

# Differential Expression Analysis for RNAseq using Poisson Mixed Models

Shiquan Sun<sup>1,2</sup>, Michelle Hood<sup>2</sup>, Laura Scott<sup>2,3</sup>, Qinke Peng<sup>1</sup>, Sayan Mukherjee<sup>4</sup>,  
Jenny Tung<sup>5,6</sup>, Xiang Zhou<sup>2,3,\*</sup>

1. Systems Engineering Institute, Xi'an Jiaotong University, Xi'an, Shaanxi  
710049, P.R.China

2. Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, USA

3. Center for Statistical Genetics, University of Michigan, Ann Arbor, MI 48109,  
USA

4. Departments of Statistical Science, Mathematics, and Computer Science,  
Duke University, Durham, NC 27708, USA

5. Departments of Evolutionary Anthropology and Biology, Duke University,  
Durham, NC 27708, USA

6. Duke University Population Research Institute, Duke University, Durham, NC  
27708, USA

\* Correspondence to: XZ ([xzhousph@umich.edu](mailto:xzhousph@umich.edu))

21 **Abstract**

22 Identifying differentially expressed (DE) genes from RNA sequencing (RNAseq)  
23 studies is among the most common analyses in genomics. However, RNAseq  
24 DE analysis presents several statistical and computational challenges, including  
25 over-dispersed read counts and, in some settings, sample non-independence.  
26 Previous count-based methods rely on simple hierarchical Poisson models (e.g.,  
27 negative binomial) to model independent over-dispersion, but do not account for  
28 sample non-independence due to relatedness, population structure and/or  
29 hidden confounders. Here, we present a Poisson mixed model with two random  
30 effects terms that account for both independent over-dispersion and sample non-  
31 independence. We also develop a scalable sampling-based inference algorithm  
32 using a latent variable representation of the Poisson distribution. With simulations,  
33 we show that our method properly controls for type I error and is generally more  
34 powerful than other widely used approaches, except in small samples ( $n < 15$ ) with  
35 other unfavorable properties (e.g., small effect sizes). We also apply our method  
36 to three real data sets that contain related individuals, population stratification, or  
37 hidden confounders. Our results show that our method increases power in all  
38 three data compared to other approaches, though the power gain is smallest in  
39 the smallest sample ( $n=6$ ). Our method is implemented in MACAU, freely  
40 available at [www.xzlab.org/software.html](http://www.xzlab.org/software.html).

41

## 42 Introduction

43 RNA sequencing (RNAseq) has emerged as a powerful tool for transcriptome  
44 analysis, thanks to its many advantages over previous microarray techniques (1-  
45 3). Compared with microarrays, RNAseq has increased dynamic range, does not  
46 rely on *a priori*-chosen probes, and can thus identify previously unknown  
47 transcripts and isoforms. It also yields allelic-specific expression estimates and  
48 genotype information inside expressed transcripts as a useful by-product (4-7).  
49 Because of these desirable features, RNAseq has been widely applied in many  
50 areas of genomics and is currently the gold standard method for genome-wide  
51 gene expression profiling.

52 One of the most common analyses of RNAseq data involves identification of  
53 differentially expressed (DE) genes. Identifying DE genes that are influenced by  
54 predictors of interest -- such as disease status, risk factors, environmental  
55 covariates, or genotype -- is an important first step towards understanding the  
56 molecular basis of disease susceptibility as well as the genetic and  
57 environmental basis of gene expression variation. Progress towards this goal  
58 requires statistical methods that can handle the complexities of the increasingly  
59 large and structurally complex RNAseq data sets that are now being collected  
60 from population and family studies (8,9). Indeed, even in classical treatment-  
61 control comparisons, the importance of larger sample sizes for maximizing power  
62 and reproducibility is increasingly well appreciated (10,11). However, identifying  
63 DE genes from such studies presents several key statistical and computational  
64 challenges, including accounting for ambiguously mapped reads (12), modeling  
65 uneven distribution of reads inside a transcript (13), and inferring transcript  
66 isoforms (14).

67 A fundamental challenge shared by all DE analyses in RNAseq, though, is  
68 accounting for the count nature of the data (3,15,16). In most RNAseq studies,  
69 the number of reads mapped to a given gene or isoform (following appropriate  
70 data processing and normalization) is often used as a simple and intuitive  
71 estimate of its expression level (13,14,17). As a result, RNAseq data display an  
72 appreciable dependence between the mean and variance of estimated gene  
73 expression levels: highly expressed genes tend to have high read counts and  
74 subsequently high between-sample variance, and vice versa (15,18). To account  
75 for the count nature of the data and the resulting mean-variance dependence,  
76 most statistical methods for DE analysis model RNAseq data using discrete

77 distributions. For example, early studies showed that gene expression variation  
78 across technical replicates can be accurately described by a Poisson distribution  
79 (19-21). More recent methods also take into account over-dispersion across  
80 biological replicates (22,23) by replacing Poisson models with negative binomial  
81 models (15,16,24-28) or other related approaches (18,29-32). While non-count  
82 based methods are also commonly used (primarily relying on transformation of  
83 the count data to more flexible, continuous distributions (33,34)), recent  
84 comparisons have highlighted the benefits of modeling RNAseq data using the  
85 original counts and accounting for the resulting mean-variance dependence (35-  
86 38), consistent with observations from many count data analyses in other  
87 statistical settings (39). Indeed, accurate modeling of mean-variance  
88 dependence is one of the keys to enable powerful DE analysis with RNAseq,  
89 especially in the presence of large sequencing depth variation across samples  
90 (25,33,40).

91 A second important feature of many RNAseq data sets, which has been largely  
92 overlooked in DE analysis thus far, is that samples often are not independent.  
93 Sample non-independence can result from individual relatedness, population  
94 stratification, or hidden confounding factors. For example, it is well known that  
95 gene expression levels are heritable. In humans, the narrow-sense heritability of  
96 gene expression levels averages from 15%-34% in peripheral blood (41-45) and  
97 is about 23% in adipose tissue (41), with a maximum heritability in both tissues  
98 as high as 90% (41,42). Similarly, in baboons, gene expression levels are about  
99 28% heritable in the peripheral blood (7). Some of these effects are attributable  
100 to nearby, putatively *cis*-acting genetic variants: indeed, recent studies have  
101 shown that the expression levels of almost all genes are influenced by *cis*-eQTLs  
102 and/or display allelic specific expression (ASE) (3,7,46-48). However, the  
103 majority of heritability is often explained by distal genetic variants (i.e., *trans*-  
104 QTLs, which account for 63%-84% of heritability in humans (41) and baboons  
105 (7)). Because gene expression levels are heritable, they will covary with kinship  
106 or population structure. Besides kinship or population structure, hidden  
107 confounding factors, commonly encountered in sequencing studies (49-52), can  
108 also induce similarity in gene expression levels across many genes even when  
109 individuals are unrelated (53-57). Failure to account for this gene expression  
110 covariance due to sample non-independence could lead to spurious associations  
111 or reduced power to detect true DE effects. This phenomenon has been  
112 extensively documented in genome-wide association studies (9,58,59) and more

113 recently, in bisulfite sequencing studies (60), but is less explored in RNAseq  
114 studies. In particular, none of the currently available count-based methods for  
115 identifying DE genes in RNAseq can appropriately control for sample non-  
116 independence. Consequently, even though count-based methods have been  
117 shown to be more powerful, recent RNAseq studies have turned to linear mixed  
118 models, which are specifically designed for quantitative traits, to deal with the  
119 confounding effects of kinship, population structure, or hidden confounders  
120 (7,42,61).

121 Here, we present a Poisson mixed model (PMM) that can explicitly model both  
122 over-dispersed count data and sample non-independence in RNAseq data for  
123 effective DE analysis. To make our model scalable to large data sets, we also  
124 develop an accompanying efficient inference algorithm based on an auxiliary  
125 variable representation of the Poisson model (62-64) and recent advances in  
126 mixed model methods (9,59,65). We refer to the combination of the statistical  
127 method and the computational algorithm developed here as MACAU (Mixed  
128 model Association for Count data via data AUgmentation), which effectively  
129 extends our previous method of the same name on the simpler binomial model  
130 (60) to the more difficult Poisson model. MACAU works directly on RNAseq count  
131 data and introduces two random effects terms to both control for sample non-  
132 independence and account for additional independent over-dispersion. As a  
133 result, MACAU properly controls for type I error in the presence of sample non-  
134 independence and, in a variety of settings, is more powerful for identifying DE  
135 genes than other commonly used methods. We illustrate the benefits of MACAU  
136 with extensive simulations and real data applications to three RNAseq studies.

137

## 138 **Methods and Materials**

### 139 **Methods for Comparison**

140 We compared the performance of seven different methods in the main text: (1)  
141 our Poisson mixed model implemented in the MACAU software package (60); (2)  
142 the linear model implemented in the *lm* function in R; (3) the linear mixed model  
143 implemented in the GEMMA software package (9,59,66); (4) the Poisson model  
144 implemented in the *glm* function in R; (5) the negative binomial model  
145 implemented in the *glm.nb* function in R; (6) edgeR implemented in the *edgeR*  
146 package in R (25); (7) DESeq2 implemented in the *DESeq2* package in R (24).  
147 All methods were used with default settings. The performance of each method in  
148 simulations was evaluated using the area under the curve (AUC) function  
149 implemented in the *pROC* package in R (67), a widely used benchmark for  
150 RNAseq method comparisons (68).

151 Both the linear model and the linear mixed model require quantitative phenotypes.  
152 Here, we considered six different transformations of count data to quantitative  
153 values, taking advantage of several methods proposed to normalize RNAseq  
154 data (e.g., (12-14,17,22,33,69) ): (1) quantile normalization (TRCQ), where we  
155 first divided the number of reads mapped to a given gene by the total number of  
156 read counts for each individual, and then for each gene, quantile normalized the  
157 resulting proportions across individuals to a standard normal distribution (7); (2)  
158 total read count normalization (TRC), where we divided the number of reads  
159 mapped to a given gene by the total number of read counts for each individual  
160 (i.e. CPM, counts per million; without further transformation to a standard normal  
161 within genes: (25)); (3) upper quantile normalization (UQ), where we divided the  
162 number of reads mapped to a given gene by the upper quantile (75-th percentile)  
163 of all genes for each individual (70); (4) relative log expression normalization  
164 (RLE) (15); (5) the trimmed mean of M-values (TMM) method (40) where we  
165 divided the number of reads mapped to a given gene by the normalization factor  
166 output from TMM; and (6) VOOM normalization (33). Simulations in a  
167 supplementary figure showed that TRCQ, VOOM and TRC worked better than  
168 the other three methods, with TRCQ showing a small advantage. Therefore, we  
169 report results using TRCQ throughout the text.

170

### 171 **Simulations**

172 To make our simulations as realistic as possible, we simulated the gene  
173 expression count data based on parameters inferred from a real baboon data set  
174 that contains 63 samples (see the next section for a detailed description of the  
175 data). We varied the sample size ( $n$ ) in the simulations ( $n = 6, 10, 14, 63, 100,$   
176  $200, 500, 800,$  or  $1000$ ). For  $n = 63$ , we used the baboon relatedness matrix  $K$   
177 (7). For sample simulations with  $n > 63$ , we constructed a new relatedness matrix

178  $K$  by filling in its off-diagonal elements with randomly drawn off-diagonal  
179 elements from the baboon relatedness matrix following (60). For sample  
180 simulations with  $n < 63$ , we constructed a new relatedness matrix  $K$  by randomly  
181 sub-sampling individuals from the baboon relatedness matrix. In cases where the  
182 resulting  $K$  was not positive definite, we used the *nearPD* function in R to find the  
183 closest positive definite matrix as the final  $K$ . In most cases, we simulated the  
184 total read count  $N_i$  for each individual from a discrete uniform distribution with a  
185 minimum (=1,770,083) and a maximum (=9,675,989) total read count (i.e.  
186 summation of read counts across all genes) equal to the minimum and maximum  
187 total read counts from the baboon data. We scaled the total read counts to  
188 ensure that the coefficient of variation was small (CV = 0.3), moderate (CV = 0.6)  
189 or high (CV = 0.9) across individuals (i.e.  $N_{new} = \bar{N} + (N - \bar{N}) CV sd(N) / \bar{N}$ ), and  
190 then discretized them. In the special case where CV = 0.3 and  $n = 63$ , we directly  
191 used the observed total read counts per individual  $i$  ( $N_i$ ) from the baboon data  
192 (which has a CV = 0.33).

193 We then repeatedly simulated a continuous predictor variable  $x$  from a standard  
194 normal distribution (without regard to the pedigree structure). We estimated the  
195 heritability of the continuous predictor using GEMMA, and retained  $x$  if the  
196 heritability ( $h_x^2$ ) estimate (with  $\pm 0.01$  tolerance) was 0, 0.4 or 0.8, representing no,  
197 moderate and highly heritable predictors. Using this procedure, approximately 30  
198 percent of  $x$  values generated were retained, with different retention percentages  
199 for different heritability values.

200 Based on the simulated sample size, total read counts and continuous predictor  
201 variable, we simulated gene expression values using the following procedure.  
202 For the expression of each gene in turn, we simulated the genetic random effects  
203  $g$  from a multivariate normal distribution with covariance  $K$ . We simulated the  
204 environmental random effects  $e$  based on independent normal distributions. We  
205 scaled the two sets of random effects to ensure a fixed value of heritability  
206 ( $h^2 = \frac{V(g)}{V(g)+V(e)}$  0 or 0.3 or 0.6) and a fixed value of over-dispersion variance  
207 ( $\sigma^2 = V(g) + V(e) = 0.1, 0.25$  or  $0.4$ , close to the lower, median and upper  
208 quantiles of the over-dispersion variance inferred from the baboon data,  
209 respectively), where the function  $V(\cdot)$  denotes the sample variance. We then  
210 generated the effect size  $\beta$  of the predictor variable on gene expression. The  
211 effect size was either 0 (for non-DE genes) or generated to explain a certain  
212 percentage of variance in  $\log(\lambda)$  (i.e.  $PVE = \frac{V(X\beta)}{V(X\beta)+\sigma^2}$ ; for DE genes). PVE values  
213 were 15%, 20%, 25%, 30% or 35% to represent different effect sizes. The  
214 predictor effects  $X\beta$ , genetic effects  $g$ , environmental effects  $e$ , and an intercept  
215 ( $= \log(\frac{100}{\bar{N}})$ ) to ensure that the expected simulated count is 100) were then  
216 summed together to yield the latent variable  $\log(\lambda) = \mu + X\beta + g + e$ . Note that  
217  $h^2$  does not include the contribution of  $X\beta$ , which in many cases represent non-



218 genetic effects. Finally, the read counts were simulated based on a Poisson  
219 distribution with rate determined by the total read counts and the latent variable  $\lambda$ ,  
220 or  $y_i \sim Poi(N_i \lambda_i)$  for the  $i$ 'th individual.

