

A peer-reviewed version of this preprint was published in PeerJ on 23 May 2017.

[View the peer-reviewed version](https://doi.org/10.7717/peerj.3334) (peerj.com/articles/3334), which is the preferred citable publication unless you specifically need to cite this preprint.

de Torrente L, Zimmerman S, Taylor D, Hasegawa Y, Wells CA, Mar JC. 2017. *pathVar*: a new method for pathway-based interpretation of gene expression variability. PeerJ 5:e3334 <https://doi.org/10.7717/peerj.3334>

***pathVar*: a new method for pathway-based interpretation of gene expression variability**

Laurence de Torrente¹, **Samuel Zimmerman**¹, **Deanne Taylor**^{2,3}, **Yu Hasegawa**⁴, **Christine A Wells**⁵, **Jessica C Mar**^{Corresp. 1, 6}

¹ Department of Systems and Computational Biology, Albert Einstein College of Medicine, Bronx, New York, USA

² Department of Biomedical and Health Informatics, Children's Hospital of Philadelphia, Philadelphia, Pennsylvania, USA

³ Pediatrics, University of Pennsylvania, Philadelphia

⁴ Department of Food Science and Technology, University of California, Davis, Davis, California, United States

⁵ Department of Anatomy and Neuroscience, University of Melbourne, Melbourne, Victoria, Australia

⁶ Department of Epidemiology and Population Health, Albert Einstein College of Medicine, Bronx, New York, USA

Corresponding Author: Jessica C Mar

Email address: jessica.mar@einstein.yu.edu

Identifying the pathways that control a cellular phenotype is the first step to building a mechanistic model. Recent examples in developmental biology, cancer genomics, and neurological disease have demonstrated how changes in the variability of gene expression can highlight important genes that are under different degrees of regulatory control. Simple statistical tests exist to identify differentially-variable genes; however, methods for investigating how changes in gene expression variability in the context of pathways and gene sets are under-explored. Here we present *pathVar*, a new method that provides functional interpretation of gene expression variability changes at the level of pathways and gene sets. *pathVar* is based on a multinomial exact test, or an asymptotic Chi-squared test as a more computationally-efficient alternative. The method can be used for gene expression studies from any technology platform in all biological settings either with a single phenotypic group, or two-group comparisons. To demonstrate its utility, we applied the method to a diverse set of diseases, species and samples. Results from *pathVar* are benchmarked against analyses based on average expression via GSEA, and demonstrate that analyses using both statistics are useful for understanding transcriptional regulation.

1 ***pathVar*: a new method for pathway-based interpretation of gene expression variability.**

2 Laurence de Torrenté^{1#}, Samuel Zimmerman¹, Deanne Taylor², Yu Hasegawa^{1%}, Christine A. Wells⁴,
3 Jessica C. Mar^{1,5*}

4 ¹Department of Systems and Computational Biology, Albert Einstein College of Medicine, Bronx, NY,
5 USA.

6 [#]Present address: New York Genome Center, New York, NY, USA.

7 ²Department of Biomedical and Health Informatics, Children's Hospital of Philadelphia, PA, USA.

8 [%]Present address: Robert Mondavi Institute for Wine and Food Science, University of California, Davis,
9 CA, USA.

10 ⁴Department of Anatomy and Neuroscience, School of Biomedical Sciences, University of Melbourne,
11 VIC, Australia.

12 ⁵Department of Epidemiology and Population Health, Albert Einstein College of Medicine, Bronx, NY,
13 USA.

14 *To whom correspondence should be addressed.

15 **Abstract**

16 Identifying the pathways that control a cellular phenotype is the first step to building a mechanistic model.
17 Recent examples in developmental biology, cancer genomics, and neurological disease have demonstrated
18 how changes in the variability of gene expression can highlight important genes that are under different
19 degrees of regulatory control. Simple statistical tests exist to identify differentially-variable genes; however,
20 methods for investigating how changes in gene expression variability in the context of pathways and gene
21 sets are under-explored. Here we present *pathVar*, a new method that provides functional interpretation of
22 gene expression variability changes at the level of pathways and gene sets. *pathVar* is based on a
23 multinomial exact test, or an asymptotic Chi-squared test as a more computationally-efficient alternative.
24 The method can be used for gene expression studies from any technology platform in all biological settings
25 either with a single phenotypic group, or two-group comparisons. To demonstrate its utility, we applied the
26 method to a diverse set of diseases, species and samples. Results from *pathVar* are benchmarked against
27 analyses based on average expression via GSEA, and demonstrate that analyses using both statistics are
28 useful for understanding transcriptional regulation.

29

30 **Availability:** The method is publicly available from the *pathVar* Bioconductor R package.

31

32 **Contact:** jessica.mar@einstein.yu.edu.

33

34

35

36 Introduction

37 Global studies of gene expression provide two quantitative parameters: a commonly-used metric is the
38 relative abundance of a transcript (and group differences in transcript abundance), likewise the expression
39 variability of that transcript provides insight into the heterogeneity of a sample group [1], and expression
40 variability changes between groups have been shown to reflect underlying changes in transcriptional
41 regulatory processes [2-5]. Patterns of variability in gene expression have provided insight into how
42 pathways are regulated in cells [6,7]; especially in the context of single cell profiling studies, where the
43 average expression of a gene in a cell population carries limited information for understanding
44 transcriptional regulation. Recent studies have identified pathways showing differential control or
45 regulatory constraint that were discovered only by modeling changes in gene expression variability and
46 were not apparent from standard analyses of average gene expression [8]. While variability is becoming
47 more prevalent as an informative metric, the current challenge lies in how to interpret these analyses to
48 maximize functional information, such as with respect to pathways and curated gene sets. Due to the
49 newness of this area, statistical methods for investigating expression variability are currently under-
50 developed, and lacking for pathway-centric approaches. It is necessary, therefore, to develop such
51 methods since information on expression variability can be used to complement analyses of average
52 expression, and improve our understanding of the transcriptional state of the cell.

