1  **MAST: A flexible statistical framework for assessing transcriptional changes and**
2  **characterizing heterogeneity in single-cell RNA-seq data.**
3
4  Greg Finak[1]\*, Andrew McDavid[1]\*, Masanao Yajima[1]\*, Jingyuan Deng[1], Vivian Gersuk[2], Alex K.
5  Shalek[3,4,5], Chloe K. Slichter[1], Hannah W. Miller[1], M. Juliana McElrath[1], Martin Prlic[1], Peter S.
6  Linsley[2], and Raphael Gottardo[1].
7  [1]Vaccine and Infectious Disease Division and [2]Public Health Sciences Division, Fred Hutchinson Cancer
8  Research Center, Seattle, WA 98109, USA
9  [2]Benaroya Research Institute at Virginia Mason, Seattle, WA 98101, USA
10  [3]Institute for Medical Engineering & Science & Department of Chemistry, MIT, Boston, MA, 01239-4307,
11  USA
12  [4]Ragon Institute of MGH, MIT, & Harvard, Boston, MA, 02139-3583, USA
13  [5]Broad Institute of MIT & Harvard, Boston, MA, 01242, USA
14
15  \*Authors contributed equally to this manuscript.
16
17  **Abstract**
18  Single-cell transcriptomic profiling enables the unprecedented interrogation of gene
19  expression heterogeneity in rare cell populations that would otherwise be obscured in
20  bulk RNA sequencing experiments. The stochastic nature of transcription is revealed in
21  the bimodality of single-cell transcriptomic data, a feature shared across single-cell
22  expression platforms. There is, however, a paucity of computational tools that take
23  advantage of this unique characteristic.  We present a new methodology to analyze
24  single-cell transcriptomic data that models this bimodality within a coherent generalized
25  linear modeling framework. We propose a two-part, generalized linear model that allows
26  one to characterize biological changes in the proportions of cells that are expressing
27  each gene, and in the positive mean expression level of that gene. We introduce the
28  *cellular detection rate*, the fraction of genes turned on in a cell, and show how it can be
29  used to simultaneously adjust for technical variation and so-called "extrinsic noise" at the
30  single-cell level without the use of control genes. Our model permits direct inference on
31  statistics formed by collections of genes, facilitating gene set enrichment analysis. The
32  residuals defined by such models can be manipulated to interrogate cellular
33  heterogeneity and gene-gene correlation across cells and conditions, providing insights
34  into the temporal evolution of networks of co-expressed genes at the single-cell level.
35  Using two single-cell RNA-seq datasets, including newly generated data from Mucosal
36  Associated Invariant T (MAIT) cells, we show how model residuals can be used to
37  identify significant changes across biologically relevant gene sets that are missed by
38  other methods and characterize cellular heterogeneity in response to stimulation.
39
40  **Introduction:**
41  Whole transcriptome expression profiling of single cells via RNA-seq (scRNA-seq) is the logical
42  apex to single cell gene expression experiments. In contrast to transcriptomic experiments on
43  mRNA derived from bulk samples, this technology provides powerful multi-parametric
44  measurements of gene co-expression at the single-cell level.  However, the development of
45  equally potent analytic tools has trailed the rapid advances in the biochemistry and molecular

46  biology, and several challenges need to be addressed to fully leverage the information in single-
47  cell expression profiles.

49  First, single-cell expression has repeatedly been shown to exhibit a characteristic bimodal
50  expression pattern, wherein the expression of otherwise abundant genes is either strongly
51  positive, or undetected within individual cells. This is due in part to low starting quantities of
52  RNA such that many genes will be below the threshold of detection, but there is also a biological
53  component to this variation (termed extrinsic noise in the literature) that is conflated with the
54  technical variability[1-3]. We and other groups[4-6] have shown that the proportion of cells with
55  detectable expression reflects both technical and biological differences between samples.
56  Results from synthetic biology also support the notion that bimodality can arise from the
57  stochastic nature of gene expression[2,3,7,8].

59  Secondly, measuring single cell gene expression might seem to obviate the need to normalize
60  for starting RNA quantities. Recent work shows that cells scale transcript copy number with cell
61  volume (a factor that affects gene expression globally) to maintain a constant mRNA
62  concentration and thus constant biochemical reaction rates[9,10]. In scRNA-seq, cells of varying
63  volume are diluted to an approximately fixed reaction volume leading to differences in detection
64  rates of various mRNA species that are driven by the initial cell volumes. Technical assay
65  variability  (e.g. mRNA quality, pre-amplification efficiency) and extrinsic biological factors (e.g.
66  nuisance biological variability due to cell size) remain, and can significantly influence expression
67  level measurements. Consequently, this may render traditional normalization strategies using
68  the expression level of a few "housekeeping" genes, like GAPDH, infeasible[10]. Recently, Shalek
69  et al[5] observed a strong relationship between average expression and detection efficiency, and
70  have proposed a computational approach to correct the estimated gene-specific probability of
71  detection. Our approach easily allows for estimation and control of the CDR simultaneously
72  while estimating treatment effects as opposed to previous approaches[5] that relied on a set of
73  control genes and could not jointly model both factors.

75  Previously, Kharchenko et al[6] developed a so-called three-component mixture model to test for
76  differential gene expression while accounting for bimodal expression. Their approach is limited
77  to two-class comparisons and cannot adjust for important biological covariates such as multiple
78  treatment groups and technical factors such as batch or time information, severely limiting its
79  utility in more complex experimental designs. On the other hand, several methods have been
80  proposed for modeling bulk RNA-seq data that permit complex modeling through linear[11] or
81  generalized linear models[12,13] but these models have not yet been adapted to single-cell data as
82  they do not properly account for the observed bimodality in expression levels. This is particularly
83  important when adjusting for covariates that might affect the expression rates. As we will
84  demonstrate later, such model mis-specification can significantly affect sensitivity and specificity
85  when detecting differentially expressed genes and gene-sets.

87  Here, we propose a Hurdle model tailored to the analysis of scRNA-seq data, providing a
88  mechanism to address the challenges noted above. It is a two-part generalized linear model that
89  simultaneously models the rate of expression over background of various transcripts, and the

90  positive expression mean. Leveraging the established theory for generalized linear modeling
91  allows us to accommodate complex experimental designs while controlling for covariates
92  (including technical factors) in both the discrete and continuous parts of the model. We
93  introduce the *cellular detection rate (CDR)*: the fraction of genes that are turned on / detected in
94  each cell, which, as discussed above, acts as a proxy for both technical (e.g. dropout,
95  amplification efficiency, etc.) and biological factors (e.g. cell volume and other extrinsic factors
96  other than treatment of interest) that can influence gene expression. As a result it represents an
97  important source of variability in scRNA-seq data that needs to be considered (Figure 1). Our
98  approach of modeling the CDR as a covariate, offers an alternative to the weight correction of
99  Shalek et al[5] that does not depend on the use of control genes and allows us to jointly estimate
100  nuisance and treatment effects. Our framework permits the analysis of complex experiments,
101  such as repeated single cell measurements under various treatments and/or longitudinal
102  sampling of single cells from multiple subjects with a variety of background characteristics (e.g.
103  gender, age, etc.) as it is easily extended to accommodate random effects. Differences between
104  treatment groups are summarized with pairs of regression coefficients whose sampling
105  distributions are available through bootstrap or asymptotic expressions, enabling us to perform
106  complementary differential gene expression and gene set enrichment analyses (GSEA). We use
107  an empirical Bayesian framework to regularize model parameters, which helps improve
108  inference for genes with sparse expression, much like what has been done for bulk gene
109  expression[14]. Our GSEA approach accounts for gene-gene correlations, which is important for
110  proper control of type I errors[15]. This GSEA framework is particularly useful for synthesizing
111  observed gene-level differences into statements about pathways or modules. Finally, our model
112  yields *single cell residuals* that can be manipulated to interrogate cellular heterogeneity and
113  gene-gene correlations across cells and conditions. We have named our approach MAST for
114  Model-based Analysis of Single-cell Transcriptomics.
115  
116  We illustrate the method on two data sets. We first apply our approach to an experiment
117  comparing primary human non-stimulated and cytokine-activated Mucosal-Associated Invariant
118  T (MAIT) cells. MAST identifies novel expression signatures of activation, and the single-cell
119  residuals produced by the model highlights a population of MAIT cells showing partial activation
120  but no induction of effector function. We then illustrate the application of MAST to a previously-
121  published complex experiment studying temporal changes in murine bone marrow-derived
122  dendritic cells subjected to LPS stimulation.  We both recapitulate the findings of the original
123  publication and describe additional coordinated gene expression changes at the single-cell level
124  across time in LPS stimulated mDC cells.
125  
126  **Results**
127  
128  **MAST can account for variation in the cellular detection rate.** As discussed previously and
129  as shown on Figure 1 by principal component analysis (PCA), the cellular detection rate (CDR,
130  see Methods for exact definition), is an important source of variability. It is highly correlated with
131  the second principal component (PC, Pearson's rho=0.76 grouped, 0.91 stimulated, 0.97 non-
132  stimulated) in the MAIT dataset and the first PC (rho=0.92 grouped, 0.97 non-stimulated, 0.92
133  LPS, 0.89 PAM, 0.92 PIC) in the mDC dataset. We observe larger CDR variability within