221 With the above procedure, we first simulated data for  $n = 63$ ,  $CV = 0.3$ ,  $h_x^2 = 0$ ,  
222  $PVE = 0.25$ ,  $h^2 = 0.3$  and  $\sigma^2 = 0.25$ . We then varied one parameter at a time to  
223 generate different scenarios for comparison. In each scenario, conditional on the  
224 sample size, total read counts, and continuous predictor variable, we performed  
225 10 simulation replicates, where “replication” is at the level described in the  
226 paragraph above. Each replicate consisted of 10,000 genes. For examining type  
227 I error control, all 10,000 genes were non-DE. For the power comparison, 1,000  
228 genes were DE while 9,000 were non-DE.

229

### 230 **RNAseq Data Sets**

231 We considered three published RNAseq data sets in this study, which include  
232 small ( $n < 15$ ), medium ( $15 \leq n \leq 100$ ), and large ( $n > 100$ ) sample sizes (based on  
233 current RNAseq sample sizes in the literature).

234 The first RNAseq data set was collected from blood samples of yellow baboons  
235 (7) from the Amboseli ecosystem of southern Kenya as part of the Amboseli  
236 Baboon Research Project (ABRP) (71). The data are publicly available on GEO  
237 with accession number GSE63788. Read counts were measured on 63 baboons  
238 and 12,018 genes after stringent quality control as in (7). As in (7), we computed  
239 pairwise relatedness values from previously collected microsatellite data (72,73)  
240 using the software COANCESTRY (74). The data contains related individuals: 16  
241 pairs of individuals have a kinship coefficient exceeding 1/8 and 48 pairs exceed  
242 1/16. We obtained sex information for each individual from GEO. Sex differences  
243 in health and survival are major topics of interest in medicine, epidemiology, and  
244 evolutionary biology (72,75). Therefore, we used this data set to identify sex-  
245 related gene expression variation. In the analysis, we included the top 5  
246 expression PCs as covariates to control for potential batch effects following the  
247 original study (7).

248 The second RNAseq data set was collected from skeletal muscle samples of  
249 Finnish individuals (61) as part of the FUSION project (76,77). The data are  
250 publicly available in dbGaP with accession code phs001068.v1.p1. Among the  
251 271 individuals in the original study, we selected 267 individuals who have both  
252 genotypes and gene expression measurements. Read counts were obtained on  
253 these 267 individuals and 21,753 genes following the same stringent quality  
254 control as in the FUSION study. For genotypes, we excluded SNPs with minor  
255 allele frequency (MAF)  $< 0.05$  and Hardy-Weinberg equilibrium  $p$ -value  $< 10^{-6}$ .  
256 We used the remaining 5,696,681 SNPs to compute the relatedness matrix using  
257 GEMMA. The data contains remotely related individuals (3 pairs of individuals



258 have a kinship coefficient exceeding 1/32 and 6 pairs exceed 1/64) and is  
259 stratified by the municipality from which samples were collected. Two predictors  
260 from the data were available to us: the oral glucose tolerance test (OGTT) which  
261 classifies  $n = 162$  individuals as either T2D patient ( $n = 66$ ) or normal glucose  
262 tolerance (NGT; i.e., control,  $n = 96$ ); and a T2D-related quantitative trait --  
263 fasting glucose levels (GL) -- measured on all  $n = 267$  individuals. We used these  
264 data to identify genes whose expression level is associated with either T2D or GL.  
265 In the analysis, we included age, sex and batch labels as covariates following the  
266 original study (61).

267 The third RNAseq data set was collected from lymphoblastoid cell lines (LCLs)  
268 derived from 69 unrelated Nigerian individuals (YRI) (3). The data are publicly  
269 available on GEO with accession number GSE19480. Following the original  
270 study (3), we aligned reads to the human reference genome (version hg19) using  
271 BWA (78). We counted the number of reads mapped to each gene on either  
272 autosomes or the X chromosome using Ensembl gene annotation information  
273 obtained from the UCSC genome browser. We then filtered out lowly expressed  
274 genes with zero counts in over 90% of individuals. In total, we obtained gene  
275 expression measurements on 13,319 genes. Sex is the only phenotype available  
276 in the data and we used sex as the predictor variable to identify sex-associated  
277 genes. To demonstrate the efficacy of MACAU in small samples, we randomly  
278 subsampled individuals from the data to create small data sets with either  $n = 6$   
279 (3 males and 3 females) or  $n = 10$  (5 males and 5 females), or  $n = 14$  individuals  
280 (7 males and 7 females). For each sample size  $n$ , we performed 20 replicates of  
281 subsampling and we evaluated method performance by averaging across these  
282 replicates. In each replicate, following previous studies (53-57), we used the  
283 gene expression covariance matrix as  $K$  (i.e.  $K = XX^T/p$ , where  $X$  is the  
284 normalized gene expression matrix and  $p$  is the number of genes) and applied  
285 MACAU to identify sex-associated genes. Note that the gene expression  
286 covariance matrix  $K$  contains information on sample non-independence caused  
287 by hidden confounding factors (53-57), and by incorporating  $K$ , MACAU can be  
288 used to control for hidden confounding factors that are commonly observed in  
289 sequencing data sets (49-52).

290 For each of these RNAseq data sets and each trait, we used a constrained  
291 permutation procedure to estimate the empirical false discovery rate (FDR) of a  
292 given analytical method. In the constrained permutation procedure, we permuted  
293 the predictor across individuals, estimated the heritability of the permuted  
294 predictor, and retained the permutation only if the permuted predictor had a  
295 heritability estimate ( $h_x^2$ ) similar to the original predictor with  $\pm 0.01$  tolerance (for  
296 the original predictors:  $h_x^2 = 0.0002$  for sex in the baboon data;  $h_x^2 = 0.0121$  for  
297 T2D and  $h_x^2 = 0.4023$  for GL in the FUSION data;  $h_x^2$  are all close to zero with  
298 small variations depending on the sub-sample size in the YRI data). We then  
299 analyzed all genes using the permuted predictor. We repeated the constrained

300 permutation procedure and analysis 10 times, and combined the  $p$ -values from  
301 these 10 constrained permutations. We used this set of  $p$ -values as a null  
302 distribution from which to estimate the empirical false discovery rate (FDR) for  
303 any given  $p$ -value threshold (60). This constrained procedure thus differs from  
304 the usual unconstrained permutation procedure (every permutation retained) (79)  
305 in that it constrains the permuted predictor to have the same  $h_x^2$  as the original  
306 predictor. We chose to use the constrained permutation procedure here because  
307 the unconstrained procedure is invalid under the mixed model assumption: the  
308 subjects are not exchangeable in the presence of sample non-independence  
309 (individual relatedness, population structure, or hidden confounders) (79,80). To  
310 validate our constrained permutation procedure and test its effectiveness in  
311 estimating FDR, we performed a simulation with 1,000 DE genes and 9,000 non-  
312 DE genes as described above. We considered three predictor variables  $x$  with  
313 different heritability:  $h_x^2 = 0$ ,  $h_x^2 = 0.4$ , and  $h_x^2 = 0.8$ . For each predictor variable  
314 and each  $p$ -value threshold, we computed the true FDR and then estimated the  
315 FDR based on either the constrained or unconstrained permutation procedures.  
316 The simulation results presented in a supplementary figure demonstrate that the  
317 constrained permutation procedure provides a much more accurate estimate of  
318 the true FDR while the unconstrained permutation procedure often under-  
319 estimates the true FDR. Therefore, we applied the constrained permutation  
320 procedure for all real data analysis.

321 Finally, we investigated whether the methods we compared were sensitive to  
322 outliers (31,81,82) in the first two data sets. To examine outlier sensitivity, we first  
323 identified genes with potential outliers using BBSeq (18). In total, we identified 8  
324 genes with potential outliers in the baboon data, 130 genes with potential outliers  
325 in the FUSION data ( $n = 267$ ) and 43 genes with potential outliers in the subset  
326 of the FUSION data for which we had T2D diagnoses ( $n = 162$ ). We counted the  
327 number of genes with potential outliers in the top 1,000 genes with strong DE  
328 association evidence. In the baboon data, 4 genes with potential outliers are in  
329 the top 1,000 genes with the strongest sex association determined by various  
330 methods: 2 of them by the negative binomial model, 3 of them by the Poisson  
331 model, but 0 of them by MACAU, linear model, or GEMMA. In the FUSION data,  
332 for T2D analysis, 9 genes with potential outliers are in the top 1,000 genes with  
333 the strongest T2D association determined by various methods: 1 by MACAU, 3  
334 by negative binomial, 6 by Poisson, 1 by linear, and 1 by GEMMA. For GL  
335 analysis, 15 genes with potential outliers are in the top 1,000 genes with the  
336 strongest GL association determined by various methods: 2 by MACAU, 7 by  
337 negative binomial, 9 by Poisson, 3 by linear, and 3 by GEMMA. All outliers are  
338 presented in supplementary figures. Therefore, the influence of outliers on DE  
339 analysis is small in the real data.

340

## 341 Results

### 342 MACAU Overview

343 Here, we provide a brief overview of the Poisson mixed model (PMM); more  
344 details are available in the Supplementary Material. To identify DE genes with  
345 RNAseq data, we examine one gene at a time. For each gene, we model the  
346 read counts with a Poisson distribution

$$y_i \sim \text{Poi}(N_i \lambda_i), \quad i = 1, 2, \dots, n,$$

347 where for the  $i$ 'th individual,  $y_i$  is the number of reads mapped to the gene (or  
348 isoform);  $N_i$  is the total read counts for that individual summing read counts  
349 across all genes; and  $\lambda_i$  is an unknown Poisson rate parameter. We model the  
350 log-transformed rate  $\lambda_i$  as a linear combination of several parameters

$$351 \quad \log(\lambda_i) = \mathbf{w}_i^T \boldsymbol{\alpha} + x_i \beta + g_i + e_i, \quad i = 1, 2, \dots, n,$$

$$\mathbf{g} = (g_1, g_2, \dots, g_n)^T \sim \text{MVN}(\mathbf{0}, \sigma^2 h^2 \mathbf{K}),$$

$$\mathbf{e} = (e_1, e_2, \dots, e_n)^T \sim \text{MVN}(\mathbf{0}, \sigma^2 (1 - h^2) \mathbf{I}),$$

352 where  $\mathbf{w}_i$  is a  $c$ -vector of covariates (including the intercept);  $\boldsymbol{\alpha}$  is a  $c$ -vector of  
353 corresponding coefficients;  $x_i$  represents the predictor variable of interest (e.g.  
354 experimental perturbation, sex, disease status, or genotype);  $\beta$  is its coefficient;  
355  $\mathbf{g}$  is an  $n$ -vector of genetic effects;  $\mathbf{e}$  is an  $n$ -vector of environmental effects;  $\mathbf{K}$  is  
356 an  $n$  by  $n$  positive semi-definite matrix that models the covariance among  
357 individuals due to individual relatedness, population structure, or hidden  
358 confounders;  $\mathbf{I}$  is an  $n$  by  $n$  identity matrix that models independent  
359 environmental variation;  $\sigma^2 h^2$  is the genetic variance component;  $\sigma^2 (1 - h^2)$  is  
360 the environmental variance component; and  $\text{MVN}$  denotes the multivariate  
361 normal distribution. In the above model, we assume that  $\mathbf{K}$  is known and can be  
362 computed based on either pedigree, genotype, or the gene expression matrix (9).  
363 For pedigree/genotype data, when  $\mathbf{K}$  is standardized to have  $\text{tr}(\mathbf{K})/n = 1$ ,  
364  $h^2 \in [0, 1]$  has the usual interpretation of heritability (9), where the  $\text{tr}(\cdot)$  denotes  
365 the trace of a matrix. Importantly, unlike several other DE methods (15,25), our  
366 model can deal with both continuous and discrete predictor variables.

367 Both of the random effects terms  $\mathbf{g}$  and  $\mathbf{e}$  model over-dispersion, the extra  
368 variance not explained by a Poisson model. However, the two terms  $\mathbf{g}$  and  $\mathbf{e}$   
369 model two different aspects of over-dispersion. Specifically,  $\mathbf{g}$  models the fraction  
370 of the extra variance that is explained by sample non-independence while  $\mathbf{e}$   
371 models the fraction of the extra variance that is independent across samples. For  
372 example, let us consider a simple case in which all samples have the same  
373 sequencing depth (i.e.  $N_i = N$ ) and there is only one intercept term  $\mu$  included as  
374 the covariate. In this case, the random effects term  $\mathbf{e}$  models the independent

375 over-dispersion: without  $\mathbf{g}$ ,  $V(y) = E(y) \left(1 + E(y)(e^{\sigma^2} - 1)\right)$  is still larger than  
376 the mean  $E(y) = Ne^{\mu + \sigma^2/2}$ , with the difference between the two increasing with  
377 increasing  $\sigma^2$ . In a similar fashion, the random effects term  $\mathbf{g}$  models the non-  
378 independent over-dispersion by accounting for the sample covariance matrix  $\mathbf{K}$ .  
379 By modeling both aspects of over-dispersion, our PMM effectively generalizes  
380 the commonly used negative binomial model -- which only models independent  
381 extra variance -- to account for sample non-independence. In addition, our PMM  
382 naturally extends the commonly used linear mixed model (LMM) (9,65,83) to  
383 modeling count data.

