1

1	
2	The Gene Expression Deconvolution Interactive Tool (GEDIT):
3	Accurate Cell Type Quantification from Gene Expression Data
4	
5	Brian B. Nadel ¹ , David Lopez ¹ , Dennis J. Montoya ¹ , Feiyang Ma ¹ , Hannah Waddel ³ , Misha M. Khan ⁴ ,
6	Serghei Mangul ⁵ , Matteo Pellegrini ^{1,2}
7	
8	¹ Bioinformatics Interdepartmental Degree Program, Molecular Biology Institute, Department of
9	Molecular Cellular and Developmental Biology, and Institute for Genomics and Proteomics, University
10	of California Los Angeles, Los Angeles, CA
11	² Department of Dermatology, David Geffen School of Medicine, University of California Los Angeles,
12	Los Angeles, CA
13	³ Department of Mathematics, University of Utah, Salt Lake City, UT
14	⁴ Departments of Biology and Computer Science, Swarthmore College, Swarthmore, PA
15	⁵ Department of Clinical Pharmacy, USC School of Pharmacy
16	

17 Abstract

18 The cell type composition of heterogeneous tissue samples can be a critical variable in both 19 clinical and laboratory settings. However, current experimental methods of cell type quantification (e.g. 20 cell flow cytometry) are costly, time consuming, and can introduce bias. Computational approaches 21 that infer cell type abundance from expression data offer an alternate solution. While these methods 22 have gained popularity, most are limited to predicting hematopoietic cell types and do not produce 23 accurate predictions for stromal cell types. Many of these methods are also limited to particular 24 platforms, whether RNA-seq or specific microarrays. We present the Gene Expression Deconvolution 25 Interactive Tool (GEDIT), a tool that overcomes these limitations, compares favorably with existing 26 methods, and provides superior versatility. Using both simulated and experimental data, we

2

extensively evaluate the performance of GEDIT and demonstrate that it returns robust results under a wide variety of conditions. These conditions include a variety of platforms (microarray and RNA-seq), tissue types (blood and stromal), and species (human and mouse). Finally, we provide reference data from eight sources spanning a wide variety of stromal and hematopoietic types in both human and mouse. This reference database allows the user to obtain estimates for a wide variety of tissue samples without having to provide their own data. GEDIT also accepts user submitted reference data, thus allowing the estimation of any cell type or subtype, provided that reference data is available.

34

35 Author Summary

36 The Gene Expression Deconvolution Interactive Tool (GEDIT) is a robust and accurate tool 37 that uses gene expression data to estimate cell type abundances. Extensive testing on a variety of 38 tissue types and technological platforms demonstrates that GEDIT provides greater versatility than 39 other cell type deconvolution tools. GEDIT utilizes reference data describing the expression profile of 40 purified cell types, and we provide in the software package a library of reference matrices from various 41 sources. GEDIT is also flexible and allows the user to supply custom reference matrices. A GUI 42 interface for GEDIT is available at http://webtools.mcdb.ucla.edu/, and source code and reference 43 matrices are available at https://github.com/BNadel/GEDIT.

44

45 Introduction

Cell type composition is an important variable in biological and medical research. In laboratory experiments, cell sample heterogeneity can act as a confounding variable. Observed changes in gene expression may result from changes in the abundance of underlying cell populations, rather than changes in expression of any particular cell type [1]. In clinical applications, the cell type composition of tissue biopsies can inform treatment. For example, in cancer, the number and type of infiltrating immune cells has been shown to correlate highly with prognosis ([2], [3], [4]). Moreover, patients with a large number of infiltrating T cells are more likely to respond positively to immunotherapy [5].

3

53

54 For many years, cell flow cytometry via FACS sorting has been the standard method of cell 55 type quantification. More recently, single cell RNA-seq methods such as 10x Chromium, Drop-Seq, 56 and Seq-Well have become available [6],[7]. However, both approaches suffer from significant 57 limitations. FACS sorting is cumbersome and expensive, and some sample types require hours of 58 highly skilled labor to generate data. Similarly, single cell RNA-seq methods remain expensive for 59 large sample studies. Additionally, cell types such as neurons, myocytes, and adipocytes are difficult 50 for these technologies to capture due to cell size and morphology.

61

62 Both FACS sorting and single cell methods have the potential to introduce bias, as these 63 technologies require that tissue samples be dissociated into single cell suspensions. Many stromal cell 64 types are tightly connected to one another in extracellular matrices. The procedures necessary to 65 create single cell suspensions can damage some cells, while others remain in larger clusters that are 66 not captured or sequenced. Consequently, subtle differences in sample preparation can produce 67 dramatically different results [8,9]. While FACS sorting and single cell methods can produce pure 68 samples of each cell type, the observed cell counts may not accurately represent the cell type 69 abundances in the original sample. Tools like SCDC and MuSiC utilize single cell reference data to 70 perform bulk deconvolution, but require that multi-subject single cell data be available for all the cell 71 types of interest, which is not always the case [10,11].

72

During the past several years, digital means of cell type quantification, often referred to as cell type deconvolution or decomposition, have become a popular complement to FACS sorting and single cell approaches. However, these methods are still developing, and often suffer from limitations. For example, tools MCP-Counter and xCell allow for deconvolution of a set of predefined cell types, but do not support the inclusion of additional cell types or subtypes in a user friendly manner [12,13]. CIBERSORT is slow to run on large datasets, particularly if signature genes are not specified, and provides reference data only for hematopoietic cell types [14].

4

80

81 To overcome some of the limitations of existing cell abundance estimation tools, we present 82 the Gene Expression Deconvolution Interactive Tool (GEDIT). GEDIT utilizes gene expression data to 83 accurately predict cell type composition of tissue samples. We have assembled a library of reference 84 data from 11 distinct sources and use these data to generate thousands of synthetic mixtures. In order 85 to produce optimal results, these synthetic mixtures are used to test and refine the approaches and 86 parameters used by GEDIT. We compare the performance of GEDIT relative to other tools using three 87 sets of mixtures containing known cell type proportions: 12 in vitro mixtures of immune cells 88 sequenced on microarrays, six RNA-seq samples collected from ovarian cancer ascites, and eight 89 RNA-seg samples collected from blood. We also use GEDIT to deconvolute two sets of human tissue 90 samples: 21 skin samples from patients with skin diseases, and 17,382 samples of varied tissues from 91 the GTEx database. Lastly, we apply GEDIT to the Mouse Body Atlas, a collection of samples 92 collected from various mouse tissues and cell types. We find that GEDIT compares favorably to other 93 cell type deconvolution tools and is effective across a broad range of datasets and conditions.

94

95 **Results**

96 **Reference Data**

97

98 Reference data profiling the expression of purified cell types is a requirement for reference-99 based deconvolution. Methods that do not directly require reference data, such as non-negative matrix 100 factorization, still require knowledge of expression profiles or marker genes in order to infer the identity 101 of the predicted components. For this study, we have assembled or downloaded a set of 11 reference 102 matrices, each containing the expression profiles of eight to 29 cell types (Table 1). These data 103 sources span multiple platforms, including bulk RNA-seq, microarray, and single-cell RNA-seq. 104 Complete details on the sources and assembly of these matrices are described in the methods [14-105 24].

	Specie			# of Cell	
Matrix	s	Reference	Platform	Types	Cell Types
			Multi-		
Human Skin Signatures	Human	(Swindell et al. 2013)	Microarray	21	Immune
			Affymetrix		
			U133A/GNF1		
Human Body Atlas	Human	(Su et al. 2004)	Н	13	Immune
			Affymetrix		
Human Primary Cell Atlas	Human	(Mabbott et al. 2013)	U133 Plus 2.0	26	Immune and Stromal
BLUEPRINT*	Human	(Martens and Stunnenberg 2013)	Bulk RNA-Seq	8	Immune
ENCODE*	Human	(ENCODE Project Consortium 2004)	Bulk RNA-Seq	29	Mostly Stroma
			Affymetrix		
LM22	Human	(Newman et al. 2015)	Microarray	22	Immune
			Single Cell		
10x Single Cell Dataset*	Human	(Zheng et al. 2017)	RNA-Seq	9	Immune
			Multi-		
ImmunoStates	Human	(Vallania et. al., 2018)	Microarray	20	Immune
			Single Cell		
Tabula Muris	Mouse	(The Tabula Muris Consortium, 2018)	RNA-seq	12	Immune and Stromal
			Affymetrix		
			Mouse		
			Genome 430		
Mouse Body Atlas	Mouse	(Lattin et al, 2008)	2.0 Array	20	Immune and Stromal
			Affymetrix		Immune with many
lmmGen	Mouse	(Heng et al, 2008)	Gene 1.0 ST	137	subtypes

106

Table 1. Library of Reference Data. Asterisk denotes matrices assembled from source data as part of
 this project. All matrices are compatible with GEDIT and available on the GitHub repository
 (<u>https://github.com/BNadel/GEDIT</u>).

110

111 GEDIT Algorithm

112

113 GEDIT requires as input two matrices of expression values. The first is expression data is

114 collected from the mixtures that will be deconvoluted; each column represents one mixture, and each

115 row corresponds to a gene. The second matrix contains reference data, with each column

116 representing a purified reference profile and each row corresponding to a gene. In a multi-step

117 process, GEDIT utilizes the reference profiles to predict the cell type proportions of each submitted

118 mixture (Figure 1).

119

Figure 1. The GEDIT pipeline. The input matrices are quantile normalized then reduced to matrices containing only signature genes. Next, a row-scaling step serves to control for the dominating effect of highly expressed genes. Lastly, linear regression is performed, and predictions of cell type

123 abundances are reported to the user.

5

- 127
- 128
- 129
- 130
- 121
- 131
- 132 133

Input	Description	Allowed Values	Default Value
		N by M matrix; N is number of	
RefMat	Matrix of purified cell types	genes, M is number of cell types	NA
	Matrix of mixtures to be	N by P matrix; N is number of	
MixMat	deconvoluted	genes, P is number of mixtures	NA
	Method of signature gene	Entropy, MeanRat, MeanDiff,	
SigMeth	selection	ZScore, fsRat, fsDiff	Entropy
NumSigs	Average number of signature genes per cell type	[1, 10,000]	50
	Minimum number of signatures		
MinSigs	per cell type	[1,NumSigs]	=NumSigs
RowScale	Extent of per-row normalization	[0.0,1.0]	0

134

135

136 Table 2. GEDIT inputs include two matrices and four parameter settings. RefMat is an expression 137 matrix documenting the expression profiles of each cell type to be estimated. MixMat is an expression 138 matrix documenting expression values for each sample to be deconvoluted. SigMeth determines the 139 method by which signature genes are selected. NumSigs determines the total number of signature genes, whereas MinSigs sets the minimum number of signature genes for each cell type. RowScale 140 141 refers to the extent to which expression vectors are transformed to lessen the dominating effect of 142 highly expressed genes, with a value of 0.0 representing the most extreme transformation. Default 143 values were determined by evaluating performance on a set of synthetic mixtures (Figure 2). 144

145 **Parameter Tuning**

146 We generated a large number of synthetic mixtures *in silico* to test the efficacy of GEDIT and

147 to assess how accuracy varies as a function of four parameter choices (SigMeth, NumSigs, MinSigs,

148 RowScale, described in Table 2). We produced a total of 10,000 simulated mixtures of known

proportions using data from four reference matrices: BLUEPRINT, The Human Primary Cell Atlas, 10x

150 Single Cell, and Skin Signatures. We then ran GEDIT on these simulated mixtures and evaluated its

151 performance while varying four parameter settings (Figure 2) and other design choices. For this

- 152 reason, these synthetic mixtures were not used to evaluate the performance of GEDIT relative to other
- 153 tools. Instead, separate datasets were used for that purpose, as described in the section

7

154 "Performance Comparison to Other Deconvolution Tools". Based on these results, we selected default

values for each parameter (SigMeth = Entropy, NumSigs = 50, MinSigs = 50, RowScale = 0.0). Full

156 details on the generation of these simulations are described in the supplementary materials.

Figure 2. Effect of GEDIT parameter choices on accuracy of predictions in simulated experiments. 10,000 simulated mixtures were generated, each using one of four reference matrices, with either four, five, six, or ten cell types being simulated. Deconvolution was performed using a separate expression matrix than the one used to generate the mixtures. When not otherwise noted, we use the following parameters: signature selection method = entropy; number of signatures = 50; row scaling = 0.0; and number of fixed genes = number of signatures.

163

164 **Preprocessing and Quantile Normalization**

165 The first step in the GEDIT pipeline is to render the two matrices comparable. This is done by

166 first excluding all genes that are not shared between the two matrices. Genes that have no detected

- 167 expression in any reference cell type are also excluded, as they contain no useful information for
- 168 deconvolution. Both matrices are then quantile normalized, such that each column follows the same
- 169 distribution as every other; this target distribution is the starting distribution of the entire reference
- 170 matrix.
- 171

172 Signature Gene Selection

173 GEDIT next identifies signature genes. Gene expression experiments can simultaneously 174 measure tens of thousands of genes, but many of these genes are uninformative for deconvolution. 175 Specifically, genes with similar expression levels across all cell types are of little use, as observed 176 expression values in the mixtures offer no insight into cell frequencies. Genes that are highly 177 expressed in a subset of cell types are more informative, and we refer to these as signature genes. By 178 using only signature genes, rather than the entire expression matrix, the problem of deconvolution 179 becomes more tractable and less computationally intensive. Moreover, identification of signature 180 genes can be valuable to researchers for other applications (e.g. cell type assignment for scRNA-seq 181 data).

In order to identify the best signature genes in a given reference matrix, GEDIT calculates a
 signature score for each gene. By default, this score is computed using the concept of information

8

184 entropy. Information entropy quantifies the amount of information in a probability distribution, with 185 highly uniform distributions having the highest entropy. The expression vector for each gene (i.e. the 186 set of expression values across all cell types in the reference) is divided by its sum, such that the 187 entries can be interpreted as probabilities. Information entropy is then calculated according to its 188 mathematical definition (see Methods), and genes with the lowest entropy are selected as signature 189 genes. Entropy is minimized when expression is detected only in a single cell type and maximized 190 when expression values are equal across all cell types. Thus, by selecting genes with low entropy, we 191 favor genes that are expressed in a cell type specific manner. By default, 50 signature genes are 192 selected for each cell type in the reference matrix. We chose 50 signature genes, and entropy as our 193 scoring method, because it returned optimal results when run on 10,000 synthetic mixtures (see 194 Figure 2a,b).

