

1 **scBASE: A Bayesian mixture model for the analysis of allelic** 2 **expression in single cells**

3 Kwangbom Choi¹, Narayanan Raghupathy¹ & Gary A. Churchill^{1,*}

4 ¹*The Jackson Laboratory, 600 Main Street, Bar Harbor, Maine, 04609*

5 **Allele-specific expression (ASE) at single-cell resolution is a critical tool for un-**
6 **derstanding the stochastic and dynamic features of gene expression. However,**
7 **low read coverage and high biological variability present challenges for analyz-**
8 **ing ASE. We propose a new method for ASE analysis from single cell RNA-Seq**
9 **data that accurately classifies allelic expression states and improves estimation**
10 **of allelic proportions by pooling information across cells.**

11 Single-cell RNA sequencing (scRNA-Seq) can reveal features of cellular gene expression
12 that cannot be observed in bulk RNA sequencing¹. Allelic imbalance is common across
13 many genes² and can range from a subtle imbalance to complete monoallelic expression as
14 in imprinted genes³ or genes under dosage compensation by X chromosome inactivation^{4,5}.
15 Allele-specific expression (ASE) in single cells can provide a rich picture of the stochastic
16 and dynamic properties of gene expression in individual cells. Analysis of single-cell ASE
17 poses unique challenges due to the low depth of sequencing coverage per cell⁶⁻¹³. In addition,
18 allelic proportions often form U-shaped or W-shaped distributions due to the occurrence of
19 monoallelic transcriptional bursts.

20 We propose a novel method for the estimation of single-cell allele proportions, **scBASE**,

21 in which we (1) disambiguate and count multi-mapping reads (multi-reads); (2) classify each
22 gene in each cell into paternal monoallelic, bi-allelic, or maternal monoallelic expression
23 states; and (3) address data sparsity by partial pooling of information across cells (Figure
24 1 and Methods). The *counting* step of **scBASE** applies an estimation-maximization (EM)
25 algorithm to count multi-reads by *weighted allocation* to estimate expected read counts^{14–16}.
26 The *classification* and *estimation* steps are iterated and together achieve *partial pooling* of
27 information among cells that are in the same allelic expression states. In the classification
28 step, we compute the posterior probabilities of paternal, bi-allelic, and maternal expression
29 states. In the estimation step we compute the posterior distributions of cell- and gene-
30 specific allelic proportions. Read counting and partial pooling can be applied together,
31 separately, or not at all. This leads to four different methods of estimating allelic propor-
32 tions. In the unique reads methods (i), we estimate allelic proportions directly from the
33 counts of uniquely mapping reads. In the weighted allocation method (ii), we apply the
34 read counting step of **scBASE** to obtain estimated expected counts. We can apply the clas-
35 sification and estimation steps to either of these counts to obtain allelic proportions from
36 unique reads with partial pooling (iii), or weighted allocation with partial pooling (iv). We
37 have implemented these methods in extensible open-source software, **scBASE**, available at
38 <https://github.com/churchill-lab/scBASE>.

39 In the following sections, we first examine the effects of weighted allocation on single-
40 cell allele expression data. Then we evaluate the effects of partial pooling on estimation of
41 allelic proportions. We then apply each of the four methods in **scBASE** to statistical testing
42 of independence of allelic bursting. Finally, we illustrate the interpretive power of allelic

43 expression by analysis of scRNA-Seq data from a development time course⁷.

44 [Figure 1 about here.]

45 Results

46 We applied scBASE to scRNA-Seq data from 286 pre-implantation mouse embryo cells from
47 an F1 hybrid mating between female CAST/EiJ (CAST) and male C57BL/6J (B6) mice⁷.
48 Cells were sampled along a time course from the zygote and early 2-cell stages through the
49 late blastocyst stage of development. We created a diploid transcriptome from CAST- and
50 B6-specific sequences of each annotated transcript (Ensembl Release 78)¹⁷ and aligned reads
51 from each cell to obtain allele-specific alignments. In order to ensure that genes had sufficient
52 polymorphic sites for ASE analysis, we restrict attention to 13,032 genes that had at least 4
53 allelic unique reads in at least 10% of cells. Where indicated below, we apply scBASE to only
54 122 cells from the blastocyst stages of development, or to only 60 cells in the mid-blastocyst
55 stage.

56 We first assessed the impact of weighted allocation of multi-reads on the estimation
57 of allelic proportions. Any read that maps to one allele of one gene is a unique read. All
58 other reads are multi-reads and they can be simple or complex. A read that maps uniquely
59 to one gene but to both allelic copies is a simple allelic multi-read. A read that maps to
60 multiple genes but only to one allele at each is a simple genomic multi-read. A read that
61 maps to multiple genes as well as to both alleles in any of those genes are complex multi-
62 reads. There are, in total, 9 patterns of simple and complex multi-read alignments for two

63 genomic loci and two alleles (Supplemental Figure S1). We estimated unique reads and
64 weighted allocation counts from each individual cell using all 286 cells to show how the
65 number of monoallelic genes changes in each cell (Figure 2a). The sequence reads from these
66 cells include 2.5% simple genomic multi-reads, 59.3% simple allelic multi-reads, and 23.3%
67 complex multi-reads. In a typical scRNA-Seq workflow for ASE, these reads are discarded
68 leaving only the unique 14.9% of the original sequence reads for analysis. This substantial
69 loss of information could lead to high variability of allelic proportions and spurious findings
70 of monoallelic gene expression. We find that using only uniquely mapping reads generates a
71 higher rate of monoallelic expression calls (Figure 2a and Supplemental Figure S1), calling
72 on average ~ 66 more genes with monoallelic expression in each cell. We also observed,
73 on average, $\sim 1,908$ genes where the unique reads method fails to call bi-allelic expression
74 compared to weighted allocation, for example, *Mtdh* (Figures 2b and 2c). These genes are
75 consistently bi-allelic in many cells according to weighted allocation, but their pattern of
76 allelic expression based on unique reads can be misinterpreted as monoallelic expression
77 and, as a result, allelic bursting appears to be more dynamic.

78 [Figure 2 about here.]

79 Next we evaluated the impact of partial pooling on the estimation of allelic proportions.
80 Since it is best to apply partial pooling to each cell type separately, we focus attention on
81 the 122 mature blastocyst cells, the largest group in Deng et al.⁷ data. These cells have the
82 coverage of ~ 14.8 M reads per cell in average, and we down-sampled these data by randomly
83 selecting 1% of reads to obtain an average coverage of ~ 148 k reads per cell. We estimated

84 allelic proportions using each of four methods: (i) unique reads, (ii) weighted allocation,
85 (iii) unique reads with partial pooling, and (iv) weighted allocation with partial pooling.
86 We compared the estimated allelic proportions from the down-sampled data to estimates
87 obtained from the full data using the corresponding unique reads or weighted allocation
88 estimates with no pooling. The full data are based on 100-fold more reads per sample and
89 provide an approximate truth standard. A similar approach to evaluation of single-cell data
90 analysis was employed by Huang et al.¹⁸. In order to assess the effects of partial pooling, we
91 computed differences in the mean squared error (MSE) of estimated allelic proportions with
92 and without partial pooling. Partial pooling applied to the unique read counts improves
93 estimation for the majority of genes (4,392 versus 1,367 out of 5,759 genes) with an average
94 MSE difference of 0.018 (Figure 3a). Partial pooling applied to the weighted allocation
95 counts improves estimation for most genes (5,078 versus 1,673 out of 6,751 genes) with an
96 average MSE difference of 0.016 (Figure 3b). In both cases, the greatest gains are seen in
97 the low expression range (<10 unique reads per gene). For the most highly expressed genes,
98 there is little or no reduction in MSE, which is consistent with our expectation that pooling
99 of information across cells is most impactful when coverage is low.

100 [Figure 3 about here.]

101 The timing of allelic bursting events is a defining feature of stochasticity in gene
102 expression¹⁹. One fundamental question is whether the occurrence of allelic bursts is coordi-
103 nated or if bursts occur independently for each allele. Statistical independence of maternal
104 and paternal bursting can be evaluated using a 2×2 table of counts of the numbers of cells for

105 which a given gene is expressed only from the maternal allele, only from the paternal allele,
106 from both, or not expressed (as in Figure 2c). If allelic bursts occur independently, the log-
107 odds ratio (logOR) computed from this 2×2 table should be close to zero. In order to relate
108 this standard approach²⁰ for testing the independence hypothesis to alternative methods^{7,21}
109 that have been proposed for scRNA-Seq data, it is helpful to consider a geometric represen-
110 tation of the proportions of cells in each allelic expression state (Figure 4a). Proportions
111 are numbers greater than or equal to zero that sum to one. They can be represented as
112 a point in a 3-dimensional tetrahedron in 4-dimensional space – the 4D simplex²². When
113 maternal and paternal bursting events occur independently, the proportions should fall near
114 the 2-dimensional surface within the simplex where the logOR is equal to zero (cross-hatched
115 region in Figure 4a). The method of testing independence used in Deng et al.⁷ and Larsson
116 et al.²¹ imposes an additional constraint on the 2×2 table proportions by assuming that the
117 frequencies of maternal and paternal bursting events are equal ($p_M = p_P$). This constraint
118 corresponds to a 2-dimensional cross-section of the simplex, indicated by the blue triangle in
119 Figure 4a. Projection of points in the 4D simplex onto this triangle produces the graphical
120 representation used by Deng et al. (e.g., Figure 4b). This illustrates how the Deng et al.
121 method is a special case of the logOR test.