384 Our goal here is to test the null hypothesis that gene expression levels are not  
385 associated with the predictor variable of interest, or  $H_0: \beta = 0$ . Testing this  
386 hypothesis requires estimating parameters in the PMM (as has previously been  
387 done in other settings (84,85), including for modeling uneven RNAseq read  
388 distribution inside transcripts (13); details in Supplementary Material). The PMM  
389 belongs to the generalized linear mixed model family, where parameter  
390 estimation is notoriously difficult because of the random effects and the resulting  
391 intractable  $n$ -dimensional integral in the likelihood. Standard estimation methods  
392 rely on numerical integration (86) or Laplace approximation (87,88), but neither  
393 strategy scales well with the increasing dimension of the integral, which in our  
394 case equals the sample size. As a consequence, standard approaches often  
395 produce biased estimates and overly narrow (i.e., anti-conservative) confidence  
396 intervals (89-95). To overcome the high-dimensionality of the integral, we instead  
397 develop a novel Markov Chain Monte Carlo (MCMC) algorithm, which, with  
398 enough iterations, can achieve high inference accuracy (96,97). We use MCMC  
399 to draw posterior samples but rely on the asymptotic normality of both the  
400 likelihood and the posterior distributions (98) to obtain the approximate maximum  
401 likelihood estimate  $\hat{\beta}_j$  and its standard error  $se(\hat{\beta}_j)$ . With  $\hat{\beta}_j$  and  $se(\hat{\beta}_j)$ , we can  
402 construct approximate Wald test statistics and  $p$ -values for hypothesis testing  
403 (Supplementary Material). Although we use MCMC, our procedure is frequentist  
404 in nature.

405 At the technical level, our MCMC algorithm is also novel, taking advantage of an  
406 auxiliary variable representation of the Poisson likelihood (62-64) and recent  
407 linear algebra innovations for fitting linear mixed models (9,59,65). Our MCMC  
408 algorithm introduces *two* continuous latent variables for each individual to replace  
409 the count observation, effectively extending our previous approach of using *one*  
410 latent variable for the simpler binomial distribution (60). Compared with a  
411 standard MCMC, our new MCMC algorithm reduces the computational  
412 complexity of each MCMC iteration from cubic to quadratic with respect to the  
413 sample size. Therefore, our method is orders of magnitude faster than the  
414 popular Bayesian software MCMCglmm (99) and can be used to analyze  
415 hundreds of samples and tens of thousands of genes with a single desktop PC

416 (Figure S1). Although our procedure is stochastic in nature, we find the MCMC  
417 errors are often small enough to ensure stable  $p$ -values across independent  
418 MCMC runs (Figure S2).

419

## 420 **Simulations: control for sample non-independence**

421 We performed a series of simulations to compare the performance of the PMM  
422 implemented in MACAU with four other commonly used methods: (1) a linear  
423 model; (2) the linear mixed model implemented in GEMMA (9,59); (3) a Poisson  
424 model; and (4) a negative binomial model. We used quantile-transformed data for  
425 linear model and GEMMA (see Methods and Materials for normalization details  
426 and a comparison between various transformations; Figure S3) and used raw  
427 count data for the other three methods. To make our simulations realistic, we use  
428 parameters inferred from a published RNAseq data set on a population of wild  
429 baboons (7,71) to perform simulations (Methods and Materials); this baboon data  
430 set contains known related individuals and hence invokes the problem of sample  
431 non-independence outlined above.

432 Our first set of simulations was performed to evaluate the effectiveness of  
433 MACAU and the other four methods in controlling for sample non-independence.  
434 To do so, we simulated expression levels for 10,000 genes in 63 individuals (the  
435 sample size from the baboon data set). Simulated gene expression levels are  
436 influenced by both independent environmental effects and correlated genetic  
437 effects, where genetic effects are simulated based on the baboon kinship matrix  
438 (estimated from microsatellite data (7)) with either zero ( $h^2 = 0.0$ ), moderate  
439 ( $h^2 = 0.3$ ), or high ( $h^2 = 0.6$ ) heritability values. We also simulated a continuous  
440 predictor variable  $x$  that is itself moderately heritable ( $h_x^2 = 0.4$ ). Because we  
441 were interested in the behavior of the null in this set of simulations, gene  
442 expression levels were not affected by the predictor variable (i.e., no genes were  
443 truly DE).

444 Figures 1, S4, and S5 show quantile-quantile plots for analyses using MACAU  
445 and the other four methods against the null (uniform) expectation, for  $h^2 = 0.6$ ,  
446  $h^2 = 0.3$ , and  $h^2 = 0.0$  respectively. When genes are heritable and the predictor  
447 variable is also correlated with individual relatedness, then the resulting  $p$ -values  
448 from the DE analysis are expected to be uniform only for a method that properly  
449 controls for sample non-independence. If a method fails to control for sample  
450 non-independence, then the  $p$ -values would be inflated, resulting in false  
451 positives.

452 Our results show that, because MACAU controls for sample non-independence,  
453 the  $p$ -values from MACAU follow the expected uniform distribution closely (and  
454 are slightly conservative) regardless of whether gene expression is moderately or  
455 highly heritable. The genomic control factors from MACAU are close to 1



456 (Figures 1 and S4). Even if we use a relatively relaxed q-value cutoff of 0.2 to  
457 identify DE genes, we do not incorrectly identify any genes as DE with MACAU.  
458 In contrast, the  $p$ -values from negative binomial are inflated and skewed towards  
459 low (significant) values, especially for gene expression levels with high heritability.  
460 With negative binomial, 27 DE genes (when  $h^2 = 0.3$ ) or 21 DE genes (when  
461  $h^2 = 0.6$ ) are erroneously detected at the q-value cutoff of 0.2. The inflation of  $p$ -  
462 values is even more acute in Poisson, presumably because the Poisson model  
463 accounts for neither individual relatedness nor over-dispersion. For non-count-  
464 based models, the  $p$ -values from a linear model are slightly skewed towards  
465 significant values, with 3 DE genes (when  $h^2 = 0.3$ ) and 1 DE gene (when  
466  $h^2 = 0.6$ ) erroneously detected at  $q < 0.2$ . In contrast, because the LMM in  
467 GEMMA also accounts for individual relatedness, it controls for sample non-  
468 independence well. Finally, when genes are not heritable, all methods except  
469 Poisson correctly control type I error (Figure S5).

470 Two important factors influence the severity of sample non-independence in  
471 RNAseq data (Figure 2). First, the inflation of  $p$ -values in the negative binomial,  
472 Poisson and linear models becomes more acute with increasing sample size. In  
473 particular, when  $h_x^2 = 0.4$ , with a sample size of  $n = 1,000$ ,  $\lambda_{gc}$  from the negative  
474 binomial, Poisson and linear models reaches 1.71, 82.28, and 1.41, respectively.  
475 In contrast, even when  $n = 1,000$ ,  $\lambda_{gc}$  from both MACAU and GEMMA remain  
476 close to 1 (0.97 and 1.01, respectively). Second, the inflation of  $p$ -values in the  
477 three models also becomes more acute when the predictor variable is more  
478 correlated with population structure. Thus, for a highly heritable predictor variable  
479 ( $h_x^2 = 0.8$ ),  $\lambda_{gc}$  (when  $n = 1,000$ ) from the negative binomial, Poisson and linear  
480 models increases to 2.13, 101.43, and 1.81, respectively, whereas  $\lambda_{gc}$  from  
481 MACAU and GEMMA remains close to 1 (1.02 and 1.05).

482 We also compared MACAU with edgeR (25) and DESeq2 (15), two commonly  
483 used methods for DE analysis (38,100). Because edgeR and DESeq2 were  
484 designed for discrete predictor variables, we discretized the continuous predictor  
485  $x$  into 0/1 based on the median predictor value across individuals. We then  
486 applied all methods to the same binarized predictor values for comparison.  
487 Results are shown in Figure S6. For the five methods compared above, the  
488 results on binarized values are comparable with those for continuous variables  
489 (i.e. Figure S6 vs Figure 1). Both edgeR and DESeq2 produce anticonservative  
490  $p$ -values and perform similarly to the negative binomial model in terms of type I  
491 error control.

492 Finally, we explored the use of principal components (PCs) from the gene  
493 expression matrix or the genotype matrix to control for sample non-independence.  
494 Genotype PCs have been used as covariates to control for population  
495 stratification in association studies (101). However, recent comparative studies  
496 have shown that using PCs is less effective than using linear mixed models



497 (83,102). Consistent with the poorer performance of PCs in association studies  
498 (83,102), using the top PCs from either the gene expression matrix or the  
499 genotype matrix does not improve type I error control for negative binomial,  
500 Poisson, linear, edgeR or DESeq2 approaches (Figures S7 and S8).

501

## 502 **Simulations: power to identify DE genes**

503 Our second set of simulations was designed to compare the power of different  
504 methods for identifying DE genes, again based on parameters inferred from real  
505 data. This time, we simulated a total of 10,000 genes, among which 1,000 genes  
506 were truly DE and 9,000 were non-DE. For the DE genes, simulated effect sizes  
507 corresponded to a fixed proportion of variance explained (PVE) in gene  
508 expression levels that ranged from 15% to 35%. For each set of parameters, we  
509 performed 10 replicate simulations and measured model performance based on  
510 the area under the curve (AUC) (as in (35,68,103)). We also examined several  
511 key factors that could influence the relative performance of the alternative  
512 methods: (1) gene expression heritability ( $h^2$ ); (2) correlation between the  
513 predictor variable  $x$  and genetic relatedness (measured by the heritability of  $x$ , or  
514  $h_x^2$ ); (3) variation of the total read counts across samples (measured by the  
515 coefficient of variation, or CV); (4) the over-dispersion parameter ( $\sigma^2$ ); (5) the  
516 effect size (PVE); and (6) sample size ( $n$ ). To do so, we first performed  
517 simulations using a default set of values ( $h^2 = 0.3$ ,  $h_x^2 = 0$ ,  $CV = 0.3$ ,  $\sigma^2 = 0.25$ ,  
518  $PVE = 0.25$ , and  $n = 63$ ) and then varied them one at a time to examine the  
519 influence of each factor on the relative performance of each method.

520 Our results show that MACAU works either as well as or better than other  
521 methods in almost all settings (Figures 3, S9-S14), probably because it both  
522 models count data directly and controls for sample non-independence. In  
523 contrast, the Poisson approach consistently fared the worst across all simulation  
524 scenarios, presumably because it fails to account for any sources of over-  
525 dispersion (Figures 3, S9-S14).

526 Among the factors that influence the relative rank of various methods, the most  
527 important factor was heritability ( $h^2$ ) (Figure 3A). While all methods perform  
528 worse with increasing gene expression heritability, heritability disproportionately  
529 affects the performance of models that do not account for relatedness (i.e.,  
530 negative binomial, Poisson and Linear), whereas when heritability is zero ( $h^2 =$   
531  $0$ ), these approaches tend to perform slightly better. Therefore, for non-heritable  
532 genes, linear models perform slightly better than GEMMA, and negative binomial  
533 models work similarly or slightly better than MACAU. This observation most likely  
534 arises because linear and negative binomial models require fewer parameters  
535 and thus have a greater number of degrees of freedom. However, even in this  
536 setting, the difference between MACAU and negative binomial is small,

537 suggesting that MACAU is robust to model misspecification and works  
538 reasonably well even for non-heritable genes. On the other hand, when  
539 heritability is moderate ( $h^2 = 0.3$ ) or high ( $h^2 = 0.6$ ), the methods that account for  
540 sample non-independence are much more powerful than the methods that do not.  
541 Because almost all genes are influenced by cis-eQTLs (47,48) and are thus likely  
542 heritable to some extent, MACAU's robustness for non-heritable genes and its  
543 high performance gain for heritable genes make it appealing.

544 The second most important factor in relative model performance was the  
545 variation of total read counts across individuals (CV; Figure 3B). While all  
546 methods perform worse with increasing CV, CV particularly affects the  
547 performance of GEMMA. Specifically, when CV is small (0.3; as the baboon  
548 data), GEMMA works well and is the second best method behind MACAU.  
549 However, when CV is moderate (0.6) or high (0.9), the performance of GEMMA  
550 quickly decays: it becomes only the fourth best method when  $CV = 0.9$ . GEMMA  
551 performs poorly in high CV settings presumably because the linear mixed model  
552 fails to account for the mean-variance dependence observed in count data, which  
553 is in agreement with previous findings (60,104).

554 The other four factors we explored had small impacts on the relative performance  
555 of the alternative methods, although they did affect their absolute performance.  
556 For example, as one would expect, power increases with large effect sizes (PVE)  
557 (Figure S9) or large sample sizes (Figure S10), and decreases with large over-  
558 dispersion  $\sigma^2$  (Figure S11) or large  $h_x^2$  (Figure S12).

559 Finally, we included comparisons with edgeR (25) and DESeq2 (15). In the basic  
560 parameter simulation setting ( $n = 63$ ,  $CV = 0.3$ ,  $h_x^2 = 0$ ,  $PVE = 0.25$ ,  $h^2 = 0.3$  and  
561  $\sigma^2 = 0.25$ ), we again discretized the continuous predictor  $x$  into a binary 0/1  
562 variable based on the median predictor value across individuals. Results for all  
563 methods are shown in Figure S13A. For the five methods also tested on a  
564 continuous predictor variable, the power on binarized values is much reduced  
565 compared with the power when the predictor variable is modeled without  
566 binarization (e.g. Figure S13A vs Figure 3). Further, neither edgeR nor DESeq2  
567 perform well, consistent with the recent move from these methods towards linear  
568 models in differential expression analysis (3,7,46-48,105). This result is not  
569 contingent on having large sample sizes. In small sample size settings ( $n=6$ ,  
570  $n=10$ , and  $n=14$ , with samples balanced between the two classes, 0 or 1),  
571 MACAU again outperforms the other methods, though the power difference is  
572 much smaller ( $n=10$  and  $n=14$ ; Figures S13C and S31D) and sometimes  
573 negligible ( $n=6$ , Figure S13B).