We also evaluated the effect of accepting more signature genes for some cell types than others, depending on how many genes have low entropy. In this scheme, on average 50 signature genes are used per cell type. However, a fourth parameter is used, which specifies the minimum number of signature genes per cell type. After these have been selected, remaining signature genes are added based only on lowest entropy, regardless of cell type of maximal expression. We found that this parameter had minimal effect on accuracy, when applied to synthetic mixtures (Figure 2c). Therefore, this option is not used by default, though it can be specified by the user.

202

203 Row Scaling and Linear Regression

One complication in the application of linear regression to gene expression data is the drastically different scale at which some genes are expressed. For example, CD14 and THEMIS (Figure 3) have both been identified as signature genes: CD14 for monocytes and THEMIS for CD4+ T cells. However, CD14 is expressed at much higher levels in most cell types and will have a larger impact on the estimation of cell type composition, relative to THEMIS. In other words, the possible penalty resulting from a poor fit of CD14 is much larger than the penalty from a poor fit of THEMIS.

210

9

Figure 3. The "row scaling" transformation, as implemented by GEDIT. CD14 and THEMIS are two examples of signature genes with drastically different magnitudes of expression. CD14 is a signature gene for monocytes, and THEMIS for CD4+ T cells. The original expression vectors are transformed, such that all values fall between 0.0 and 1.0, equalizing the effect of genes with varying magnitudes of expression.

216

217 In order to equalize the effect of each signature gene on the linear regression, we implement a 218 transformation that we term row scaling. Specifically, the range of all observed values for a particular 219 gene (including reference cell types and samples) is adjusted such that the maximum value is 1.0 and 220 the minimum value is 0.0. As a result, all genes have a comparable influence on the calculation of the 221 linear regression solution, regardless of overall magnitude of expression. This transformation can be 222 modulated by adjusting the row scaling parameter. By default, the value of this parameter is 0.0, and 223 the transformation is applied as described above. Values between 0.0 and 1.0 are also allowed, which 224 reduces the extent of the transformation (see Methods for details). Linear regression is then performed 225 in R using the glmnet package, as described in the methods.

226

227 **Performance Comparison to Other Deconvolution Tools**

In order to assess the performance of GEDIT relative to other tools, we perform an experiment
comparing GEDIT to 4 other deconvolution tools on datasets of known cell-type content
(CIBERSORT, DeconRNASeq, dtangle and xCell; [13,14,25,26]). Non-deconvolution tools like MCPcounter, SAVANT, and the DCQ algorithm are excluded from this study because they do not predict
cell type fractions [12,27,28]. Tools that require single cell data, such as MuSiC and CPM, are also
excluded, as this study is limited to tools that operate on bulk expression data [11,29]. See Table 3 for
a summary of current bulk deconvolution methods.

					Reference Data Provided with Tool				
					Number of				
Tool	Publication	Custom Reference	Approach	Output	Datasets	Cell Types	Species		
				Predicted		Immune	Human,		
GEDIT	Nadel et. al., 2020	Yes	Deconvolution	Fractions	11	and Stromal	Mouse		
				Predicted					
CIBERSORT	Newman et. al., 2015	Yes	Deconvolution	Fractions	1	Immune	Human		
				Predicted		Immune			
xCell	Aran et. al., 2017	No	Marker Genes	Fractions	5	and Stromal	Human		
		Yes, if marker		Predicted					
dtangle	Hunt et. al., 2018	genes specified	Deconvolution	Fractions	0	N/A	N/A		
DeconRNASeq	Gong et. al., 2013	Yes	Deconvolution	Predicted	0	N/A	N/A		

				Fractions			
	Altboum et. al., 2014;						Human,
DCQ/ImmQuant	Frishberg et. al., 2016	Yes	Deconvolution	Scores	3	Immune	Mouse
CIBERSORT							
(absolute mode)	Newman et. al., 2015	Yes	Deconvolution	Scores	1	Immune	Human
		Yes, if marker				Immune	Human,
SaVant	Lopez et. al., 2017	genes specified	Marker Genes	Scores	12	and Stromal	Mouse
						Immune	
MCP-Counter	Becht et. al., 2016	No	Marker Genes	Scores	N/A	and Stroma	Human

235

Table 3. High level characteristics of current cell type estimation tools. Some tools accept custom references, which allows the tool to estimate the abundance of cell types not present in the default reference. Tools listed here take one of two approaches: they either perform deconvolution (most commonly regression) or calculate a score based on intensity of marker gene expression. Depending on the tool, the output can be interpreted as fractions corresponding to the abundance of each cell type, or as scores for each cell type that cannot necessarily be compared in an inter-cellular manner.

To perform this study, we utilize three datasets for which cell type fractions have been estimated using orthogonal methods. Two of these datasets were used in a recent benchmarking study [30]. Both are profiled using RNA-seq, and represent samples collected either from human cancer ascites or human blood [31,32]. In both cases, cell type fractions have been evaluated by FACS sorting. The final dataset was prepared *in vitro* and consists of six cell types that were physically mixed together (in known proportions) to prepare 12 mixtures. These mixtures were then

profiled using an Illumina HT12 BeadChip microarray. Adding to the previous benchmarking study, we

also explore the effect of using four separate reference datasets: The Human Primary Cell Atlas,

LM22, ImmunoStates, and a reference constructed from BLUEPRINT data. For each dataset, all tools

252 (except xCell) were run four times, each time using a different reference matrix.

253 The optimal choice of reference matrix varies greatly depending on the exact combination of tool,

dataset, and cell type. While using LM22 often produces the most accurate results, there are many

255 exceptions. For instance, DeconRNASeq and GEDIT produce their best results for the blood dataset

when using the BLUEPRINT reference. For the ascites data, several tools prefer ImmunoStates as the

257 optimal reference choice. The best choice of reference is highly dependent on the nature of the input

data and on the tool being used. In practice, researchers may wish to perform deconvolution multiple

times--in each case using a separate reference matrix--and compare results for consistency.

260

- 11
- 261 Compared to the other tools, GEDIT produces robust and consistently accurate results 262 (Figures 4.5). For many tools, the quality of predictions varies greatly depending on the cell type, 263 dataset, or choice of reference matrix. When results are averaged across the four possible reference 264 choices, GEDIT produces the minimum error and maximum correlation for all three datasets. This 265 result suggests that GEDIT is a strong choice when researchers are using novel references matrices 266 that have not been curated or tested. 267 268 Figure 4. Performance of five deconvolution tools when applied to a set of 26 physical samples from 269 three sources. Actual cell type fractions are either known due to controlled cell mixing (Cell Mix) or 270 estimated by FACS sorting (Ascites and Blood). In each instance, we calculate the correlation 271 between actual cell type fractions and those predicted by deconvolution; deeper blues represent 272 higher correlations. We test four different reference datasets for each tool, and averaged correlations 273 across these 5 cases are shown in boxes. We calculate correlations for each cell type (right 5 274 columns), for each of the 3 mixtures (middle 3 columns), and for all predictions regardless of cell type 275 or data source. 276 277 Figure 5. Distribution of errors, calculated as the difference between predicted and actual cell type 278 fractions. Each point on the graph represents the percentage of predictions (y-axis) that are accurate 279 within a particular error range (x-axis). 280 281 We also perform 2 additional comparisons between GEDIT and other deconvolution tools. 282 Firstly, we create 100 simulated mixtures of pancreatic cells (alpha, beta, gamma, delta) using single 283 cell data from a recent single cell experiment (details in supplementary materials). We evaluate the 284 accuracy of each tool when used to predict the cell type content of these synthetic mixtures, and 285 GEDIT provides the lowest overall error (Supplementary Figure 3). 286 Lastly, we perform an evaluation of runtime required for each tool. We randomly select batches 287 of 100, 200, 500, 1000, and 2000 samples from the GTEx database, and measure CPU time required 288 to deconvolute these batches for each tool. The runtime of GEDIT, dtangle, and DeconRNASeq 289 scales well with growing input size, taking at most 20 minutes on average (Supplementary Figure 4). 290 Skin Expression Data 291 292 We further validate GEDIT by using it to deconvolute a set of skin biopsies from humans with a 293 variety of skin diseases [13]. The exact cell type composition of these samples is unknown, but we

12

294	have reasonable expectations based on skin and disease biology. For example, macrophages are
295	known to be abundant in granulomas of leprosy legions, and Steven-Johnson Syndrome produces
296	blisters that fill with large numbers of monocytes [33,34]. We find that, in all cases, predictions made
297	by GEDIT conform well with these biological expectations. Keratinocytes are highly predicted in most
298	cases, as one would expect with skin samples (Figure 6). Deviations from this pattern correspond with
299	disease biology. Monocytes are highly predicted in Stevens-Johnson syndrome, as are macrophages
300	in the three leprosy samples, and T cells in the Mycosis Fungoides (T cell lymphoma) sample.
301 302 303 304 305 306 307 308 309 310	Figure 6. GEDIT predictions for 21 samples of various skin diseases. GEDIT correctly identifies keratinocytes and subcutaneous adipose as the most common cell. Deviations from this pattern correspond to disease biology. SJS represents blister fluid from Steven Johnson Syndrome, and is predominantly immune cells. LL and RR represent two forms of leprosy, which result in large numbers of macrophages. MF is a T Cell Lymphoma.
311	Application of GEDIT to Mouse Data
312	GEDIT can be used to decompose data from any organism for which reference data is
313	available. Here, we demonstrate the efficacy of GEDIT when applied to the Mouse Body Atlas, a
314	collection of tissue and cell type samples collected from mice [23]. As reference data, we assembled a
315	matrix of 12 cell types using single cell data from the Tabula Muris [20]. GEDIT correctly infers the
316	identity of purified cell types, including six samples that consist of either pure NK cells, B cells, T cells,
317	or granulocytes (Figure 7). An entry for macrophages is not available in the reference used, but most
318	macrophage complex are identified as manageted, which is the most similar call type present in the
	macrophage samples are identified as monocytes, which is the most similar cell type present in the
319	reference matrix. For more complex tissues, GEDIT predicts cell type fractions that correspond to the
319 320	
	reference matrix. For more complex tissues, GEDIT predicts cell type fractions that correspond to the
320	reference matrix. For more complex tissues, GEDIT predicts cell type fractions that correspond to the biology of the samples. Hepatocytes are predicted to be highly prevalent in the liver sample (84%) and
320 321	reference matrix. For more complex tissues, GEDIT predicts cell type fractions that correspond to the biology of the samples. Hepatocytes are predicted to be highly prevalent in the liver sample (84%) and are not predicted in any other sample (less than 5% in all cases). Similar patterns hold for
320 321 322	reference matrix. For more complex tissues, GEDIT predicts cell type fractions that correspond to the biology of the samples. Hepatocytes are predicted to be highly prevalent in the liver sample (84%) and are not predicted in any other sample (less than 5% in all cases). Similar patterns hold for keratinocytes in the epidermis, epithelial cells in two intestinal samples and cardiac muscle cells in

325

- 326
- 327 Figure 7. GEDIT predictions on 30 samples collected from various mouse tissues and cell types
- 328 (mouse body atlas [23]). Predictions largely conform with tissue and cell biology.
- 329

Deconvolution of GTEx Database

- To assess the use of GEDIT across very large datasets, we applied the tool to 17,382 GTEx RNA-seq samples collected from various tissues. However, no single reference contained all relevant cell types. For example, none of the available references contain both myocytes and adipocytes (Supplementary Figure 1). Therefore, we predicted proportions three times using three separate references (BlueCode, Human Primary Cell Atlas, Skin Signatures). We then combined these outputs by taking their median value. This allowed us to produce predictions spanning a larger number of cell types than are present in any one reference matrix (Figure 8).
- 338

Figure 8. GEDIT cell type predictions when applied to 17,382 samples from the GTEx database. Here,
 predictions have been averaged for each tissue of origin.

341

These predictions largely conform to biological expectations. For example, immune cells are predicted to have high abundance in blood and spleen, adipocytes in adipose tissue, Shwann cells in nerve and heart, and keratinocytes in skin. Each of these patterns matches expectations of which cell types should be present in these tissues. Neither cardiac myocytes nor smooth muscle are highly abundant in GTEx muscle samples. This is likely because the GTEx samples are collected from skeletal muscle, which is known to have an expression profile that is distinct from that of cardiac and smooth muscle.

349

350 **GEDIT Availability**

351 GEDIT can be run online at <u>http://webtools.mcdb.ucla.edu/.</u> Source code, associated data, and

352	relevant files are available on GitHub at https://github.com/BNadel/GEDIT. We provide access to the
353	tool, a set of varied reference data, and two sample mixture matrices. The website automatically
354	produces a heatmap of predicted proportions for the user, as well as a .tsv file. The user also has
355	access to the parameter choices of GEDIT (signature gene selection method, number of signature
356	genes, row scaling).

- 357
- 358 Methods
- 359 **GEDIT Algorithm**

360 Signature Gene Selection

361

362 During signature gene selection, we automatically exclude genes with zero detected 363 expression in half or more of cell types. Observed expression values of exactly zero are often the 364 result of either technical artefacts or resolution issues. Using such genes as signatures can result in 365 inaccurate and highly unstable results, particularly when working with scRNA-seg derived data. As an 366 additional safeguard, we treat all remaining expression values of zero as the lowest observed non-367 zero value in the matrix. Implementing this change has minimal effect on most genes but prevents 368 genes with resolution issues from achieving artificially high scores. We consider this transformation 369 valid, since values of zero generally do not represent zero expression, but rather an expression level 370 below the detection limit of the technology used.

371

For any given gene, a scoring method takes as input the vector of the expression values across all reference cell types, and outputs a score. A gene is considered a potential signature gene in cell type X if it is expressed more highly in X than any other cell type. For each cell type, we keep only the N genes with the highest scores, where N is the NumSigs parameter.