122 We evaluated bursting independence on the 122 mature blastocyst cells as was done in
123 Jiang et al.²³. We first simulated data under the assumption of independent allelic bursting
124 (Methods) and plotted the results to illustrate how points will be distributed in this diagram
125 when the pure independence model is true with and without the constraint of $p_M = p_P$
126 (Figure 4b). Next we estimated the 2×2 table proportions of allelic expression states using

127 each of the four methods (i~iv) implemented in **scBASE**. The appearance of the data in
128 Figures 4c is qualitatively distinct from the simulated data (Figure 4b). Moreover, the
129 null hypothesis of independence is rejected by the method used in Jiang et al.²³ for the
130 majority of genes regardless of the method used to estimate the allelic state proportions
131 (Supplemental Figure S2a). **SCALE** reports 3,381 genes that are non-independent using the
132 results of unique-reads method, 4,815 genes using weighted allocation, 6,068 genes by unique
133 reads with partial pooling, and 6,761 genes based on weighted allocation with partial pooling
134 at the FDR level of 5%. Similarly the logOR are away from 0 for thousands of genes. For
135 example, 2,845 and 3,763 out of 8,290 genes had $|\log\text{OR}|>2$ using unique reads and weighted
136 allocation. More genes have $|\log\text{OR}|>2$ after partial pooling with **scBASE**: 5,622 and 6,209
137 respectively. The majority of genes had positive logOR, indicating a tendency for bursting to
138 occur more in synchrony than chance would predict (Supplemental Figure S2b). We repeated
139 this analysis using three additional data sets^{21,24,25} and arrive at similar conclusions in each
140 case (Supplemental Figures S3, S4, and S5). The evidence for statistical dependence of
141 bursting is strong and application of weighted allocation and partial pooling strengthens
142 this conclusion.

143 [Figure 4 about here.]

144 The **scBASE** classification step provides a novel way to characterize allelic imbalance
145 across a population of cells by estimating the expected proportions of cells in different tran-
146 scriptional states. Using **scBASE**, we can compute the posterior probability of allelic ex-
147 pression states of genes in each cell. This probabilistic classification allows for uncertainty

148 associated with statistical sampling from the pool of transcripts that are present in the cell
149 including the occurrence of zero read counts. Based on the posterior probabilities, we can
150 derive the expected proportions of cells in states P, B, and M, which can be represented as
151 point in a triangular simplex diagram. (Note that this representation is a projection of points
152 in the 4D simplex onto the bottom triangular region, Figure 4a.) The classification step of
153 scBASE assumes that all genes are expressed at some level, which may be very low for some
154 genes. This allows us to classify the allelic expression of cells that may have zero read counts
155 due to statistical sampling. To interpret the distribution of allelic expression across cells,
156 we designate seven patterns of allelic expression (Figure 5a). Genes that are predominantly
157 expressed as P, B, or M will appear near the corresponding vertex of the triangle (**P**, **B** or
158 **M** region). Genes with mixed allelic states will appear along the edges (**PB**, **BM**, or **MP**
159 region) or near the center of the triangle (all three states, **PBM** region). For example, the
160 gene *Pacs2*, which is expressed from either the maternal or the paternal allele but rarely
161 both, is classified as an **MP** gene. The bi-allelic region (**B**) includes genes that are consis-
162 tently expressed from both alleles e.g., *Mtdh*. The **PB** and **BM** regions include genes that
163 show a mixture of bi-allelic and monoallelic expression with a strong allelic imbalance, e.g.,
164 *Timm23* and *Tulp3*. The majority of genes (56.9%) in the blastocyst stages of development
165 are in the **PBM** region (Supplemental Figure S6). These genes display a mix of mono- and
166 bi-allelic expression states (e.g., *Akr1b3*) that is consistent with dynamic allele-specific gene
167 expression with a low bursting rate relative to mRNA half life.

168

[Figure 5 about here.]

169 We applied **scBASE** (with weighted allocation and partial pooling) to track changes in
170 the ASE patterns of cells sampled over a developmental time course (Figure 5b, Supplemental
171 Figure S6 and S7). Our aim is to classify allelic state distributions within subpopulations of
172 cells defined by developmental stages. To achieve this, we first ran **scBASE** MCMC algorithm
173 on all 286 cells to estimate the prior parameters, α_g^s and β_g^s (Figure 1 and Supplemental
174 Methods). These parameters describe the distribution of allelic proportions in each allelic
175 state. According to our diagnostic criteria, **scBASE** MCMC algorithm produced reliable
176 parameter estimation for 10,017 out of 13,032 genes. We then ran **scBASE** EM algorithm
177 (with the prior parameters fixed) on each subpopulation of cells to estimate developmental
178 stage-specific parameters (Details are provided in Methods.). In the zygote and early 2-cell
179 stages, essentially all genes show monoallelic maternal expression. At this stage, the hybrid
180 embryo genome is not being transcribed and the mRNA present is derived from the mother
181 (inbred CAST genome). At the mid 2-cell stage the hybrid embryo is being transcribed and
182 we start to see expression of the paternal allele for some genes. Many genes exhibit the **M**
183 and **BM** patterns through the 8- or 16-cell stages perhaps due to the persistence of long-lived
184 mRNA species that were present at the 2-cell stage. The bi-allelic class **B** dominates the late
185 2-cell and 4-cell stages indicating high levels of expression at rates that exceed the half-life
186 of most mRNA species. In the later stages of development, 8-cell through late blastocyst,
187 most genes transition into the **PBM** pattern.

188 There are ~ 400 genes that make dramatic transitions across allelic expression states.
189 For example, *Akr1b3* (Figure 5c) starts in the zygote and early 2-cell stage with only ma-
190 ternal alleles present. It transitions to bi-allelic expression by the mid 2-cell stage indicating

191 the onset of transcription of the paternal allele. It then transitions through the paternal
192 monoallelic state. Our interpretation is that the early maternally derived transcripts were
193 present prior to fertilization and these transcripts are still present when the paternal allele
194 in the hybrid embryo gene starts to express. The early maternal transcripts are largely de-
195 graded by the 4- to 8-cell stages where we see only expression from the paternal allele. In
196 the early blastocyst stages, we start to see embryonic expression of maternal alleles resulting
197 in a bi-allelic expression pattern by the late blastocyst stage.

198 Discussion

199 Allelic expression in single cells has provided new insights into the dynamic regulation of gene
200 expression²⁴. However, estimates of allelic proportions can display high statistical variation
201 due to low depth of sequencing coverage per cell. The common practice of discarding multi-
202 mapping reads exacerbates this problem. The scBASE algorithm reduces statistical variability
203 by retaining and disambiguating multi-read data. It further improves estimation of allelic
204 proportions by partial pooling of information across cells in the same ASE states. As a
205 result we can obtain a more precise and accurate picture of gene expression dynamics in
206 which biological stochasticity is revealed by reducing statistical variation.

207 Weighted allocation has been demonstrated to improve gene expression estimation in
208 whole-tissue RNA-Seq¹⁴⁻¹⁶. When estimating total gene expression with weighted allocation,
209 only genomic multi-reads need to be resolved and these typically represent a small proportion
210 of all reads. When estimating allele-specific expression, however, depending on the levels

211 of nucleotide heterozygosity, the majority of reads may lack distinguishing polymorphisms
212 and will be allelic multi-reads. Complex multi-reads with ambiguity in both genomic and
213 allelic alignment can carry useful information about allele-specific expression, as illustrated
214 in Supplemental Figure S1.

