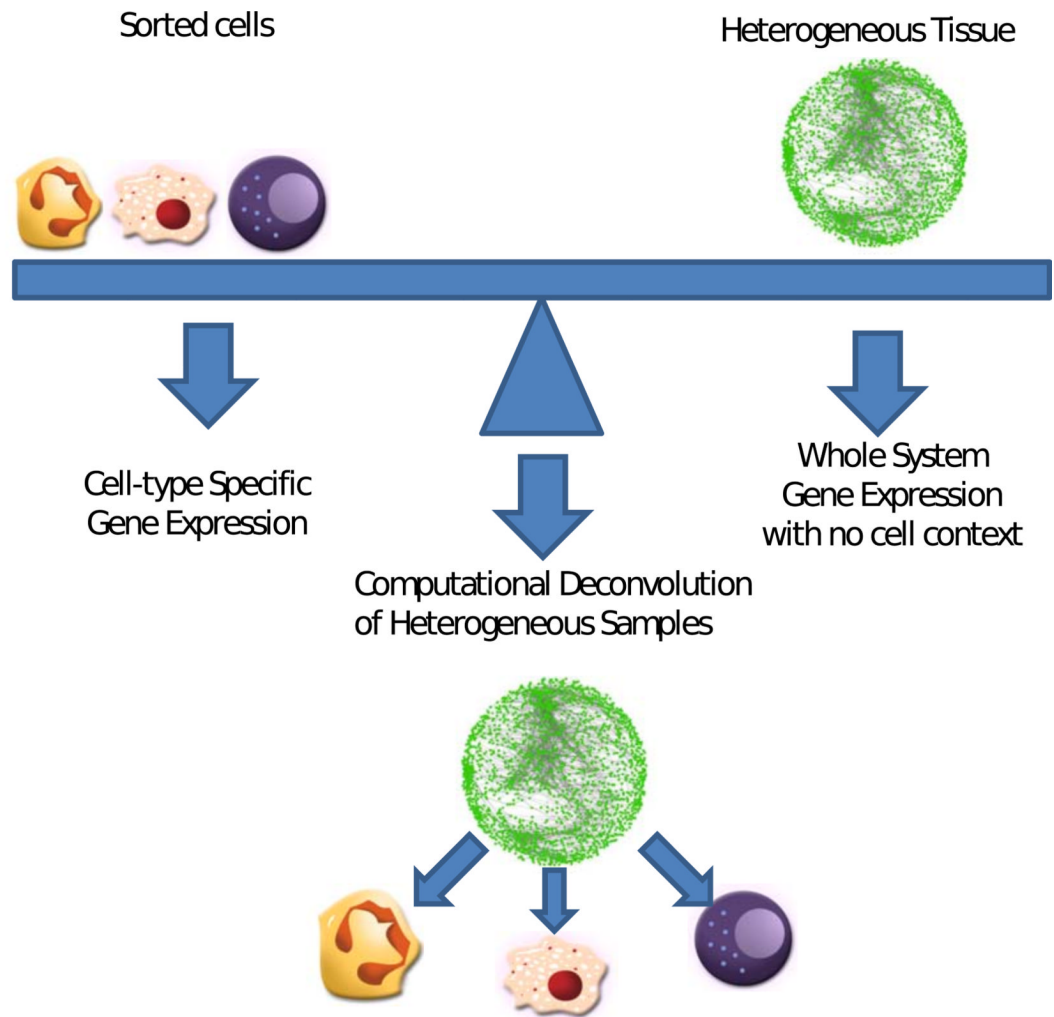


**Figure 1. Biological samples are heterogeneous with respect to underlying cell subsets, with strong implications on downstream analysis**

A) Most tissue samples are composed of multiple cell subsets, and different samples show high variance between one another in relative cell subset proportions, especially under pathological conditions. B) This implies that the total measured transcript abundance of a gene (as well as many other molecular species) is strongly affected by the cell subset composition of the sample and may be decoupled into three Abundance Components. Implications of sample heterogeneity include (C) an inability to identify whether increased total expression is due to the over-expression of a gene or to merely having more cells of a given subset in the sample, as well as (D) a difficulty in interpreting results and identifying the cell subset of origin of any detected differences.