376

377 Information entropy (H) is calculated using the following formula:

15

378 $H = -\sum_{1}^{i} [p_{i} * (p_{i})]$ (1) 379 380 where p_i is the probability of the *i*th observation. To apply this to expression values, we convert 381 382 the vector of expression values into a vector of probabilities by dividing by its sum. In an equal mixture of each cell type, the *i*th probability can be interpreted as the fraction of transcripts originating from the 383 *i*th cell type. 384 385 386 387 **Row Scaling** 388 389 During this step, we apply a transformation on the expression values for each gene. Each gene 390 has measured expression in N purified cell types and M samples. Each of these values, X_{old}, is 391 transformed according to the following formula: 392 $X_{new} = (X_{old} - Min)/(Max - Min) * Max^{p}$ (2) 393 394 395 Where Min is the minimum of all M + N original values, Max is the maximum of those values, 396 and p is a tunable parameter with natural range p \in [0.0,1.0]. This procedure produces values 397 between the range of 0.0 and Max^p. 398 Linear Regression: 399 400 Non-negative linear regression was performed using the glmnet package in R. The glmnet 401 function is used with lower.limits=0, alpha=0, lambda=0, intercept=FALSE. These settings perform a 402 linear regression where all coefficients are non-negative, and with no regularization and no intercept 403 term.

404

405 **Reference Data**

406 **BLUEPRINT Reference Dataset**

407 35 gene counts files were downloaded from the BLUEPRINT database, all collected from 408 venous blood [18]. This included entries for CD14-positive, CD16-negative classical monocytes (5 409 samples), CD38-negative naive B cells (1), CD4-positive, alpha-beta T cell (8), central memory CD4positive, alpha-beta T cell (2), cytotoxic CD56-dim natural killer cell (2), macrophage (4). mature 410 411 neutrophil (10), and memory B Cell (1). When two or more transcripts appeared for a single gene, the 412 transcript with the highest average expression was selected and others were excluded. Genes with no 413 detected expression in any sample were also excluded, and then each sample was quantile 414 normalized. Samples generally clustered by cell type, but we excluded one CD4-positive alpha-beta T 415 cell. Replicates for each cell type were then collapsed into a single entry by taking the median value 416 for each gene.

417

418 **ENCODE Reference Dataset**

106 transcript quantification files were downloaded from the ENCODE database [19]. These
included all RNA-seq experiments collected from adult primary cells, excluding four with warnings.
Warnings indicated that three samples suffered from low replicate concordance and one sample from
low read depth, and these samples were excluded. All samples were processed by the Gingeras Lab
at Cold Spring Harbor and mapped to GRCH38.

The samples were quantile normalized and clustered. In cases where multiple transcripts were measured for a single gene, the expression of that gene was calculated as the sum of all transcripts. At this time, 18 additional samples were excluded as they did not cluster with their replicates. Based on sample descriptions and data clustering, we found that the remaining 88 samples represented 28 unique cell types. We produced an expression profile for each cell type by merging all samples of that cell type via median average. For example, a cluster of 19 samples were labelled as endothelial cells

17

430	(collected from	various body	locations) a	and were n	nerged into a	a single entr	y termed	canonical

431 endothelial cells. This dataset spans a wide range of stromal cell types (e.g. smooth muscle,

fibroblast, epithelial), but contains only a single entry for blood cells, which are labelled mononuclearcells.

434 We also combined the ENCODE and BLUEPRINT reference matrices into a single reference

435 matrix, which we call BlueCode. We combined, then quantile normalized, the columns of both

436 matrices. Possible batch effects in this combined matrix have not been fully evaluated.

437

438 **10x Reference Dataset**

We obtained single cell expression data for nine varieties of immune cells from the 10x website [20]. This included at least 2446 cells for each cell type, and at least 7566 cells for all cells other than CD14 monocytes. For each cell type, expression values for all cells were mean averaged to form an expression profile.

443

444 **Tabula Muris Reference Dataset**

We downloaded from the Tabula Muris single cell data for 12 clusters of mouse cell types. For each cluster, we averaged all cells of that cluster to produce a reference profile for the corresponding cell type.

448 **Other Reference Datasets**

449 Other datasets used in this project were obtained from their corresponding publications or

450 GEO repositories. This includes a reference matrix of human skin signatures, the Human Body Atlas,

451 the Human Primary Cell Atlas, LM22, ImmunoStates, the Mouse Body Atlas, and ImmGen [14–

452 17,21,23,24].

453

454 Skin Diseases Data

455 We obtained expression data from 21 skin biopsies, collected from human patients with a 456 variety of skin diseases. These data originally came from a wide range of sources and platforms, and

- 457 were compiled into a single dataset by previous work [35].
- 458
- 459 **GTEx Data**
- 460 GTExX data for 17,382 samples were obtained from the GTExX database
- 461 (<u>https://gtexportal.org/</u>). We ran GEDIT on all samples three times, each time using a different
- 462 reference matrix (BlueCode, the Human Primary Cell Atlas, and Skin Signatures). For each cell type,
- 463 we calculated our initial estimate as the median estimate across the three sets of predictions (or fewer,
- if that cell type is missing from one to two of the reference matrices). Lastly, for each sample we
- 465 divided the vector of predictions by its sum, such that the final predictions sum to 100%.
- 466

467 Multi-Tool Performance Evaluation

468 In Vitro Immune Cell Mixture

469 Combinations of six immune cells (Neutrophils, Monocytes, Natural Killer Cells, B cells, and

470 CD4 and CD8 T Cells) were mixed together and sequenced using an affymetrix array. Whole blood

471 from healthy human donors was supplied with informed consent through a sample sharing agreement

472 with the UCLA/CFAR Virology Core Lab (grant number 5P30 Al028697). CD4+ T cells, CD8+ T cells,

- 473 B cells, and NK cells were isolated using Stem Cell Technologies (Vancouver, BC, Canada)
- 474 RosetteSep negative selection. Neutrophils were positively selected through the EasySep approach,
- 475 according to the manufacturer's specifications. Cells were then counted by hemocytometer and added
- 476 at defined percentages to a total cell count of two million cells to create six different mixtures.
- 477 Subsequently cells were processed for RNA isolation by AllPrep DNA/RNA. Illumina HT12 BeadChip
- 478 microarray was performed by the UCLA Neuroscience Genomics Core. Data was normalized by
- 479 quantile normalization through R 'normalize.quantiles' function (R Core Team, 2013).
- 480

481 **RNA-seq Mixtures Used for Tool Evaluation**

482 We also obtained two datasets used in a recent benchmarking study [30]. The first dataset is

19

483 composed of three RNA-seg samples, each with two technical replicates that represent biopsies of 484 ovarian cancer ascites [32]. The second dataset is composed of RNA-seq collected from the blood of 485 healthy individuals, some of whom recently received an influenza vaccine [31]. These data were 486 downloaded from the GitHub for the benchmarking paper, which also contained FACS estimates for 487 six cell types for the ascites data (B cells, dendritic cells, NK cells, T cells, macrophages, neutrophils) 488 and five cell types for the blood data (B cells, dendritic cells, T cells, monocytes, natural killer cells). 489 However, since dendritic cells were never present at more than 3.5% abundance, we did not evaluate 490 performance for this cell type.

- 491
- 492 **Tools**

We installed and ran GEDIT, CIBERSORT, DeconRNASeq and dtangle on the hoffman2 computational cluster at UCLA. xCell was run using the online interface at <u>https://xcell.ucsf.edu/</u>. The default choice for genes signatures (xCell =64) was used. The RNA-seq option was selected for the 2 RNA-seq datasets (blood and ascites), but not for the *in vitro* dataset, which was sequenced on microarray.

498 xCell produces 67 output scores, seven of which were used in this study. These were the 499 entries labelled "B-Cells", "Macrophages", "Monocytes", "NK cells", "Neutrophils", "CD4+ T cells" and 500 "CD8+ T Cells". As suggested by the xCell authors, the outputs for CD4 and CD8 T cell subtypes were 501 summed to produce a final output for total T cells.

502

503 **Reference Data**

We evaluated the performance of the four reference-based tools (GEDIT, CIBERSORT, DeconRNASeq and dtangle) using each of four choices of reference matrix (LM22, ImmunoStates, BLUEPRINT, and the Human Primary Cell Atlas). The BLUEPRINT and Human Primary Cell Atlas reference matrices differ from ImmunoStates and LM22 in that they contain tens of thousands of genes, many of which should not be considered signature genes. This contrasts to ImmunoStates and LM22; each reference matrix contains fewer than 600 genes, which have been specifically identified

20

as signature genes by previous work [14,21]. We include both forms of reference matrices in order to
evaluate the input requirements of the tools studied.

512 Depending on the choice of reference matrix, reference-based tools often produce multiple 513 outputs for some cell types, each representing a cell sub-type. This includes B cells (naïve and 514 memory), Monocytes (CD14 and CD16), NK cells (resting and active) and T cells (many subtypes 515 including varieties of CD4 and CD8). In each case, the outputs for each sub-type were summed in 516 order to produce a total score for each greater cell type.

517

518 **Discussion:**

519 GEDIT is an expression-based cell type quantification tool that offers unprecedented flexibility 520 and accuracy in a wide variety of contexts. Using both simulated and experimental data, we 521 demonstrate that GEDIT produces high-quality predictions for multiple platforms, species, and a 522 diverse range of cell types, outperforming other tools in many cases. We include in the software 523 package a comprehensive library of reference data, which facilitates application of GEDIT to a wide 524 range of tissue types in both human and mouse. GEDIT can also accept reference data supplied by 525 the user, which can be derived from bulk RNA-seq, scRNA-seq, or microarray experiments. GEDIT 526 represents a competitive addition to the suite of existing tissue decomposition tools while maintaining 527 flexibility and performance robustness.

528 As part of this project, we perform a study in which we compare the performance of several 529 deconvolution tools using multiple metrics. Unlike previous evaluation studies, we explore the effect of 530 reference choice by running tools multiple times with reference data from different sources. Choice of 531 optimal reference has a large impact on the accuracy of many tools, but GEDIT provides robust 532 performance and accurate estimates for many possible reference choices. While all efforts were taken 533 to perform this comparison in an unbiased manner, the authors note that development of the tool was 534 still underway when the first comparisons were made. All code and inputs used to reproduce this study 535 are included in the github (https://github.com/BNadel/GEDIT), with the exception of CIBERSORT

536 code, which is limited by copyright.

537 The high performance of GEDIT is due to two key innovations. Firstly, signature gene selection 538 by information entropy serves to select genes that are the most informative for deconvolution. 539 Secondly, the row scaling step, which aims to equally weight all signature genes into the final 540 estimate, even those with comparatively low expression. In addition, the flexibility of GEDIT and the 541 diverse set of reference matrices we provide allows GEDIT to be easily applied in a wide range of 542 circumstances.

The output of GEDIT represents the fraction of mRNA originating from each cell type. This is an effective measure of the transcriptional contribution of each cell type in a mixture. However, in cases where some cell types consistently produce more or less mRNA per cell, this measure may not represent cell counts. Data capturing the average mRNA content per cell is becoming more widely available in the form of single cell experiments and could in principle be used to convert our fractions into cell counts.

When extensively applied to several large public datasets, GEDIT produces predicted cell type fractions that conform with biological expectations. When used to decompose skin biopsies, keratinocytes are found to be the most abundant cell type. Variations in the abundance of other cell types conform to expected immune responses across diseases. Similarly, cell type predictions of GTEx samples are concordant with our expectations of the dominant cell types across tissues. Schwann cells, keratinocytes, adipose cells, and immune cells are found to be most abundant in nerve, skin, adipose tissue, and blood, respectively.

Single cell RNA-seq is an emerging approach to study the composition of cell types within a sample. Due to biases associated with the capture of different cell types, these methods are not always capable of accurately quantifying cell type populations [8]. However, the pure reference profiles produced by existing methods can be used by GEDIT to generate accurate estimates of cell type populations. Thus, GEDIT circumvents some of the biases associated with the preparation of samples for both scRNA-seq and FACS. GEDIT is freely available, and therefore an extremely economical option for researchers, particularly those who profile expression data for other purposes.

- 22
- 563 GEDIT produces accurate results when tested on mixtures of human immune cells. Compared
- to other tools, GEDIT produces the lowest error in majority of scenarios in the studied mixtures.
- 565 GEDIT provides increased flexibility over previously developed tools, as we provide a set of reference
- 566 matrices for varied cell types for both mouse and human datasets.
- 567 GEDIT provides unique advantages to researchers, especially in terms of cell type, species
- and platform flexibility, and constitutes a useful addition to the existing set of tools for tissue
- 569 decomposition. Our efficient decomposition methodology has been extensively optimized and we find
- 570 that it performs robustly across a broad range of tissues in both mouse and human datasets. Our
- 571 future work will extend reference matrices to facilitate application of GEDIT on varied bulk gene
- 572 expression datasets.