215 **scBASE** uses partial pooling in the context of a mixture model with three allelic expres-
216 sion states (paternal monoallelic, bi-allelic, and maternal monoallelic) to preserve cell-to-cell
217 heterogeneity by pooling information across cells that are in the same state. Combining
218 information across cells, therefore, does not weaken the signals of strong allelic imbalance.
219 We applied **scBASE** to X chromosome genes in female cells of three different data sets^{7,24,25}.
220 In the Reinius et al. fibroblast data, partial pooling corrected the allelic proportions of Xist
221 gene expression towards either maternal or paternal monoallelic expression for both unique
222 reads and weighted allocation counts (Supplemental Figure S9a). Looking at expression of
223 all X chromosome genes in these same cells, we observe that partial pooling strengthens the
224 expected pattern of expression due to X chromosome inactivation (XCI) consistent with Xist
225 allele expression (Supplemental Figure S9b). We observe that XCI is often incomplete and
226 not uniform across cells. In the Chen et al. and Deng et al. data sets, Xist is clearly in
227 the bi-allelic expression state in many of mouse embryo cells, epistem cells, or motor neu-
228 ron cells and this is preserved after partial pooling. We also observe that XCI is not fully
229 established for these cells (Supplemental Figure S10, and S11). In addition, for genes that
230 are reported to be imprinted²⁶⁻²⁸ we examined their allelic expression. Irrespective of the
231 estimation method applied, many of these genes do not appear to be fully imprinted in these
232 three data sets (Supplementary Figure S12 and S13). However, for those genes that do show

233 evidence of imprinting, i.e., appear in **M**- or **P**-class, partial pooling improves the evidence
234 for monoallelic expression for both unique reads and weighted allocation counts.

235 The **scBASE** analysis incorporates statistical uncertainty in both the classification of
236 allelic expression state and the estimated allelic proportions of a gene. To evaluate the pre-
237 cision of the estimated parameters, we have computed the posterior standard deviation of
238 allele proportions across a range of total read counts and with varying numbers of cells (286
239 cells versus 60 cells). The trends are as expected, deeper read coverage or more cells im-
240 proves the precision of estimation (Supplemental Figure **S8**). Our probabilistic classification
241 accounts for uncertainty and can estimate the allelic expression state of a gene even when
242 few or no reads are sampled from a given cell based on the behavior of other cells. The
243 **scBASE** model is still reliable with degenerate inputs, for example, in the most extreme case
244 of a single cell and a gene with zero total reads, the algorithm provides a sensible answer:
245 class probabilities are $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ and a nearly uniform distribution for allelic proportion (mean
246 at 0.5 with standard deviation of 0.2), indicating that the data does not contain any infor-
247 mation. As the number of cells or the read depth increases, the class probabilities become
248 more concentrated and the posterior distribution for the allelic proportion gets narrower.
249 Partial pooling has the biggest impact when read coverage is low and the number of cells is
250 large (Figure **3** and Supplemental Figure **S8**).

251 **scBASE** software can be implemented as part of a scRNA-Seq analysis pipeline. For
252 example, we applied **SCALE** software using counts based on four methods: (i) unique reads,
253 (ii) weighted allocation, (iii) unique reads with partial pooling, and (iv) weighted allocation

254 with partial pooling implemented in `scBASE`. We found that more genes appeared to be
255 non-independent when weighted allocation-based counts are used in `SCALE`. Even more
256 genes were identified as non-independent using counts based on partial pooling (Results and
257 Supplemental Figure [S2a](#)). Although it is not mentioned in Jiang et al.²³, a substantial
258 number (3,485 at FDR=5%) of genes were identified as non-independent using the allelic
259 counts (unique reads) reported by Deng et al. Our findings suggest that running `SCALE` with
260 `scBASE` estimated read counts as input will result in more accurate estimates of bursting
261 kinetics and reduced levels of monoallelic gene expression when compared to results obtained
262 using unique read counts.

263 The statistical properties of allelic bursting shed light on the nature of gene expres-
264 sion regulation. If expression bursts are statistically independent, this would imply that
265 the regulation of allelic expression is local and acting autonomously at each allele. Under
266 the perfect independence model, there would be no shared regulation of expression across
267 alleles and the counts of cells in each allelic state will satisfy statistical criteria for inde-
268 pendence. Under an alternative model, perfect dependence, bursting would be precisely
269 coordinated across alleles and bursts would occur synchronously. All cells would be in ei-
270 ther the bi-allelic or not expressed states. Our analysis of published scRNA-Seq data from
271 four different experiments^{7,21,24,25} indicates that neither of these extremes is true (Figure 4
272 and Supplemental Figure [S2](#), [S3](#), [S4](#), and [S5](#)). We observed that the pattern of bursting is
273 statistically dependent and positively correlated ($\log\text{OR} > 0$) for the majority of genes. It
274 is neither statistically independent nor perfectly synchronous. This suggests that regulation
275 of allelic expression has both shared and locally autonomous components. While our statis-

276 tical analysis cannot identify the mechanisms of regulation, it seems plausible that diffusible
277 transcription factors could be responsible for the coordinated component of regulation. Lo-
278 cal control is likely to be cis-acting and may involve stochastic variation in the activation
279 of the transcriptional machinery. Additional experimental work would be required to test
280 these hypotheses and to identify the cis-acting molecular events that trigger bursting of gene
281 expression. However, the available data are sufficient to reject both hypotheses of perfect
282 independence and of perfect dependence of allelic bursting.

283 When estimating parameters associated with many genomic features in each of many
284 individual cells, one can improve the estimated parameters by pooling information across
285 cells. The motivation behind partial pooling is that the individual estimates are unbiased but
286 lack precision whereas the average provides a precise but biased estimate for individual cells
287 and also masks cell to cell heterogeneity entirely. Weighted allocation of multi-mapping reads
288 is not just to avoid information loss but is effective to prevent possible bias due to the genomic
289 multi-reads that contain allele information. For these reasons, we generally recommend the
290 strategy (iv) weighted allocation with partial pooling. But we provide all four options in
291 **scBASE** so more evaluation could be performed in other contexts. These general principles
292 — retention of multi-mapping sequence reads and partial pooling of information across cells
293 — apply broadly to analysis of genomic sequencing data but they are especially critical in
294 single cell applications where the observed numbers of reads for each gene in each cell may
295 be very small.

296 Methods

297 **Data.** Deng et al.⁷ sampled 286 pre-implantation embryo cells from an F1 hybrid of CAST×B6
298 along the stages of prenatal development. Embryos were manually dissociated into single cells
299 using Invitrogen TrypLE and single-end RNA-Seq sequencing was performed using Illumina
300 HiSeq 2000 (Platform GPL12112). We downloaded the data, Series GSE45719 from Gene Ex-
301 pression Omnibus (GEO) at <http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE45719>.
302 There were fastq-format read files for 4 single-cell samples from zygote stage, 8 from early
303 2-cell, 12 from mid 2-cell, 10 from late 2-cell, 14 from 4-cell, 47 from 8-cell, 30 from 16-cell,
304 43 from early blastocyst, 60 from mid blastocyst, and 58 from late blastocyst stage. The
305 Reinius et al. data²⁴ consist of primary mouse fibroblast cells from the F1 reciprocal crosses
306 of CAST×B6 (125 cells, sex-typed) and B6×CAST (113 cells, sex-typed), available from
307 GEO at <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE75659>. The Chen et al.
308 data²⁵ are from mouse embryonic stem cells (mESCs) from an F1 hybrid of B6×CAST: 111
309 mESCs cultured in 2i and LIF, 120 mESCs cultured in serum and LIF, 183 mouse Epistem
310 cells (mEpiSCs), and 74 post-mitotic neuron cells. The samples are sex-typed. We down-
311 loaded SRA format files available at <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE74155>.
312 Larsson et al.²¹ generated 224 individual primary mouse fibroblast cells from the F1 hybrid of
313 CAST×B6. As the data are from non-standard SMART-Seq2 platform, we downloaded the
314 allele-specific UMI counts from <https://github.com/sandberg-lab/txburst/tree/master/data>
315 (as of April 19th, 2019), and we were unable to apply weighted allocation to these data.

316 **scRNA-Seq read alignment.** For the F1 hybrid mouse we aligned reads to a phase-known
317 diploid transcriptome – this is a best-case scenario for phasing. When dealing with more
318 complex genomes, phasing should be performed beforehand if haplotype-specific transcrip-
319 tomes are not available and `scphaser`²⁹ is one possible approach. We reconstructed the
320 CAST genome by incorporating known SNPs and short indels (Sanger REL-1505) into the
321 reference mouse genome sequence (Genome Reference Consortium Mouse Reference 38) using
322 `g2gtools` (<http://churchill-lab.github.io/g2gtools/>). We lifted the reference gene annotation
323 (Ensembl Release 78) over to the CAST genome coordinates, and derived a CAST-specific
324 transcriptome. The B6 transcriptome is based on the mouse reference genome. We con-
325 structed a bowtie (v1.0.0) index to represent the diploid transcriptome with two alleles of
326 each transcript. We aligned reads using bowtie with parameters ‘-all’, ‘-best’, and ‘-strata’,
327 allowing for 3 mismatches (‘-v 3’). These settings enable us to find all of the best alignments
328 for each read. For example, if there is a zero-mismatch alignment for a read, all alignments
329 with zero mismatch will be accepted.