574 In summary, the power of MACAU and other methods, as well as the power  
575 difference between methods, is influenced in a continuous fashion by multiple  
576 factors. Larger sample sizes, larger effect sizes, lower read depth variation, lower  
577 gene expression heritability, lower predictor variable heritability, and lower over-

578 dispersion all increase power. However, MACAU's power is less diminished by  
579 high gene expression heritability and high read depth variability than the non-  
580 mixed model methods, while retaining the advantage of modeling the count data  
581 directly. In challenging data analysis settings (e.g., when sample size is low *and*  
582 effect size is low: Figure S13B for  $n=6$ ), no method stands out, and using  
583 MACAU results in no or negligible gains in power relative to other methods.  
584 When the sample size is low ( $n=6$ ) and effect sizes are large, however, MACAU  
585 consistently outperforms the other methods ( $n=6$ , Figure S14).

586

## 587 **Real Data Applications**

588 To gain insight beyond simulation, we applied MACAU and the other six methods  
589 to three recently published RNAseq data sets.

590 The first data set we considered is the baboon RNAseq study (7) used to  
591 parameterize the simulations above. Expression measurements on 12,018 blood-  
592 expressed genes were collected by the Amboseli Baboon Research Project  
593 (ABRP) (71) for 63 adult baboons (26 females and 37 males), among which  
594 some were relatives. Here, we applied MACAU and the six other methods to  
595 identify genes with sex-biased expression patterns. Sex-associated genes are  
596 known to be enriched on sex chromosomes (106,107), and we use this  
597 enrichment as one of the criteria to compare method performance, as in (18).  
598 Because the same nominal  $p$ -value from different methods may correspond to  
599 different type I errors, we compared methods based on empirical false discovery  
600 rate (FDR). In particular, we permuted the data to construct an empirical null,  
601 estimated the FDR at any given  $p$ -value threshold, and counted the number of  
602 discoveries at a given FDR cutoff (see Methods and Materials for permutation  
603 details and a comparison between two different permutation procedures; Figure  
604 S15).

605 In agreement with our simulations, MACAU was the most powerful method of  
606 those we considered. Specifically, at an empirical FDR of 5%, MACAU identified  
607 105 genes with sex-biased expression patterns, 40% more than that identified by  
608 the linear model, the second best method at this FDR cutoff (Figure 4A). At a  
609 more relaxed FDR of 10%, MACAU identified 234 sex-associated genes, 47%  
610 more than that identified by the negative binomial model, the second best  
611 method at this FDR cutoff (Figure 4A). Further, as expected, the sex-associated  
612 genes detected by MACAU are enriched on the X chromosome (the Y  
613 chromosome is not assembled in baboons and is thus ignored), and this  
614 enrichment is stronger for the genes identified by MACAU than by the other  
615 methods (Figure 4B). Of the remaining approaches, the negative binomial, linear  
616 model, and GEMMA all performed similarly and are ranked right after MACAU.

617 The Poisson model performs the worst, and edgeR and DESeq2 fall between the  
618 Poisson model and the other methods (Figures 4A and 4B).

619 The second data set we considered is an RNAseq study on type II diabetes (T2D)  
620 collected as part of the Finland-United States Investigation of NIDDM Genetics  
621 (FUSION) Study (61). Here, the data were collected from skeletal muscle  
622 samples from 267 individuals with expression measurements on 21,753 genes.  
623 Individuals are from three municipalities (Helsinki, Savitaipale, and Kuopio) in  
624 Finland. Individuals within each municipality are more closely related than  
625 individuals between municipalities (e.g., the top genotype principal components  
626 generally correspond to the three municipalities; Figure S16). Two related  
627 phenotypes were available to us: 162 individuals with T2D or NGT (normal  
628 glucose tolerance) status (i.e., case/control) based on the oral glucose tolerance  
629 test (OGTT) and 267 individuals with the quantitative trait fasting glucose level  
630 (GL), a biologically relevant trait of T2D.

631 We performed analyses to identify genes associated with T2D status as well as  
632 genes associated with GL. To accommodate edgeR and DESeq2, we also  
633 discretized the continuous GL values into binary 0/1 categories based on the  
634 median GL value across individuals. We refer to the resulting values as GL01.  
635 Therefore, we performed two sets of analyses for GL: one on the continuous GL  
636 values and the other on the discretized GL01 values. Consistent with simulations  
637 and the baboon data analysis, MACAU identified more T2D-associated genes  
638 and GL-associated genes than other methods across a range of empirical FDR  
639 values. For the T2D analysis, MACAU identified 23 T2D-associated genes at an  
640 FDR of 5%, while GEMMA and the linear model, the second best methods at this  
641 FDR cutoff, identified only 1 T2D-associated gene (Figure 4C). Similarly, at an  
642 FDR of 10%, MACAU identified 123 T2D-associated genes, 51% more than that  
643 identified by the linear model, the second best method at this FDR cutoff (Figure  
644 4C). For GL analysis, based on an FDR of 5%, MACAU detected 12 DE genes,  
645 while the other methods did not identify any DE genes at this FDR cutoff. At an  
646 FDR of 10%, MACAU identified 100 GL associated genes, while the second best  
647 methods -- the linear model and GEMMA -- identified 12 DE genes (Figure 4E).  
648 For the dichotomized GL01, none of the methods detected any DE genes even at  
649 a relaxed FDR cutoff of 20%, highlighting the importance of modeling the original  
650 continuous predictor variable in DE analysis.

651 Several lines of evidence support the biological validity of the genes detected by  
652 MACAU. First, we performed Gene Ontology (GO) analysis using LRpath (108)  
653 on T2D and GL associated genes identified by MACAU, as in the FUSION study  
654 (61) (Figure S17). The GO analysis results for T2D and GL are consistent with  
655 previous studies (61,109) and are also similar to each other, as expected given  
656 the biological relationship between the two traits. In particular, T2D status and  
657 high GL are associated with decreased expression of cellular respiratory pathway

658 genes, consistent with previous observations (61,109). T2D status and GL are  
659 also associated with several pathways that are related to mTOR, including  
660 generation of precursor metabolites, poly-ubiquitination and vesicle trafficking, in  
661 agreement with a prominent role of mTOR pathway in T2D etiology (110-113).

662 Second, we performed overlap analyses between T2D and GL associated genes.  
663 We reasoned that T2D-associated genes are likely associated with GL because  
664 T2D shares a common genetic basis with GL (114-116) and T2D status is  
665 determined in part by fasting glucose levels. Therefore, we used the overlap  
666 between genes associated with T2D and genes associated with GL as a  
667 measure of method performance. In the overlap analysis, genes with the  
668 strongest T2D association identified by MACAU show a larger overlap with the  
669 top 1,000 genes that have the strongest GL association than did genes identified  
670 by other methods (Figure 4D). For instance, among the top 100 genes with the  
671 strongest T2D-association evidence from MACAU, 63 of them also show strong  
672 association evidence with GL. In contrast, only 55 of the top 100 genes with the  
673 strongest T2D-association identified by GEMMA, the second best method, show  
674 strong association evidence with GL. We observed similar results, with MACAU  
675 performing the best, when performing the reciprocal analysis (overlap between  
676 genes with the strongest GL-association and the top 1,000 genes that have the  
677 strongest T2D-association: Figure 4F). To include the comparison with edgeR  
678 and DESeq2, we further examined the overlap between T2D associated genes  
679 and GL01 associated genes for all methods (Figure S18). Again, MACAU  
680 performs the best, followed by GEMMA and the linear model, and neither edgeR  
681 nor DESeq2 perform well in this context (Figure S18). Therefore, MACAU  
682 appears to both confer more power to identify biologically relevant DE genes and  
683 be more consistent across analyses of related phenotypes.

684 To assess the type I error rate of various methods, we permuted the trait data  
685 from the baboon and the FUSION studies. Consistent with our simulation results,  
686 the  $p$ -values from MACAU and GEMMA under the permuted null were close to  
687 uniformly distributed (slightly conservative) in both data sets, whereas the other  
688 methods were not (Figures S19 and S20). In addition, none of the methods  
689 compared here are sensitive to outliers in the two data sets (Figures S21-S23).

690 Finally, although large, population-based RNAseq data sets are becoming more  
691 common, MACAU's flexible PMM modeling framework allows it to be applied to  
692 DE analysis in small data sets with unrelated individuals as well. In this setting,  
693 MACAU can use the gene expression covariance matrix as the  $K$  matrix to  
694 control for hidden confounding effects that are commonly observed in  
695 sequencing studies (49-52). Hidden confounders can induce similarity in gene  
696 expression levels across many genes even though individuals are unrelated (53-  
697 57), similar to the effects of kinship or population structure. Therefore, by  
698 defining  $K$  using a gene expression (instead of genetic) covariance matrix,



699 MACAU can effectively control for sample non-independence induced by hidden  
700 confounders, thus extending the linear mixed model widely used to control for  
701 hidden confounders in array based studies (53-57) to sequencing count data.

702 To illustrate this application, we analyzed a third data set on lymphoblastoid cell  
703 lines (LCLs) derived from 69 unrelated Nigerian individuals (YRI) (3) from the  
704 HapMap project (117), with expression measurements on 13,319 genes. We also  
705 aimed to identify sex-associated genes in this data set. To demonstrate the  
706 effectiveness of MACAU in small samples, we randomly subsampled individuals  
707 from the data to create small data sets with either  $n = 6$  (3 males and 3 females),  
708  $n = 10$  (5 males and 5 females), or  $n = 14$  individuals (7 males and 7 females).  
709 For each sample size  $n$ , we performed 20 replicates of random subsampling and  
710 then evaluated method performance by averaging across replicates. In each  
711 replicate, we used the gene expression covariance matrix as  $K$  and compared  
712 MACAU's performance against other methods. Because of the small sample size,  
713 none of the methods were able to identify DE genes at an FDR cutoff of 10%,  
714 consistent with recent arguments that at least 6-12 biological replicates are  
715 needed to ensure sufficient power and replicability in DE analysis (11). We  
716 therefore used enrichment of genes on the sex chromosomes to compare the  
717 performance of different methods (Figure S24). The enrichment of top ranked  
718 sex-associated genes on sex chromosomes has previously been used for  
719 method comparison and is especially suitable for comparing methods in the  
720 presence of batch effects and other hidden confounding factors (118).

721 In this comparison, MACAU performs the best of all methods when the sample  
722 size is either  $n = 10$  or  $n = 14$ , and is ranked among the best (together with the  
723 negative binomial model) when  $n = 6$ . For instance, when  $n = 6$ , among the top  
724 50 genes identified by each method, the number of genes on the sex  
725 chromosomes for MACAU, negative binomial, Poisson, edgeR, DESeq2,  
726 GEMMA, and Linear are 3.3, 2.7, 3.1, 1.8, 3.0, 2.0, and 2.4, respectively. The  
727 advantage of MACAU becomes larger when the sample size increases: for  
728 example, when  $n = 14$ , an average of 10.6 genes in the top 50 genes from  
729 MACAU are on the sex chromosomes, which is again larger than that from the  
730 negative binomial (8.3), Poisson (6.0), edgeR (6.65), DESeq2 (8.8), GEMMA  
731 (9.8), or Linear (8.05). These results suggest that MACAU can also perform  
732 better than existing methods in relatively small sample study designs with  
733 unrelated individuals by controlling for hidden confounders. However, MACAU's  
734 power gain is much smaller in this setting than in the first two data sets we  
735 considered (the baboon and Fusion data). In addition, MACAU's power gain is  
736 negligible in the case of  $n=6$  when compared with the second best method,  
737 though its power gain over the commonly used edgeR and DESeq2 is still  
738 substantial. MACAU's small power gain in this data presumably stems from both  
739 the small sample size and the small effect size of sex in the data, consistent with  
740 previous reports for blood cell-derived gene expression (3,7,119).



741

742

## 743 Discussion

744 Here, we present an effective Poisson mixed effects model, together with a  
745 computationally efficient inference method and software implementation in  
746 MACAU, for identifying DE genes in RNAseq studies. MACAU directly models  
747 count data and, using two random effects terms, controls for both independent  
748 over-dispersion and sample non-independence. Because of its flexible modeling  
749 framework, MACAU controls for type I error in the presence of individual  
750 relatedness, population structure, and hidden confounders, and MACAU  
751 achieves higher power than several other methods for DE analysis across a  
752 range of settings. In addition, MACAU can easily accommodate continuous  
753 predictor variables and biological or technical covariates. We have demonstrated  
754 the benefits of MACAU using both simulations and applications to three recently  
755 published RNAseq data sets.

756 MACAU is particularly well-suited to data sets that contain related individuals or  
757 population structure. Several major population genomic resources contain  
758 structure of these kinds. For example, the HapMap population (117), the Human  
759 Genome Diversity Panel (120), the 1000 Genomes Project in humans (121) as  
760 well as the 1001 Genomes Project in Arabidopsis (122) all contain data from  
761 multiple populations or related individuals. Several recent large-scale RNAseq  
762 projects also collected individuals from genetically differentiated populations (46).  
763 MACAU is also well-suited to analyzing genes with moderate to high heritability.  
764 Previous studies in humans have shown that, while heritability varies across  
765 genes, many genes are moderately or highly heritable, and almost all genes  
766 have detectable eQTL (47,123). Analyzing these data with MACAU can reduce  
767 false positives and increase power. Notably, even when genes exhibit zero  
768 heritability, our results show that MACAU incurs minimal loss of power compared  
769 with other approaches.

770 While we have mainly focused on illustrating the benefits of MACAU for  
771 controlling for individual relatedness and population stratification, we note that  
772 MACAU can be used to control for sample non-independence occurred in other  
773 settings. For example, cell type heterogeneity (55) or other hidden confounding  
774 factors (53) are commonly observed in sequencing studies and can induce gene  
775 expression similarity even when individuals are unrelated (49-52). Because the  
776 gene expression covariance matrix  $K$  contains information on sample non-  
777 independence caused by hidden confounding factors (53-57), MACAU could be  
778 applied to control for hidden confounding effects by using the gene expression  
779 covariance as the  $K$  matrix. Therefore, MACAU provides a natural avenue for  
780 extending the commonly used mixed effects model for controlling for hidden  
781 confounding effects (53-56) in array-based studies to sequencing studies. In  
782 addition, although we have designed MACAU for differential expression analysis,  
783 we note that MACAU may also be effective in other common settings. For

784 example, MACAU could be readily applied in QTL mapping studies to identify  
785 genetic variants that are associated with gene expression levels estimated using  
786 RNAseq or related high-throughput sequencing methods.