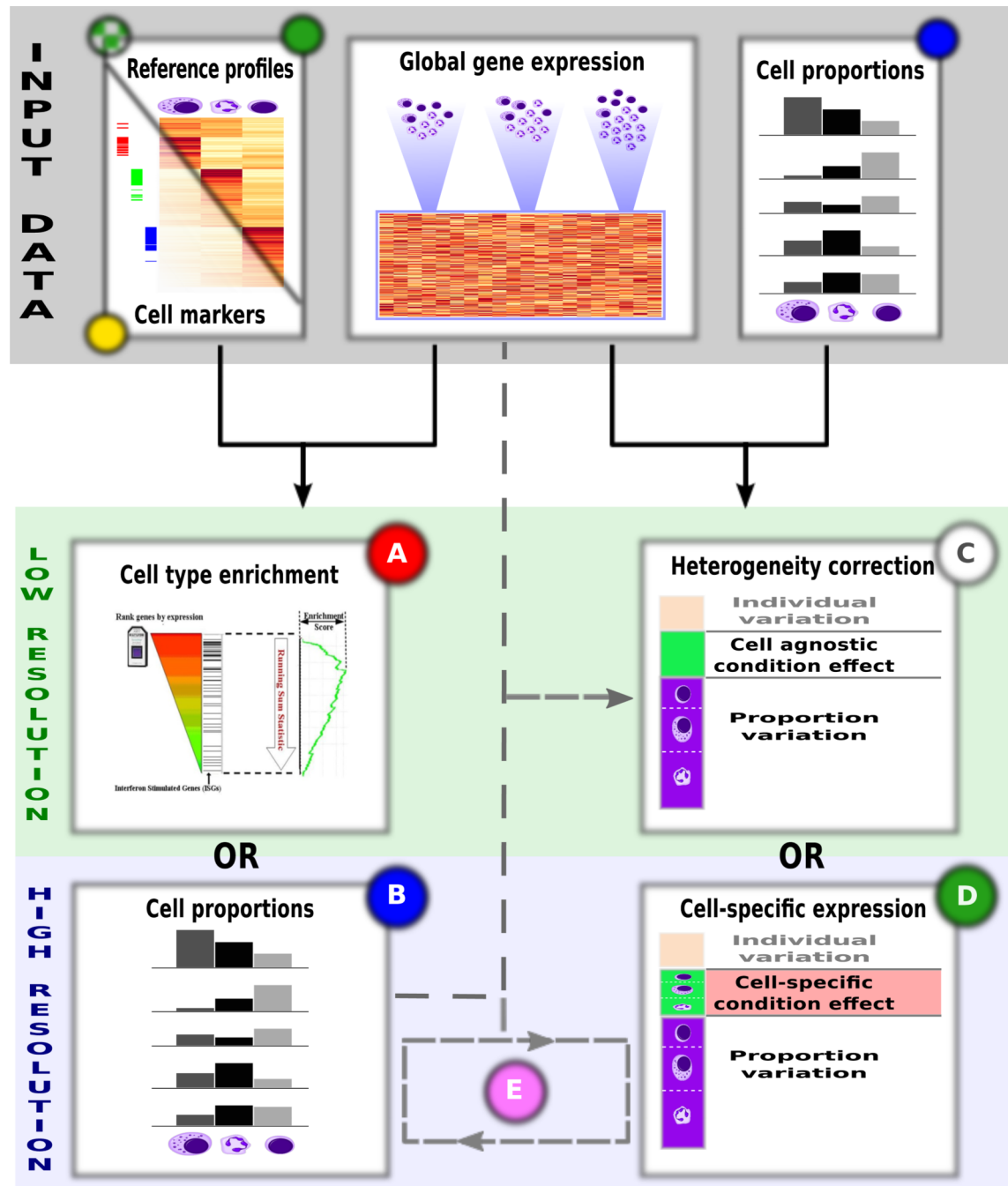
## Profiling Methodology

## Knowledge Gained



**Figure 2. Computational deconvolution methodologies enable capturing both cell-centered and system wide information**

Experimental methodologies for dealing with sample heterogeneity require either to isolate cells of interest, which perturbs the cells, entails a loss of perspective on the whole system, and is biased towards prior knowledge on which cells are of interest. The alternative, that of profiling the heterogeneous sample directly, provides a whole system view, which, however, lacks any cellular context. Computational deconvolution methodologies offer an intermediate alternative and allow to capture system level information in a cell-centered manner, a model proper to immunology, namely cells interacting with one another.