330 **Overview of the scBASE model.** The scBASE algorithm is composed of three steps:
331 *read counting*, *classification*, and *estimation* (Figure 1). The read counting step is applied
332 first to resolve read mapping ambiguity due to multi-reads and to estimate expected read
333 counts. The read counting step is not a requirement since the following steps are applicable
334 to any allele-specific count estimates. The classification and estimation steps are executed
335 iteratively to classify the allelic expression state and to estimate the allelic proportions for
336 each gene in each cell using a hierarchical mixture model. We have implemented scBASE as a
337 Monte Carlo Markov chain (MCMC) algorithm³⁰, which randomly samples parameter values

338 from their conditional posterior distributions. We have also implemented the classification
339 and estimation steps as an Expectation-Maximization (EM) algorithm³¹ that converges to
340 the maximum a posteriori parameter estimates (Supplemental Methods). MCMC is flexible,
341 and the sampling distributions and priors are easy to change in the MCMC code. MCMC
342 provides the full posterior distribution of allelic proportions and thus provides useful in-
343 formation about the uncertainty of estimated parameters. We also found that MCMC is
344 more stable when fitting allelic proportion of monoallelic classes. The EM algorithm is much
345 faster, but it provides only point estimation. We provide a brief description of the algorithm
346 here and provide additional details in Supplemental Methods.

347 ***Read counting:*** In order to count all of the available sequence reads for each gene and
348 allele, we have to resolve read mapping ambiguity that occur when aligning reads to a diploid
349 genome. Genomic multi-reads align with equal quality to more than one gene. Allelic multi-
350 reads align with equal quality to both alleles of a gene. In `scBASE`, multi-reads are resolved
351 by computing a weighted allocation based on the estimated probability of each alignment.
352 We use an EM algorithm implemented in `EMASE` software for this step¹⁶. Alternatively, read
353 counting could be performed using similar methods implemented in `RSEM`¹⁴ or `kallisto`¹⁵
354 software. The estimated maternal read count (x_{gk}) for each gene (g) in each cell (k) is the
355 weighted sum of all reads that align to the maternal allele, where the weights are proportional
356 to the probability of the read alignment. Similarly, the estimated paternal read count (y_{gk})
357 is the weighted sum of all reads that align to the paternal allele. The total read count is
358 the sum of the allele-specific counts ($n_{gk} = x_{gk} + y_{gk}$). A parameter of interest is the allelic
359 proportion p_{gk} . The read counting step provides an initial estimate $\hat{p}_{gk} = x_{gk}/n_{gk}$, which we

360 refer to as the weighted allocation estimated counts (ii).

361 **Classification:** In the classification step, we estimate the allelic expression state
362 (z_{gk}) for each gene in each cell. The allelic expression state is a latent variable with three
363 possible values $z_{gk} \in \{P, B, M\}$ representing paternal monoallelic, bi-allelic, and maternal
364 monoallelic expression, respectively. Uncertainty about the allelic expression state derives
365 from sampling variation that can produce zero counts for one or both alleles even when
366 the allele-specific transcripts may be present in the cell. We account for this uncertainty by
367 computing a probabilistic classification based on a mixture model in which the maternal read
368 counts x_{gk} are drawn from one of three beta-binomial distributions (given n_{gk}) according to
369 the allelic expression state z_{gk} . For a gene in the bi-allelic expression state the maternal
370 allelic proportion is denoted p_{gk}^B and, as suggested by the notation, it may vary from cell to
371 cell following a beta distribution. For a gene in the paternal monoallelic expression state,
372 the allelic proportion p_g^P follows a beta distribution with a high concentration of mass near
373 zero. Similarly, for a gene in the maternal monoallelic expression state, we model p_g^M using a
374 beta distribution with the concentration of mass near one. The beta distribution parameters
375 for the maternal and paternal states are gene-specific but are constant across cells.

376 **Estimation:** The classification step assumes that the mixture model parameters are
377 known. This model describes gene-specific allelic proportions for each cell and thus it has
378 a very large number of parameters. In the scRNA-Seq setting where thousands of genes
379 are measured but low read counts and sampling zeros are prevalent, we may have limited
380 data to support their reliable estimation. Bayesian analysis of the hierarchical model treats

381 parameters as random variables and is well suited for this type of estimation. In this con-
382 text, the hierarchical model improves the precision of estimation by borrowing information
383 across cells for each gene, giving more weight to cells that are in the same allelic expression
384 state. This estimation technique is referred to as *partial pooling*. Specifically, we sample the
385 mixture weights $(\pi_{g^*}^P, \pi_{g^*}^B, \pi_{g^*}^M)$ and the class-specific allele proportions (p_g^P, p_{gk}^B, p_g^M) ; gener-
386 ate classification probabilities $(\pi_{gk}^P, \pi_{gk}^B, \pi_{gk}^M)$; and then estimate the allelic proportions as a
387 weighted average

$$p_{gk} = \pi_{gk}^P p_g^P + \pi_{gk}^B p_{gk}^B + \pi_{gk}^M p_g^M \quad (1)$$

388 The average value across many iterations is \tilde{p}_{gk} , the partial pooling estimator.

389 **Estimating allelic proportions in subpopulations of cells or genes.** The scBASE
390 algorithm is designed to model heterogeneous ASE states in any population of cells. In some
391 cases, as in the developmental series of Deng et al., it is of interest to focus on different
392 subpopulations. When subpopulations of cells or groups of genes, e.g., X chromosome genes,
393 are expected to have different distributions of allelic states, we recommend two options. The
394 first option is to run the MCMC implementation of scBASE separately for each group. The
395 strength of this approach is that it provides the posterior distribution of group-specific allelic
396 proportions. However the level of uncertainty could increase for estimated parameters when
397 the number of cells in any group is limited. The second option is to first run MCMC with
398 all the available cells and estimate the prior parameters, α_g^s and β_g^s . These prior parameters
399 describe how allelic proportions are distributed in the monoallelic and bi-allelic states, and

400 therefore, are common across all groups. Then using estimated values for these parameters,
401 re-estimate the remaining parameters, π_g^s , π_{gk}^s , and p_{gk} , within each cell type using the EM
402 algorithm. In the restricted version of EM, we iteratively update π_{gk}^s (E-step) and π_g^s (M-
403 step) for cells within each subpopulation. Once π_g^s and π_{gk}^s converge, we can compute \tilde{p}_{gk}
404 using Equation (1). We applied this second approach to Deng et al. time series data along
405 mouse embryo development (n=286 cells). Genes on the X chromosome present another
406 example where it makes sense to run scBASE separately, in this case on two subpopulations
407 of genes. Our analyses of female X chromosome genes used this strategy (Supplemental
408 Figures S9, S10, and S11).

409 **Assigning allelic expression states from estimated counts.** Unique read counts are
410 obtained directly from counting reads after discarding all genomic and allelic multi-reads.
411 Weighted allocation counts are derived from the EM algorithm as described above. To
412 estimate counts after partial pooling, we multiply \tilde{p}_{gk} by the total gene expression counts.
413 We note that estimated counts are not integers and may be non-zero but less than one.
414 Classification of allelic expression states for each gene in each cell directly from observed or
415 estimated counts requires setting a threshold for monoallelic expression. For each allele, we
416 regarded it as expressed if its estimated abundance is greater than one reads (or one UMI
417 as in Larsson et al²¹).

418 **Classification of a gene according to its ASE profile across many cells.** We classify
419 a gene according to the proportion of cells in P-, B-, and M-states, $(\pi_g^P, \pi_g^B, \pi_g^M)$, that are
420 estimated by the partial pooling model. If a majority of cells ($\pi_g^s > 0.7$) are in a particular

421 ASE state, $s \in \{P, B, M\}$, then we will assign the gene to the class **P** (monoallelic paternal;
422 blue), **B** (bi-allelic; yellow), or **M** (monoallelic maternal; red) respectively. When a majority
423 of cells are a mixture of two of those classes ($\pi_g^{s_1} + \pi_g^{s_2} > 0.9$ where $s_1, s_2 \in \{P, B, M\}$),
424 we classify it into either of **PB** (mixture of monoallelic paternal and bi-allelic; green), **BM**
425 (mixture of monoallelic maternal and bi-allelic; orange), or **MP** (a mixture of monoallelic
426 maternal and paternal; purple). Otherwise, genes that present all three ASE states are
427 classified as **PBM** (mixture of all; gray). We specified these seven classes in a ternary simplex
428 diagram (Figure 5a)³². The class boundaries are arbitrary but the aim of this classification is
429 to provide a simple descriptive summary of the gene expression states present in a population
430 of cells.