787 In the present study, we have focused on demonstrating the performance of  
788 MACAU in three published RNAseq data sets with sample sizes ranging from  
789 small ( $n=6$ ) to medium ( $n=63$ ) to large ( $n=267$ ), relative to the size of most  
790 current RNAseq studies. Compared with small sample studies, RNAseq studies  
791 with medium or large sample sizes are better powered and more reproducible,  
792 and are thus becoming increasingly common in genomics (10,11). For example,  
793 a recent comparative study makes explicit calls for medium to large sample  
794 RNAseq studies performed with at least 12 replicates per condition (i.e.  $n \geq 24$ )  
795 (11). However, we recognize that many RNAseq studies are still carried out with  
796 a small number of samples (e.g. 3 replicates per condition). As our simulations  
797 make clear, the power of all analysis methods is dramatically reduced with  
798 decreasing sample size, conditional on fixed values of other factors that influence  
799 power (e.g., effect size, gene expression heritability). Thus, MACAU's advantage  
800 is no longer obvious in simulated data with only 3 replicates per condition when  
801 the effect size is also small (although its advantage becomes apparent when the  
802 simulated effect size increases: Figures S13B and S14). In addition, MACAU's  
803 advantage is much smaller and sometimes negligible in the small real data set  
804 when compared with the medium and large data sets analyzed here.  
805 Furthermore, because MACAU requires estimating one more parameter than  
806 other existing methods, MACAU requires at least five samples to run while  
807 existing DE methods require at least four. Therefore, MACAU may not confer  
808 benefits to power in some settings, and is especially likely (like all methods) to be  
809 underpowered in very small sample sizes with small effect sizes. Future  
810 extensions of MACAU are likely needed to ensure its robust performance in small  
811 as well as moderate to large samples. For example, further power improvements  
812 could be achieved by borrowing information across genes to estimate the over-  
813 dispersion parameter (15,22,25) or building in a hierarchical structure to model  
814 many genes at once.

815 Like other DE methods (24,25), MACAU requires data pre-processing to obtain  
816 gene expression measurements from raw sequencing read files. This data pre-  
817 processing step may include read alignment, transcript assembly, alternative  
818 transcript quantification, transcript measurement, and normalization. Many  
819 methods are available to perform these tasks (12,14,68,124-129) and different  
820 methods can be differentially advantageous across settings (68,124,130).  
821 Importantly, MACAU can be paired with any pre-processing method that retains  
822 the count nature of the data. While we provide a preliminary comparison of  
823 several methods here (see Materials and Methods; Figure S3), a full analysis of  
824 how different data pre-processing choices affect MACAU's performance in  
825 alternative study designs is beyond the scope of this paper. Notably, recent

826 results suggest that a recommended approach is to incorporate data pre-  
827 processing and DE analysis into the same, joint statistical framework (131),  
828 which represents an important next step for the MACAU software package.

829 We note that, like many other DE methods (15,25), we did not model gene length  
830 in MACAU. Because gene length does not change from sample to sample, it  
831 does not affect differential expression analysis on any particular gene (15,25).  
832 However, gene length will affect the power of DE analysis across different genes:  
833 genes with longer length tend to have a larger number of mapped reads and  
834 more accurate expression measurements, and as a consequence, DE analysis  
835 on these genes tends to have higher statistical power (2,70,132). Gene length  
836 may also introduce sample-specific effects in certain data sets (133). Therefore,  
837 understanding the impact of, and taking into account gene length effects, in  
838 MACAU DE analysis represents another possible future extension.

839 Currently, despite the newly developed computationally efficient algorithm,  
840 applications of MACAU can still be limited by its relatively heavy computational  
841 cost. The MCMC algorithm in MACAU scales quadratically with the number of  
842 individuals/samples and linearly with the number of genes. Although MACAU is  
843 two orders of magnitude faster than the standard software MCMCglmm for fitting  
844 Poisson mixed effects models (Table S1), it can still take close to 20 hours to  
845 analyze a data set of the size of the FUSION data we considered here (267  
846 individuals and 21,753 genes). Therefore, new algorithms will be needed to use  
847 MACAU for data sets that are orders of magnitude larger.

## 848 **URLs**

849 The software implementation of MACAU is freely available at:  
850 [www.xzlab.org/software.html](http://www.xzlab.org/software.html).

851

## 852 **Competing Interests**

853 The authors declare that they have no competing interests.

854

## 855 **Funding**

856 This study was supported by NIH fund R01HG009124. XZ is also supported by  
857 R01HL117626 (PI Abecasis), R21ES024834 (PI Pierce), R01HL133221 (PI  
858 Smith), and a grant from the Foundation for the National Institutes of Health  
859 through the Accelerating Medicines Partnership (BOEH15AMP, co-PIs Boehnke  
860 and Abecasis). JT is supported by 1R01GM102562 and R21AG049936. LS is  
861 supported by U01DK062370 (PI Boehnke). SS is supported by a scholarship  
862 from the China Scholarship Council.

863

## 864 **Acknowledgements**

865 We thank Matthew Stephens for insight and support on previous versions of  
866 MACAU. We thank Baylor College of Medicine Human Genome Sequencing  
867 Center for access to the current version of the baboon genome assembly (Panu  
868 2.0). We thank FUSION investigators for access to the FUSION expression data.

869

## 870 References

- 871 1. Nagalakshmi, U., Wang, Z., Waern, K., Shou, C., Raha, D., Gerstein, M. and Snyder, M.  
872 (2008) The transcriptional landscape of the yeast genome defined by RNA sequencing.  
873 *Science*, **320**, 1344-1349.
- 874 2. Mortazavi, A., Williams, B.A., Mccue, K., Schaeffer, L. and Wold, B. (2008) Mapping and  
875 quantifying mammalian transcriptomes by RNA-Seq. *Nat Methods*, **5**, 621-628.
- 876 3. Pickrell, J.K., Marioni, J.C., Pai, A.A., Degner, J.F., Engelhardt, B.E., Nkadori, E., Veyrieras,  
877 J.B., Stephens, M., Gilad, Y. and Pritchard, J.K. (2010) Understanding mechanisms  
878 underlying human gene expression variation with RNA sequencing. *Nature*, **464**, 768-  
879 772.
- 880 4. Wang, Z., Gerstein, M. and Snyder, M. (2009) RNA-Seq: a revolutionary tool for  
881 transcriptomics. *Nat Rev Genet*, **10**, 57-63.
- 882 5. Oshlack, A., Robinson, M.D. and Young, M.D. (2010) From RNA-seq reads to differential  
883 expression results. *Genome Biol*, **11**, 220.
- 884 6. Ozsolak, F. and Milos, P.M. (2011) RNA sequencing: advances, challenges and  
885 opportunities. *Nat Rev Genet*, **12**, 87-98.
- 886 7. Tung, J., Zhou, X., Alberts, S.C., Stephens, M. and Gilad, Y. (2015) The genetic  
887 architecture of gene expression levels in wild baboons. *Elife*, **4**, e04729.
- 888 8. Bennett, B.J., Farber, C.R., Orozco, L., Kang, H.M., Ghazalpour, A., Siemers, N., Neubauer,  
889 M., Neuhaus, I., Yordanova, R., Guan, B. *et al.* (2010) A high-resolution association  
890 mapping panel for the dissection of complex traits in mice. *Genome Res*, **20**, 281-290.
- 891 9. Altshuler, D.M., Durbin, R.M., Abecasis, G.R., Bentley, D.R., Chakravarti, A., Clark, A.G.,  
892 Donnelly, P., Eichler, E.E., Flicek, P., Gabriel, S.B. *et al.* (2012) An integrated map of  
893 genetic variation from 1,092 human genomes. *Nature*, **491**, 56-65.
- 894 10. Wood, A.R., Esko, T., Yang, J., Vedantam, S., Pers, T.H., Gustafsson, S., Chu, A.Y., Estrada,  
895 K., Luan, J., Kutalik, Z. *et al.* (2014) Defining the role of common variation in the genomic  
896 and biological architecture of adult human height. *Nature genetics*, **46**, 1173-1186.
- 897 11. Schurch, N.J., Schofield, P., Gierlinski, M., Cole, C., Sherstnev, A., Singh, V., Wrobel, N.,  
898 Gharbi, K., Simpson, G.G., Owen-Hughes, T. *et al.* (2016) How many biological replicates  
899 are needed in an RNA-seq experiment and which differential expression tool should you  
900 use? *Rna*, **22**, 839-851.
- 901 12. Li, B., Ruotti, V., Stewart, R.M., Thomson, J.A. and Dewey, C.N. (2010) RNA-Seq gene  
902 expression estimation with read mapping uncertainty. *Bioinformatics*, **26**, 493-500.
- 903 13. Hu, M., Zhu, Y., Taylor, J.M.G., Liu, J.S. and Qin, Z.H.S. (2012) Using Poisson mixed-  
904 effects model to quantify transcript-level gene expression in RNA-Seq. *Bioinformatics*,  
905 **28**, 63-68.
- 906 14. Wu, M.C., Lee, S., Cai, T., Li, Y., Boehnke, M. and Lin, X. (2011) Rare-variant association  
907 testing for sequencing data with the sequence kernel association test. *American journal*  
908 *of human genetics*, **89**, 82-93.
- 909 15. Anders, S. and Huber, W. (2010) Differential expression analysis for sequence count data.  
910 *Genome Biol*, **11**, R106.
- 911 16. Trapnell, C., Roberts, A., Goff, L., Pertea, G., Kim, D., Kelley, D.R., Pimentel, H., Salzberg,  
912 S.L., Rinn, J.L. and Pachter, L. (2012) Differential gene and transcript expression analysis  
913 of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc*, **7**, 562-578.
- 914 17. Li, J., Jiang, H. and Wong, W.H. (2010) Modeling non-uniformity in short-read rates in  
915 RNA-Seq data. *Genome Biol*, **11**.



- 916 18. Zhou, Y.H., Xia, K. and Wright, F.A. (2011) A powerful and flexible approach to the  
917 analysis of RNA sequence count data. *Bioinformatics*, **27**, 2672-2678.
- 918 19. Marioni, J.C., Mason, C.E., Mane, S.M., Stephens, M. and Gilad, Y. (2008) RNA-seq: An  
919 assessment of technical reproducibility and comparison with gene expression arrays.  
920 *Genome Res*, **18**, 1509-1517.
- 921 20. Wang, L.K., Feng, Z.X., Wang, X., Wang, X.W. and Zhang, X.G. (2010) DEGseq: an R  
922 package for identifying differentially expressed genes from RNA-seq data. *Bioinformatics*,  
923 **26**, 136-138.
- 924 21. Langmead, B., Hansen, K.D. and Leek, J.T. (2010) Cloud-scale RNA-sequencing  
925 differential expression analysis with Myrna. *Genome Biol*, **11**, R83.
- 926 22. Li, J., Witten, D.M., Johnstone, I.M. and Tibshirani, R. (2012) Normalization, testing, and  
927 false discovery rate estimation for RNA-sequencing data. *Biostatistics*, **13**, 523-538.
- 928 23. Auer, P.L. and Doerge, R.W. (2011) A Two-Stage Poisson Model for Testing RNA-Seq  
929 Data. *Stat Appl Genet Mol*, **10**.
- 930 24. Love, M.I., Huber, W. and Anders, S. (2014) Moderated estimation of fold change and  
931 dispersion for RNA-seq data with DESeq2. *Genome Biol*, **15**, 1--21.
- 932 25. Robinson, M.D., McCarthy, D.J. and Smyth, G.K. (2010) edgeR: a Bioconductor package  
933 for differential expression analysis of digital gene expression data. *Bioinformatics*, **26**,  
934 139-140.
- 935 26. McCarthy, D.J., Chen, Y.S. and Smyth, G.K. (2012) Differential expression analysis of  
936 multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Res*,  
937 **40**, 4288-4297.
- 938 27. Di, Y.M., Schafer, D.W., Cumbie, J.S. and Chang, J.H. (2011) The NBP Negative Binomial  
939 Model for Assessing Differential Gene Expression from RNA-Seq. *Stat Appl Genet Mol*,  
940 **10**.
- 941 28. Wu, H., Wang, C. and Wu, Z.J. (2013) A new shrinkage estimator for dispersion improves  
942 differential expression detection in RNA-seq data. *Biostatistics*, **14**, 232-243.
- 943 29. Van De Wiel, M.A., Leday, G.G.R., Pardo, L., Rue, H., Van Der Vaart, A.W. and Van  
944 Wieringen, W.N. (2013) Bayesian analysis of RNA sequencing data by estimating  
945 multiple shrinkage priors. *Biostatistics*, **14**, 113-128.
- 946 30. Hardcastle, T.J. and Kelly, K.A. (2010) baySeq: Empirical Bayesian methods for identifying  
947 differential expression in sequence count data. *Bmc Bioinformatics*, **11**, 1.
- 948 31. Li, J. and Tibshirani, R. (2013) Finding consistent patterns: A nonparametric approach for  
949 identifying differential expression in RNA-Seq data. *Stat Methods Med Res*, **22**, 519-536.
- 950 32. Tarazona, S., Garcia-Alcalde, F., Dopazo, J., Ferrer, A. and Conesa, A. (2011) Differential  
951 expression in RNA-seq: A matter of depth. *Genome Res*, **21**, 2213-2223.
- 952 33. Law, C.W., Chen, Y.S., Shi, W. and Smyth, G.K. (2014) voom: precision weights unlock  
953 linear model analysis tools for RNA-seq read counts. *Genome Biol*, **15**, R29.
- 954 34. Zwiener, I., Frisch, B. and Binder, H. (2014) Transforming RNA-Seq Data to Improve the  
955 Performance of Prognostic Gene Signatures. *Plos One*, **9**, e85150.
- 956 35. Soneson, C. and Delorenzi, M. (2013) A comparison of methods for differential  
957 expression analysis of RNA-seq data. *Bmc Bioinformatics*, **14**, 1.
- 958 36. Kvam, V.M., Lu, P. and Si, Y.Q. (2012) A Comparison Of Statistical Methods for Detecting  
959 Differentially Expressed Genes From Rna-Seq Data. *Am J Bot*, **99**, 248-256.
- 960 37. Zhang, Z.H., Jhaveri, D.J., Marshall, V.M., Bauer, D.C., Edson, J., Narayanan, R.K.,  
961 Robinson, G.J., Lundberg, A.E., Bartlett, P.F., Wray, N.R. *et al.* (2014) A Comparative  
962 Study of Techniques for Differential Expression Analysis on RNA-Seq Data. *Plos One*, **9**.