573 Availability of Source Code and Requirements

- 574 Project name: GEDIT
- 575 Project Home Page: <u>https://github.com/BNadel/GEDIT</u>
- 576 Programming Languages: Python 2.0, R
- 577 Other requirements: numpy, glmnet
- 578 Ø Operating Systems: Linux
- 579 🛛 License: MIT
- 580 Availability of Data and Materials
- 581 All data used in this paper are freely available on GitHub
- 582 (https://github.com/purebrawn/GEDIT), as well as their original sources. Code for DeconRNASeq was
- 583 obtained as an R package from the CRAN repository. Code for CIBERSORT was obtained by
- 584 requesting it via the web portal (<u>https://cibersort.stanford.edu/download.php</u>), and code for dtangle
- 585 from the project's GitHub page (<u>https://github.com/gjhunt/dtangle</u>).
- 586 Reference data is also available from their original sources. Most datasets can be found on
- 587 project website pages or from public databases. These include BLUEPRINT (<u>http://www.blueprint-</u>
- 588 <u>epigenome.eu/</u>), ENCODE (<u>https://www.encodeproject.org</u>), the Human Primary Cell Atlas

23

- 589 (<u>http://biogps.org/dataset/BDS_00013/primary-cell-atlas/</u>), LM22
- 590 (http://cibersort.stanford.edu/%20or%20GEO:GSE65136), 10x Genomics
- 591 (https://support.10xgenomics.com/single-cell-gene-expression/datasets), Tabula Muris (https://tabula-
- 592 <u>muris.ds.czbiohub.org/</u>), the Mouse Body Atlas (GEO:GSE10246), and ImmGen
- 593 (http://www.immgen.org/Databrowser19/DatabrowserPage.html). Some reference matrices were
- obtained as supplementary files from the publications listed in Table 1.
- 595 Expression values for the blood and ascites RNA-seq datasets were obtained from the GitHub
- 596 repository <u>https://github.com/grst/immune_deconvolution_benchmark</u>, and are also available at at
- 597 <u>https://figshare.com/s/711d3fb2bd3288c8483a</u> and GEO: GSE64655). The *in vitro* mixture of immune
- 598 cells was prepared by our lab, and available on our GitHub page.
- 599

600 Acknowledgments

- 601 We acknowledge the Biomedical Big Data Grant (5T32LM012424-03) for supporting Brian
- Nadel during the course of this research. We also acknowledge the Bruins-in-Genomics Summer
- 603 Undergraduate Research Program for supporting Hannah Waddel and Misha Khan during the summer
- of 2017, when they contributed to this work. We also thank Lana Martin for her help with editing and
- 605 proofreading the manuscript.
- 606 References

607 **Competing Interests**

- The authors declare that they have no competing interests.
- Bolen CR, Uduman M, Kleinstein SH. Cell subset prediction for blood genomic studies. BMC
 Bioinformatics. 2011;12: 258.
- 611 2. Gentles AJ, Newman AM, Liu CL, Bratman SV, Feng W, Kim D, et al. The prognostic landscape
 612 of genes and infiltrating immune cells across human cancers. Nat Med. 2015;21: 938–945.
- Li B, Severson E, Pignon J-C, Zhao H, Li T, Novak J, et al. Comprehensive analyses of tumor
 immunity: implications for cancer immunotherapy. Genome Biol. 2016;17: 174.
- Fridman WH, Pagès F, Sautès-Fridman C, Galon J. The immune contexture in human tumours:
 impact on clinical outcome. Nat Rev Cancer. 2012;12: 298–306.

- 24
- 5. Şenbabaoğlu Y, Gejman RS, Winer AG, Liu M, Van Allen EM, de Velasco G, et al. Tumor
 immune microenvironment characterization in clear cell renal cell carcinoma identifies prognostic
 and immunotherapeutically relevant messenger RNA signatures. Genome Biol. 2016;17: 231.
- 6. Gierahn TM, Wadsworth MH, Hughes TK, Bryson BD, Butler A, Satija R, et al. Seq-Well: portable,
 low-cost RNA sequencing of single cells at high throughput. Nature Methods. 2017. pp. 395–398.
 doi:10.1038/nmeth.4179
- Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, et al. Highly Parallel Genome wide Expression Profiling of Individual Cells Using Nanoliter Droplets. Cell. 2015;161: 1202–1214.
- 8. Hines WC, Su Y, Kuhn I, Polyak K, Bissell MJ. Sorting out the FACS: a devil in the details. Cell
 Rep. 2014;6: 779–781.
- 627 9. Denisenko E, Guo BB, Jones M, Hou R, de Kock L, Lassmann T, et al. Systematic assessment of
 628 tissue dissociation and storage biases in single-cell and single-nucleus RNA-seq workflows.
 629 Genome Biol. 2020;21: 130.
- 10. Dong M, Thennavan A, Urrutia E, Li Y, Perou CM, Zou F, et al. SCDC: bulk gene expression
 deconvolution by multiple single-cell RNA sequencing references. Brief Bioinform. 2020.
 doi:10.1093/bib/bbz166
- Mang X, Park J, Susztak K, Zhang NR, Li M. Bulk tissue cell type deconvolution with multi-subject
 single-cell expression reference. Nat Commun. 2019;10: 380.
- Becht E, Giraldo NA, Lacroix L, Buttard B, Elarouci N, Petitprez F, et al. Estimating the population
 abundance of tissue-infiltrating immune and stromal cell populations using gene expression.
 Genome Biol. 2016;17: 218.
- Aran D, Hu Z, Butte AJ. xCell: digitally portraying the tissue cellular heterogeneity landscape.
 Genome Biol. 2017;18: 220.
- 14. Newman AM, Liu CL, Green MR, Gentles AJ, Feng W, Xu Y, et al. Robust enumeration of cell
 subsets from tissue expression profiles. Nat Methods. 2015;12: 453–457.
- 5. Swindell WR, Johnston A, Voorhees JJ, Elder JT, Gudjonsson JE. Dissecting the psoriasis
 transcriptome: inflammatory- and cytokine-driven gene expression in lesions from 163 patients.
 BMC Genomics. 2013;14: 527.
- Su AI, Wiltshire T, Batalov S, Lapp H, Ching KA, Block D, et al. A gene atlas of the mouse and
 human protein-encoding transcriptomes. Proc Natl Acad Sci U S A. 2004;101: 6062–6067.
- Mabbott NA, Baillie JK, Brown H, Freeman TC, Hume DA. An expression atlas of human primary
 cells: inference of gene function from coexpression networks. BMC Genomics. 2013;14: 632.
- 649 18. Martens JHA, Stunnenberg HG. BLUEPRINT: mapping human blood cell epigenomes.
 650 Haematologica. 2013;98: 1487–1489.
- 651 19. ENCODE Project Consortium. The ENCODE (ENCyclopedia Of DNA Elements) Project. Science.
 652 2004;306: 636–640.
- 20. Zheng GXY, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al. Massively parallel digital
 transcriptional profiling of single cells. Nat Commun. 2017;8: 14049.
- 21. Vallania F, Tam A, Lofgren S, Schaffert S, Azad TD, Bongen E, et al. Leveraging heterogeneity

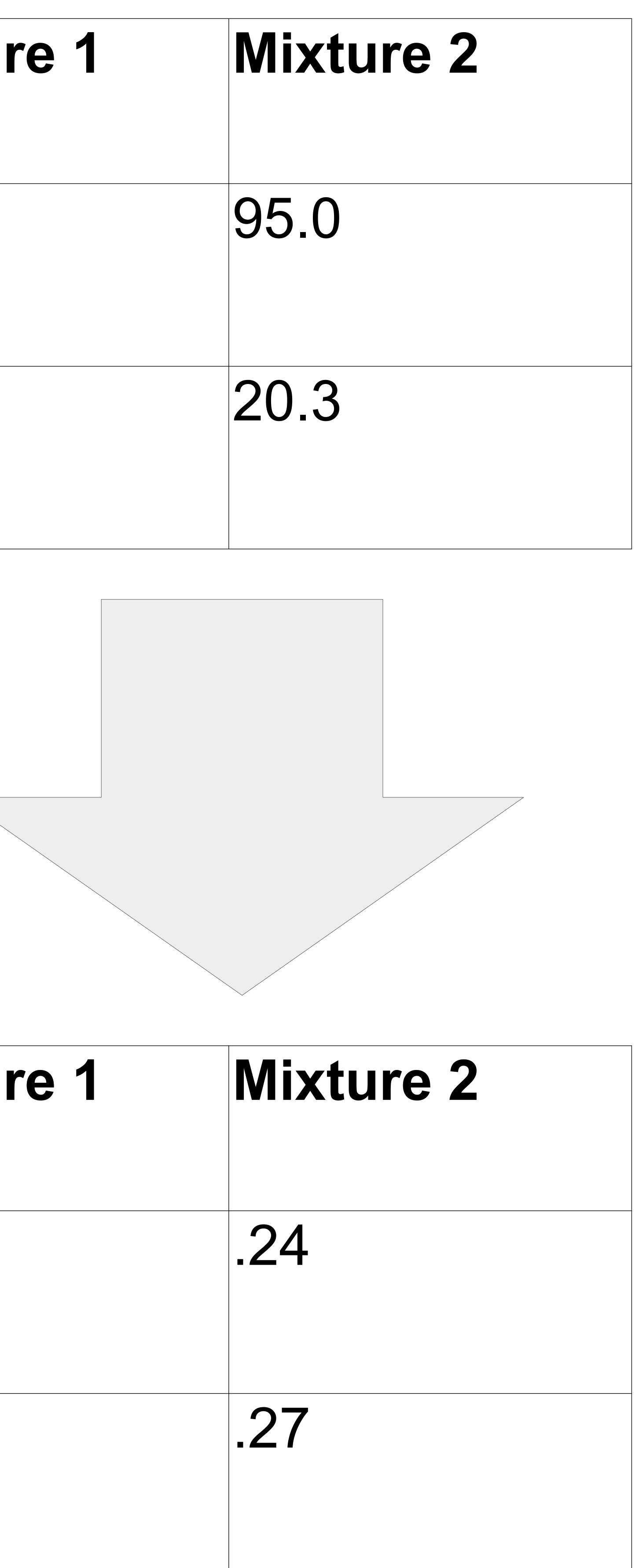
- 25
- 656across multiple datasets increases cell-mixture deconvolution accuracy and reduces biological657and technical biases. Nat Commun. 2018;9: 4735.
- Consortium TTM, The Tabula Muris Consortium, Coordination O, Coordination L, Organ collection
 and processing, Library preparation and sequencing, et al. Single-cell transcriptomics of 20
 mouse organs creates a Tabula Muris. Nature. 2018. pp. 367–372. doi:10.1038/s41586-0180590-4
- 23. Lattin JE, Schroder K, Su AI, Walker JR, Zhang J, Wiltshire T, et al. Expression analysis of G
 Protein-Coupled Receptors in mouse macrophages. Immunome Res. 2008;4: 5.
- 4. Heng TSP, Painter MW, Immunological Genome Project Consortium. The Immunological
 Genome Project: networks of gene expression in immune cells. Nat Immunol. 2008;9: 1091–
 1094.
- 667 25. Gong T, Szustakowski JD. DeconRNASeq: a statistical framework for deconvolution of
 668 heterogeneous tissue samples based on mRNA-Seq data. Bioinformatics. 2013;29: 1083–1085.
- 469 26. Hunt GJ, Freytag S, Bahlo M, Gagnon-Bartsch JA. dtangle: accurate and robust cell type
 470 deconvolution. Bioinformatics. 2018 [cited 17 Jan 2019]. doi:10.1093/bioinformatics/bty926
- 27. Lopez D, Montoya D, Ambrose M, Lam L, Briscoe L, Adams C, et al. SaVanT: a web-based tool
 for the sample-level visualization of molecular signatures in gene expression profiles. BMC
 Genomics. 2017;18: 824.
- Altboum Z, Steuerman Y, David E, Barnett-Itzhaki Z, Valadarsky L, Keren-Shaul H, et al. Digital
 cell quantification identifies global immune cell dynamics during influenza infection. Mol Syst Biol.
 2014;10: 720.
- Frishberg A, Peshes-Yaloz N, Cohn O, Rosentul D, Steuerman Y, Valadarsky L, et al. Cell
 composition analysis of bulk genomics using single-cell data. Nat Methods. 2019;16: 327–332.
- Sturm G, Finotello F, Petitprez F, Zhang JD, Baumbach J, Fridman WH, et al. Comprehensive
 evaluation of transcriptome-based cell-type quantification methods for immuno-oncology.
 Bioinformatics. 2019;35: i436–i445.
- 31. Hoek KL, Samir P, Howard LM, Niu X, Prasad N, Galassie A, et al. A Cell-Based Systems Biology
 Assessment of Human Blood to Monitor Immune Responses after Influenza Vaccination. PLOS
 ONE. 2015. p. e0118528. doi:10.1371/journal.pone.0118528
- 32. Schelker M, Feau S, Du J, Ranu N, Klipp E, MacBeath G, et al. Estimation of immune cell content
 in tumour tissue using single-cell RNA-seq data. Nat Commun. 2017;8: 2032.
- de Sousa JR, Lucena Neto FD, Sotto MN, Quaresma JAS. Immunohistochemical characterization
 of the M4 macrophage population in leprosy skin lesions. BMC Infect Dis. 2018;18: 576.
- 4. Lin C-C, Chen C-B, Wang C-W, Hung S-I, Chung W-H. Stevens-Johnson syndrome and toxic
 epidermal necrolysis: risk factors, causality assessment and potential prevention strategies.
 Expert Rev Clin Immunol. 2020;16: 373–387.
- 35. Inkeles MS, Scumpia PO, Swindell WR, Lopez D, Teles RMB, Graeber TG, et al. Comparison of
 molecular signatures from multiple skin diseases identifies mechanisms of immunopathogenesis.
 J Invest Dermatol. 2015;135: 151–159.
- 695

Cell Type	Monocytes	Neutrophils	B Cells	NK Cells	CD4+ T cells	Macrophages	Mixture
CD14	338.4	163.9	18.9	16.9	19.2	105.9	22.3
THEMIS	9.7	11.6	8.4	13.2	52.0	8.7	50.3

Cell Type	Monocytes	Neutrophils	BCells	NK Cells	CD4+ T cells	Macrophages	Mixture
CD14	1.0	.46	.01	0.0	.01	.28	.02
THEMIS	.03	.07	0.0	.11	1.0	.01	.96
728493; this version posted September 14, 2020. The copyright holder for this preprint (which was under, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.							