53 Intuitively, the distribution of gene expression variability in a pathway highlights the subset of genes with
54 different degrees of regulatory control (Fig 1). In the case of a one-group design, where multiple profiles
55 represent replicates of the same phenotype, e.g. different embryonic stem cell (ESC) lines, identifying
56 pathways that have an unexpected proportion of low variability genes may point to those that contribute
57 integral roles for stem maintenance or regulation [1]. To appreciate this, consider two previous studies
58 that provided evidence linking criticality of genes and their decreased variability in expression. One study
59 [8] identified genes with decreased expression variability in tumors relative to normal tissue; this gene set,
60 termed the Posed Gene Cassette, included key genes whose expression impacted metastasis and patient
61 survival as demonstrated through *in vivo* and *in vitro* experimental approaches. More recently, a second
62 study [9] showed that genes with decreased variability in expression for four stages of early embryonic
63 development (4-cell, 8-cell, morula and blastocyst) were more likely to be associated with essentiality,
64 haploinsufficiency or ubiquitous expression, suggesting that these stably-expressed genes contribute to
65 cell survival.

66 In the two-group design, where profiles are compared between two contrasting phenotypes, e.g. ESCs
67 versus induced pluripotent stem cells (iPSCs), identifying pathways associated with different patterns of
68 expression variability may highlight those pathways that contribute to group-specific differences.
69 Previous studies have analyzed the enrichment of genes with different levels of expression variability for
70 specific pathways [9,10]; however, these analyses are based on gene lists defined by an arbitrary cut-off
71 and do not take into account the expression distribution of genes in the pathway. One would expect that
72 more informative results could be obtained by focusing on the shape of the expression distribution in a
73 statistically rigorous manner, much like a gene set enrichment analysis (GSEA) [11] analogue for
74 variability instead of relying only on average expression, or over-representation (OR) analyses [12].
75 Computational methods to implement these kinds of approaches are currently lacking for expression
76 variability.

77 Our method, *pathVar*, addresses this gap by providing a pathway-based analysis of gene expression
78 variability where pathways are assessed based on deviations of their gene expression variability distribution
79 relative to a reference. In the one-group setting, the reference can be the global distribution constructed
80 from all genes. In a two-group setting, one of the groups serves as the reference or control group. For each
81 pathway, our method also identifies which genes in a pathway show aberrant levels of gene expression
82 variability (Fig 1).

83

84 **Methods**

85 The *pathVar* method can be summarized in three main steps (Fig. 2).

86 ***Step 1: Selecting a statistical measure to estimate gene expression variability.***

87 Variability is defined as the amount of dispersion in a given distribution [13]. Different statistical measures
 88 are available to estimate gene expression variability, and in genomics, the estimators that are most often
 89 employed are the standard deviation (SD) (Equation 1) [9], the coefficient of variation (CV) (Equation 2)
 90 [1,10], and the median absolute deviation (MAD (Equation 3) [14]. Conceptually, these statistics share
 91 similarities in their mathematical definition, and each one comes with their own advantages, as is often the
 92 case with any estimator that is applied to data.

$$SD = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \quad (\text{Equation 1})$$

$$CV = \frac{SD}{\bar{x}} \quad (\text{Equation 2})$$

$$MAD = \text{median}_i(|x_i - \text{median}_j(x_j)|) \quad (\text{Equation 3})$$

93 A consensus on which estimator should be adopted for variability analysis remains unclear, and this is
 94 partly because performance of the estimators appears to be data-specific. The SD is often the preferred
 95 estimator for measuring gene expression variability since it is on the same scale as the average and
 96 therefore easy to interpret. The variance (i.e. SD²) is also characterized by the second central moment of a
 97 distribution, and hence the SD is directly linked to one of the fundamental metrics characterizing the
 98 probability distribution. A criticism of the SD however is that it may be dependent on average expression
 99 and therefore it is necessary to investigate the association between these two measures. To address this
 100 concern, the CV, which represents the ratio of the SD and the average, is often used however it also has
 101 its own drawbacks. Most significantly, it can be affected by zero-inflation which occurs when very small
 102 levels of average expression result in extremely large values of CV that do not necessarily reflect a large
 103 degree of overall variability. The MAD is a robust measurement of dispersion that behaves well in the
 104 presence of outlier data points. It has previously been used to study DNA methylation variability [14].

105 Data simulations suggest that in general, the SD shows stronger performance as the variability estimator
 106 compared to MAD and CV in the *pathVar* method (see Text S1). The simulations were conducted to test
 107 the performance of each estimator under a wide variety of conditions. The choice of an estimator can also
 108 be motivated by the expectation that an ideal estimator of gene expression variability will be uncorrelated
 109 with average gene expression. If the variability estimator is highly correlated with the average expression,
 110 then trends observed for expression variability may simply be recapitulated by those observed for average
 111 expression. We therefore desire an estimator of variability that is the least correlated with average gene
 112 expression for an analysis of gene expression variability to be maximally informative. Our suggestion is to
 113 use the SD since overall this estimator displayed stronger performance in simulations or the estimator with
 114 the lowest correlation with average expression; however, ultimately, the final choice of the estimator is left
 115 to the user and can be easily specified in the software.