134     treatment groups than across groups, suggesting that it is likely to be a nuisance factor. This is
135     further supported by the fact that the CDR calculated within control (e.g. housekeeping) genes
136     is highly correlated with the CDR calculated over all genes (Supplementary Figure 1). Its role as
137     a principal source of variation persists across experiments (Figure 1).
138

139     We thus conjecture that CDR is a proxy for unobserved nuisance factors that should be
140     explicitly modeled. In particular, it is not unreasonable to suggest that the CDR captures
141     variation in global transcription rate due to variations in cell size (among other factors)[10], as well
142     as technical variation such as dropout, with dropout rates possibly correlated with cell-size.
143     Fortunately, MAST easily accommodates covariates, such as the CDR, and more importantly
144     allows joint, additive modeling of them with other biological variables of interest, with the effect
145     of each covariate decomposed into its discrete and continuous parts. This two-part modeling is
146     key to account for the CDR that directly reflects the gene-level transcription rates. Applying an
147     analysis of deviance with MAST (see Methods), we quantified the amount of variability that
148     could be attributed to CDR. The CDR accounts for 5.2% of the deviance in the MAIT data set
149     and 4.8% in the mDC data set for the average gene, and often times much more than that: it
150     comprises more than 9% of the deviance in over 10% of genes in both data sets, particularly for
151     the discrete component of the model (Supplementary Figure 2). It should also be noted that the
152     CDR deviance estimates for many of the genes are comparable (if not greater) to the treatment
153     deviance estimates showing that it.
154

155     That CDR predicts expression levels contradicts the model of independent expression between
156     genes, since the level of expression (averaged across many genes) would not affect the level in
157     any given gene were expression independent. This pervasiveness suggests latent factors are
158     creating coordinated changes in expression across genes. In light of the work of Padovan-
159     Merhar et al[10], we conjecture the latent factor relates to differences in cell volumes, since cells
160     of different volumes compensate to conserve mRNA species molarity, which implies higher copy
161     numbers of all transcripts in larger cells. Higher copy numbers result in higher scRNA-seq
162     detection rates globally across transcripts.
163

164     Finally, we have investigated the relationship between our approach and the weight correction
165     of Shalek et al[5] (Supplementary Figure 3). We observe a strong linear relationship between the
166     CDR and the weights of Shalek et al[5]. Thus, use of the CDR as a covariate can be seen as a
167     statistically rigorous way to correct for the dropout biases of Shalek et al[5], without the need to
168     use control genes,, and more importantly with the ability to control for these while estimating
169     treatment effects.
170

171     **Single-cell sequencing identifies a transcriptional profile of MAIT cell activation**
172     We applied MAST to our MAIT dataset to identify genes up- or down-regulated by cytokine
173     stimulation while accounting for variation in the CDR (see Methods). We detected 291
174     differentially expressed genes, as opposed to 1413 when excluding CDR. To determine whether
175     this was due to a change in ranking or a simply a shift in significance, we compared the overlap
176     between the top $n$ genes in both models (varying $n$ from 100 to 1413), and found that, on
177     average, 35% (range 32% - 38%) of genes are excluded when CDR is modeled, suggesting that

178    inclusion of this variable allows global changes in expression, manifest in the CDR, to be
179    decomposed from local changes in expression.  This is supported by gene ontology enrichment
180    analysis (Supplementary Figure 4) of these CDR-specific genes (n=539), where we see no
181    enrichment for modules associated with treatment of interest.
182
183    In order to assess the type-I error rate of our approach, we also applied MAST to identify
184    differentially expressed genes across random splits of the non-stimulated MAIT cells. As
185    expected, MAST did not detect any significant differences (Supplementary Figure 5A), whereas
186    DEseq and edgeR, designed for bulk RNA-seq, detected large number of differentially
187    expressed genes even at very low FDR thresholds. We examined the GO enrichment of genes
188    detected by limma or edgeR or DESeq but not MAST and found that these sets lacked
189    significant enrichment for modules related to the treatment of interest (Supplementary Figures
190    5B and 6-8). MAST's testing framework evidently has better specificity than these approaches.
191
192    Figure 2A shows the single-cell expression (log$_2$-TPM) of the top 100 genes identified as
193    differentially expressed between cytokine (IL18, IL15, IL12) stimulated (purple) and non-
194    stimulated (pink) MAIT cells using MAST.  Following stimulation with IL12/15/18, we observe
195    increased expression in genes with effector function including Interferon$-\gamma$ (IFN$-\gamma$), granzyme-
196    B (GZMB) as has been reported in NK, NKT and memory T cells, and a concomitant
197    downregulation of the AP-1 transcription factor network. CD69 is an early and only transient
198    marker of activation that can be induced by stimulation of the T cell receptor or by cytokine
199    signals. Its downregulation at the mRNA level after 24h is likely preceding subsequent protein-
200    level downregulation[16-18].
201
202    We used these lists of up- and down-regulated genes to define a MAIT activation score that
203    differentiates between stimulated and non-stimulated MAITs as shown in Figure 2B. This score
204    (see Methods), for each cell, is defined as the expected expression level across genes in a
205    module (based on the model fit) corrected for nuisance factors (such as CDR, see Methods).
206    The score enables us to cleanly differentiate stimulated and non-stimulated cells, and
207    demonstrates that the stimulated MAIT population is much more heterogeneous in its
208    expression phenotype. In particular, a few stimulated MAIT cells (SC08, SC54, SC48, SC15,
209    SC46, and SC61 in Figure 2A) exhibit low expression of IFN$-\gamma$ response genes, suggesting
210    these cells did not fully activate despite stimulation. Post-sort experiments via FCM show that
211    the sorted populations were over 99% pure MAITs (Supplementary Figure 9A), and exhibited a
212    change in cell size upon stimulation (Supplementary Figure 9B), and that up to 26% of
213    stimulated MAITs didn't express IFN-$\gamma$ or GZMB following cytokine stimulation (Supplementary
214    Figure 9C). The non-responding cells in the RNA-seq experiment likely correspond to these
215    non-responding cells from the flow cytometry experiment, and the observed frequencies of
216    these cells in the RNA-seq and flow populations are consistent with each other ( Pr(observing 6
217    or fewer non-responding cells) = 0.16 under binomial sampling). We discuss this heterogeneity
218    in a further section. Importantly, the lists of up- and down-regulated genes can be used to define
219    gene sets for gene set enrichment analysis in order to identify transcriptional changes related to
220    MAIT activation in bulk experiments.
221