**Figure 3. Five classes of computational approaches that extract cell type-specific information from heterogeneous sample data**


Different classes of deconvolution methods defined according to the combination of the input data they require and the type and resolution of output they offer. All methods use data from heterogeneous samples, combined with either markers, signatures or proportions to (A) detect cell presence or implication of cell types, (B) estimate cell proportions, (C) correct for heterogeneity, or (D) estimate cell type-specific expression profiles. Dotted line indicates a possibility of using the output of one class of methods as input for another. Complete deconvolution methods (E) alternately estimate proportions from cell type-specific

expression and vice-versa, starting with some limited prior knowledge on proportions or expression profiles (signatures, markers).



**Table 1**

Deconvolution methods with an available user interface. Methods are grouped by classes, which are identified using the labels from Figure 3. Methods are ordered by class. Methods matching more than one class are classified by the highest resolution they provide. For each method, the type of input required and output generated is listed. The color bullets match the labels of relevant blocks in Figure 3.

In	Out	Name	Description	Availability	Tissue
		SPEC*	Predicts the most likely cellular source of a given gene expression signature [11]	<a href="http://clip.med.yale.edu/SPEC">clip.med.yale.edu/SPEC</a> R:SPEC	Blood
		CTEN*	Identifies enriched cell types in heterogeneous microarray data [12]	<a href="http://www.inuenza-x.org/jshoemaker/cten">www.inuenza-x.org/jshoemaker/cten</a>	-
		ssGSEA	Single sample GSEA [10]	ssGSEAProjection in <a href="#">GenePattern</a> [9]	Cancer
		collapseRows	Aggregates/selects a proportion proxy within cell type-specific co-expression modules [37]	R:WGCNA	-
		Abbas	Estimates proportions of 17 immune cell subsets using IRIS-based signatures [15]	CellMix	Blood
		DeconRNAseq	Similar to <i>Abbas</i> but uses quadratic programming instead of standard regression [38]	R + CellMix	-
		PERT	Perturbation model that estimates proportions and a global condition effect the reflects deviance from reference pure profiles [26]	Octave code	Blood
		methylSpectrum	Estimate proportions from DNA methylation reference profiles [28]	R:methylSpectrum	Blood
		qpure	Estimate tumor cellularity SNP microarray data from paired (tumor and normal) samples [17]	R:qpure	Cancer
		ABSOLUTE	Infers tumor purity and malignant cell ploidy directly Copy-Number-Variation data and precomputed models of recurrent cancer karyotypes [18]	R	Cancer
		csSAM*	Estimates cell/tissue specific signatures from known proportions using SAM [25]	R:csSAM + XLS plugin + CellMix	-
		PSEA*	Population-Specific Expression Analysis, using a regression model selection schema [30]	R function(s)	-
		DeMix*	Estimate tumor fraction and individual purified tumor profiles using normal tissue profiles [35]	R function(s)	Cancer
		ISOpure	Estimate tumor fraction and individual purified sample profiles using normal tissue profiles [34]	Matlab	Cancer
		DSection*	Bayesian MCMC-based estimation from priors on cell proportions [31]	<a href="http://informatics.systemsbiology.net/DSection">informatics.systemsbiology.net/DSection</a> Matlab + CellMix	-
		DSA	Digital Sorting Algorithm: complete deconvolution using a set of linear equations and quadratic programming [19]	R:dsa + CellMix	-
		deconf	Alternate least-square NMF method, using heuristic constraints [33]	R:deconf + CellMix	Blood
		ssNMF	Semi-supervised NMF algorithms, enforcing marker gene expression patterns [32]	CellMix	-

\* denotes methods that provide built-in capabilities for statistical testing or confidence interval estimation. For each method, the type of software implementation is indicated. For  $R$ , the name of the package implementing the method is listed, if available. [CellMix](#), is an  $R$  package compiling together in a standardized interface many of the published computational deconvolution methodologies for gene expression data. The tissue for which a method was developed and may be expected to perform best is also listed.