- 963 38. Schurch, N.J., Schofield, P., Gierliński, M., Cole, C., Sherstnev, A., Singh, V., Wrobel, N.,  
964 Gharbi, K., Simpson, G.G. and Owen-Hughes, T. (2016) Evaluation of tools for differential  
965 gene expression analysis by RNA-seq on a 48 biological replicate experiment. *RNA &*  
966 *Bioinformatics*.
- 967 39. P. McCullagh, J.A.N.F. (1989) *Generalized Linear Models*. Springer US USA.
- 968 40. Robinson, M.D. and Oshlack, A. (2010) A scaling normalization method for differential  
969 expression analysis of RNA-seq data. *Genome Biol*, **11**, R25.
- 970 41. Price, A.L., Helgason, A., Thorleifsson, G., McCarroll, S.A., Kong, A. and Stefansson, K.  
971 (2011) Single-Tissue and Cross-Tissue Heritability of Gene Expression Via Identity-by-  
972 Descent in Related or Unrelated Individuals. *Plos Genet*, **7**.
- 973 42. Wright, F.A., Sullivan, P.F., Brooks, A.I., Zou, F., Sun, W., Xia, K., Madar, V., Jansen, R.,  
974 Chung, W.I., Zhou, Y.H. *et al.* (2014) Heritability and genomics of gene expression in  
975 peripheral blood. *Nat Genet*, **46**, 430-437.
- 976 43. Monks, S.A., Leonardson, A., Zhu, H., Cundiff, P., Pietrusiak, P., Edwards, S., Phillips, J.W.,  
977 Sachs, A. and Schadt, E.E. (2004) Genetic inheritance of gene expression in human cell  
978 lines. *Am J Hum Genet*, **75**, 1094-1105.
- 979 44. Emilsson, V., Thorleifsson, G., Zhang, B., Leonardson, A.S., Zink, F., Zhu, J., Carlson, S.,  
980 Helgason, A., Walters, G.B., Gunnarsdottir, S. *et al.* (2008) Genetics of gene expression  
981 and its effect on disease. *Nature*, **452**, 423-U422.
- 982 45. Yang, S.J., Liu, Y.Y., Jiang, N., Chen, J., Leach, L., Luo, Z.W. and Wang, M.H. (2014)  
983 Genome-wide eQTLs and heritability for gene expression traits in unrelated individuals.  
984 *Bmc Genomics*, **15**, 1.
- 985 46. Lappalainen, T., Sammeth, M., Friedlander, M.R., 't Hoen, P.A.C., Monlong, J., Rivas,  
986 M.A., Gonzalez-Porta, M., Kurbatova, N., Griebel, T., Ferreira, P.G. *et al.* (2013)  
987 Transcriptome and genome sequencing uncovers functional variation in humans. *Nature*,  
988 **501**, 506-511.
- 989 47. Ardlie, K.G., DeLuca, D.S., Segre, A.V., Sullivan, T.J., Young, T.R., Gelfand, E.T.,  
990 Trowbridge, C.A., Maller, J.B., Tukiainen, T., Lek, M. *et al.* (2015) The Genotype-Tissue  
991 Expression (GTEx) pilot analysis: Multitissue gene regulation in humans. *Science*, **348**,  
992 648-660.
- 993 48. Battle, A., Mostafavi, S., Zhu, X.W., Potash, J.B., Weissman, M.M., McCormick, C.,  
994 Haudenschild, C.D., Beckman, K.B., Shi, J.X., Mei, R. *et al.* (2014) Characterizing the  
995 genetic basis of transcriptome diversity through RNA-sequencing of 922 individuals.  
996 *Genome Res*, **24**, 14-24.
- 997 49. Stegle, O., Parts, L., Piipari, M., Winn, J. and Durbin, R. (2012) Using probabilistic  
998 estimation of expression residuals (PEER) to obtain increased power and interpretability  
999 of gene expression analyses. *Nat Protoc*, **7**, 500-507.
- 1000 50. Leek, J.T. (2014) svaseq: removing batch effects and other unwanted noise from  
1001 sequencing data. *Nucleic Acids Res*, **42**.
- 1002 51. Leek, J.T. and Storey, J.D. (2007) Capturing heterogeneity in gene expression studies by  
1003 surrogate variable analysis. *Plos Genet*, **3**, 1724-1735.
- 1004 52. Risso, D., Ngai, J., Speed, T.P. and Dudoit, S. (2014) Normalization of RNA-seq data using  
1005 factor analysis of control genes or samples. *Nat Biotechnol*, **32**, 896-902.
- 1006 53. Kang, H.M., Ye, C. and Eskin, E. (2008) Accurate Discovery of Expression Quantitative  
1007 Trait Loci Under Confounding From Spurious and Genuine Regulatory Hotspots. *Genetics*,  
1008 **180**, 1909-1925.

- 1009 54. Listgarten, J., Kadie, C., Schadt, E.E. and Heckerman, D. (2010) Correction for hidden  
1010 confounders in the genetic analysis of gene expression. *P Natl Acad Sci USA*, **107**, 16465-  
1011 16470.
- 1012 55. Zou, J., Lippert, C., Heckerman, D., Aryee, M. and Listgarten, J. (2014) Epigenome-wide  
1013 association studies without the need for cell-type composition. *Nat Methods*, **11**, 309-  
1014 U283.
- 1015 56. Rahmani, E., Zaitlen, N., Baran, Y., Eng, C., Hu, D.L., Galanter, J., Oh, S., Burchard, E.G.,  
1016 Eskin, E., Zou, J. *et al.* (2016) Sparse PCA corrects for cell type heterogeneity in  
1017 epigenome-wide association studies. *Nat Methods*, **13**, 443-+.
- 1018 57. McGregor, K., Bernatsky, S., Colmegna, I., Hudson, M., Pastinen, T., Labbe, A. and  
1019 Greenwood, C.M.T. (2016) An evaluation of methods correcting for cell-type  
1020 heterogeneity in DNA methylation studies. *Genome Biol*, **17**.
- 1021 58. Price, A.L., Patterson, N.J., Plenge, R.M., Weinblatt, M.E., Shadick, N.A. and Reich, D.  
1022 (2006) Principal components analysis corrects for stratification in genome-wide  
1023 association studies. *Nat Genet*, **38**, 904-909.
- 1024 59. Zhou, X. and Stephens, M. (2014) Efficient multivariate linear mixed model algorithms  
1025 for genome-wide association studies. *Nat Methods*, **11**, 407-409.
- 1026 60. Lea, A.J., Alberts, S.C., Tung, J. and Zhou, X. (2015) A flexible, efficient binomial mixed  
1027 model for identifying differential DNA methylation in bisulfite sequencing data. *Plos*  
1028 *Genet*, **11**, e1005650.
- 1029 61. Scott, L.J., Erdos, M.R., Huyghe, J.R., Welch, R.P., Beck, A.T., Boehnke, M., Collins, F.S.  
1030 and Parker, S.C.J. (2016) The genetic regulatory signature of type 2 diabetes in human  
1031 skeletal muscle. *Nature Communications*, DOI: 10.1038/ncomms11764.
- 1032 62. Fruhwirth-Schnatter, S. and Wagner, H. (2006) Auxiliary mixture sampling for  
1033 parameter-driven models of time series of counts with applications to state space  
1034 modelling. *Biometrika*, **93**, 827-841.
- 1035 63. Scott, S.L. (2011) Data augmentation, frequentist estimation, and the Bayesian analysis  
1036 of multinomial logit models. *Stat Pap*, **52**, 87-109.
- 1037 64. Fruhwirth-Schnatter, S. and Fruhwirth, R. (2010) *Data Augmentation and MCMC for*  
1038 *Binary and Multinomial Logit Models*. Springer, New York.
- 1039 65. Lippert, C., Listgarten, J., Liu, Y., Kadie, C.M., Davidson, R.I. and Heckerman, D. (2011)  
1040 FaST linear mixed models for genome-wide association studies. *Nat Methods*, **8**, 833-  
1041 U894.
- 1042 66. Zhou, X., Carbonetto, P. and Stephens, M. (2013) Polygenic modeling with bayesian  
1043 sparse linear mixed models. *Plos Genet*, **9**, e1003264.
- 1044 67. Robin, X., Turck, N., Hainard, A., Tiberti, N., Lisacek, F., Sanchez, J.C. and Muller, M.  
1045 (2011) pROC: an open-source package for R and S plus to analyze and compare ROC  
1046 curves. *Bmc Bioinformatics*, **12**, 74.
- 1047 68. Teng, M., Love, M.I., Davis, C.A., Djebali, S., Dobin, A., Graveley, B.R., Li, S., Mason, C.E.,  
1048 Olson, S., Pervouchine, D. *et al.* (2016) A benchmark for RNA-seq quantification  
1049 pipelines. *Genome Biol*, **17**, 74.
- 1050 69. Ritchie, M.E., Phipson, B., Wu, D., Hu, Y.F., Law, C.W., Shi, W. and Smyth, G.K. (2015)  
1051 limma powers differential expression analyses for RNA-sequencing and microarray  
1052 studies. *Nucleic Acids Res*, **43**.
- 1053 70. Bullard, J.H., Purdom, E., Hansen, K.D. and Dudoit, S. (2010) Evaluation of statistical  
1054 methods for normalization and differential expression in mRNA-Seq experiments. *Bmc*  
1055 *Bioinformatics*, **11**, 94.

- 1056 71. Alberts, S. and Altmann, J. (2012) In Kappeler, P. M. and Watts, D. P. (eds.), *Long-Term*  
1057 *Field Studies of Primates*. Springer Berlin Heidelberg, pp. 261-287.
- 1058 72. Alberts, S.C., Buchan, J.C. and Altmann, J. (2006) Sexual selection in wild baboons: from  
1059 mating opportunities to paternity success. *Anim Behav*, **72**, 1177-1196.
- 1060 73. Buchan, J.C., Alberts, S.C., Silk, J.B. and Altmann, J. (2003) True paternal care in a multi-  
1061 male primate society. *Nature*, **425**, 179-181.
- 1062 74. Altmann, J., Altmann, S. and Hausfater, G. (1981) Physical Maturation And Age Estimates  
1063 Of Yellow Baboons, *Papio-Cynocephalus*, In Amboseli National-Park, Kenya. *Am J*  
1064 *Primatol*, **1**, 389-399.
- 1065 75. Archie, E.A., Tung, J., Clark, M., Altmann, J. and Alberts, S.C. (2014) Social affiliation  
1066 matters: both same-sex and opposite-sex relationships predict survival in wild female  
1067 baboons. *P Roy Soc B-Biol Sci*, **281**.
- 1068 76. Valle, T., Ehnholm, C., Tuomilehto, J., Blaschak, J., Bergman, R.N., Langefeld, C.D., Ghosh,  
1069 S., Watanabe, R.M., Hauser, E.R., Magnuson, V. *et al.* (1998) Mapping genes for NIDDM -  
1070 Design of the Finland United States Investigation of NIDDM Genetics (FUSION) study.  
1071 *Diabetes Care*, **21**, 949-958.
- 1072 77. Vaatainen, S., Keinanen-Kiukaanniemi, S., Saramies, J., Uusitalo, H., Tuomilehto, J. and  
1073 Martikainen, J. (2014) Quality of life along the diabetes continuum: a cross-sectional  
1074 view of health-related quality of life and general health status in middle-aged and older  
1075 Finns. *Qual Life Res*, **23**, 1935-1944.
- 1076 78. Li, H. and Durbin, R. (2009) Fast and accurate short read alignment with Burrows-  
1077 Wheeler transform. *Bioinformatics*, **25**, 1754-1760.
- 1078 79. Churchill, G.A. and Doerge, R.W. (2008) Naive application of permutation testing leads  
1079 to inflated type I error rates. *Genetics*, **178**, 609-610.
- 1080 80. Abney, M. (2015) Permutation Testing in the Presence of Polygenic Variation. *Genet*  
1081 *Epidemiol*, **39**, 249-258.
- 1082 81. Zhou, X.B., Lindsay, H. and Robinson, M.D. (2014) Robustly detecting differential  
1083 expression in RNA sequencing data using observation weights. *Nucleic Acids Res*, **42**, e91.
- 1084 82. George, N.I., Bowyer, J.F., Crabtree, N.M. and Chang, C.W. (2015) An Iterative Leave-  
1085 One-Out Approach to Outlier Detection in RNA-Seq Data. *Plos One*, **10**.
- 1086 83. Kang, H.M., Sul, J.H., Service, S.K., Zaitlen, N.A., Kong, S.Y., Freimer, N.B., Sabatti, C. and  
1087 Eskin, E. (2010) Variance component model to account for sample structure in genome-  
1088 wide association studies. *Nat Genet*, **42**, 348-U110.
- 1089 84. Tempelman, R.J. and Gianola, D. (1996) A mixed effects model for overdispersed count  
1090 data in animal breeding. *Biometrics*, **52**, 265-279.
- 1091 85. Tempelman, R.J. (1998) Generalized linear mixed models in dairy cattle breeding. *J Dairy*  
1092 *Sci*, **81**, 1428-1444.
- 1093 86. Pinheiro, J.C. and Chao, E.C. (2006) Efficient Laplacian and adaptive Gaussian quadrature  
1094 algorithms for multilevel generalized linear mixed models. *J Comput Graph Stat*, **15**, 58-  
1095 81.
- 1096 87. Goldstein, H. (1991) Nonlinear Multilevel Models, with an Application To Discrete  
1097 Response Data. *Biometrika*, **78**, 45-51.
- 1098 88. Breslow, N.E. and Clayton, D.G. (1993) Approximate Inference in Generalized Linear  
1099 Mixed Models. *J Am Stat Assoc*, **88**, 9-25.
- 1100 89. Breslow, N.E. and Lin, X.H. (1995) Bias Correction In Generalized Linear Mixed Models  
1101 with a Single-Component Of Dispersion. *Biometrika*, **82**, 81-91.
- 1102 90. Browne, W.J. and Draper, D. (2006) A comparison of Bayesian and likelihood-based  
1103 methods for fitting multilevel models. *Bayesian Anal*, **1**, 473-513.