222 **Gene set enrichment analysis highlights pathways implicated in MAIT cell activation.**
223 We used MAST to perform gene set enrichment analysis (GSEA, see methods) in the MAIT
224 data using the blood transcriptional modules of Li et al[19]. The cell-level scores for the top 9
225 enriched modules (Figure 3A) continue to show significant heterogeneity in the stimulated cells,
226 particularly for modules related to T-cell signaling, protein folding, proteasome function, and the
227 AP-1 transcription factor network.  Enrichment in stimulated cells (green) and non-stimulated
228 cells (pink) is displayed for each module for the discrete and continuous components of the
229 model (Figure 3B, see Methods), as well as a Z-score combining the discrete and continuous
230 parts. The enrichment in the T-cell signaling module is driven by the increased expression of
231 IFN-$\gamma$, GZMB, IL2RA, IL2RB, and TNFRSF9, 5 of the 6 genes in the module.  Stimulated cells
232 also exhibit increased energy usage, translation and protein synthesis, while down-regulating
233 genes involved in cell cycle growth and arrest (and other cell cycle related modules). The down-
234 regulation of cell cycle growth inhibition genes indicates that IL-12/15/18 signals are sufficient to
235 prepare MAIT cells for cell proliferation. Interestingly, we observe down-regulation of mRNA
236 transcripts from genes in the AP-1 transcription factor network. This has been previously
237 described in dendritic cells in response to LPS stimulation[20] and, indeed, we observe this effect
238 in the mDC data set analyzed here (Supplementary Figure 10).
239
240 Our GSEA approach is more powerful than existing methods for bulk RNA-seq data
241 (Supplementary Figure 11), and we discover significantly enriched modules with clear patterns
242 of stimulation-induced changes that other methods omit (Supplementary Figure 12). Two such
243 modules include the "T-cell surface signature" and "chaperonin mediated protein folding, whose
244 component genes show elevated expression in response to stimulation (Supplementary Figure
245 12A-D). These additional discoveries are not solely due to greater permissiveness in MAST.
246 We applied MAST to identify differentially expressed gene sets across random partitions of the
247 non-stimulated cells, to examine its false discovery rate. As expected, MAST did not detect any
248 significant differences, which suggests that it has good type I error control.
249
250 **Residual analysis identifies networks of co-expressed genes implicated in MAIT cell**
251 **activation.** Much of the heterogeneity between the non-responding and responding stimulated
252 cells remains even after removal of marginal (gene level) stimulation effects. Since, MAST
253 models the expected expression value for each cell, we can compute residuals adjusted for
254 known sources of variability (See Methods). The residuals can be compared across genes to
255 characterize cellular heterogeneity and correlation. We observe co-expression in the residuals
256 from stimulated cells that is not evident in the non-stimulated group (Figure 4A,B).  Since the
257 residuals have removed any marginal changes due to stimulation in each gene, the average
258 residual in the two groups is comparable. The co-expression observed, meanwhile, is due to
259 individual cells expressing these genes *dependently*, where pairs of genes appear together
260 more often than expected under a model of independent expression.
261
262 Two clusters of co-expressed genes stand out in the residuals of the stimulated cells (Figure 4
263 B). These clusters show coordinated, early up-regulation of GZMB and IFN-$\gamma$ in response to
264 stimulation in MAIT cells and a concomitant decrease in CD69 expression, an early and

6

265    transient activation marker. PCA of the model residuals highlights the non-responsive stimulated
266    MAIT cells (Figure 4C).
267
268    Accounting for the CDR reduces the background correlation observed between genes
269    (Supplementary Figure 13) where nearly 25% of pairwise correlations decrease after CDR
270    correction. When the CDR is included in the model, the number of differentially expressed
271    genes with significant correlations across cells (FDR adjusted p-value < 1%) decreases from 73
272    to 61 in the stimulated cells, and from 808 to 15 in non-stimulated cells. This shows that
273    adjusting for CDR is also important for co-expression analyses as it reduces background co-
274    expression attributable to cell volume, which otherwise results in dense, un-interpretable gene
275    networks.
276
277

278    **MAST on complex experimental designs: temporal expression patterns of mouse**
279    **dendritic cell maturation**
280    Shalek et al[5] analyzed murine bone-marrow derived dendritic cells simulated using three
281    pathogenic components over the course of six hours and estimated the proportion of cells that
282    expressed a gene and the expression level of expressing cells.  We compared results from
283    applying our model to those obtained by Shalek et al[5] when analyzing their lipopolysaccharide
284    (LPS) stimulated cells. As with the MAIT analysis, we used MAST adjusting for the CDR. MAST
285    identified a total of 1359 differentially expressed genes (1996 omitting the CDR), and the CDR
286    accounted for 5.2% of the model deviance in the average gene.
287    The most significantly elevated genes at 6h include CCL5, CD40, IL12B, and Interferon-
288    inducible (IFIT) gene family members, while down-regulation was observed for EGR1 and
289    EGR2, transcription factors that are known to negatively regulate dendritic cell
290    immunogenicity[21].
291

292    **GSEA of mouse bone marrow-derived dendritic cells**
293    We performed GSEA with the Mouse GO modules and three modules Shalek et al[5] identified.
294    The blood transcriptional modules of Li et al[19] are shown in Supplementary Figure 10. Figure 5
295    shows module scores for significant GSEA modules for the LPS stimulated cells where the
296    heatmap represents Z values (see methods for details). Besides finding signatures consistent
297    with the modules from Shalek et. al. (Figure 5A), we identify modules that show similar
298    annotation and overlap significantly with the *core antiviral* and *sustained inflammatory*
299    signatures, including several modules linked to type 1 interferon response and antiviral
300    signatures (Figure 5B).  The "cellular response to interferon- beta" signature (n = 22) overlaps
301    with the original core antiviral signature (n = 99) by 13 genes (hypergeometric p = $1.24 \times 10^{-23}$).
302    The *response* and *defense response to virus* signatures overlap with the core antiviral signature
303    by 17 of 43 and 22 of 74 genes (hypergeometric p=$3.64 \times 10^{-26}$ and $4.08 \times 10^{-29}$, respectively),
304    suggesting the core antiviral signature captures elements of these known signatures. The
305    *chemokine* (n=16) and *cytokine activity* (n=51) modules overlap with the sustained inflammatory
306    (n = 95) module by 5 and 12 genes, respectively (hypergeometric p=$5.10 \times 10^{-9}$ and $9.53 \times 10^{-16}$).
307    Our modeling approach identifies the two "early marcher" cells in the core antiviral module
308    (marked with triangles on Figure 5A) corresponding to the same cells highlighted in Figure 4b of

7

309   Shalek et al[5]. Other modules exhibiting significant time-dependent trends include a module of
310   genes involved in the AP-1 transcription factor network that is down-regulated (Supplementary
311   Figure 10), a finding which has been previously shown in human monocytes following LPS
312   stimulation[20]. As with the MAITs, GSEA permutation analysis to evaluate type I error rates did
313   not identify any significant modules (data not shown). These results further confirm the original
314   findings and demonstrate the increased sensitivity of our approach. GSEA heatmaps for the
315   other stimulations can be found in Supplementary Figure 14.
316
317   **Residual analysis of mouse bone marrow-derived dendritic cells identifies sets of co-**
318   **expressed genes.**
319   We also explored stimulation-driven correlation patterns. Principal component analysis (Figure
320   6A) of the model residuals demonstrates a clear time trend associated with PC1, as cells
321   increase co-expression of interferon-activated genes. After removing the marginal stimulation
322   and adjusting for the CDR, we observe correlation between chemokines CCL5, TNF receptor
323   CD40, and interferon-inducible (IFIT) genes (Figure 6B). A principal finding of the original
324   publication was the identification of a subset of cells that exhibited an early temporal response
325   to LPS stimulation. Recapitulating the original results here, when we examine the PCA of the
326   residuals using the genes in the core antiviral module, we can identify the "early marcher" cells
327   at the 1h time-point (Supplementary Figure 15).  The co-expression plot for other stimulations
328   can be found in the supplementary material (Supplementary Figures 16 and 17).
329
330   **Discussion**
331   We have presented MAST, a flexible statistical framework for the analysis of scRNA-seq data.
332   MAST is suitable for supervised analyses about differential expression of genes and gene-
333   modules, as well as unsupervised analyses of model residuals, to generate hypotheses
334   regarding co-expression of genes. MAST accounts for the bimodality of single-cell data by
335   jointly modeling rates of expression (discrete) and positive mean expression (continuous)
336   values. Information from the discrete and continuous parts is combined to perform inference
337   about changes in expression levels using gene or gene-set based statistics. Because our
338   approach uses a generalized linear framework, it can be used to jointly estimate nuisance
339   variation from biological and technical sources, as well as biological effects of interest. In
340   particular, we have shown that it is important to control for the proportion of genes detected in
341   each cell, which we refer to as the cellular detection rate (CDR), as this factor can single-
342   handedly explain 13% of the variability in the 90% percentile gene. Adjusting for CDR at least
343   partially controls for differences in abundance due to cell size (and other extrinsic biological and
344   technical effects), while omitting it would lead to overestimated effects of the treatment on the
345   system. Using several scRNA-seq datasets, we showed that our approach provides a
346   statistically rigorous improvement to methods proposed by other groups in this context[5].
347
348   Because our approach is regression-based, it can be used to compute residuals to explore
349   cellular heterogeneity and gene-gene correlations after selected technical and/or biological
350   effects have been removed. In particular, using this approach, we identify MAIT cells that do not
351   have a typical activated expression profile in response to stimulation (Figures 2 and 3). The
352   proportion of these cells detected in the scRNASeq experiment is consistent with what was