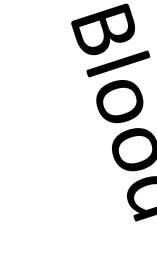
Row Scaling Transformation



						Ρ	earson	Corre	lations	
Reference Used:	.2 0	.25 .5	.75	.8 .9	•	1				
BLUEPRINT	0.64	-0.15	0.76	0.74		0.38	0.33	0.6	-0.23	0.63
HPCA	0.83	0.65	0.9	0.76	(0.01	0.79	0.64	0.75	0.95
ImmunoStates	0.89	0.72	0.95	0.78	(0.79	0.84	0.92	0.83	0.97
LM22	0.86	0.68	0.96	0.61	(0.89	0.78	0.8	0.91	0.92
Average	0.8	0.47	0.89	0.72	(0.33	0.69	0.74	0.56	0.87
BLUEPRINT	0.46	-0.19	0.44	0.77		-0.1	0.05	-0.06	-0.23	0.65
HPCA	0.16	-0.07	0.1	0.54		-0.3	0	-0.39	0.46	0.69
ImmunoStates	0.5	0.3	0.87	0	_	0.17	0.33	0.45	0.43	0.9
LM22	0.65	0.81	0.91	0.25	(0.54	0.38	0.74	0.65	0.88
Average	0.44	0.21	0.58	0.39	_	0.01	0.19	0.18	0.32	0.78
BLUEPRINT	0.57	0.18	0.54	0.98		0.1	0.33	0.57	0.41	0.84
HPCA	0.58	0.38	0.51	0.94	(0.12	0.18	0.55	0.43	0.45
ImmunoStates	0.35	0.67	0.57	0.41	(0.31	0.25	-0.19	0.39	0.48
LM22	0.6	0.68	0.55	0.96		0.5	0.27	-0.13	0.41	0.62
Average	0.53	0.48	0.54	0.82	(0.26	0.26	0.19	0.41	0.6
BLUEPRINT	0.84	0.13	0.94	0.99	(0.57	0.78	0.93	0.91	0.89
HPCA	0.93	0.82	0.97	0.87	(0.63	0.81	0.89	0.91	0.98
ImmunoStates	0.9	0.81	0.96	0.83	(0.22	0.74	0.95	0.82	0.97
LM22	0.9	0.71	0.97	0.8	(0.75	0.7	0.89	0.94	0.96
Average	0.89	0.62	0.96	0.87	(0.55	0.76	0.92	0.89	0.95
xCell Signatures	0.64	0.82	0.88	0.06	(0.65	0.55	0.76	0.71	0.7
bioRxiv preprint doi: https://doi.org/10.1101/728493; this version posted September 14, 2020. The copyright holder for this preprint (which was										