431 **Sampling reads.** We randomly sampled 1% of reads in each of 122 cells at the early, mid,
432 and late blastocyst stages to obtain an average read count of ~ 148 k reads per cell. We chose
433 the blastocyst cell types because, unlike cells in earlier developmental stages, they show
434 the widest range of different states of allelic expression. The original analysis of SCALE²³
435 also used the same 122 cells. We applied the unique-reads method and weighted allocation
436 algorithm to the full set of ~ 14.8 M reads and also applied each of four estimation methods
437 (unique reads, weighted allocation counts, unique reads with partial pooling, and weighted
438 allocation with partial pooling) to the down-sampled data. We compared estimates obtained
439 from the down-sampled data to the full data estimates and computed the mean squared error
440 of estimation across cells for each gene.

441 **Simulation of counts under perfect independence model.** We randomly sampled the
442 marginal probabilities of maternal and paternal allelic expression, p_M and p_P from uniform
443 distribution between 0 and 1. Then we generated 2×2 tables by sampling counts from
444 multinomial distribution with probability $\{p_M p_P, p_M(1-p_P), (1-p_M)p_P, (1-p_M)(1-p_P)\}$
445 for bi-allelic, maternal monoallelic, paternal monoallelic, and silent cells respectively.

446 **Funding Statement**

447 This work has been supported by the National Institutes of Health (NIH) grant R01-
448 GM070683.

449 **Acknowledgments**

450 We would like to thank Steven C. Munger and Daniel A. Skelly for their helpful comments
451 on this manuscript.

452 **Author Contributions**

453 KC, NR, and GC conceived and planned the study. KC performed the model implementation
454 and analyses. KC and GC interpreted the scientific findings. KC and GC discussed, and
455 wrote the manuscript.

456 References

- 457 1. Linnarsson, S. & Teichmann, S. A. Single-cell genomics: coming of age. *Genome Biol.*
458 **17**, 97 (2016).
- 459 2. Crowley, J. J. *et al.* Analyses of allele-specific gene expression in highly divergent mouse
460 crosses identifies pervasive allelic imbalance. *Nat Genet* **47**, 353–360 (2015).
- 461 3. Santoni, F. A. *et al.* Detection of imprinted genes by single-cell allele-specific gene
462 expression. *The American Journal of Human Genetics* **100**, 444–453 (2017).
- 463 4. Tukiainen, T. *et al.* Landscape of x chromosome inactivation across human tissues.
464 *Nature* **550**, 244–248 (2017).
- 465 5. Garieri, M. *et al.* Extensive cellular heterogeneity of x inactivation revealed by single-cell
466 allele-specific expression in human fibroblasts. *Proceedings of the National Academy of*
467 *Sciences* **115**, 13015–13020 (2018).
- 468 6. Brennecke, P. *et al.* Accounting for technical noise in single-cell RNA-seq experiments.
469 *Nat. Methods* **10**, 1093–1095 (2013).
- 470 7. Deng, Q., Ramsköld, D., Reinius, B. & Sandberg, R. Single-cell rna-seq reveals dynamic,
471 random monoallelic gene expression in mammalian cells. *Science* **343**, 193–196 (2014).
- 472 8. Kim, J. K. *et al.* Characterizing noise structure in single-cell RNA-seq distinguishes
473 genuine from technical stochastic allelic expression. *Nat Commun* **6**, 8687 (2015).
- 474 9. Stegle, O., Teichmann, S. A. & Marioni, J. C. Computational and analytical challenges
475 in single-cell transcriptomics. *Nat. Rev. Genet.* **16**, 133–145 (2015).
- 476 10. Bacher, R. & Kendziorski, C. Design and computational analysis of single-cell RNA-
477 sequencing experiments. *Genome Biol.* **17**, 63 (2016).
- 478 11. Rostom, R., Svensson, V., Teichmann, S. A. & Kar, G. Computational approaches for
479 interpreting scRNA-seq data. *FEBS Lett.* **591**, 2213–2225 (2017).
- 480 12. Qiu, X. *et al.* Single-cell mRNA quantification and differential analysis with Census.
481 *Nat. Methods* **14**, 309–315 (2017).
- 482 13. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell
483 transcriptomic data across different conditions, technologies, and species. *Nat. Biotech-*
484 *nol.* **36**, 411–420 (2018).
- 485 14. Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data
486 with or without a reference genome. *BMC Bioinformatics* **12**, 323 (2011).
- 487 15. Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic rna-seq
488 quantification. *Nature Biotechnology* **34**, 525–527 (2016).

- 489 16. Raghupathy, N. *et al.* Hierarchical analysis of rna-seq reads improves the accuracy of
490 allele-specific expression. *Bioinformatics* **34**, 2177–2184 (2018).
- 491 17. Keane, T. M. *et al.* Mouse genomic variation and its effect on phenotypes and gene
492 regulation. *Nature* **477**, 289–294 (2011).
- 493 18. Huang, M. *et al.* Saver: gene expression recovery for single-cell rna sequencing. *Nature*
494 *Methods* **15**, 539–542 (2018).
- 495 19. Reinius, B. & Sandberg, R. Random monoallelic expression of autosomal genes: stochas-
496 tic transcription and allele-level regulation. *Nat Rev Genet* **16**, 653–664 (2015). Review.
- 497 20. Agresti, A. *Contingency Tables* (John Wiley & Sons, Inc., 2007), 2 edn.
- 498 21. Larsson, A. J. M. *et al.* Genomic encoding of transcriptional burst kinetics. *Nature* **565**,
499 251–254 (2019).
- 500 22. Slavković, A. & Fienberg, S. Algebraic geometry of 2x2 contingency table. In Gibilisco,
501 P., Riccomagno, E., Rogantin, M. P. & Wynn, H. P. (eds.) *Algebraic and Geometric*
502 *Methods in Statistics*, chap. 3, 63–81 (Cambridge University Press, 2009).
- 503 23. Jiang, Y., Zhang, N. R. & Li, M. Scale: modeling allele-specific gene expression by
504 single-cell rna sequencing. *Genome Biology* **18**, 74 (2017).
- 505 24. Reinius, B. *et al.* Analysis of allelic expression patterns in clonal somatic cells by single-
506 cell rna-seq. *Nature Genetics* **48**, 1430–1435 (2016).
- 507 25. Chen, G. *et al.* Single-cell analyses of x chromosome inactivation dynamics and pluripo-
508 tency during differentiation. *Genome research* **26**, 1342–1354 (2016).
- 509 26. Babak, T. *et al.* Global survey of genomic imprinting by transcriptome sequencing.
510 *Current Biology* **18**, 1735 – 1741 (2008).
- 511 27. The Jackson Laboratory. Mouse genome informatics (2019). URL [http://www.informatics.jax.org/searchtool/Search.do?query=genetic+](http://www.informatics.jax.org/searchtool/Search.do?query=genetic+imprinting&submit=Quick%0D%0ASearch)
512 [imprinting&submit=Quick%0D%0ASearch](http://www.informatics.jax.org/searchtool/Search.do?query=genetic+imprinting&submit=Quick%0D%0ASearch).
513
- 514 28. Jirtle, R. L. Imprinted genes: by species (2012). URL <http://www.geneimprint.com/site/genes-by-species.Mus+musculus>.
515
- 516 29. Edsgård, D., Reinius, B. & Sandberg, R. scphaser: haplotype inference using single-cell
517 RNA-seq data. *Bioinformatics* **32**, 3038–3040 (2016).
- 518 30. Carpenter, B. *et al.* Stan: A probabilistic programming language. *Journal of Statistical*
519 *Software, Articles* **76**, 1–32 (2017).
- 520 31. Kleinman, J. C. Proportions with extraneous variance: Single and independent sample.
521 *Journal of the American Statistical Association* **68**, 46–54 (1973).
- 522 32. Harper, M. *et al.* python-ternary: Ternary plots in python (2015).

523 Figure Captions

524 **Figure 1: Overview of the scBASE algorithm.** We summarize the three steps of
525 the scBASE algorithm. The **Counting** step estimates the expected read counts using an
526 EM algorithm to compute a weighted allocation of multi-reads. Each read is represented
527 as an incidence matrix that summarizes all best-quality alignments to genes and alleles ①.
528 Weighted allocation of multi-reads uses a current estimate of allele-specific gene expression
529 to compute weights equal to the probability of each possible alignment ②. The weights are
530 summed across reads to obtain the expected read counts for each gene and allele ③. Steps
531 ② and ③ are repeated until the read counts converge. The weighted allocation estimates
532 of maternal allelic proportion (\hat{p}_{gk}) are obtained at this step. The **Classification** step
533 computes the posterior probability of paternal monoallelic (P), bi-allelic (B), or maternal
534 monoallelic (M) expression (π_{gk}^s) using current estimates of the model parameters (Equation 3
535 in Supplemental Methods). The classification model is a beta-binomial mixture model with
536 three components. The model parameters are initialized to non-informative values and are
537 obtained from the estimation step in subsequent iterations. The **Estimation** step uses the
538 classification results to re-estimate the weights of mixture components ($\pi_{g.}^s$) and parameters
539 of the Beta densities (α_g^s, β_g^s) that define the distribution of the within-class the maternal
540 allelic proportions (p_g^s). The partial pooling estimate of the maternal allelic proportions (\tilde{p}_{gk})
541 is obtained as an average of the class-specific proportions weighted by the class membership
542 probabilities (Equation 1 in Methods).