116 ***Step 2: Identifying genes that belong to categories of high, medium or low levels of expression variability.***

117 Using a specified estimator, all genes are assigned to a discrete level of expression variability. In the one-
 118 group case, assignments are based on clustering the data using Normal mixture models via the *mclust*
 119 algorithm [15]. The number of clusters or mixtures corresponds to the discrete levels of variability and is a
 120 parameter inferred by *mclust*. The *mclust* algorithm considers a finite range of values (starting with a
 121 minimum of one level to a maximum of four by default) and chooses the number that is most appropriate
 122 for the data using the Bayesian Information Criterion. The upper limit of four levels is recommended out
 123 of simplicity, where it is more useful to model a handful of variability levels, e.g. low, medium, high and

124 very high, whereas for much larger numbers, the interpretation ceases to be as intuitive. In the *pathVar*
125 package, the user is, however, free to use whatever upper limit is appropriate for the analysis.

126 For the two-group case, assignments are based instead on the 33rd and 66th percentiles that are computed
127 from the combined gene expression variability distribution of all genes from both groups in the dataset.
128 Low variability genes are defined as those with values falling between 0 and the 33rd percentile, medium
129 variability genes are those between the 33rd percentile and the 66th percentile, and high variability genes
130 being greater than the 66th percentile. The variability levels are defined by a fixed number of standardized
131 percentiles, instead of inferred as in the one-group case, because it is possible that a different number of
132 variability levels might be inferred for the two groups. To ensure a straightforward and balanced
133 comparison between the two groups, the number of discrete categories are fixed and the boundaries for
134 variability levels are based on percentiles calculated from all the data. Under both the one-group and two-
135 group cases, the outcome of Step 2 is to identify the fixed boundaries that define each of the discrete levels
136 of expression variability.

137 ***Step 3: Testing pathways for aberrant gene expression variability signatures.***

138 The *pathVar* method decomposes each pathway into a set of counts corresponding to the number of genes
139 in each discrete level of expression variability. By default, pathways from the Kyoto Encyclopedia of Genes
140 and Genomes (KEGG) [16] and REACTOME [17] are used, however users may also import their own
141 definitions. For a specific pathway, let O_i denote the observed count of genes annotated to this pathway
142 with expression variability in the i -th level, where $i = 1, \dots, m$ discrete levels. The total number of genes n
143 is defined as $n = \sum_{i=1}^m O_i$. A statistical test is used to evaluate whether the set of counts (O_1, \dots, O_m)
144 associated with a pathway deviates significantly from either a reference count distribution in the one-group
145 case, or between the two phenotypes in the two-group case. In the one-group case, the reference distribution
146 is obtained by counting the total number of genes in each level of gene expression variability (Fig. S1).

147 The null hypothesis that *pathVar* evaluates is:

- 148 • $H_{0(1)}$: the expression variability counts observed for a specific pathway were generated under the same
149 distribution as the reference counts for the one-group case.
- 150 • $H_{0(2)}$: the variability-based counts for both groups were drawn from the same underlying distribution,
151 for the two-group case.

152 To assess the deviation observed between the variability count distribution and an expected distribution for
153 a specific pathway, we provide two statistical tests as options available for this analysis. Option 1 is the
154 multinomial exact test, and Option 2 is the Chi-squared test.

155 For Option 1, the exact test models the counts as a multinomial random variable $(O_1, \dots, O_m) \sim \text{Multi}(n, p)$
156 where $p = (p_1, \dots, p_m)$ and p_i is the probability that a gene belongs to the i -th variability level for $i = 1, \dots, m$
157 in the reference distribution. Under the exact test, the P-value is obtained by summing over all possible
158 events that are less likely than the set of counts observed. If the P-value falls below a specified significance
159 threshold (e.g. P-value < 0.05), sufficient evidence exists to reject the null hypothesis and we conclude that
160 the pathway has an aberrant distribution of expression variability counts. In the one-group case, this means
161 that the variability count distribution of the specific pathway deviates from the reference distribution i.e. all
162 genes surveyed. In the two-group case, this result means that the variability is not identically distributed
163 between the two contrasting phenotypes.

164 While the exact test is attractive because it calculates the P-value exactly, from a practical perspective, these
165 kinds of tests can often be time and memory-intensive for genomic data, especially as the number of levels
166 m and the number of genes in the pathway n grows. For example, consider a pathway with 30 genes, where
167 the number of possible sets of counts to consider with three variability levels for the calculation of the P-
168 value is 496. If the size of the pathway increases to 100, then the number of possibilities to consider grows

169 to 5151. From this basic example, it can be seen how increasing n or m can lead to some very extensive
170 calculations.