- 1104 91. Lin, X.H. and Breslow, N.E. (1996) Bias correction in generalized linear mixed models  
1105 with multiple components of dispersion. *J Am Stat Assoc*, **91**, 1007-1016.
- 1106 92. Goldstein, H. and Rasbash, J. (1996) Improved approximations for multilevel models  
1107 with binary responses. *J Roy Stat Soc a Sta*, **159**, 505-513.
- 1108 93. Rodriguez, G. and Goldman, N. (2001) Improved estimation procedures for multilevel  
1109 models with binary response: a case-study. *J Roy Stat Soc a Sta*, **164**, 339-355.
- 1110 94. Jang, W. and Lim, J. (2009) A Numerical Study of PQL Estimation Biases in Generalized  
1111 Linear Mixed Models Under Heterogeneity of Random Effects. *Commun Stat-Simul C*, **38**,  
1112 692-702.
- 1113 95. Fong, Y.Y., Rue, H. and Wakefield, J. (2010) Bayesian inference for generalized linear  
1114 mixed models. *Biostatistics*, **11**, 397-412.
- 1115 96. Smith, A.F.M. and Roberts, G.O. (1993) Bayesian Computation Via the Gibbs Sampler  
1116 And Related Markov-Chain Monte-Carlo Methods. *J Roy Stat Soc B Met*, **55**, 3-23.
- 1117 97. Gelman, A. and Shirley, K. (2011) Inference from Simulations and Monitoring  
1118 Convergence. *Ch Crc Handb Mod Sta*, 163-174.
- 1119 98. Schwartz, L. (1965) On Bayes procedures. *Zeitschrift f{"u}r Wahrscheinlichkeitstheorie  
1120 und Verwandte Gebiete*, **4**, 10--26.
- 1121 99. Hadfield, J.D. (2010) MCMC Methods for Multi-Response Generalized Linear Mixed  
1122 Models: The MCMCglmm R Package. *J Stat Softw*, **33**, 1-22.
- 1123 100. Seyednasrollah, F., Laiho, A. and Elo, L.L. (2015) Comparison of software packages for  
1124 detecting differential expression in RNA-seq studies. *Brief Bioinform*, **16**, 59-70.
- 1125 101. Patterson, N., Price, A.L. and Reich, D. (2006) Population structure and eigenanalysis.  
1126 *Plos Genet*, **2**, 2074-2093.
- 1127 102. Yang, J., Zaitlen, N.A., Goddard, M.E., Visscher, P.M. and Price, A.L. (2014) Advantages  
1128 and pitfalls in the application of mixed-model association methods. *Nat Genet*, **46**, 100-  
1129 106.
- 1130 103. Rapaport, F., Khanin, R., Liang, Y.P., Pirun, M., Krek, A., Zumbo, P., Mason, C.E., Socci,  
1131 N.D. and Betel, D. (2013) Comprehensive evaluation of differential gene expression  
1132 analysis methods for RNA-seq datas. *Genome Biol*, **14**, R95.
- 1133 104. Chen, H.a.W., Chaolong and Conomos, Matthew P and Stilp, Adrienne M and Li, Zilin and  
1134 Sofer, Tamar and Szpiro, Adam A and Chen, Wei and Brehm, John M and Celed{\o}n,  
1135 Juan C and others. (2016) Control for Population Structure and Relatedness for Binary  
1136 Traits in Genetic Association Studies via Logistic Mixed Models. *The American Journal of  
1137 Human Genetics*, **98**, 653--666.
- 1138 105. Zhou, X., Cain, C.E., Myrthil, M., Lewellen, N., Michelini, K., Davenport, E.R., Stephens,  
1139 M., Pritchard, J.K. and Gilad, Y. (2014) Epigenetic modifications are associated with  
1140 inter-species gene expression variation in primates. *Genome Biol*, **15**, 547.
- 1141 106. Vawter, M.P., Evans, S., Choudary, P., Tomita, H., Meador-Woodruff, J., Molnar, M., Li, J.,  
1142 Lopez, J.F., Myers, R., Cox, D. *et al.* (2004) Gender-specific gene expression in post-  
1143 mortem human brain: Localization to sex chromosomes. *Neuropsychopharmacol*, **29**,  
1144 373-384.
- 1145 107. Lemos, B., Branco, A.T., Jiang, P.P., Hartl, D.L. and Meiklejohn, C.D. (2014) Genome-Wide  
1146 Gene Expression Effects of Sex Chromosome Imprinting in Drosophila. *G3-Genes Genom  
1147 Genet*, **4**, 1-10.
- 1148 108. Kim, J.H., Karnovsky, A., Mahavisno, V., Weymouth, T., Pande, M., Dolinoy, D.C., Rozek,  
1149 L.S. and Sartor, M.A. (2012) LRpath analysis reveals common pathways dysregulated via  
1150 DNA methylation across cancer types. *Bmc Genomics*, **13**.

- 1151 109. Vamsi K Mootha, C.M.L., Karl-Fredrik Eriksson, Aravind Subramanian, Smita Sihag,  
1152 Joseph Lehar, Pere Puigserver, Emma Carlsson, Martin Ridderstråle, Esa Laurila, Nicholas  
1153 Houstis, Mark J Daly, Nick Patterson, Jill P Mesirov, and Todd R Golub. (2003) PGC-  
1154 1alpha-responsive genes involved in oxidative phosphorylation are coordinately  
1155 downregulated in human diabetes. *Nat Genet*, **34**, 267 - 273.
- 1156 110. Leibowitz, G., Cerasi, E. and Ketzinel-Gilad, A. (2008) The role of mTOR in the adaptation  
1157 and failure of beta-cells in type 2 diabetes. *Diabetes Obes Metab*, **10**, 157-169.
- 1158 111. Ost, A., Svensson, K., Ruishalme, I., Brannmark, C., Franck, N., Krook, H., Sandstrom, P.,  
1159 Kjolhede, P. and Stralfors, P. (2010) Attenuated mTOR Signaling and Enhanced  
1160 Autophagy in Adipocytes from Obese Patients with Type 2 Diabetes. *Mol Med*, **16**, 235-  
1161 246.
- 1162 112. Laplante, M. and Sabatini, D.M. (2012) mTOR Signaling in Growth Control and Disease.  
1163 *Cell*, **149**, 274-293.
- 1164 113. Zoncu, R., Efeyan, A. and Sabatini, D.M. (2011) mTOR: from growth signal integration to  
1165 cancer, diabetes and ageing. *Nat Rev Mol Cell Bio*, **12**, 21-35.
- 1166 114. Matthews, D.R., Hosker, J.P., Rudenski, A.S., Naylor, B.A., Treacher, D.F. and Turner, R.C.  
1167 (1985) Homeostasis Model Assessment - Insulin Resistance And Beta-Cell Function From  
1168 Fasting Plasma-Glucose And Insulin Concentrations In Man. *Diabetologia*, **28**, 412-419.
- 1169 115. Lyssenko, V., Nagorny, C.L.F., Erdos, M.R., Wierup, N., Jonsson, A., Spiegel, P., Bugliani,  
1170 M., Saxena, R., Fex, M., Pulizzi, N. *et al.* (2009) Common variant in MTNR1B associated  
1171 with increased risk of type 2 diabetes and impaired early insulin secretion. *Nat Genet*, **41**,  
1172 82-88.
- 1173 116. Dupuis, J., Langenberg, C., Prokopenko, I., Saxena, R., Soranzo, N., Jackson, A.U.,  
1174 Wheeler, E., Glazer, N.L., Bouatia-Naji, N., Gloyn, A.L. *et al.* (2010) New genetic loci  
1175 implicated in fasting glucose homeostasis and their impact on type 2 diabetes risk. *Nat*  
1176 *Genet*, **42**, 105-U132.
- 1177 117. Gibbs, R.A., Belmont, J.W., Hardenbol, P., Willis, T.D., Yu, F.L., Yang, H.M., Ch'ang, L.Y.,  
1178 Huang, W., Liu, B., Shen, Y. *et al.* (2003) The International HapMap Project. *Nature*, **426**,  
1179 789-796.
- 1180 118. Gagnon-Bartsch, J.A. and Speed, T.P. (2012) Using control genes to correct for unwanted  
1181 variation in microarray data. *Biostatistics*, **13**, 539-552.
- 1182 119. Powell, J.E., Henders, A.K., McRae, A.F., Wright, M.J., Martin, N.G., Dermitzakis, E.T.,  
1183 Montgomery, G.W. and Visscher, P.M. (2012) Genetic control of gene expression in  
1184 whole blood and lymphoblastoid cell lines is largely independent. *Genome Res*, **22**, 456-  
1185 466.
- 1186 120. Cann, H.M., de Toma, C., Cazes, L., Legrand, M.F., Morel, V., Piouffre, L., Bodmer, J.,  
1187 Bodmer, W.F., Bonne-Tamir, B., Cambon-Thomsen, A. *et al.* (2002) A human genome  
1188 diversity cell line panel. *Science*, **296**, 261-262.
- 1189 121. Landi, M.T., Wang, Y.F., Mckay, J.D., Rafnar, T., Wang, Z.M., Timofeeva, M., Broderick, P.,  
1190 Stefansson, K., Risch, A., Chanock, S.J. *et al.* (2014) Imputation from The 1000 Genomes  
1191 Project identifies rare large effect variants of BRCA2-K3326X and CHEK2-I157T as risk  
1192 factors for lung cancer; a study from the TRICL consortium. *Cancer Res*, **74**, 942--942.
- 1193 122. Weigel, D. and Mott, R. (2009) The 1001 Genomes Project for Arabidopsis thaliana.  
1194 *Genome Biol*, **10**, 107.
- 1195 123. Lonsdale, J., Thomas, J., Salvatore, M., Phillips, R., Lo, E., Shad, S., Hasz, R., Walters, G.,  
1196 Garcia, F., Young, N. *et al.* (2013) The Genotype-Tissue Expression (GTEx) project. *Nat*  
1197 *Genet*, **45**, 580-585.



- 1198 124. Kanitz, A., Gypas, F., Gruber, A.J., Gruber, A.R., Martin, G. and Zavolan, M. (2015)  
1199 Comparative assessment of methods for the computational inference of transcript  
1200 isoform abundance from RNA-seq data. *Genome Biol*, **16**.
- 1201 125. Langmead, B. and Salzberg, S.L. (2012) Fast gapped-read alignment with Bowtie 2. *Nat*  
1202 *Methods*, **9**, 357-U354.
- 1203 126. Trapnell, C., Williams, B.A., Pertea, G., Mortazavi, A., Kwan, G., van Baren, M.J., Salzberg,  
1204 S.L., Wold, B.J. and Pachter, L. (2010) Transcript assembly and quantification by RNA-Seq  
1205 reveals unannotated transcripts and isoform switching during cell differentiation. *Nat*  
1206 *Biotechnol*, **28**, 511-U174.
- 1207 127. Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson,  
1208 M. and Gingeras, T.R. (2013) STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*,  
1209 **29**, 15-21.
- 1210 128. Kim, D., Pertea, G., Trapnell, C., Pimentel, H., Kelley, R. and Salzberg, S.L. (2013) TopHat2:  
1211 accurate alignment of transcriptomes in the presence of insertions, deletions and gene  
1212 fusions. *Genome Biol*, **14**.
- 1213 129. Bray, N.L., Pimentel, H., Melsted, P. and Pachter, L. (2016) Near-optimal probabilistic  
1214 RNA-seq quantification. *Nat Biotechnol*, **34**, 525-527.
- 1215 130. Conesa, A., Madrigal, P., Tarazona, S., Gomez-Cabrero, D., Cervera, A., McPherson, A.,  
1216 Szczesniak, M.W., Gaffney, D.J., Elo, L.L., Zhang, X. *et al.* (2016) A survey of best  
1217 practices for RNA-seq data analysis. *Genome Biol*, **17**, 13.
- 1218 131. Trapnell, C., Hendrickson, D.G., Sauvageau, M., Goff, L., Rinn, J.L. and Pachter, L. (2013)  
1219 Differential analysis of gene regulation at transcript resolution with RNA-seq. *Nat*  
1220 *Biotechnol*, **31**, 46-+.
- 1221 132. Oshlack, A. and Wakefield, M.J. (2009) Transcript length bias in RNA-seq data confounds  
1222 systems biology. *Biology direct*, **4**, 14.
- 1223 133. Hansen, K.D., Irizarry, R.A. and Wu, Z.J. (2012) Removing technical variability in RNA-seq  
1224 data using conditional quantile normalization. *Biostatistics*, **13**, 204-216.
- 1225 134. Kang, H.M., Zaitlen, N.A., Wade, C.M., Kirby, A., Heckerman, D., Daly, M.J. and Eskin, E.  
1226 (2008) Efficient control of population structure in model organism association mapping.  
1227 *Genetics*, **178**, 1709-1723.
- 1228 135. Venables, W.N.a.R., B. D. (2002) *Modern Applied Statistics with S*. Springer, New York.
- 1229