543 **Figure 2: Weighted allocation of multi-reads reduces monoallelic expression calls.**
544 **(a)** For each of 13,032 genes, we obtained the allele-specific read counts by unique reads and
545 by weighted allocation. We counted the numbers of genes in each cell that showed either
546 maternal or paternal monoallelic expression and display the results as points (one per cell)
547 overlaid on boxplots. Each data point in this figure represents a cell and we are showing all
548 286 cells including zygote and 2-cells (highlighted in red). The zygote and 2-cell stage cells
549 have large numbers of genes with maternal monoallelic expression. On average there are ~ 66
550 fewer monoallelic calls per cell with the weighted allocation counts. The outlier cell with
551 high levels of paternal monoallelic expression was noted in Deng et al.⁷. **(b)** We selected one
552 gene (*Mtdh*) to illustrate the distribution of maternal (X-axis) and paternal (Y-axis) counts
553 across 286 cells. The weighted allocation counts (green) are connected to their corresponding
554 unique counts by a line in the scatter plot. **(c)** Cross-tabulation (2×2 table) of maternal and
555 paternal allelic expression for *Mtdh* gene with unique reads and weighted allocation counts.
556 The unique counts resulted in 88 cells with monoallelic expression while only 7 monoallelic
557 calls were seen with weighted allocation.

558 **Figure 3: Partial pooling improves the accuracy of estimated allelic proportions.**
559 We randomly sampled 1% of reads from the full data of 122 mature blastocyst cells to
560 obtain a sub-sample of 147,538 reads per cell, on average. We estimated gene- and cell-
561 specific allelic proportions from the sub-sampled data, and computed mean squared error
562 (MSE) between the estimated allelic proportions from the full data versus the sub-sampled
563 data. We compared the MSE based on partial pooling versus the MSE from no pooling

564 estimates, and display the difference on the y-axis along the expression level in unique-read
565 counts on the x-axis. We made this comparison for **(a)** unique reads and for **(b)** weighted
566 allocation. Points representing individual genes are shown as a density heatmap.

567 **Figure 4: Independence of allelic bursting.** **(a)** The geometry of the 2x2 table pro-
568 portions can be represented as a simplex, a 3D tetrahedral region of 4D space in which
569 proportions are all non-negative and sum to one. The vertices of the simplex correspond
570 to genes where all cells are in the same allelic expression state as indicated by labels. The
571 distance from a vertex is inversely related (1-x) to the proportion of cells in that state. The
572 shaded surface inside the simplex represents proportions corresponding to the perfect inde-
573 pendence model, i.e., the logOR equals zero. The blue triangle indicates proportions with
574 equal maternal and paternal expression $p_M = p_P$. **(b)** We simulated data under the perfect
575 independence model without assuming $p_M = p_P$ and plotted the proportions of bi-allelic and
576 silent cells as in Deng et al.⁷. **(c)** Four panels illustrate the proportions of bi-allelic and
577 silent cells as estimated from (i) unique reads, (ii) weighted allocation, (iii) unique reads
578 with partial pooling, and (iv) weighted allocation with partial pooling. Points representing
579 individual genes are shown as a density heat map.

580 **Figure 5: Classification of allele-specific expression patterns across cells.** **(a)** For
581 each gene in each cell, the classification step of scBASE estimates allelic state probabilities
582 π_{gk}^s , where s indicates paternal monoallelic (P), bi-allelic (B), or maternal monoallelic (M)
583 expression. The average proportions of cells in each allelic state (π_g^s) can be represented as
584 a point in a triangular diagram which is a 3D simplex corresponding to the projection of
585 points onto the bottom triangular region of the 4D simplex in Figure 4a. A gene that is
586 predominantly paternal, bi-allelic, or maternal across the cell population will be plotted near
587 the corresponding vertex. Points representing genes with mixed classification states across
588 the cell population will appear along the edges or in the center of the triangle. We delineate
589 seven patterns of allelic expression for a gene as indicated by the different colored regions in
590 the diagram: **P** (blue), **B** (yellow), **M** (red), **PB** (green), **BM** (orange), **MP** (purple), and
591 **PBM** (gray). Examples of genes from each pattern are shown as scatter plots of maternal
592 and paternal read counts (log10 scale). Each point in the scatter plot corresponds to one cell
593 (n=286 embryo cells). For example, the gene *Pacs2* is expressed from either the maternal or
594 the paternal allele but rarely both and is classified as an **MP** gene. The bi-allelic region (**B**)
595 includes genes that may show allelic imbalance ($p_{gk} \neq \frac{1}{2}$) across many cells but consistently
596 express both alleles (e.g., *Mtdh*). The **PB** and **BM** regions will include genes that show
597 a mixture of bi-allelic expression and monoallelic expression. Many of the genes in these
598 regions have strong allelic imbalance and cells with monoallelic expression could be due
599 to statistical sampling zeros in the lower expressed allele (e.g., *Tmim23* and *Tulp3*). The
600 expression pattern in blastocyst cells for the majority of genes (57%) fall in the **PBM** region
601 and display a pattern that is a mix of mono- and bi-allelic expression states across cells (e.g.,
602 *Akr1b3*). **(b)** Cells were divided into nine developmental stages as indicated on the X-axis.
603 The cell types and numbers of expressed genes at each stage are indicated in parentheses on
604 the X-axis. For each stage, we counted the proportion of expressed genes that fall into each of
605 the seven allelic expression patterns (Y-axis), indicated by lines using the same color coding

606 used in Figure 5a. In the zygote and early 2-cell stage, most genes show purely maternal
607 expression (**M**). The proportion of maternally expressed genes decreases through subsequent
608 stages of development. The numbers of genes showing purely paternal expression (**P**) is low
609 across all developmental stages. The **M** and **P** classes become equally represented in the
610 later stages of development. The 2- and 4-cell stages show high levels of bi-allelic expression
611 (**B**) and the mixed class (**PBM**) proportion becomes highest by the 8-cell stage. (c) The
612 expected proportions of cells in each allelic state (π_g^s) for one gene *Akr1b3* at each stage of
613 the developmental time course is shown as a trajectory in the 3D simplex. Yellow to blue
614 color line segments indicates the transitions between developmental stages. This gene starts
615 in the maternal monoallelic state (**M**), it transitions through **PBM** to a paternal expression
616 state (**P**), and then transitions to bi-allelic expression (**B**) in the blastocyst stages.