8

353 detected in the flow cytometry experiment. These cells do not produce IFN-$\gamma$ or GZMB upon to

354 cytokine stimulation and exhibit expression profiles intermediate to non-stimulated and

355 stimulated cells (Supplementary Figure 18C). The cells exhibit lower levels of IFN-$\gamma$ and GZMB

356 than activated cells (Supplementary Figure 18A), but also exhibit decreased expression of AP-1

357 component genes Fos and FosB, consistent with other stimulated cells (Supplementary Figure

358 18B).

359

360 As discussed by Padovan-Merhar et al[10], care must be taken when interpreting experiments

361 where the system shows global changes in CDR across treatment groups, as this could result in

362 confounding treatment effect with differences in cell volume, which are not necessarily of

363 biological interest. Our approach addresses this issue as MAST allows joint modeling of CDR

364 and treatment effects, so the interpretation of the treatment effect is that the cell volume/CDR

365 has been held constant. It is also possible to only use CDR as a precision variable by centering

366 the CDR within each treatment groups, which makes the CDR measurement orthogonal to

367 treatment.  This would implicitly assume that the observed changes are treatment induced,

368 while still modeling the heterogeneity in cell volume within each treatment group. An alternative

369 approach would be to estimate the CDR coefficient using a set of control genes assumed to be

370 treatment invariant, such as housekeeping or ERCC spike-ins[22,23] and including it as an offset to

371 the linear predictors in the regression. An analogous approach is undertaken by Buettner et.

372 al.[22], however it does not account for bimodality and does not jointly model technical and

373 biological effects.

374

375 MAST is available as an R package (http://www.github.com/RGLab/MAST, doi:

376 10.5281/zenodo.18539). All data and results presented in this paper – including code to

377 reproduce the results – are available at:

378 (http://github.com/RGLab/MASTdata/archive/v1.0.0.tar.gz, doi: 10.5281/zenodo.18540). It

379 should also be noted that while most of the methodology presented here was developed for

380 scRNA-seq, it should be applicable to other single-cell gene expression platforms.

381

382 **Figure Captions**

383 **Figure 1.** The fraction of genes expressed, or cellular detection rate (CDR), explains the

384 principal components of variation in MAIT and DC data sets.

385

386 **Figure 2.**  Single-cell expression (log$_2$-TPM) of the top 100 genes identified as differentially

387 expressed between cytokine (IL18, IL15, IL12) stimulated (purple) and non-stimulated (pink)

388 MAIT cells using MAST (A).  Partial residuals for up- and down- regulated genes are

389 accumulated to yield an activation score (B), and this score suggests that the stimulated cells

390 have a more heterogeneous response to stimulation than do the non-stimulated cells.

391

392 **Figure 3.** Module scores for individual cells for the top 9 enriched modules (A) and decomposed

393 Z-scores (B) for single-cell gene set enrichment analysis in MAIT data set, using the blood

394 transcription modules (BTM) database. The distribution of module scores suggests

395 heterogeneity among individual cells with respect to different biological processes. Enrichment

396 of modules in stimulated and non-stimulated cells is due to a combination of differences in the

397 discrete (proportion) and continuous (mean conditional expression) of genes in modules. The
398 combined Z-score reflects the enrichment due to differences in the continuous and discrete
399 components.
400
401 **Figure 4.** Gene-gene correlation (Pearson's rho) of model residuals in non-stimlated (A) and
402 stimulated (B) cells, and principal components analysis biplot of model residuals (C) on both
403 populations using the top 50 marginally differentially expressed genes. As marginal changes in
404 the genes attributable to stimulation and CDR have been removed, clustering of subpopulations
405 in (C) indicates co-expression of the indicated genes on a cellular basis.
406
407 **Figure 5.** Module scores (A) and decomposed Z-scores (B) for single-cell gene set enrichment
408 analysis for LPS stimulated cells, mDC data set, using the mouse GO biological process
409 database. The change in single-cell module scores over time for the nine most significantly
410 enriched modules in response to LPS stimulation are shown in A. The *core antiviral*, *peaked*
411 *inflammatory* and *sustained inflammatory* modules are among the top enriched modules,
412 consistent with the original publication. Additionally we identify GO modules *cellular response to*
413 *interferon-beta* and *response to virus*, which behave analogously to the core antiviral and
414 sustained inflammatory modules. No GO analog for the *peaked inflammatory* module was
415 detected. The majority of modules detected exhibit enrichment relative to the 1h time point (thus
416 increasing with time). The "early marcher" cells identified in the original publication are
417 highlighted here with triangles. We show the top 50 most significant modules (B). The combined
418 Z-score summarizes the changes in the discrete and continuous components of expression.
419
420 **Figure 6.** Principal components analysis biplot of model residuals (A) and Gene-gene
421 correlation (Pearson's R) of model residuals (B) by time point for LPS cells, mDC experiment
422 using 20 genes with largest log-fold changes, given significant (FDR q <.01) marginal changes
423 in expression. PC1 is correlated with change over time. The two "early marcher" cells are
424 highlighted by an asterisk at the 1h time-point. Correlation structure in the residuals is
425 increasingly evident over time and can be clearly observed at the 6h time-point compared to the
426 earlier time-points.
427

428 ## METHODS

429
430 ### Data Sets
431 Data for the MAIT study were derived from a single donor who provided written informed
432 consent for immune response exploratory analyses. The study was approved by the
433 relevant institutional review boards.
434
435 ### MAIT cell isolation and stimulation
436 Cryopreserved PBMC were thawed and stained with Aqua Live/Dead Fixable Dead Cell Stain
437 and the following antibodies: CD3, CD8, CD4, CD161, V$\alpha$7.2, CD56 and CD16. CD8$^+$ MAIT
438 cells were sorted as live CD3$^+$CD8$^+$ CD4$^-$CD161$^{hi}$V$\alpha$7.2$^+$ cells and purity was confirmed by post-
439 sort FACS analysis. Sorted MAIT cells were divided into aliquots and immediately processed on

10

440    a C1 Fluidigm machine or treated with a combination of IL-12 (eBioscience), IL-15
441    (eBioscience), and IL-18 (MBL) at 100ng/mL for 24 hours followed by C1 processing.
442
443