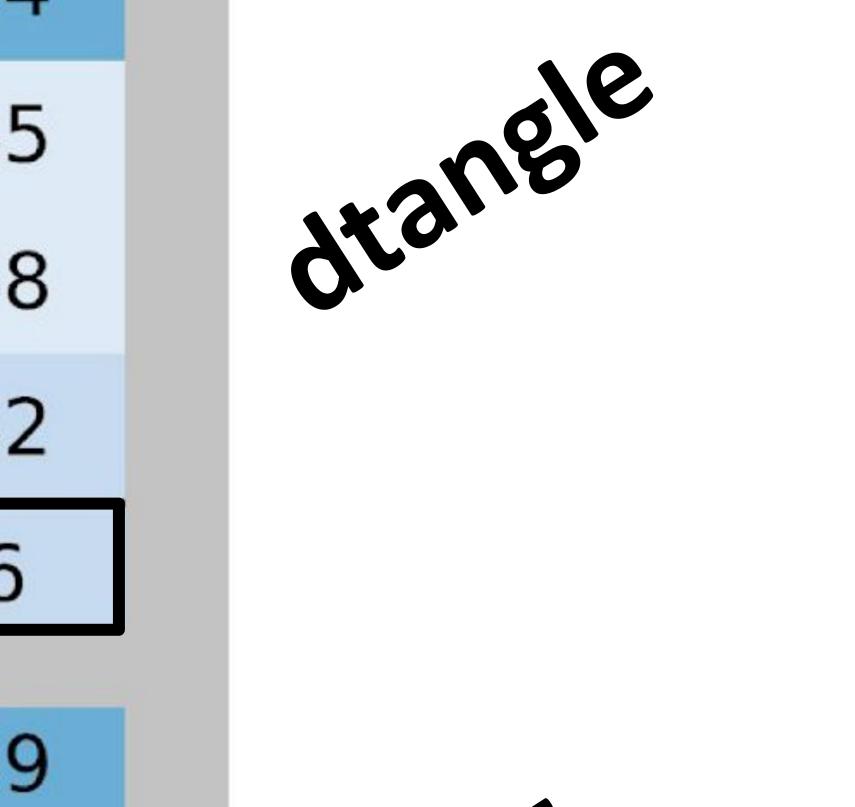




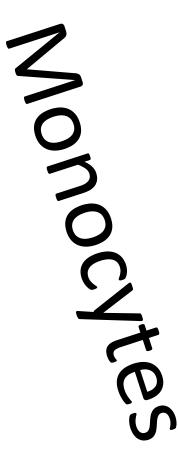


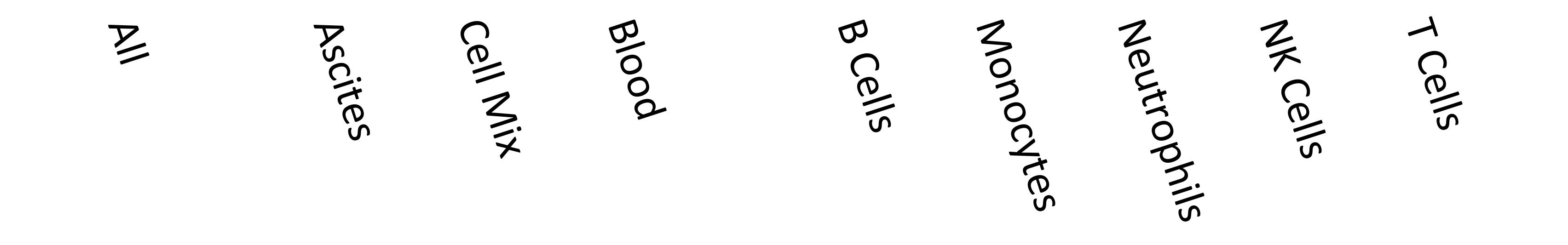
All Predictions

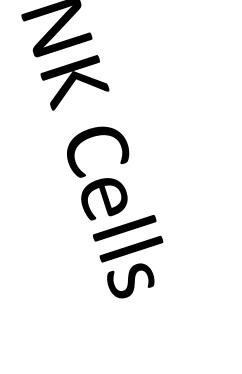
Datasets

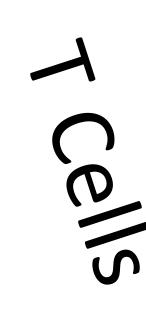


Kel



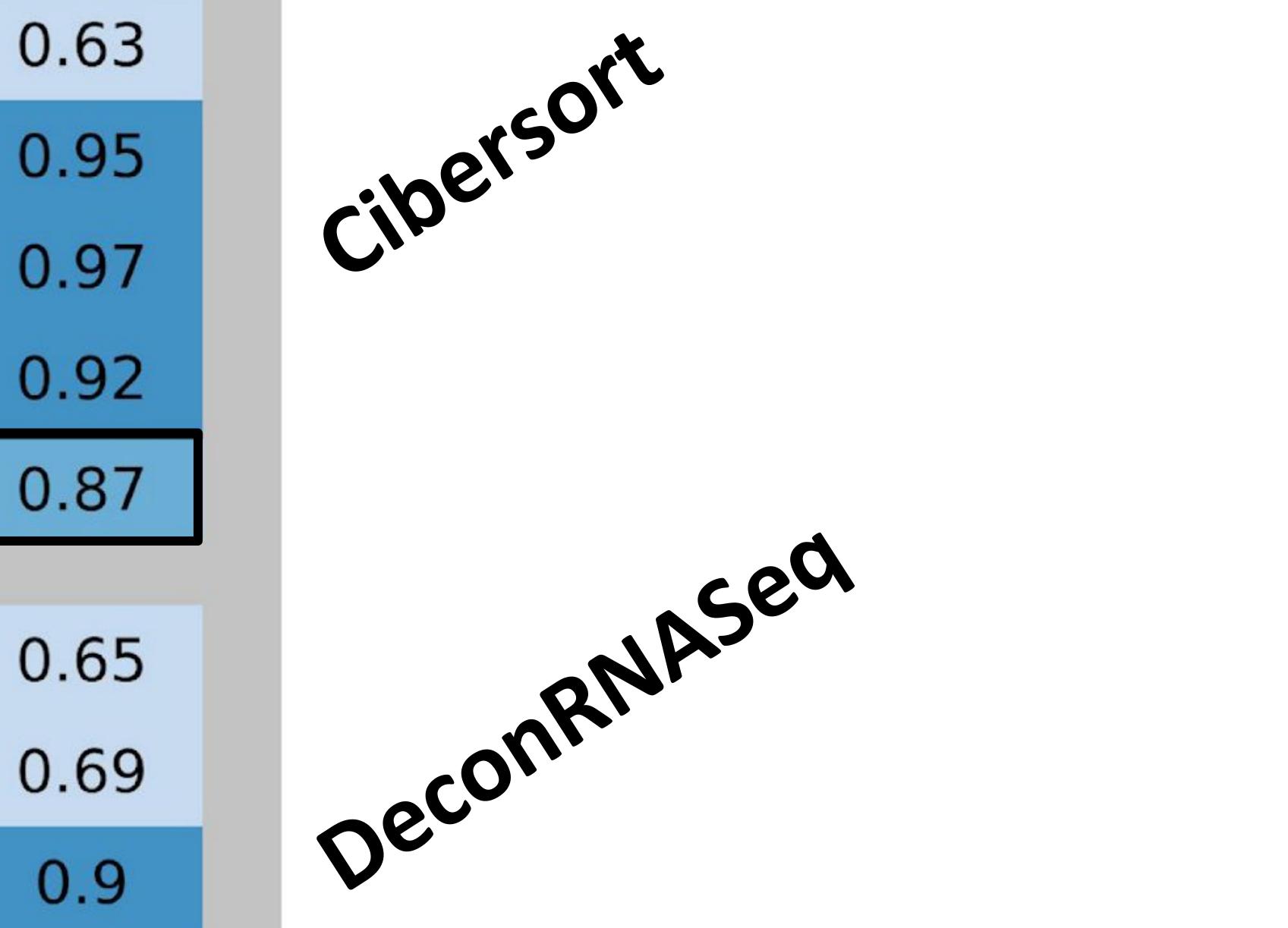




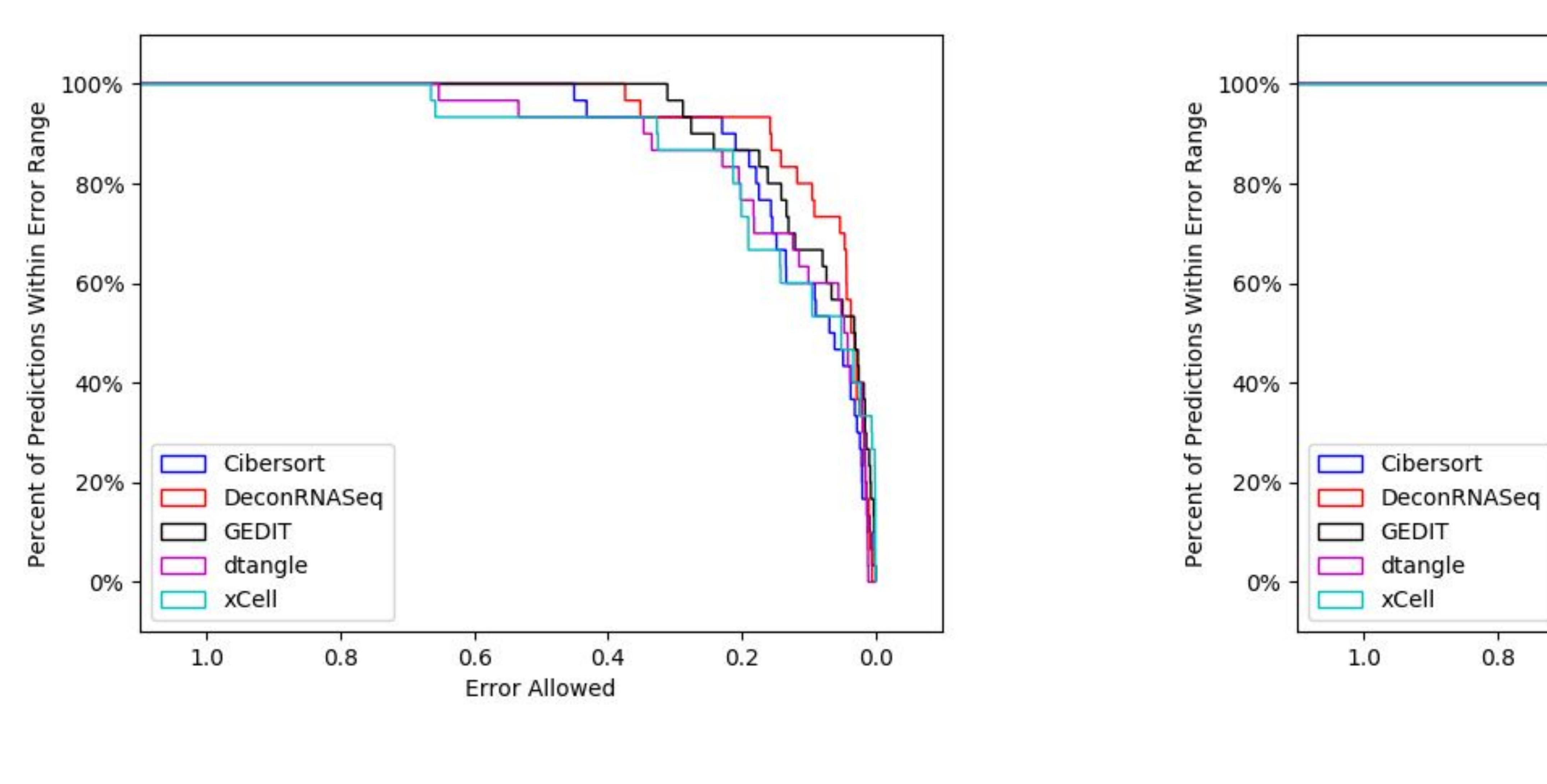




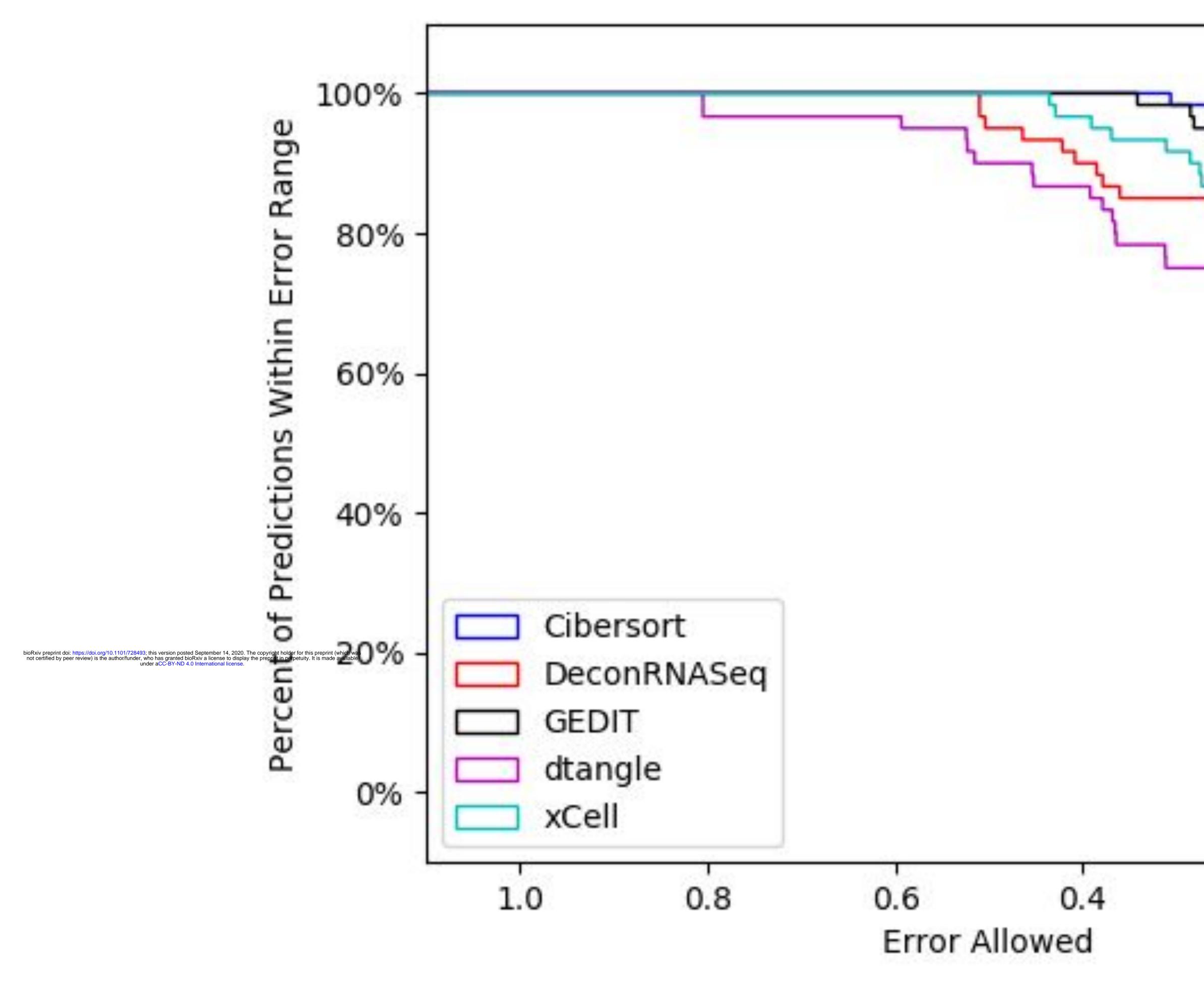




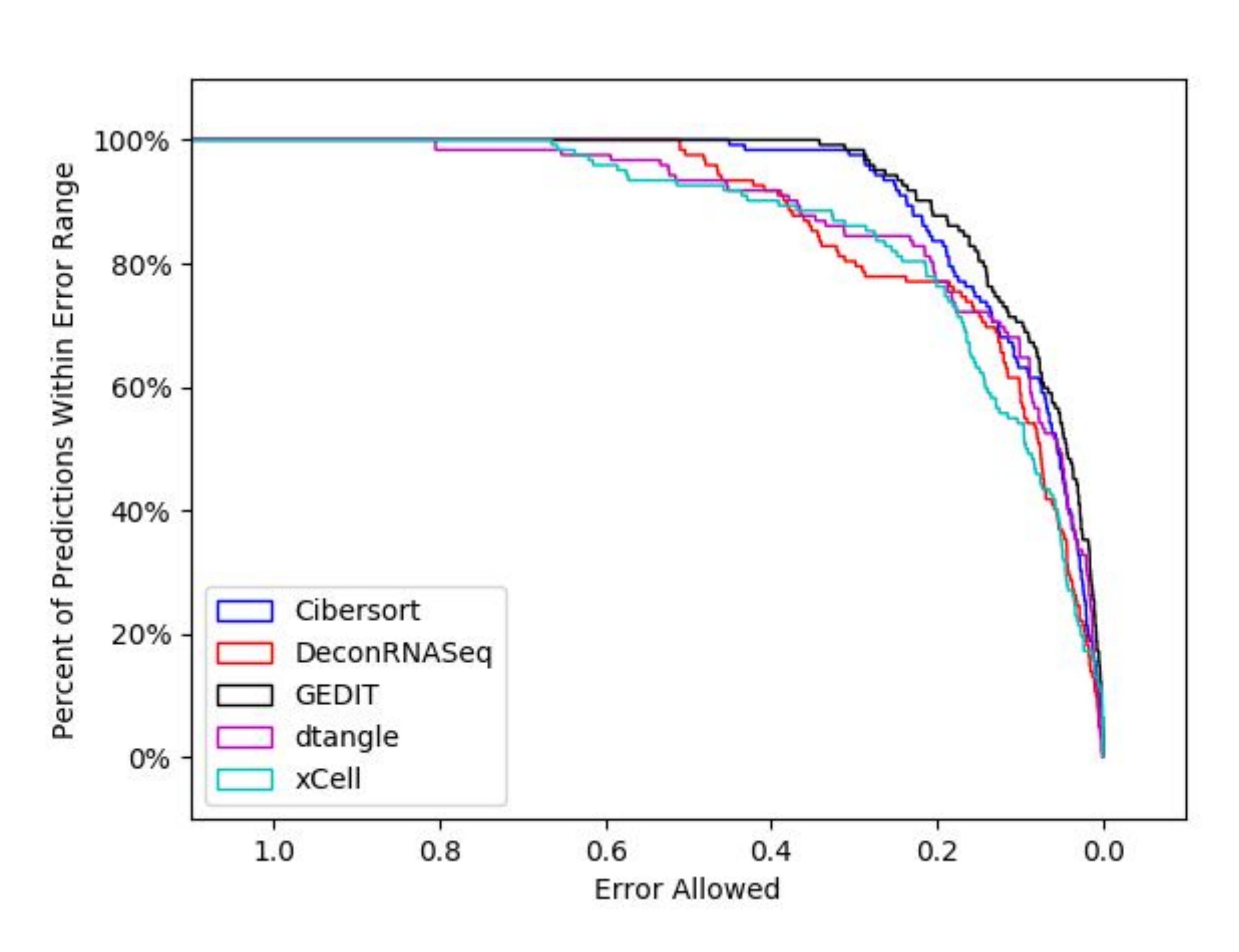
Cancer Ascites



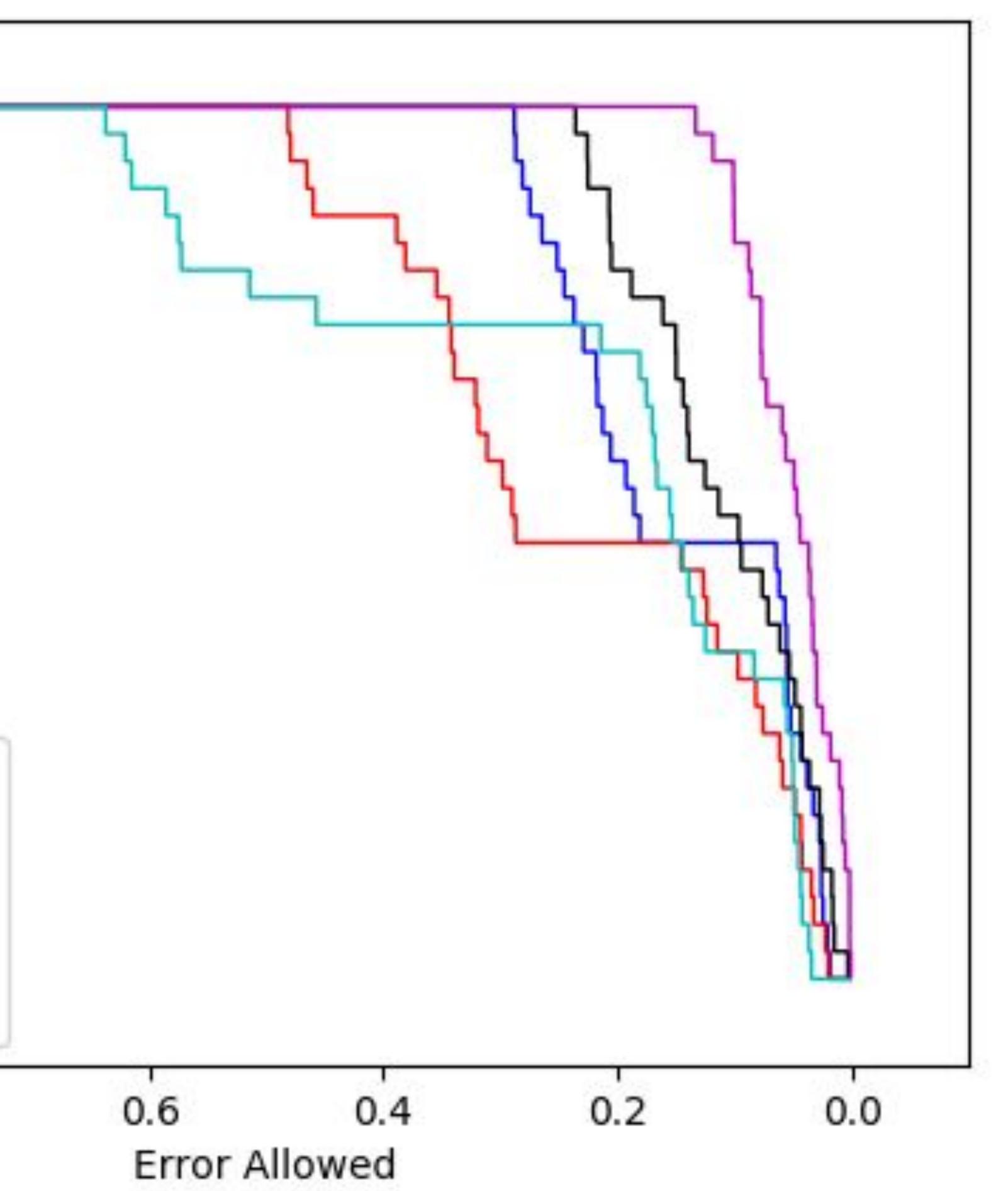
In Vitro Cell Mixtures



0.0 0.2

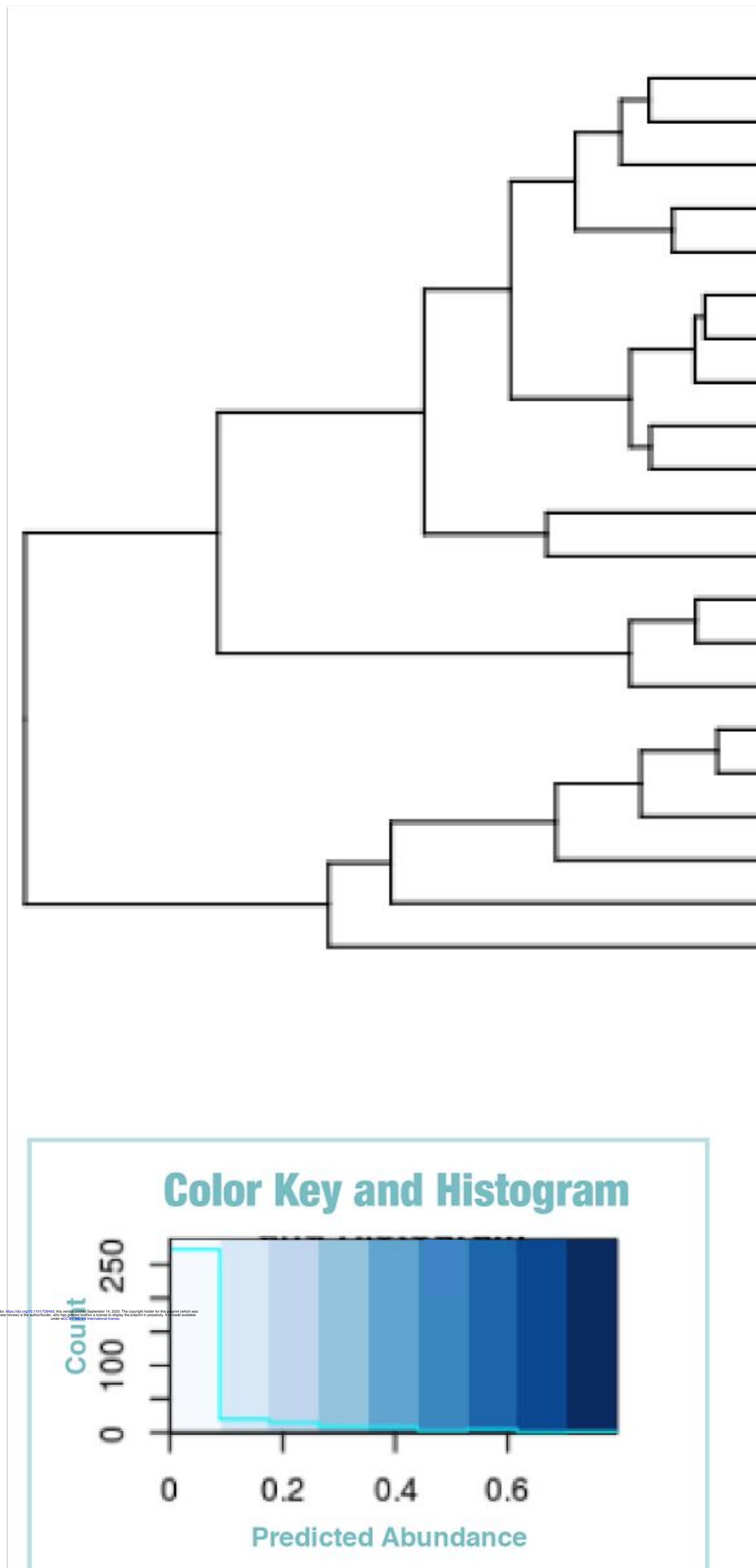


BIOOd



All Datasets

23.46



 -					_			
-								
_								
-								
	 	11 I						
-		-						
_								

ŝ

S - 10 Berlin -UU.