171 Option 2 overcomes this limitation, as it is less computationally intensive, and tests the same null
172 hypotheses using the Chi-squared test as an alternative to the exact test. The test statistic $\chi^2 = \sum_{i=1}^m \frac{(O_i - E_i)^2}{E_i}$

173 follows a Chi-squared distribution $\chi^2 \sim \chi_{m-1}^2$ with $m-1$ degrees of freedom: $\chi^2 \sim \chi_{m-1}^2$. The expected
174 counts $E_i = n * p_i$ are the expected number of genes in each level of expression variability within a specific
175 pathway. The Chi-squared test achieves its computational efficiency because it is based on an asymptotic
176 approximation, where the test becomes more accurate as n , the total number of genes in a pathway increases.

177 Both the exact test and the Chi-squared test assess whether a pathway has a significant change in gene
178 expression variability. The resulting P-values from all pathways tested are adjusted using the Benjamini-
179 Hochberg method which is based on controlling the false discovery rate [18]. Finally, a pathway of interest
180 can be investigated further using a Binomial test to assess within each level of variability, whether a
181 significant deviation exists between either the pathway and the reference in the one-group case, or between
182 the two phenotypic groups. This pairwise difference for each independent level provides a means to
183 pinpoint the subset of genes within a pathway showing deviation in expression variability.

184 **Power calculations indicate that *pathVar* has greater power to detect changes for pathways with
185 skewed expression variability distributions than symmetric ones.**

186 Simulations were designed to investigate how experimental parameters influence the power of *pathVar*
187 (Text S1). The results demonstrate that for a fixed effect size, the power of the Chi-squared test increases
188 when the number of genes in a pathway also increases. Similarly, for a pathway of fixed size, the power of
189 the test is higher for a larger effect size. The simulations also showed that the Chi-squared test had more
190 power to detect differences between the reference distribution and a pathway that has a skewed variability
191 distribution compared to a symmetric one. As an example, consider pathway p_1 where an effect size of 0.24
192 results in 80% power for approximately 150 genes in the pathway. For the symmetric distribution in
193 pathway p_2 , in order to obtain the same level of power, a pathway size of at least 200 genes is required (see
194 Text S1).

195

196 Results

197 *Application of pathVar to human embryonic stem cell datasets identify significant pathways with distinct*
198 *profiles of gene expression variability.*

199 To demonstrate the utility of *pathVar* in practice, the method was applied to three gene expression datasets
200 of human ESCs (Text S2). The Bock dataset [19] had twenty ESC samples that were generated using
201 microarray profiling, and the Yan dataset [20] had hESC samples profiled using single cell RNA-
202 sequencing (RNA-seq) for eight cells at passage 0 (p0), and 26 cells at passage 10 (p10). *pathVar* was run
203 independently on the three datasets, and pathways with a statistically significant deviation in their gene
204 expression variability profile relative to the reference distribution were detected (Table S1, P-value < 0.01).
205 Significant KEGG pathways reflected aberrant gene expression counts in ribosomes, metabolism (oxidative
206 phosphorylation), the spliceosome and neurodegenerative pathways (Alzheimer, Parkinson and
207 Huntington) (Table S2A-S4A). Significant REACTOME pathways fell into three main classes representing
208 cell cycle, metabolism and infectious disease (Table S2B-S4B). Considerable overlap was observed in the
209 significant results obtained between the three datasets (Fig. S2), where cell cycle was the most highly
210 represented REACTOME category. Similar percentage distributions were observed (Fig. S2),
211 demonstrating consistency in the results obtained from *pathVar* despite differences in technology platforms,
212 ESC lines and passage number.

213 We next inspect the variability count distributions for individual pathways of interest. As an example,
214 consider the Bock dataset, where the significant KEGG pathways for spliceosome, oxidative
215 phosphorylation and ECM-receptor interaction each cover different aspects of hESC regulation. For both
216 the spliceosome and oxidative phosphorylation pathways, greater transcriptional stability was manifested
217 through a significantly higher number of low variability genes, in addition to significantly fewer medium
218 variability genes compared to the reference distribution (Fig 3A-C). A significant reduction in genes with
219 high variability in the ESCs was also observed for the oxidative phosphorylation pathway compared to the
220 reference (Fig 3C). The opposite trend was observed for genes in the ECM-receptor interaction where there
221 was a significant increase in genes with medium, high and very high levels of expression variability (Fig
222 3D). This pathway also had a concurrent reduction in low variability genes compared to the reference
223 distribution.

224 *Highlighting differences between ESC and iPSC usage of global gene expression programs using*
225 *pathVar.*

226 Human iPSCs were also profiled by Bock using microarrays, and we use this data to investigate how ESCs
227 and iPSCs differ with respect to expression variability. *pathVar* identified five KEGG and thirty
228 REACTOME statistically significant pathways between the ESCs and iPSCs (Table S5, P-value < 0.01).
229 The significant KEGG pathways reflected aberrant gene expression activity in ribosome, oxidative
230 phosphorylation, DNA replication and disease processes (Huntington's disease and Parkinson's disease,
231 Table S6). The significant REACTOME terms were associated with cell cycle, splicing, and metabolic
232 processes as well as DNA replication and repair, namely homologous recombination pathways (Table S6).

233 We then compared the variability count distributions in ESCs versus iPSCs for two significant KEGG
234 pathways, oxidative phosphorylation and DNA replication (Fig. 4). Both pathways showed increased gene
235 expression stability in ESCs compared to iPSCs. For the oxidative phosphorylation pathway, a significantly
236 higher number of low variability genes were in ESCs versus iPSCs. The same trend was observed for the
237 DNA replication pathway where there was a significant reduction in the number of highly variable genes
238 in ESCs compared to iPSCs.