444    **C1 processing, Sequencing, and Alignment**
445    After flow sorting, single cells were captured on the Fluidigm$^{TM}$ C1 Single-Cell Auto Prep
446    System (C1), lysed on chip and subjected to reverse transcription and cDNA amplification using
447    the SMARTer® Ultra™ Low Input RNA Kit for C1 System (Clontech).  Sequencing libraries were
448    prepared using the Nextera XT DNA Library Preparation Kit (Illumina) according to C1 protocols
449    (Fluidigm).  Barcoded libraries were pooled and quantified using a Qubit® Fluorometer (Life
450    Technologies). Single-read sequencing of the pooled libraries was carried out either on a
451    HiScanSQ or a HiSeq2500 sequencer (Illumina) with 100-base reads, using TruSeq v3 Cluster
452    and SBS kits (Illumina) with a target depth of >2.5M reads. Sequences were aligned to the
453    UCSC Human genome assembly version 19 and gene expression levels quantified using
454    RSEM[25] and TPM values were loaded into R[26]  for analyses. See supplement for more details
455    on data processing procedures.
456
457

458    **Time-series stimulation of mouse bone-marrow derived dendritic cells (mDC)**
459    Processed RNA-seq data (transcripts-per-million, TPM) were downloaded from GEO under
460    accession number GSE41265. Alignment, pre-processing and filtering steps have been
461    previously described[5]. Low quality cells were filtered as described in Shalek et al[5].
462

463    <u>**Single Cell RNA Seq Hurdle model**</u>
464    We model the log$_2$(TPM+1) expression matrix as a two part generalized regression model. The
465    cell expression rate given a design is modeled using logistic regression and the expression level
466    is modeled as conditionally Gaussian given that they are expressed.
467

468    Given normalized, possibly thresholded (see supplementary material), scRNA-seq expression
469    $Y = [y_{ig}]$, the rate of expression and the level of expression for the expressed cells are modeled
470    conditionally independent for each gene *g*.  Define the indicator $Z = [z_{ig}]$ indicating whether
471    gene *g* is expressed in cell *i*, i.e. $z_{ig} = 0$ if $y_{ig} = 0$ and $z_{ig} = 1$ if $y_{ig} > 0$.  We fit logistic
472    regression models for the discrete variable $Z$ and Gaussian linear model for the continuous
473    variable $(Y \mid Z = 1)$ independently, as follows,
474

$$logit\left(P(Z_{ig} = 1)\right) = X_i \beta_g^D$$
$$\Pr(Y_{ig} = y | Z_{ig} = 1) = N(X_i\beta_g^C, \sigma_g^2)$$

475
476    The regression coefficients of the discrete component are regularized using a Bayesian
477    approach as implemented in the *bayesglm* function of the *arm* R package, which uses weakly
478    informative priors[27] to provide sensible estimates under linear separation (See supplementary
479    material for details). We also perform regularization of the continuous model variance

480    parameter, as described below, which helps increases robustness of gene-level differential
481    expression analysis when a gene is only expressed in a few cells.
482
483    We define the *cellular detection rate* (CDR) as the proportion of genes detected in each cell.
484    The CDR for cell $i$ is:

$$\mathrm{CDR}_i = 1/N \sum_{g=1}^{N} z_{ig}$$

485    An advantage of our approach is that it is straightforward to account for CDR variability by
486    adding the variable as a covariate in the discrete and continuous models (column of the design
487    matrix, $X$, defined above). In the context of our hurdle model, inclusion the CDR covariate can
488    be thought of as the discrete analog of global normalization, and as we show in the examples,
489    this normalization yields more interpretable results and helps decrease background correlation
490    between genes, which is desirable for detecting genuine gene co-expression.
491

492    **Shrinkage of the continuous variance**
493    As the number of expressed cells varies from gene to gene, so does the amount of information
494    available to estimate the residual variance of the gene.  On the other hand, many genes can be
495    expected to have similar variances. To accommodate this feature of the assay, we shrink the
496    gene-specific variances estimates to a global estimate of the variance using an empirical Bayes
497    method.  Let $\tau_g^2$ be the precision (1/variance) for $Y_g|Z_g = 1$ in gene g.  We suppose
498    $\tau_g^2 \sim Gamma(\alpha, \beta)$, find the joint likelihood (across genes) and integrate out the gene-specific
499    inverse variances. Then maximum likelihood is used to estimate $\alpha$ and $\beta$.  Due to conjugacy,
500    these parameters are interpretable providing $2\alpha$ pseud-observations with precision $\beta/\alpha$.  This
501    leads to a simple procedure where the shrunken gene-specific precision is a convex
502    combination of its MLE and the common precision. This approach accounts for the fact that the
503    number of cells expressing a gene varies from gene to gene. Genes with fewer expressed cells
504    end up with proportionally stronger shrinkage, as the ratio of pseudo observations to actual
505    observations is greater. Further details are available in the supplement.
506

507    **Testing for differential expression**
508    Because $Z_g$ and $Y_g$ are defined conditionally independent for each gene, tests with asymptotic
509    $\chi^2$ null distributions, such as the likelihood ratio or Wald tests can be summed and remain
510    asymptotically $\chi^2$, with the degrees of freedom of the component tests added. For the
511    continuous part, we use the shrunken variance estimates derived through our empirical Bayes
512    approach described above. The test results across genes can be combined and adjusted for
513    multiplicity using the false discovery rate (FDR) adjustment[28]. In this paper, we declare a gene
514    differentially expressed if the FDR adjusted p-value is less than 0.01 and the estimated fold-
515    change is greater than 1.5 (on $\log_2$ scale).
516

517    **Gene Set Enrichment Analysis (GSEA)**
518    Our competitive GSEA compares the average model coefficient in the *test* set (gene set of
519    interest) to the average model coefficient in the *null* set (everything else) with a Z-test.  Suppose
520    the genes are sorted so that the first $G_0$ genes are in the null set, and the last $G - G_0$ genes are

521     in the test set.  Then, for example, to test the continuous coefficients in the gene set, the sample

522     means of the coefficients in the test and null sets are calculated, that is, calculate $\hat{\theta} =$

523     $1/(G - G_0)\sum_{g=G_0+1}^{G}\hat{\beta}_g$ and $\hat{\theta}_0 = 1/G_0 \sum_{g=1}^{G_0}\hat{\beta}_g$.  The sampling variance of $\hat{\theta}_0$, in principle, is

524     equal to $1/G_0\left(\sum_{g=1}^{n} Var(\hat{\beta}_g) + 2\sum_{1\le g<h<G_0} Cov(\hat{\beta}_g, \hat{\beta}_h)\right)$, and similarly for $\hat{\theta}$.

525     Given this sampling variance, a Z test can be formed by comparing $Z = \dfrac{\hat{\theta}-\hat{\theta}_0}{\sqrt{\widehat{Var}(\hat{\theta})+\widehat{Var}(\hat{\theta}_0)}}$.

526

527     We estimate $Var(\hat{\beta}_g)$ and $Cov(\hat{\beta}_g, \hat{\beta}_h)$ via bootstrap, to avoid relying on asymptotic

528     approximations. In practice, we find only a few (<100) bootstrap replicates are necessary to

529     provide stable variance-covariance estimates, however even this modest requirement can be

530     relaxed for exploratory analysis by assuming independence across genes and using model-

531     based (asymptotic) estimates.

532

533     Z scores are formed and calculated equivalently for the logistic regression coefficients. GSEA

534     tests are done separately on the two components of the hurdle model and the results from the

535     two components are combined using the Stouffer's method[29], which favors consensus in the two

536     components[30] (see supplement for details).  The approach is similar to that used by CAMERA[15]

537     for bulk experiments in its accounting for inter-gene correlation that is known to inflate the false

538     significance (type-I error) in permutation-based GSEA protocols[15], although it differs in that it

539     uses the sampling variance of each model coefficient to find the variance of the average

540     coefficient, whereas CAMERA uses the empirical variance of the model coefficients. In our

541     analyses we used the Emory blood transcriptional modules[19] as well as mouse gene ontology

542     annotations available from the Mouse Genome Informatics web site[32].

543

544     **GO Enrichment Analysis**

545     Testing for enriched Gene Ontology terms based on list of genes was performed with the

546     GOrilla online tool using the approach of comparing an unranked target list against a

547     background list[33].