5

S

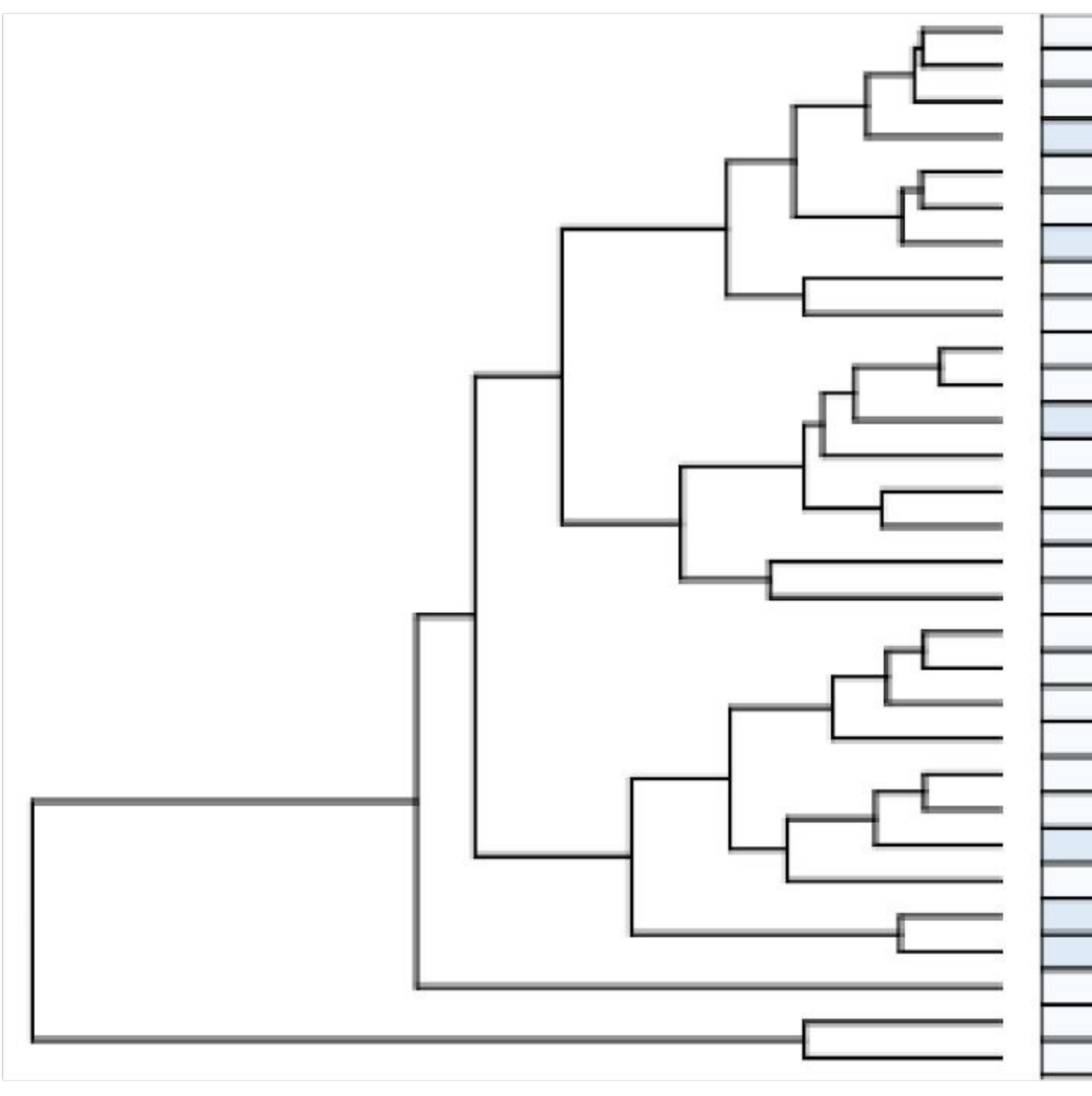
20

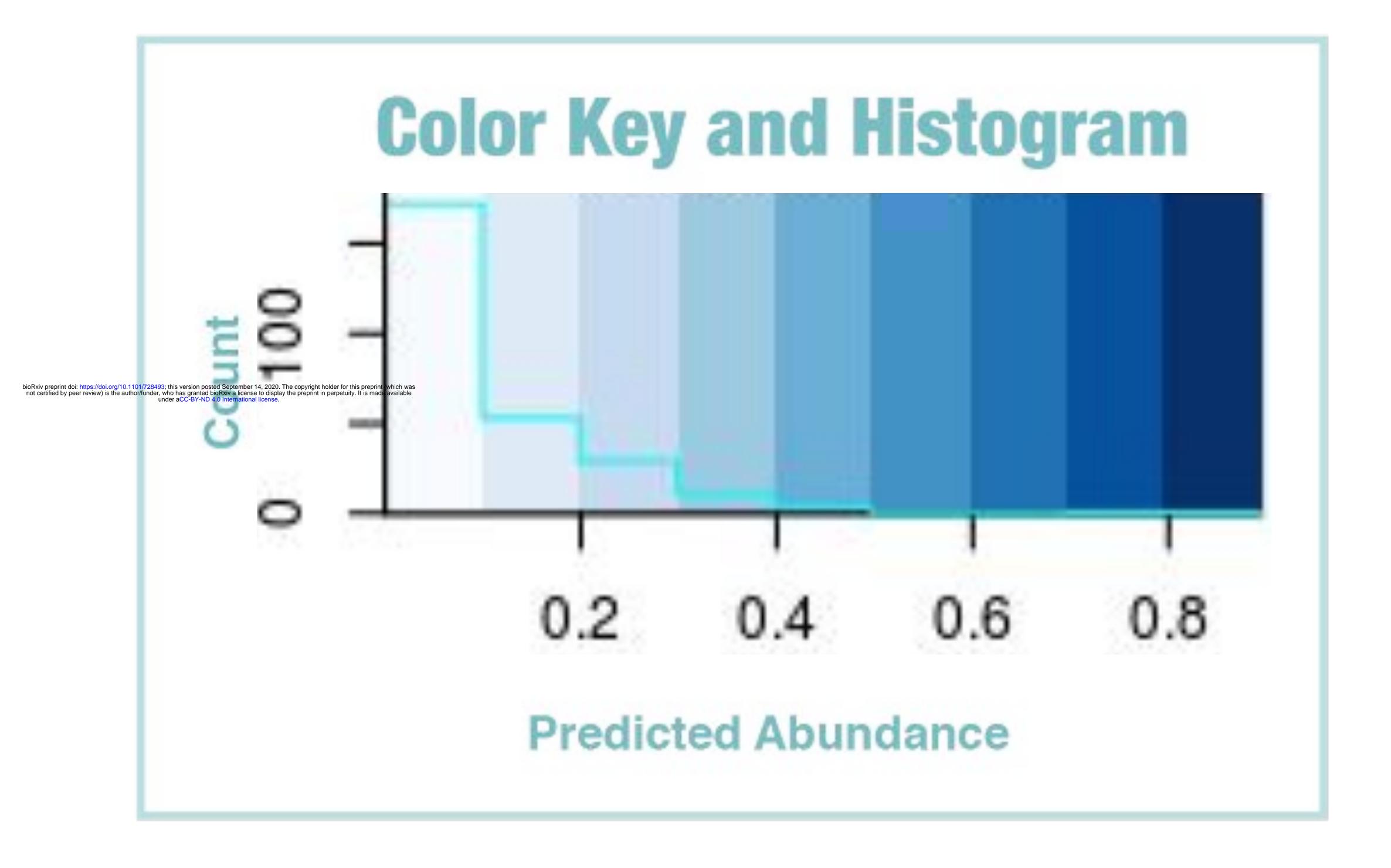
$\tilde{\omega}$ ч.v $\Box 0$

Burn Chancroid Discoid Lupus Erythematosus Wound- Post Operative Melanoma Alopecia Areata Atopic Dermatitis Psoriasis Normal Skin Acne Allergic Contact Dermititis Acute Wound Squamous Cell Carcinoma Irritant Contact Dermatitis Basal Cell Carcinoma Cutaneous Sarcoidosis Leprosy Reversal Reaction Erythema Nodosum Leprosum L-Leprosy Mycosis Fungoides Stevens-Johnson Syndrome

	0.01	0.06	0.55	0.01	0.01	0.1	0	0.01	0.12	0.01	0	0.12	m
	0.01	0.06	0.57	0.01	0	0.05	0	0.01	0.18	0.01	0.01	0.08	m
	0.01	0.06	0.58	0.01	0	0	0	0	0.08	0	0	0.25	m
	0.03	0.21	0.51	0.04	0.01	0	0	0.01	0.07	0.02	0.01	0.09	m
	0.03	0.2	0.56	0.01	0	0	0	0.01	0.06	0.03	0	0.1	m
					0								
	0.01	0.07	0.78	0.01	0	0	0	0	0.04	0	0	0.07	m
					0								
					0.02								
	0.03	0.06	0.31	0.03	0	0.08	0.01	0.02	0.22	0.02	0	0.2	g
					0.01								
	0.01	0.04	0.11	0.01	0	0.01	0.01	0.01	0.75	0.01	0	0.05	g
	0.02	0.03	0.06	0.03	0.6	0	0	0	0	0.13	0.11	0.01	s
	0.01	0.01	0.05	0.01	0.72	0	0	0	0	0.04	0.11	0.04	h
	0.01	0	0.03	0	0.05	0	0	0.82	0.02	0.04	0.03	0.01	e
<u>.</u>	0.01	0.01	0.04	0.01	0	0	0.84	0	0	0.01	0.01	0.06	liv
	0.8	0.03	0.02	0.01	0	0.11	0.01	0.01	0	0	0	0.02	tł
	0.8	0.04	0.03	0	0	0.08	0.01	0.01	0	0	0	0.02	tł
	0.73	0.07	0.02	0.02	0	0.12	0.01	0.01	0	0	0	0.02	T.
	0.78	0.09	0.03	0.01	0	0.04	0	0.01	0	0	0	0.03	T.
	0.9	0.02	0.01	0.01	0	0	0	0	0.01	0	0	0.03	tł
	0.63	0.18	0.04	0.01	0	0.07	0	0.01	0.01	0	0	0.04	T.
	0	0.01	0.03	0.01	0	0.9	0	0	0	0	0	0.04	Ν
	0.21	0.4	0.27	0	0	0.03	0	0	0	0.03	0.06	0	ly
	0.1	0.41	0.11	0.06	0	0.06	0.01	0.02	0.06	0.02	0.07	0.07	st
	0.1	0.62	0.16	0	0	0.02	0	0	0.01	0.02	0.03	0.04	S
Γ	0	0.93	0.02	0	0	0	0	0	0	0	0	0.05	fc
	0.01	0.93	0.01	0	0	0	0.01	0	0.01	0	0	0.03	В
8493; this version posted September 14, 2020. The copyright holder for this preprint (which was er, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.	0	0	0.01	0.96	0	0	0	0	0	0.02	0	0	ir
	0.02	0.02	0.02	0.83	0	0.01	0.05	0	0.01	0.01	0.03	0	ir
	T.cell	B.cell	Monocyte	Epithelial.cell	Cardiac.Muscle	NK.cell	Hepatocyte	Epidermis	Granulocyte	Stromal.cell	Endothelial.cell	Erythrocyte	

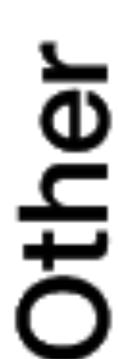
macrophage_bone_marrow_24h_LPS macrophage_bone_marrow_6hr_LPS macrophage_bone_marrow_2hr_LPS macrophage_peri_LPS_thio_1hrs macrophage_peri_LPS_thio_0hrs microglia macrophage_bone_marrow_0hr macrophage_peri_LPS_thio_7hrs mega_erythrocyte_progenitor granulo_mono_progenitor lung granulocytes_mac1.gr1. skeletal_muscle heart epidermis liver thymocyte_SP_CD8. thymocyte_SP_CD4. T.cells_CD8. T.cells_CD4. thymocyte_DP_CD4.CD8. T.cells_foxP3. NK_cells lymph_nodes stem_cells__HSC spleen follicular_B.cells B.cells_marginal_zone intestine_large intestine_small