239 *Benchmarking pathVar results using gene expression variability versus those based on average*
240 *expression via GSEA.*

241 *pathVar* results were benchmarked against those obtained using an average expression statistic to
242 investigate the utility of expression variability analyses. Under this average-based setting, genes were first
243 classified into discrete levels using average expression, corresponding to low, medium and high levels of
244 absolute expression. This average-based implementation defaults to the generic version of GSEA [21] and
245 allows for direct comparison with the results that are obtained when studying gene expression variability.
246 Four statistically significant pathways (P-value < 0.01, Table S7) were identified as having differences in
247 average expression between iPSCs and ESCs (three REACTOME terms: Heme biosynthesis, Sphingolipid
248 metabolism, Metabolism of porphyrins; and one KEGG pathway African *trypanosomiasis*). These results
249 do not seem very informative or relevant to stem cell regulation suggesting that average expression alone
250 does not always identify the pathways involved in transcriptional control of a phenotype.

251 *Regulatory insights from other datasets confirm utility of looking at pathways with changes in both gene*
252 *expression variability and average gene expression using pathVar and GSEA.*

253 Examples using stem cells illustrate how *pathVar* works in practice, however the method can be applied
254 to virtually any gene expression data set. To highlight the generalizability of *pathVar*, we selected ten
255 other data sets that cover a variety of biological and experimental variables. Collectively, these ten data
256 sets were generated from multiple technology platforms that featured samples from human, mouse and
257 parasite which represent a range of different disease phenotypes (see Text S3).

258 Three cancer RNA-seq datasets from the Cancer Genome Atlas (TCGA) were selected; these were the
259 ovarian serous cystadenocarcinoma (OVC) [22], acute myeloid leukemia (AML) [23], and glioblastoma
260 multiforme (GBM) [24] cohorts. An infectious disease was included where transcriptomes from patients
261 infected with cerebral malaria were profiled using microarrays [25], as well as the *Plasmodium falciparum*
262 parasites that the patients were infected with [26]. A genetic disorder was featured where patient-derived
263 iPSCs were collected from Down syndrome (DS) donors and profiled using microarrays, with a set of
264 matched controls from healthy subjects [27]. A microarray data set from a normal human population via
265 the Geuvadis study (1000 Genomes Project) was used [28], as well as two mouse data sets that profiled
266 tissues from two different regions of the brain, the hippocampus and the striatum using microarrays [29].
267 *pathVar* identified statistically significant pathways from KEGG and REACTOME pathway terms for the
268 ten different one-group analyses (Table S8), and five independent two-group comparisons (Table S9) of all
269 ten data sets.

270 The results from the different analyses were used to investigate the uniqueness of analyses based on GSEA
271 versus gene expression variability. The number of significant pathways that had changes in both mean and
272 variability via GSEA and *pathVar* respectively showed that for all cases, the amount of overlap in
273 significant pathways differed depending on the data sets that were used. This suggests that average and
274 variability-based statistics reflect different ways in which cells may use their transcription programs
275 depending on the biological context (Table S8, S9). It is interesting to note that for the KEGG pathways,
276 *pathVar* results for the DS versus WT iPSC comparison, and the mouse hippocampus versus striatum
277 comparison both had zero overlap between average-based and variability-based significant pathways (Table
278 S9A). In fact, both comparisons also yielded no significant pathways with a difference in average, whereas
279 pathways were found to have a difference in gene expression variability. These two comparisons are
280 extreme examples where the analyses of gene expression variability identify changes in the transcriptional
281 program, whereas average-based analyses do not yield significant results.

282 Overall, it was apparent that the transcriptional features responsible for distinguishing one phenotype from
283 another are exerted through changes in average expression or variability in expression for key pathways.
284 To further investigate the relevance of these different modes, we focused only on the ten most significant
285 pathways from the *pathVar* results obtained for the three cancer versus normal comparisons (Table S10,
286 S11). For all three cancers, the KEGG DNA replication pathway and REACTOME “DNA strand
287 elongation term” had significant changes in both average and variability of expression. Other terms with
288 changes in both average and variability were related to DNA damage response pathways, such as “base
289 excision repair” (Table S10) for AML versus normal, and “non-homologous end-joining” for OVC versus
290 normal (Table S11).

291 Five KEGG pathways had changes in variability only that were consistent in all three cancer comparisons;
292 these were the pathways involved in Epstein-Bar virus infection, cell cycle, Fanconi anemia, lysosome and
293 apoptosis (Table S10). The Epstein-Bar virus is associated with certain kinds of cancer like lymphoma or
294 carcinoma. Apoptosis is also an important pathway for tumors because its inactivation is central in the
295 development of cancer. Similarly, for the REACTOME terms, those unique to changes in variability were
296 related to DNA repair and replication (SLBP dependent processing of replication-dependent histone pre-
297 mRNAs) for the AML and GBM comparisons. For OVC, several terms were related to the cell cycle, e.g.
298 G1 phase, cyclin D associated events in G1, cyclin A/B1 associated events during G2/M transition (Table
299 S11C).