548

549     **Residual Analysis**

550     The hurdle model, in general, provides two residuals: one for the discrete component and one

551     for the continuous component. Standardized deviance residuals are calculated for the discrete

552     and continuous component separately, and then we combine the residuals by averaging them.

553     If a cell is unexpressed, then its residual is missing and it is omitted from the average.  See the

554     supplement for details.

555

556     **Module Scores**

557     In order to assess the degree to which each cell exhibits enrichment for each gene module, we

558     use quantities available through our model to define module "scores", which are defined as the

559     observed expression corrected for CDR effect, analogous to those defined by Shalek et al[5]. The

560     score $s_{ij}$ for cell $i$ and gene $j$ is defined as the observed expression corrected for the CDR

561     effect: $s_{ij} = y_{ij} - \tilde{y}_{ij}$ where $\tilde{y}_{ij}$ is the predicted effect from the fitted model that excludes thre

562     treatment effects of interest. This can be interpreted as correcting the observed expression of

13

563    gene $j$ in cell $i$ by subtracting the conditional expectation of nuisance effects. In our two part

564    model, $\tilde{y}_{ij} = \hat{z}_{ij}\hat{y}_{ij}$ where $\hat{z}_{ij}$ and $\hat{y}_{ij}$ are the predicted values from the discrete and continuous

565    components of our hurdle model.

566    A gene module score for cell I is the average of the scores for the genes contained in the

567    module, i.e. $\sum_{\{j \in module\}} s_{ij} / |module|$

568

569    **Author Contributions:**

570    GF, AM, MY and RG developed the statistical methods, and wrote the manuscript. AM, MY,

571    and GF wrote the R package and performed data analysis. CKS, HWM, and MP designed and

572    performed the MAIT cell experiments and contributed to data interpretation and provided

573    manuscript feedback. VG and PL coordinated the collection of the single cell sequencing data

574    and contributed manuscript preparation and feedback and data interpretation. AKS contributed

575    the mDC data and contributed manuscript feedback and to data interpretation. JD contributed to

576    data analysis of the single-cell expression data. MJM contributed samples and to study design.
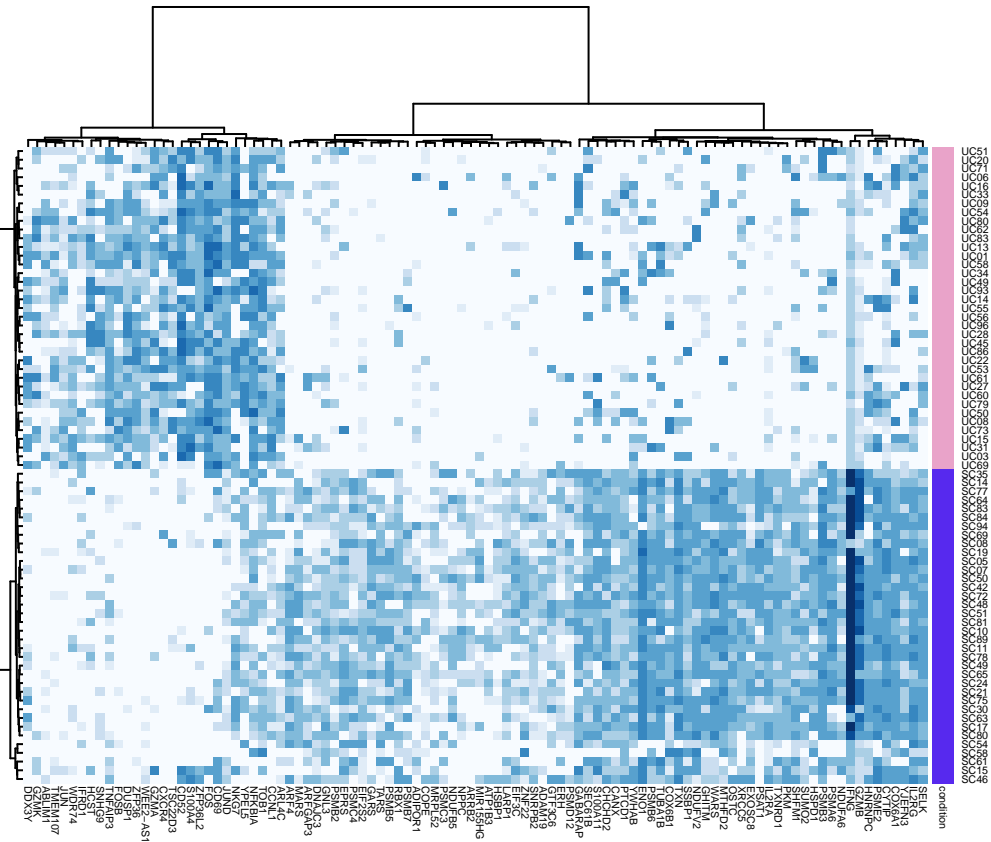
577

583

584    1.    Elowitz, M. B., Levine, A. J., Siggia, E. D. & Swain, P. S. Stochastic gene expression in a single cell.
585      *Science* **297,** 1183–1186 (2002).

586    2.    Raj, A., van den Bogaard, P., Rifkin, S. A., van Oudenaarden, A. & Tyagi, S. Imaging individual mRNA
587      molecules using multiple singly labeled probes. *Nature methods* **5,** 877–879 (2008).

588    3.    Sanchez, A. & Golding, I. Genetic determinants and cellular constraints in noisy gene expression. *Science*
589      **342,** 1188–1193 (2013).

590    4.    McDavid, A. *et al.* Data exploration, quality control and testing in single-cell qPCR-based gene expression
591      experiments. *Bioinformatics* **29,** 461–467 (2013).

592    5.    Shalek, A. K. *et al.* Single-cell RNA-seq reveals dynamic paracrine control of cellular variation. *Nature* **510,**
593      263–269 (2014).

594    6.    Kharchenko, P. V., Silberstein, L. & Scadden, D. T. Bayesian approach to single-cell differential expression
595      analysis. *Nature methods* 1–5 (2014). doi:10.1038/nmeth.2967

596    7.    Kaufmann, B. B. & van Oudenaarden, A. Stochastic gene expression: from single molecules to the
597      proteome. *Current opinion in genetics \& development* **17,** 107–112 (2007).

598    8.    Marinov, G. K. *et al.* From single-cell to cell-pool transcriptomes: stochasticity in gene expression and RNA
599      splicing. *Genome research* **24,** 496–510 (2014).

600    9.    Marguerat, S. *et al.* Quantitative analysis of fission yeast transcriptomes and proteomes in proliferating and
601      quiescent cells. *Cell* **151,** 671–683 (2012).

602    10.    Single Mammalian Cells Compensate for Differences in Cellular Volume and DNA Copy Number through
603      Independent Global Transcriptional Mechanisms. **58,** 339–352 (2015).

604    11.    Voom: precision weights unlock linear model analysis tools for RNA-seq read counts. **15,** R29 (2014).

605    12.    Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for differential
606      expression analysis of digital gene expression data. *Bioinformatics* **26,** 139–140 (2010).

607    13.    Anders, S. & Huber, W. Differential expression analysis for sequence count data. *Genome biology* **11,** R106
608      (2010).

609    14.    Smyth, G. K. Linear models and empirical bayes methods for assessing differential expression in microarray
610      experiments. **3,** Article3 (2004).