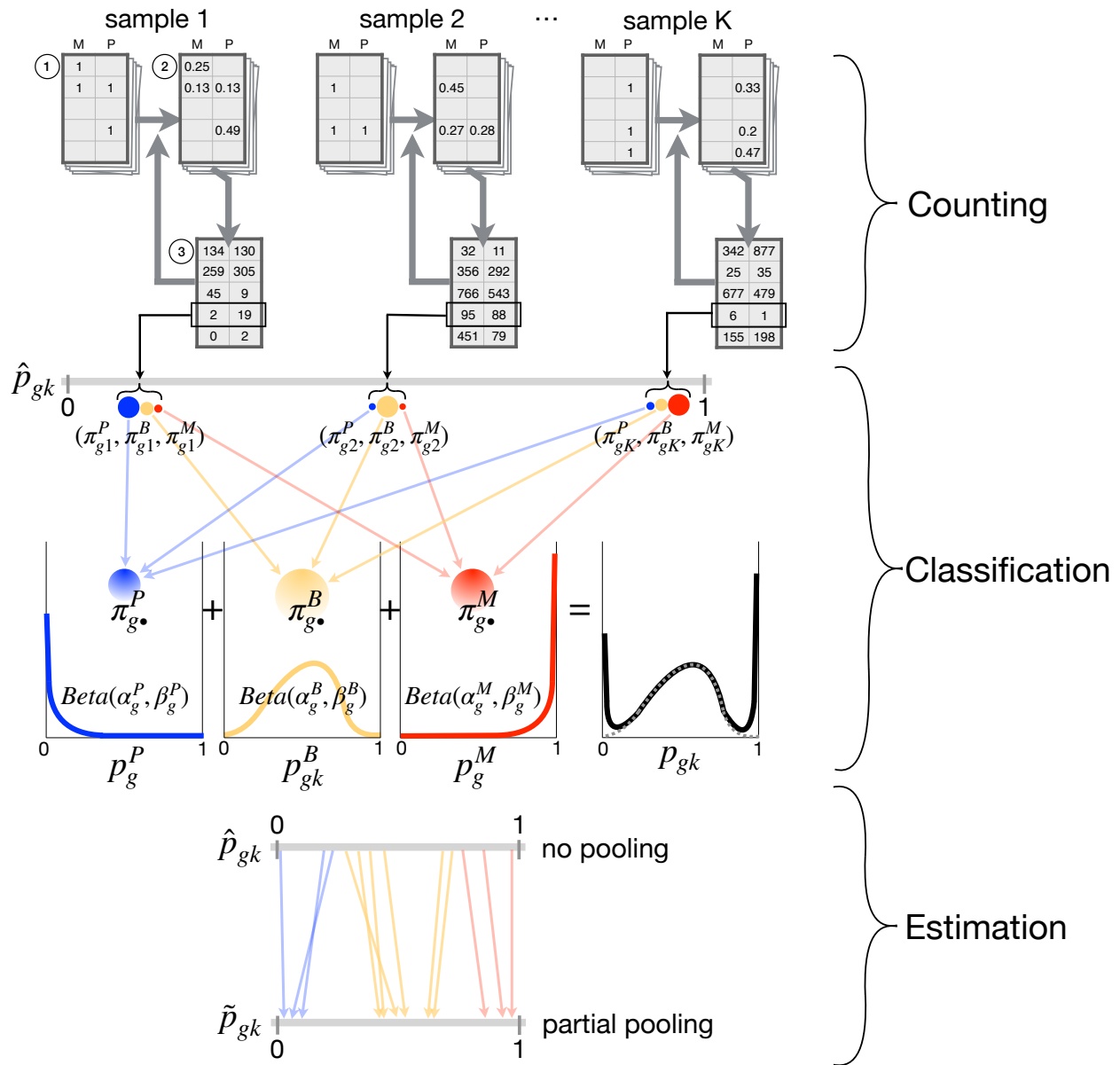
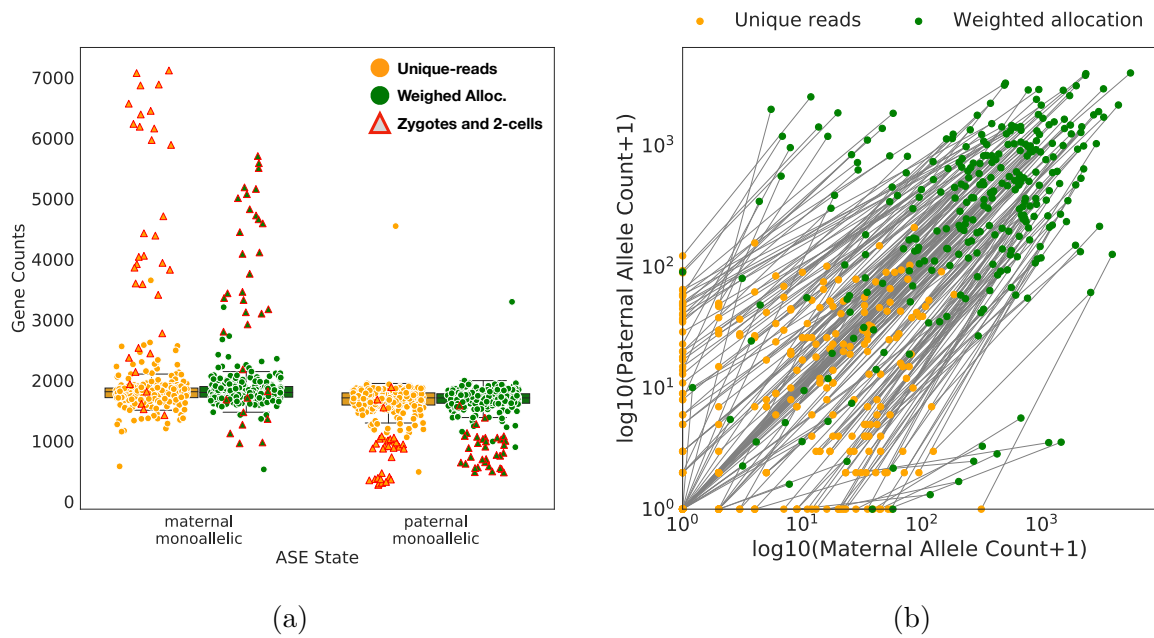


Figure 1

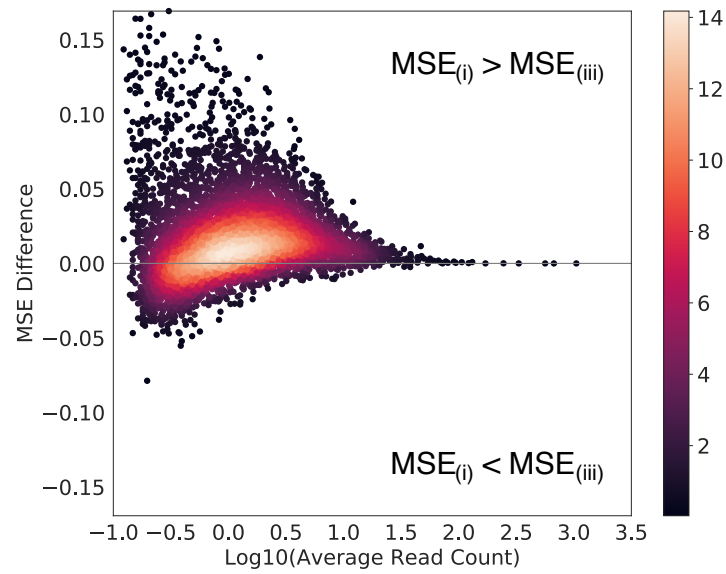


Unique reads		Maternal allele	
		Not expressed	Expressed
Paternal allele	Not expressed	56	39
	Expressed	49	142

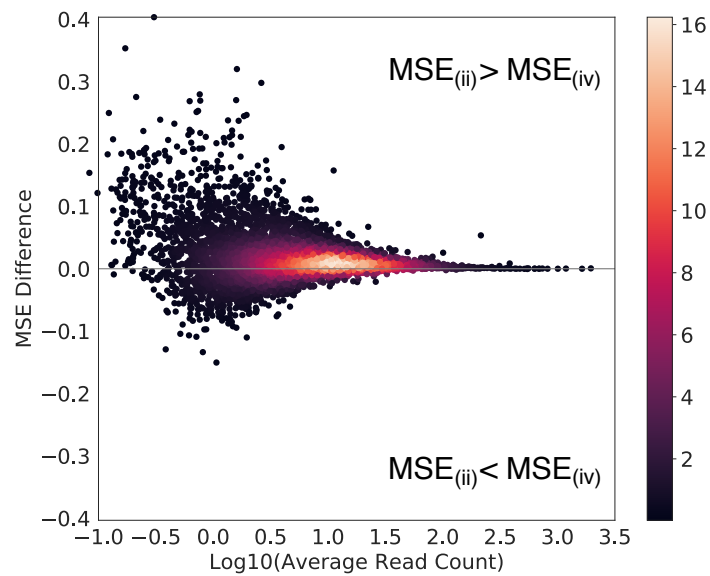
Weighted allocation		Maternal allele	
		Not expressed	Expressed
Paternal allele	Not expressed	0	5
	Expressed	2	279

(c)