300 Discussion

301 With pathway-centric approaches like GSEA and OR now such ubiquitous features of transcriptomic
302 analyses, *pathVar* represents a natural adjunct to this kind of analysis. Our results from analyses of ESCs
303 and other datasets have demonstrated that it is not uncommon for phenotypes to be regulated by pathways
304 that have altered levels in both average expression and expression variability, as well as pathways unique

305 to either statistic. Therefore, to derive more accurate insights into transcriptional control, our results suggest
306 that pathway-based analyses should include the detection of changes in both population statistics. *pathVar*
307 may also be used to investigate the regulatory control associated with common targets of transcription
308 factors [30,31], microRNA [32,33], or lncRNAs [34,35], genes with common variants identified from
309 genome-wide association studies [36], or other regulatory features [37] that may benefit from further study
310 of gene expression variability patterns.

311 In the analysis of ESCs versus iPSCs, *pathVar* identified very few significant pathways relative to the other
312 two-group comparisons conducted (Table S9). This result likely reflects the high degree of similarity that
313 exists between iPSC and ESC transcriptional programs. While the two cell populations have identical
314 developmental capabilities, in some instances, iPSCs retain a limited memory of the gene expression
315 program of the cell of origin. Some of the significant pathways identified by *pathVar* may point to different
316 usage of metabolic processes or the cell cycle by iPSCs and ESCs. Overall, we see more variability in the
317 iPSCs than the ESCs and the increased heterogeneity for these pathways could reflect underlying
318 differences due to donor variability, or experimental factors associated with their generation.

319 Of the six two-group comparisons performed, it is interesting to note that the ESC and iPSC comparison
320 also had the least number of significant pathways (Table S9) and this may have been due to the fact that all
321 other comparisons were between a disease and normal group, or in the case of the mouse data, between two
322 distinct regions of the brain (Text S3). This result suggests that the degree of perturbation to a transcriptome
323 in the presence of a tumor, or extra chromosome, or even a different anatomical region of the same organ,
324 is greater globally, than how iPSCs differ from ESCs.

325 The observation that pathways were significant for changes in both average expression and gene expression
326 variability reflects the different modes in which cells are using pathways to regulate transcriptional signals.
327 For the cancer-based comparisons, common themes were observed across cancer types where pathways
328 involved in DNA replication and DNA damage response had significant changes in average and variability
329 (AML versus 1000 Genomes, OVC versus 1000 Genomes, Table S10, S11). The reliance of DNA
330 replication pathways may be to facilitate the proliferative nature of tumor cells, while the pathways that
331 control DNA damage response are important for tumor cells to remain viable in the presence of increased
332 rates of mutation. This result suggests that a critical factor to understanding how cancer subverts cellular
333 pathways to promote growth and evade apoptosis more accurately may lie in focusing on how gene
334 expression is being regulated based on average expression and expression variability from cell to cell, or
335 from patient to patient.

336 Single cell heterogeneity, or inter-cellular variation is a common reality of all cell populations since even
337 isogenic cells have some degree of stochastic gene expression. Across the transcriptome, gene expression
338 variability is not distributed uniformly, and its functional contribution of transcriptional regulation at the
339 single cell-level remains largely unknown. Genes with decreased variability may be useful as potential
340 markers since they have a higher degree of generalizability, where it is easier to predict the expression state
341 for such a gene in any cell in the population. Although the *pathVar* method is applicable for both single cell
342 and bulk cell datasets, the interpretation of gene expression variability in the context of single cells would
343 provide even more precise insights into how cells are controlled by the transcriptional regulation of certain
344 pathways.

345 **Conclusion**

346 The *pathVar* method identifies pathways with aberrant distributions in gene expression variability relative
347 to either a reference distribution, or a contrasting control group. The method is based on an intuitive
348 framework where either a multinomial exact test or Chi-squared test is employed to assess the differences
349 in variability distributions for each pathway using definitions from any standard or custom annotation

350 system. A Binomial test is then used to identify genes within a specific pathway that show differences in
351 gene expression variability. Comparisons benchmarking results from *pathVar* applied to a variety of gene
352 expression data sets against those obtained using GSEA identified significant pathways showing changes
353 either in average expression, expression variability or both. These results indicate that both population
354 statistics are useful for interpreting significant alterations of pathways and gene sets that underlie
355 transcriptional regulation. The implications of these results suggest that future studies may benefit from
356 analyses of gene expression variability to complement standard analyses of average expression.

357 Acknowledgements

358 We thank Drs. Barbra Birshstein, Maureen Charron, and Deepa Rastogi for valuable feedback. We also
359 thank Raymund Bueno for testing the method.