611    15.    Wu, D. & Smyth, G. K. Camera: a competitive gene set test accounting for inter-gene correlation. *Nucleic*
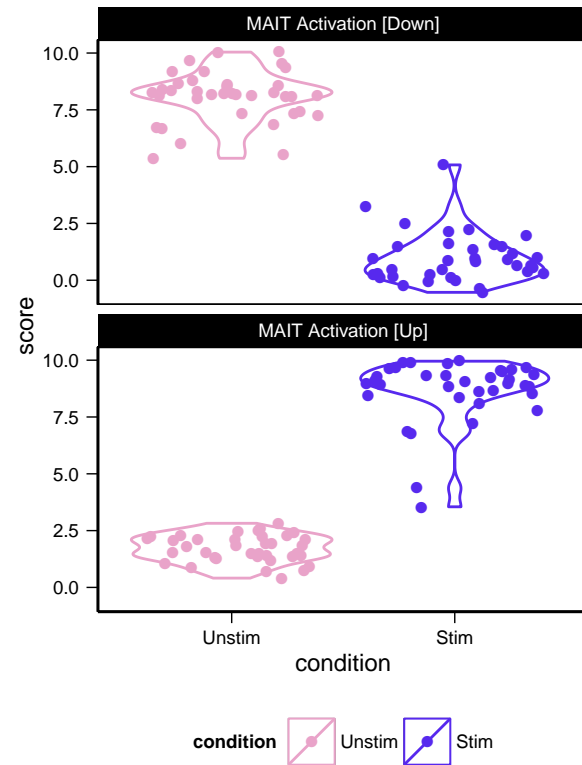
612          *Acids Research* **40,** e133 (2012).
613   16.   Chu, T. *et al.* Bystander-activated memory CD8 T cells control early pathogen load in an innate-like,
614          NKG2D-dependent manner. *Cell Rep* **3,** 701–708 (2013).
615   17.   Tyznik, A. J., Verma, S., Wang, Q., Kronenberg, M. & Benedict, C. A. Distinct requirements for activation
616          of NKT and NK cells during viral infection. *Journal of immunology (Baltimore, Md. : 1950)* **192,** 3676–
617          3685 (2014).
618   18.   Smeltz, R. B. Profound enhancement of the IL-12/IL-18 pathway of IFN-gamma secretion in human CD8+
619          memory T cell subsets via IL-15. *Journal of immunology (Baltimore, Md. : 1950)* **178,** 4786–4792 (2007).
620   19.   Li, S. *et al.* Molecular signatures of antibody responses derived from a systems biology study of five human
621          vaccines. *Nat Immunol* **15,** 195–204 (2014).
622   20.   Interferon-gamma modulates the lipopolysaccharide-induced expression of AP-1 and NF-kappa B at the
623          mRNA and protein level in human monocytes. **24,** 228–235 (1996).
624   21.   Egr2 induced during DC development acts as an intrinsic negative regulator of DC immunogenicity. **43,**
625          2484–2496 (2013).
626   22.   Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden
627          subpopulations of cells. **33,** 155–160 (2015).
628   23.   Brennecke, P. *et al.* Accounting for technical noise in single-cell RNA-seq experiments. *Nature methods*
629          **advance on,** (2013).