Endothelia

Smooth Muscle



~~		

Fibroblasts

Schwann Cells

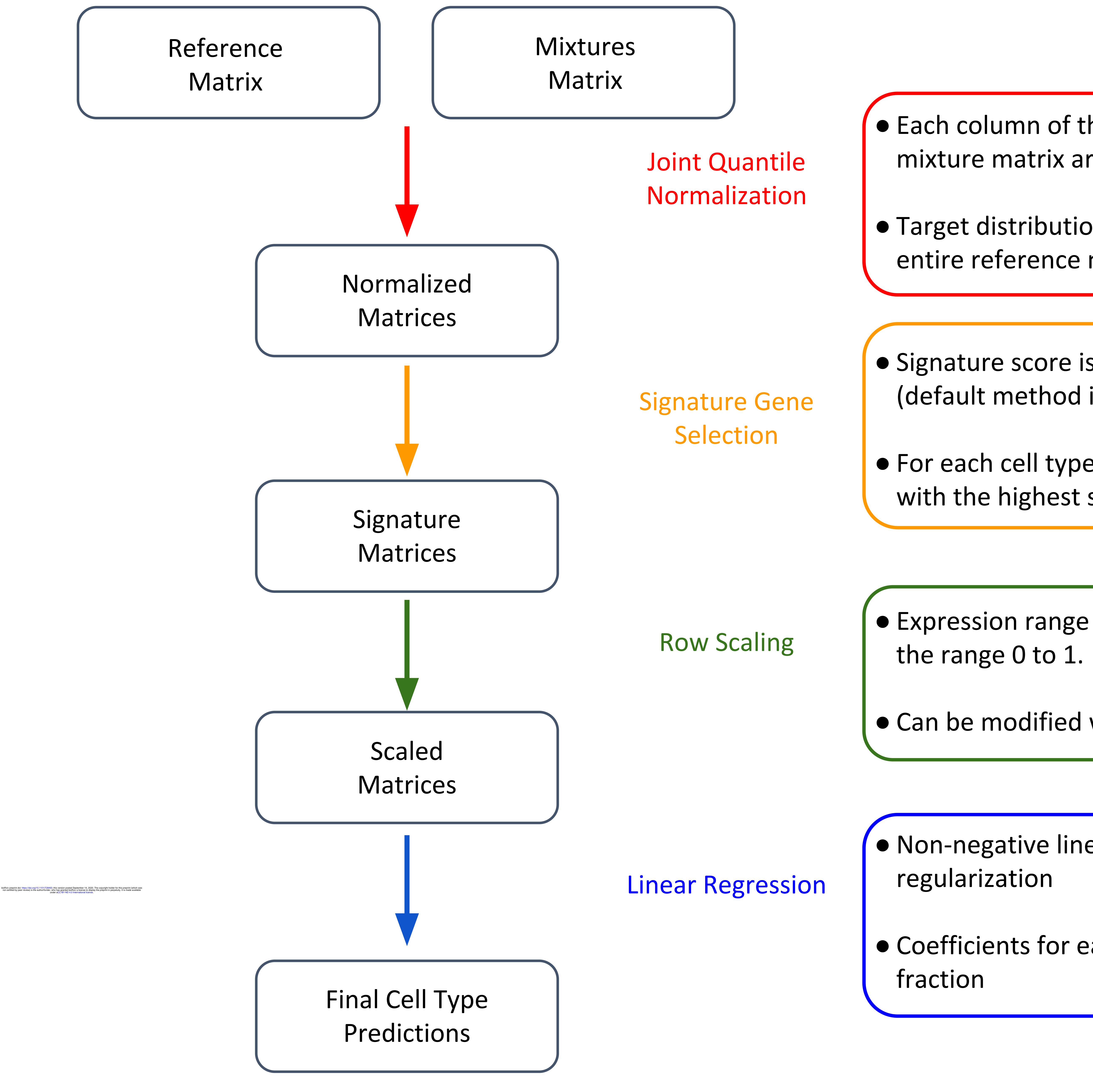
Keratinocyte

Epithelia

	PROSTATE
	BLADDER
	ESOPHAGUS
	THYROID
	PANCREAS
	STOMACH
	KIDNEY
	SALIVARY GLAND
	VAGINA
	MUSCLE
	ADRENAL GLAND
	LUNG
	LIVER
	SMALL INTESTINE
	PITUITARY
	TESTIS
	BRAIN
	BLOOD VESSEL
	UTERUS
	OVARY
	CERVIX UTERI
	FALLOPIAN TUBE
	COLON
	HEART
	NERVE
	ADIPOSE TISSUE
	BREAST
	SKIN
	SPLEEN
	BLOOD

Adipocyte

Total Imme



• Each column of the reference matrix and mixture matrix are quantile normalized

 Target distribution is starting distribution of entire reference matrix

 Signature score is calculated for each gene (default method is negative information entropy)

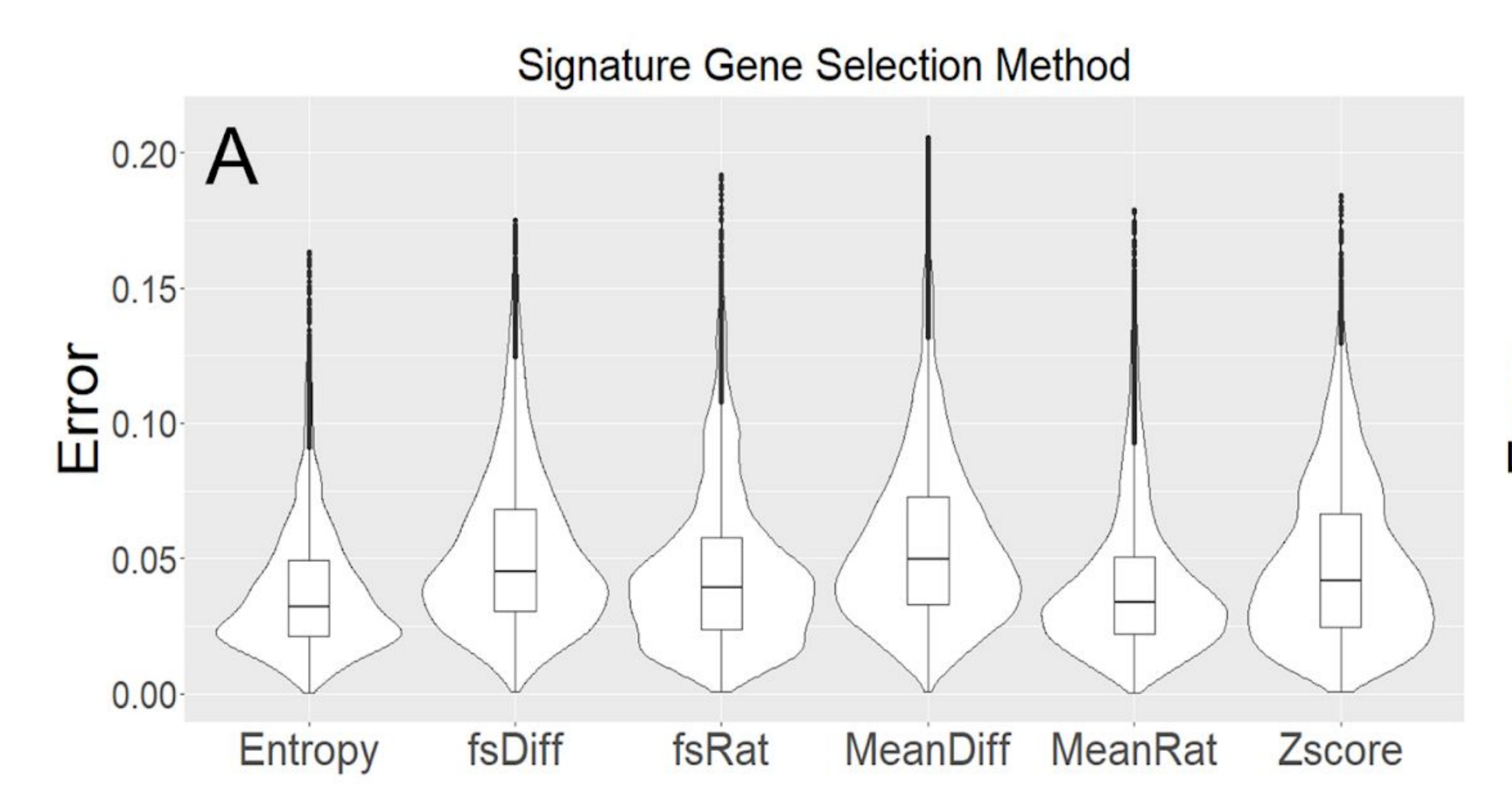
 For each cell type, a number of genes (NumSigs) with the highest scores are selected

 Expression range of each gene is adjusted to the range 0 to 1.

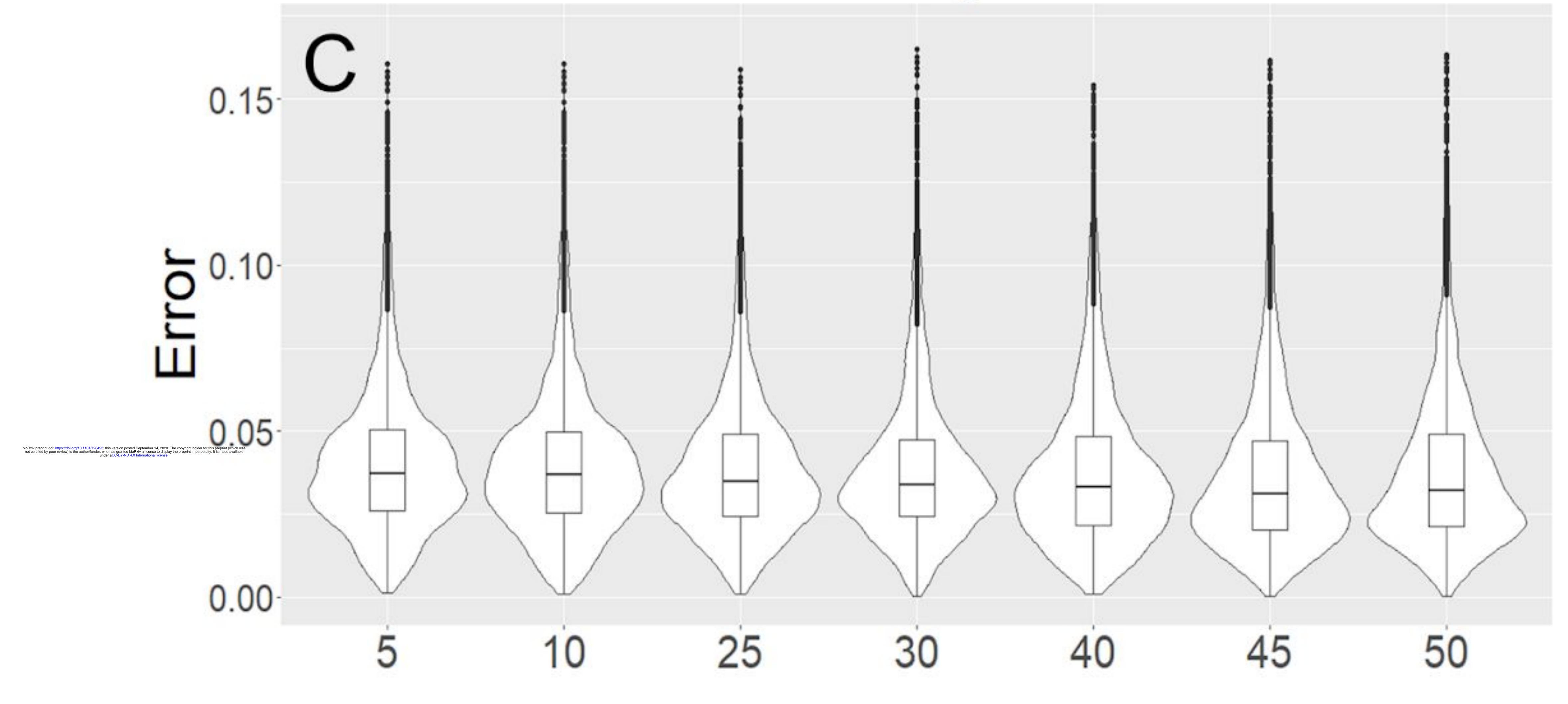
Can be modified with the RowScale parameter

Non-negative linear regression with no regularization

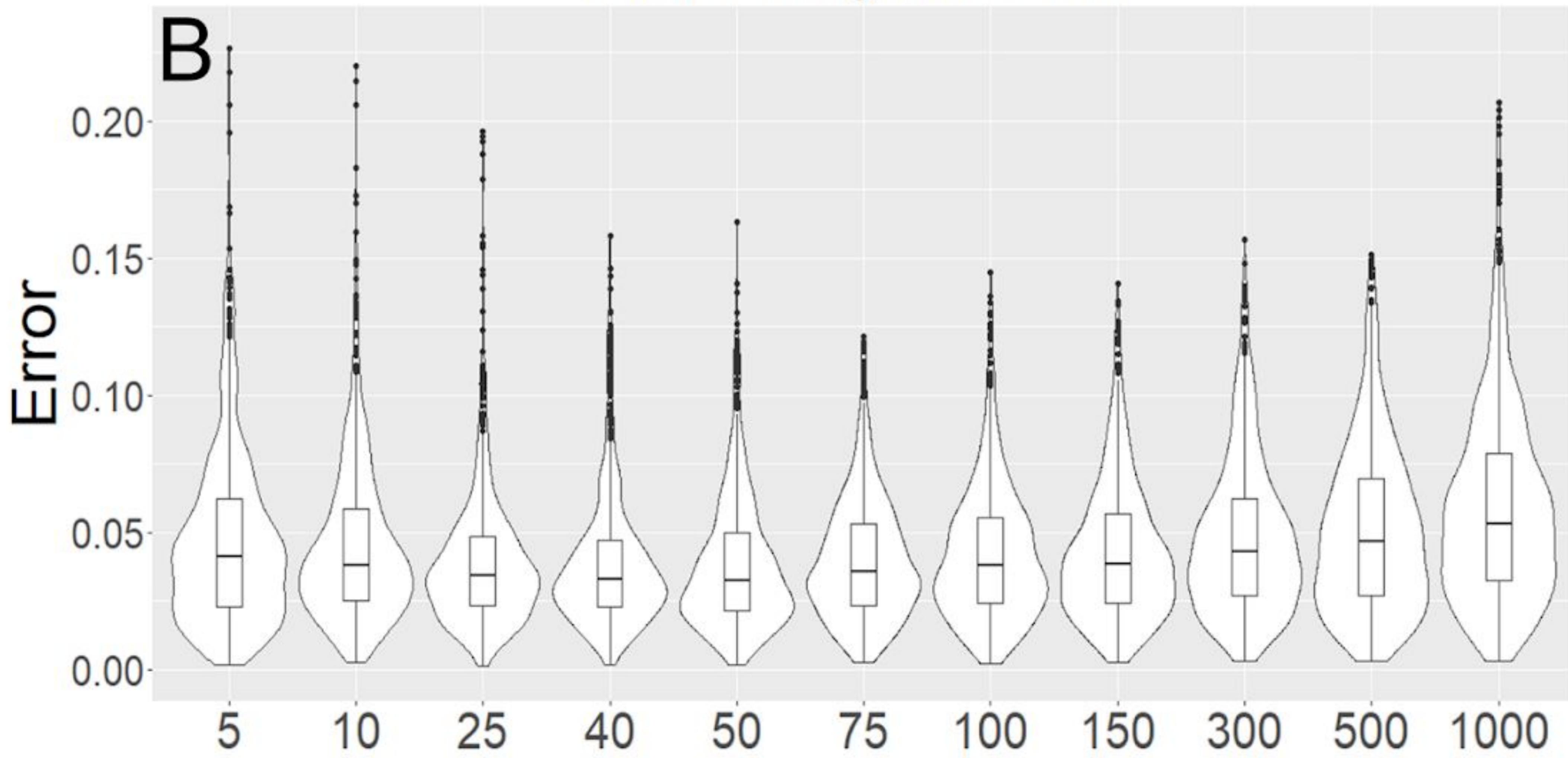
Coefficients for each cell type refer to predicted



Number of Fixed Signature Genes



Number of Signature Genes



Row Scaling Value