360

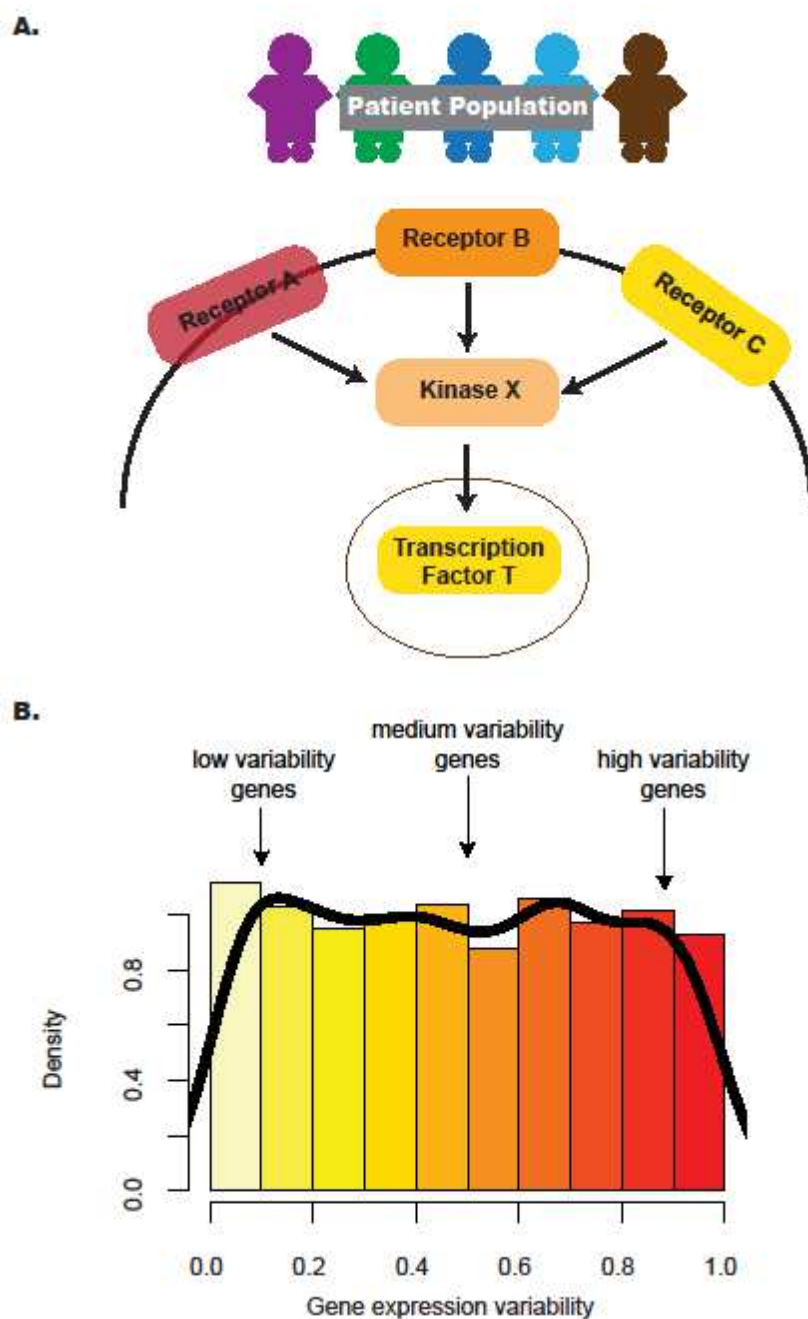
361 References

- 362 1. Mason EA, Mar JC, Laslett AL, Pera MF, Quackenbush J, et al. (2014) Gene expression variability as a
363 unifying element of the pluripotency network. *Stem Cell Reports* 3: 365-377.
- 364 2. Chalancon G, Ravarani CN, Balaji S, Martinez-Arias A, Aravind L, et al. (2012) Interplay between gene
365 expression noise and regulatory network architecture. *Trends Genet* 28: 221-232.
- 366 3. Munsky B, Neuert G, van Oudenaarden A (2012) Using gene expression noise to understand gene
367 regulation. *Science* 336: 183-187.
- 368 4. Blake WJ, M KA, Cantor CR, Collins JJ (2003) Noise in eukaryotic gene expression. *Nature* 422: 633-
369 637.
- 370 5. Raser JM, O'Shea EK (2004) Control of stochasticity in eukaryotic gene expression. *Science* 304: 1811-
371 1814.
- 372 6. Raj A, Rifkin SA, Andersen E, van Oudenaarden A (2010) Variability in gene expression underlies
373 incomplete penetrance. *Nature* 463: 913-918.
- 374 7. Burga A, Casanueva MO, Lehner B (2011) Predicting mutation outcome from early stochastic variation
375 in genetic interaction partners. *Nature* 480: 250-253.
- 376 8. Yu K, Ganesan K, Tan LK, Laban M, Wu J, et al. (2008) A precisely regulated gene expression cassette
377 potently modulates metastasis and survival in multiple solid cancers. *PLoS Genet* 4: e1000129.
- 378 9. Hasegawa Y, Taylor D, Ovchinnikov DA, Wolvetang EJ, de Torrente L, et al. (2015) Variability of Gene
379 Expression Identifies Transcriptional Regulators of Early Human Embryonic Development. *PLoS*
380 *Genet* 11: e1005428.
- 381 10. Mar JC, Matigian NA, Mackay-Sim A, Mellick GD, Sue CM, et al. (2011) Variance of gene expression
382 identifies altered network constraints in neurological disease. *PLoS Genet* 7: e1002207.
- 383 11. Mootha VK, Lindgren CM, Eriksson KF, Subramanian A, Sihag S, et al. (2003) PGC-1alpha-responsive
384 genes involved in oxidative phosphorylation are coordinately downregulated in human diabetes.
385 *Nat Genet* 34: 267-273.
- 386 12. Falcon S, Gentleman R (2007) Using GOstats to test gene lists for GO term association. *Bioinformatics*
387 23: 257-258.
- 388 13. Larsen RJ, Marx ML (2012) An introduction to mathematical statistics and its applications.: Prentice
389 Hall.
- 390 14. Wijetunga NA, Delahaye F, Zhao YM, Golden A, Mar JC, et al. (2014) The meta-epigenomic structure
391 of purified human stem cell populations is defined at cis-regulatory sequences. *Nat Commun* 5:
392 5195.
- 393 15. Fraley C, Raftery AE (2002) Model-based clustering, discriminant analysis and density estimation.
394 *Journal of the American Statistical Association* 97: 611-631.
- 395 16. Kanehisa M, Goto S (2000) KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res* 28:
396 27-30.