630   24.   Normalization of RNA-seq data using factor analysis of control genes or samples. **32,** 896–902 (2014).
631   25.   Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with or without a
632          reference genome. *BMC Bioinformatics* **12,** 323 (2011).
633   26.   Gentleman, R. C. *et al.* Bioconductor: open software development for computational biology and
634          bioinformatics. *Genome biology* **5,** R80 (2004).
635   27.   Gelman, A., Jakulin, A., Pittau, M. G. & Su, Y.-S. A Weakly Informative Default Prior Distribution for
636          Logistic and Other Regression Models. *The annals of applied statistics* **2,** 1360–1383 (2008).
637   28.   Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to
638          multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)* **57,** 289–300 (1995).
639   29.   Annotated Bibliography of Some Papers on Combining Significances or p-values. **physics.data-an,** (2007).
640   30.   Combining independent p values: extensions of the Stouffer and binomial methods. **5,** 496–515 (2000).
641   31.   Molecular signatures database (MSigDB) 3.0. **27,** 1739–1740 (2011).
642   32.   Blake, J. A. *et al.* The Mouse Genome Database (MGD): premier model organism resource for mammalian
643          genomics and genetics. *Nucleic Acids Research* **39,** D842–8 (2011).
644   33.   Eden, E., Navon, R., Steinfeld, I., Lipson, D. & Yakhini, Z. GOrilla: a tool for discovery and visualization of
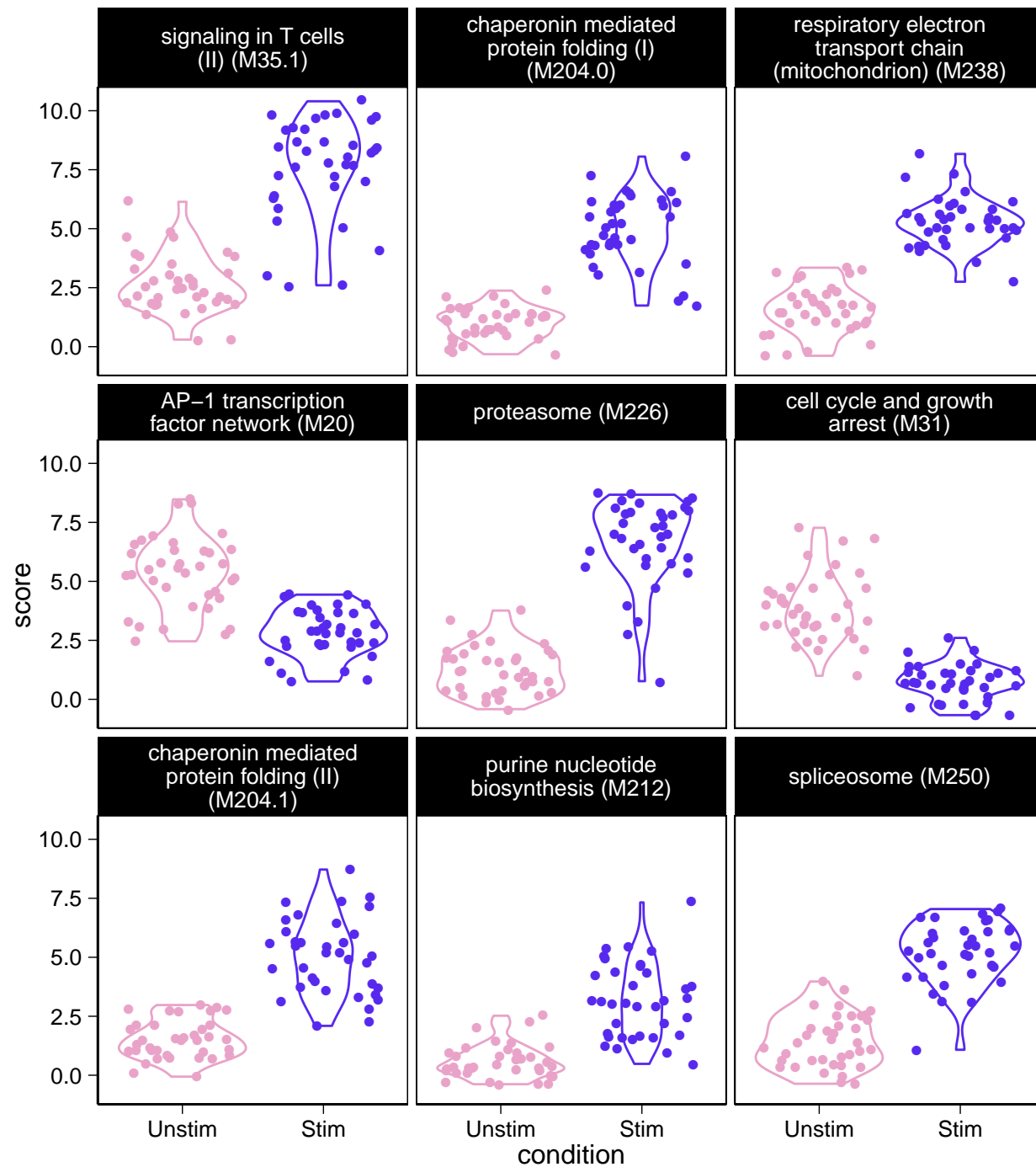645          enriched GO terms in ranked gene lists. *BMC Bioinformatics* **10,** 48 (2009).
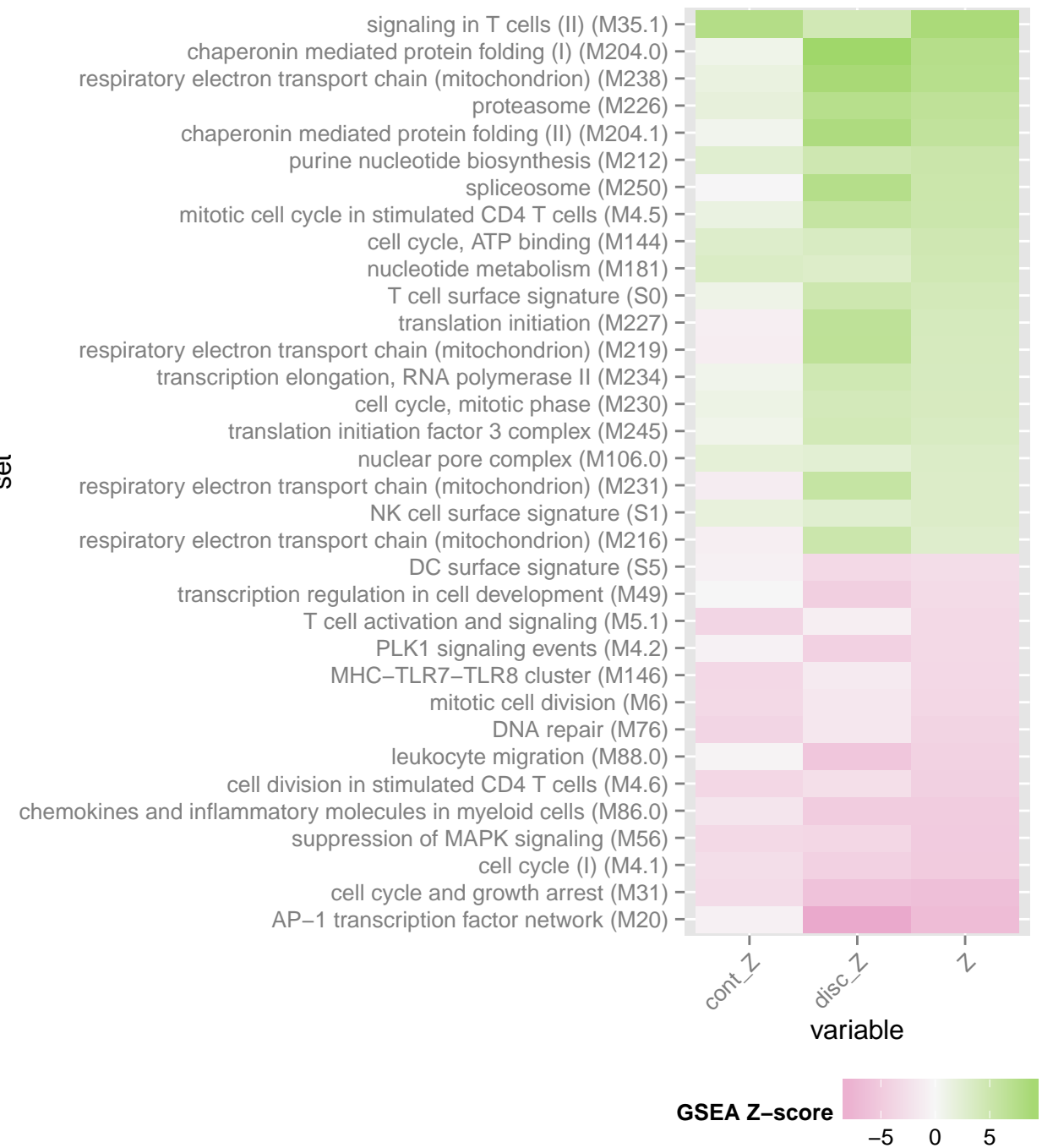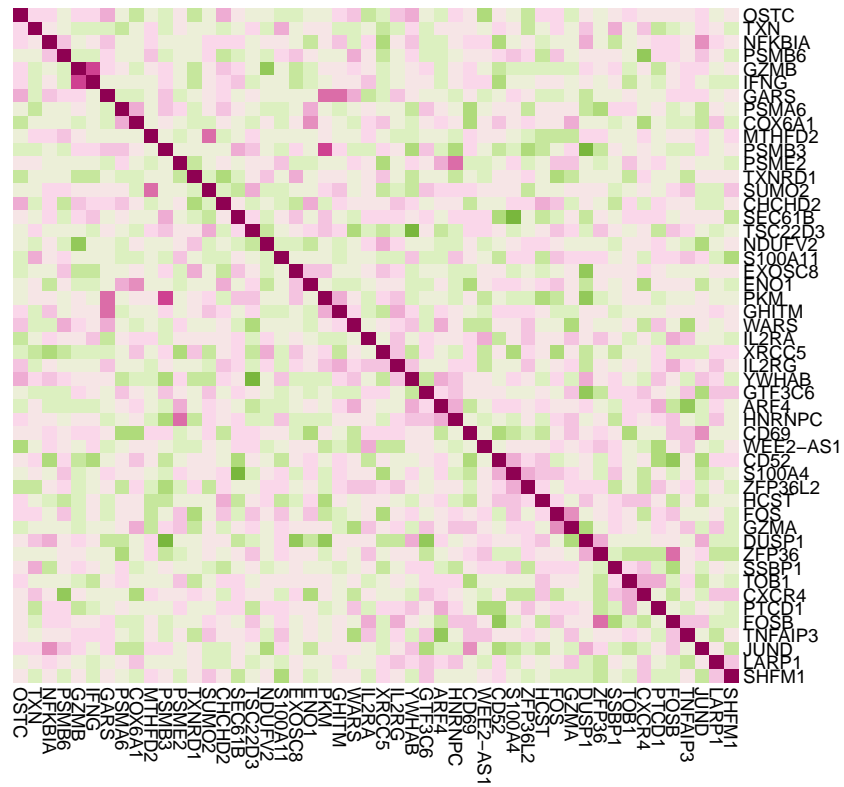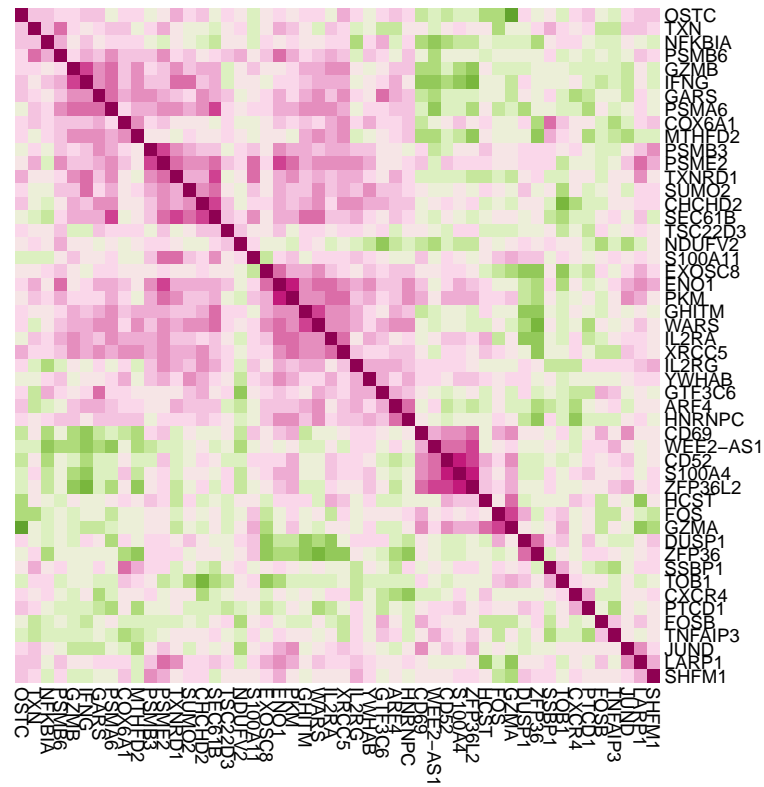646

15

A

B

A

B

C

**A**



**B**