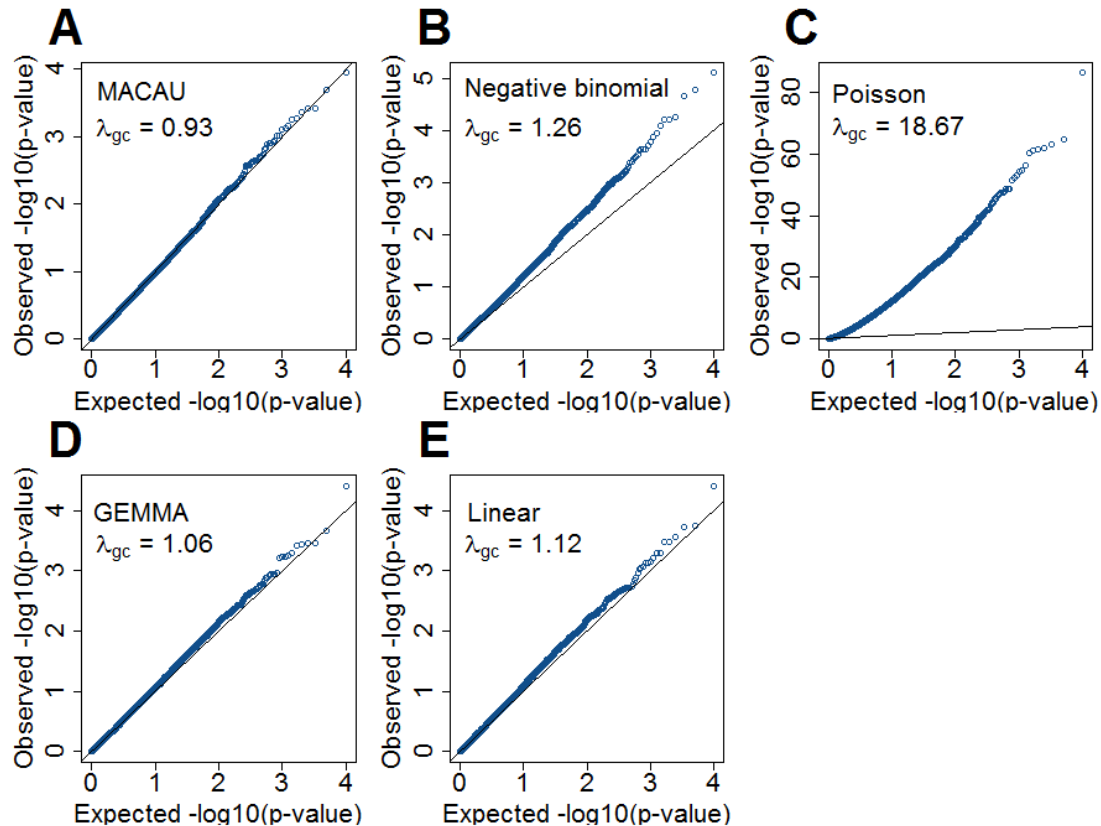
1230 **Table 1. Current approaches for identifying differentially expressed genes**  
1231 **in RNAseq.**

<b>Statistical method</b>	<b>Directly models counts?</b>	<b>Controls for biological covariates?</b>	<b>Controls for sample non-independence?</b>	<b>Example software that implements the method</b>
<b>Linear regression</b>	No	Yes	No	R and many others
<b>Linear mixed model</b>	No	Yes	Yes	GEMMA (9) and EMMA (134)
<b>Poisson model</b>	Yes	Some methods do	No	GLMP (135) and DEGseq (20)
<b>Negative binomial model</b>	Yes	Some methods do	No	edgeR (25), DESeq (15) and GLMNB(135)
<b>Poisson mixed model</b>	Yes	Yes	Yes	MACAU

1232

1233

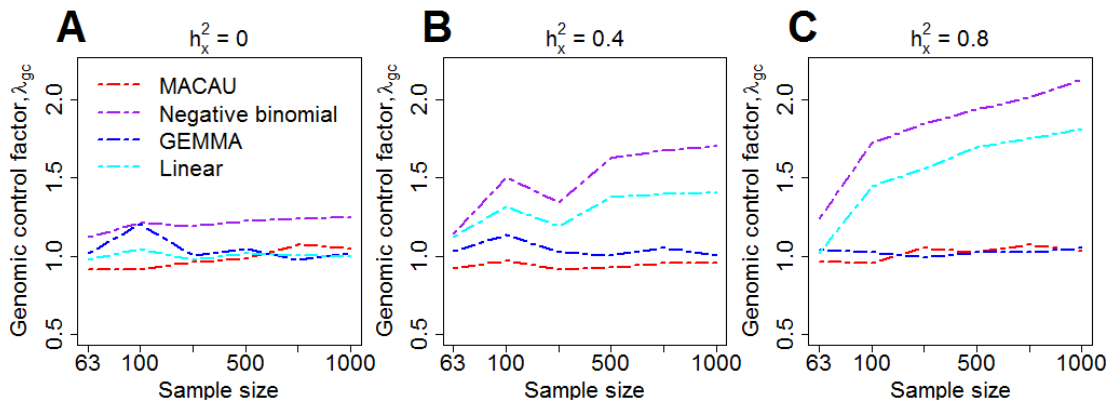
1234 **Figure 1. QQ-plots comparing expected and observed  $p$ -value distributions**  
1235 **generated by different methods for the null simulations in the presence of**  
1236 **sample non-independence.** In each case, 10,000 non-DE genes were  
1237 simulated with  $n = 63$ ,  $CV = 0.3$ ,  $\sigma^2 = 0.25$ ,  $h^2 = 0.6$  and  $h_x^2 = 0.4$ . Methods for  
1238 comparison include MACAU (A), Negative binomial (B), Poisson (C), GEMMA  
1239 (D), and Linear (E). Both MACAU and GEMMA properly control for type I error  
1240 well in the presence of sample non-independence.  $\lambda_{gc}$  is the genomic control  
1241 factor.



1242

1243

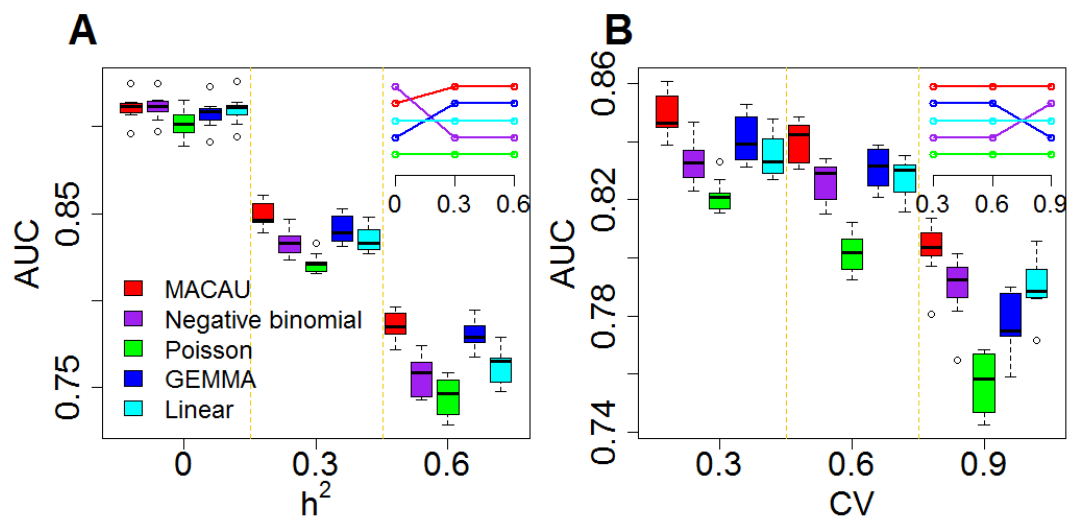
1244 **Figure 2. Comparison of the genomic control factor  $\lambda_{gc}$  from different**  
1245 **methods for the null simulations in the presence of sample non-**  
1246 **independence.** 10,000 null genes were simulated with  $CV = 0.3$ ,  $\sigma^2 = 0.25$ ,  $h^2 =$   
1247  $0.6$ , and (A)  $h_x^2 = 0$ ; (B)  $h_x^2 = 0.4$ ; or (C)  $h_x^2 = 0.8$ .  $\lambda_{gc}$  (y-axis) changes with  
1248 sample size  $n$  (x-axis). Methods for comparison were MACAU (red), Negative  
1249 binomial (purple), GEMMA (blue), and Linear (cyan). Both MACAU and GEMMA  
1250 provide calibrated test statistics in the presence of sample non-independence  
1251 across a range of settings.  $\lambda_{gc}$  from Poisson exceeds 10 in all settings and is  
1252 thus not shown.



1253

1254

1255 **Figure 3. MACAU exhibits increased power to detect true positive DE genes**  
1256 **across a range of simulation settings.** Area under the curve (AUC) is shown  
1257 as a measure of performance for MACAU (red), Negative binomial (purple),  
1258 Poisson (green), GEMMA (blue), and Linear (cyan). Each simulation setting  
1259 consists of 10 simulation replicates, and each replicate includes 10,000  
1260 simulated genes, with 1,000 DE and 9,000 non-DE. We used  $n = 63$ ,  $h_x^2 = 0.0$ ,  
1261  $PVE = 0.25$ ,  $\sigma^2 = 0.25$ . In (A) we increased  $h^2$  while maintaining  $CV = 0.3$  and in  
1262 (B) we increased  $CV$  while maintaining  $h^2 = 0.3$ . Boxplots of AUC across  
1263 replicates for different methods show that (A) heritability ( $h^2$ ) influences the  
1264 relative performance of the methods that account for sample non-independence  
1265 (MACAU and GEMMA) compared to the methods that do not (negative binomial,  
1266 Poisson, linear); (B) variation in total read counts across individuals, measured  
1267 by the coefficient of variation ( $CV$ ), influences the relative performance of  
1268 GEMMA and negative binomial. Insets in the two figures show the rank of  
1269 different methods, where the top row represents the highest rank.



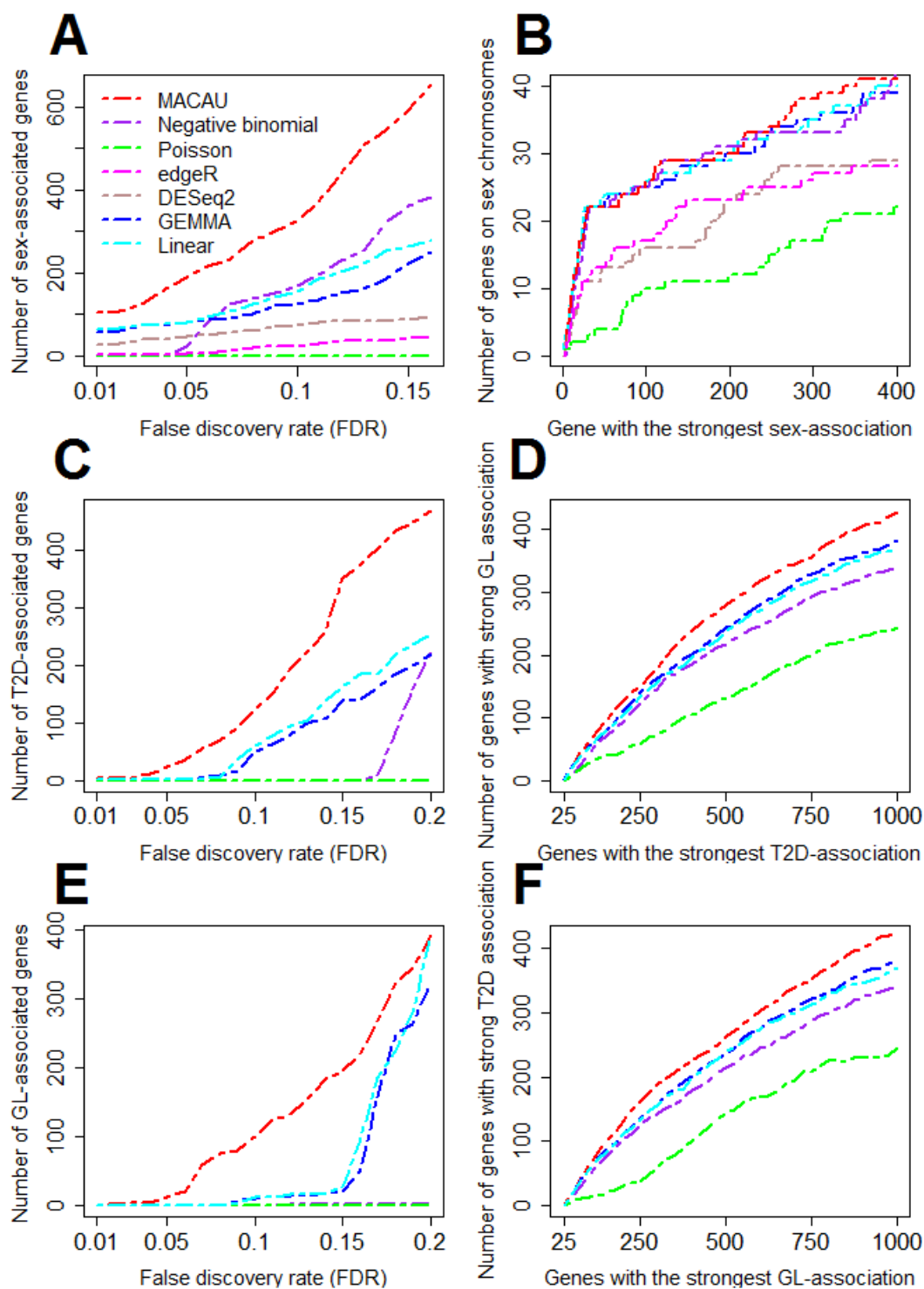
1270

1271

1272 **Figure 4. MACAU identifies more differentially expressed genes than other**  
1273 **methods in the baboon (panels A and B) and FUSION (panels C, D, E, and F)**  
1274 **data sets.** Methods for comparison include MACAU (red), Negative binomial  
1275 (purple), Poisson (green), edgeR (magenta), DESeq2 (rosybrown), GEMMA  
1276 (blue), and Linear (cyan). (A) shows the number of sex-associated genes  
1277 identified by different methods at a range of empirical false discovery rates  
1278 (FDRs). (B) shows the number of genes that are on the X chromosome out of the  
1279 genes that have the strongest sex association for each method (note that the Y  
1280 chromosome is not assembled in baboons and is thus ignored). For instance, in  
1281 the top 400 genes identified by MACAU, 41 of them are also on the X  
1282 chromosome. (C) shows the number of T2D-associated genes identified by  
1283 different methods at a range of empirical false discovery rates (FDRs). (D) shows  
1284 the number of genes that are in the list of top 1,000 genes most significantly  
1285 associated with GL out of the genes that have the strongest association for T2D  
1286 for each method. For instance, in the top 1,000 genes with the strongest T2D  
1287 association identified by MACAU, 428 of them are also in the list of top 1,000  
1288 genes with the strongest GL association identified by the same method. (E)  
1289 shows the number of GL-associated genes identified by different methods at a  
1290 range of FDRs. (F) shows the number of genes that are in the list of top 1,000  
1291 genes most significantly associated with T2D out of the genes that have the  
1292 strongest association for GL for each method. T2D: type II diabetes; GL: fasting  
1293 glucose level.

1294





1295

1296