Figure 2



(a) $MSE_{(i)} - MSE_{(iii)}$ along the expression level



(b) $MSE_{(ii)} - MSE_{(iv)}$ along the expression level

Figure 3

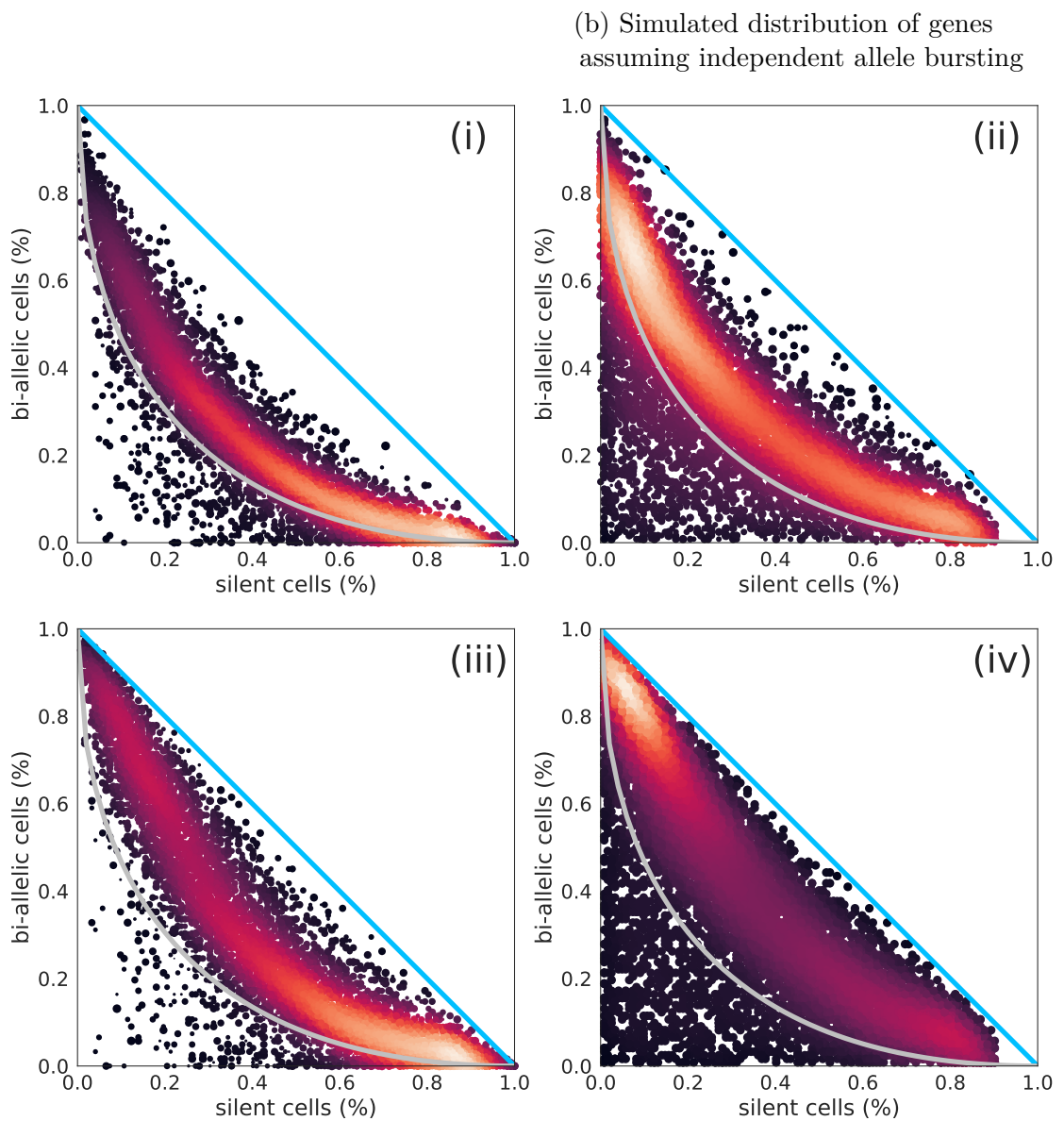
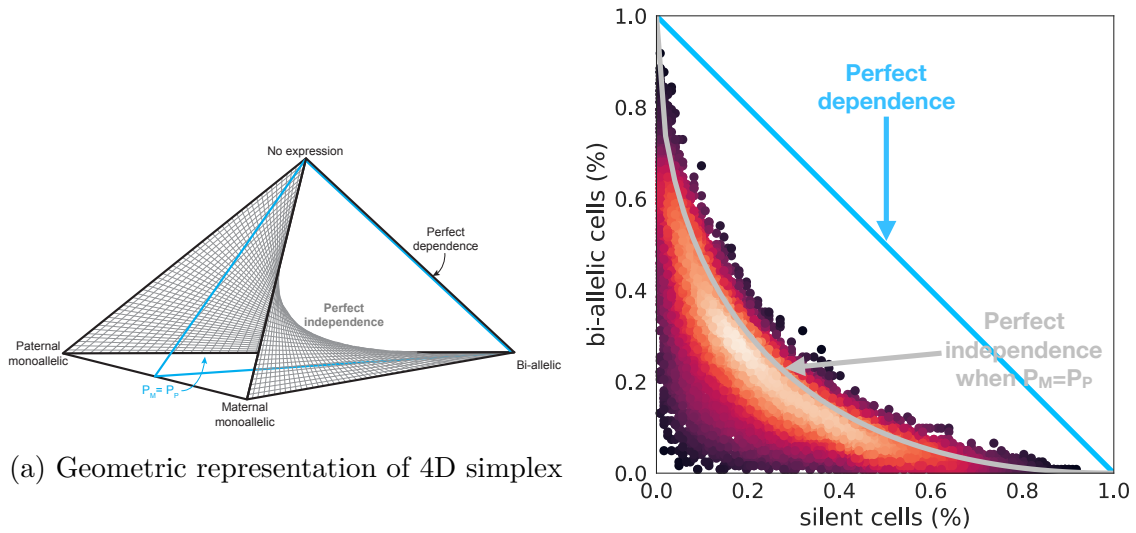


Figure 4

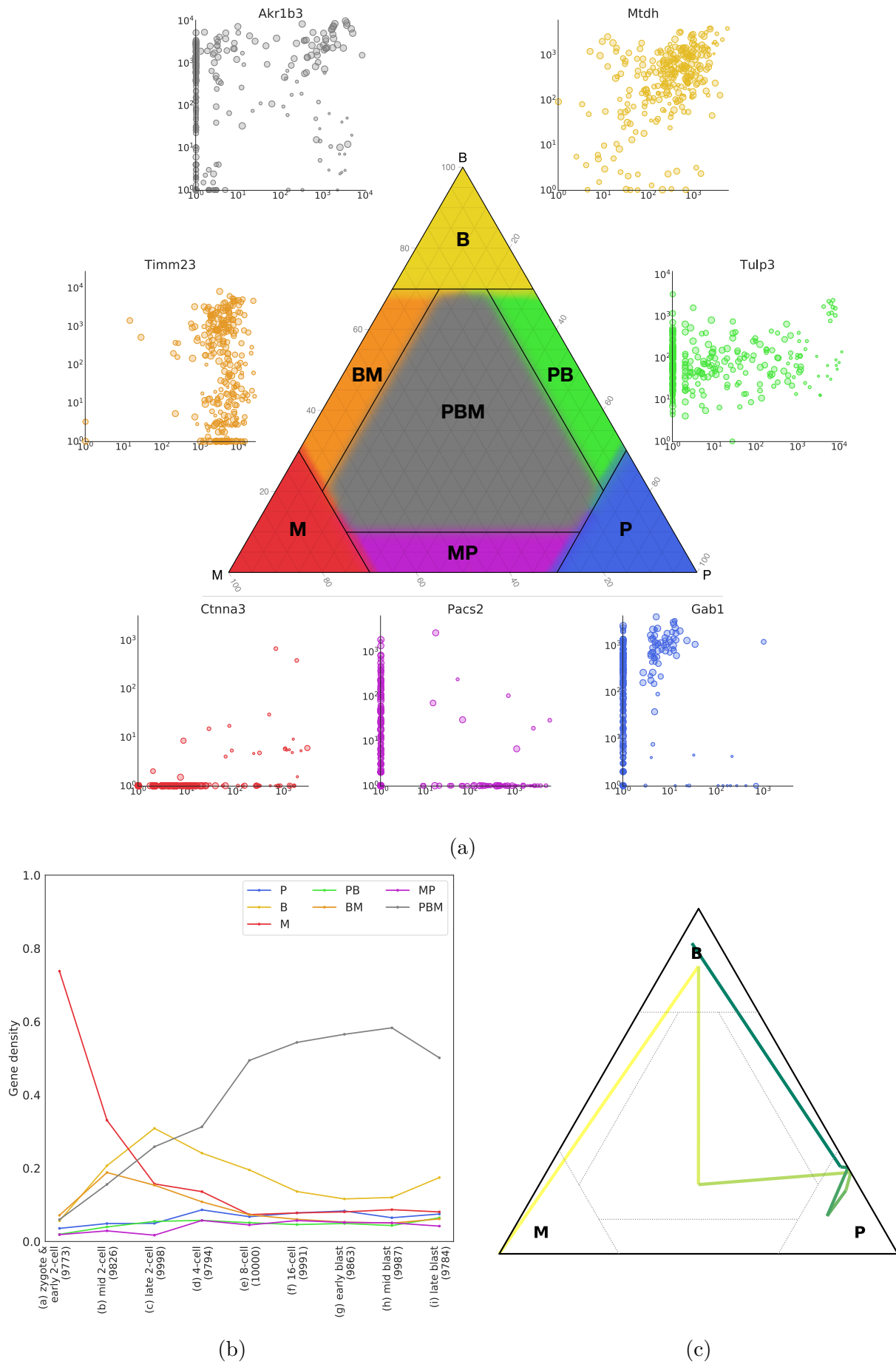


Figure 5