- 397 17. Croft D, Mundo AF, Haw R, Milacic M, Weiser J, et al. (2014) The Reactome pathway knowledgebase.
398 Nucleic Acids Res 42: D472-477.
- 399 18. Hochberg Y, Benjamini Y (1990) More powerful procedures for multiple significance testing. Stat Med
400 9: 811-818.
- 401 19. Bock C, Kiskinis E, Verstappen G, Gu H, Boulting G, et al. (2011) Reference Maps of human ES and
402 iPS cell variation enable high-throughput characterization of pluripotent cell lines. Cell 144: 439-
403 452.
- 404 20. Yan L, Yang M, Guo H, Yang L, Wu J, et al. (2013) Single-cell RNA-Seq profiling of human
405 preimplantation embryos and embryonic stem cells. Nat Struct Mol Biol 20: 1131-1139.
- 406 21. Irizarry RA, Wang C, Zhou Y, Speed TP (2009) Gene set enrichment analysis made simple. Stat
407 Methods Med Res 18: 565-575.
- 408 22. Cancer Genome Atlas Research N (2011) Integrated genomic analyses of ovarian carcinoma. Nature
409 474: 609-615.
- 410 23. Cancer Genome Atlas Research N (2013) Genomic and epigenomic landscapes of adult de novo acute
411 myeloid leukemia. N Engl J Med 368: 2059-2074.
- 412 24. Brennan CW, Verhaak RG, McKenna A, Campos B, Nounshmehr H, et al. (2013) The somatic genomic
413 landscape of glioblastoma. Cell 155: 462-477.
- 414 25. Daily JP, Scanfeld D, Pochet N, Le Roch K, Plouffe D, et al. (2007) Distinct physiological states of
415 Plasmodium falciparum in malaria-infected patients. Nature 450: 1091-1095.
- 416 26. Feintuch CM, Saidi A, Seydel K, Chen G, Goldman-Yassen A, et al. (2016) Activated Neutrophils Are
417 Associated with Pediatric Cerebral Malaria Vasculopathy in Malawian Children. MBio 7: e01300-
418 01315.
- 419 27. Briggs JA, Sun J, Shepherd J, Ovchinnikov DA, Chung TL, et al. (2013) Integration-free induced
420 pluripotent stem cells model genetic and neural developmental features of down syndrome etiology.
421 Stem Cells 31: 467-478.
- 422 28. Genomes Project C, Abecasis GR, Altshuler D, Auton A, Brooks LD, et al. (2010) A map of human
423 genome variation from population-scale sequencing. Nature 467: 1061-1073.
- 424 29. Park CC, Gale GD, de Jong S, Ghazalpour A, Bennett BJ, et al. (2011) Gene networks associated with
425 conditional fear in mice identified using a systems genetics approach. BMC Syst Biol 5: 43.
- 426 30. Lachmann A, Xu H, Krishnan J, Berger SI, Mazloom AR, et al. (2010) ChEA: transcription factor
427 regulation inferred from integrating genome-wide ChIP-X experiments. Bioinformatics 26: 2438-
428 2444.
- 429 31. Matys V, Kel-Margoulis OV, Fricke E, Liebich I, Land S, et al. (2006) TRANSFAC and its module
430 TRANSCompel: transcriptional gene regulation in eukaryotes. Nucleic Acids Res 34: D108-110.
- 431 32. Wong N, Wang X (2015) miRDB: an online resource for microRNA target prediction and functional
432 annotations. Nucleic Acids Res 43: D146-152.
- 433 33. Chou CH, Chang NW, Shrestha S, Hsu SD, Lin YL, et al. (2016) miRTarBase 2016: updates to the
434 experimentally validated miRNA-target interactions database. Nucleic Acids Res 44: D239-247.
- 435 34. Quek XC, Thomson DW, Maag JL, Bartonicek N, Signal B, et al. (2015) lncRNADB v2.0: expanding
436 the reference database for functional long noncoding RNAs. Nucleic Acids Res 43: D168-173.
- 437 35. Jiang Q, Wang J, Wu X, Ma R, Zhang T, et al. (2015) lncRNA2Target: a database for differentially
438 expressed genes after lncRNA knockdown or overexpression. Nucleic Acids Res 43: D193-196.
- 439 36. Welter D, MacArthur J, Morales J, Burdett T, Hall P, et al. (2014) The NHGRI GWAS Catalog, a
440 curated resource of SNP-trait associations. Nucleic Acids Res 42: D1001-1006.
- 441 37. Guo L, Du Y, Qu S, Wang J (2016) rVarBase: an updated database for regulatory features of human
442 variants. Nucleic Acids Res 44: D888-893.

443

444



445

446

447 **Fig 1. The distribution of gene expression variability highlights the regulatory control that different**
 448 **genes in the pathway are subjected to. A.** Absolute gene expression is a proxy for how genes are
 449 transcriptionally regulated between samples. Studying the consistency of how genes are expressed can
 450 also add information on pathway control e.g. lower levels of inter-individual variability may reflect
 451 increased regulatory control. **B.** By considering the distribution of gene expression variability, we may be
 452 able to understand transcriptional regulation in a more comprehensive manner – this is the premise of the
 453 *pathVar* method.

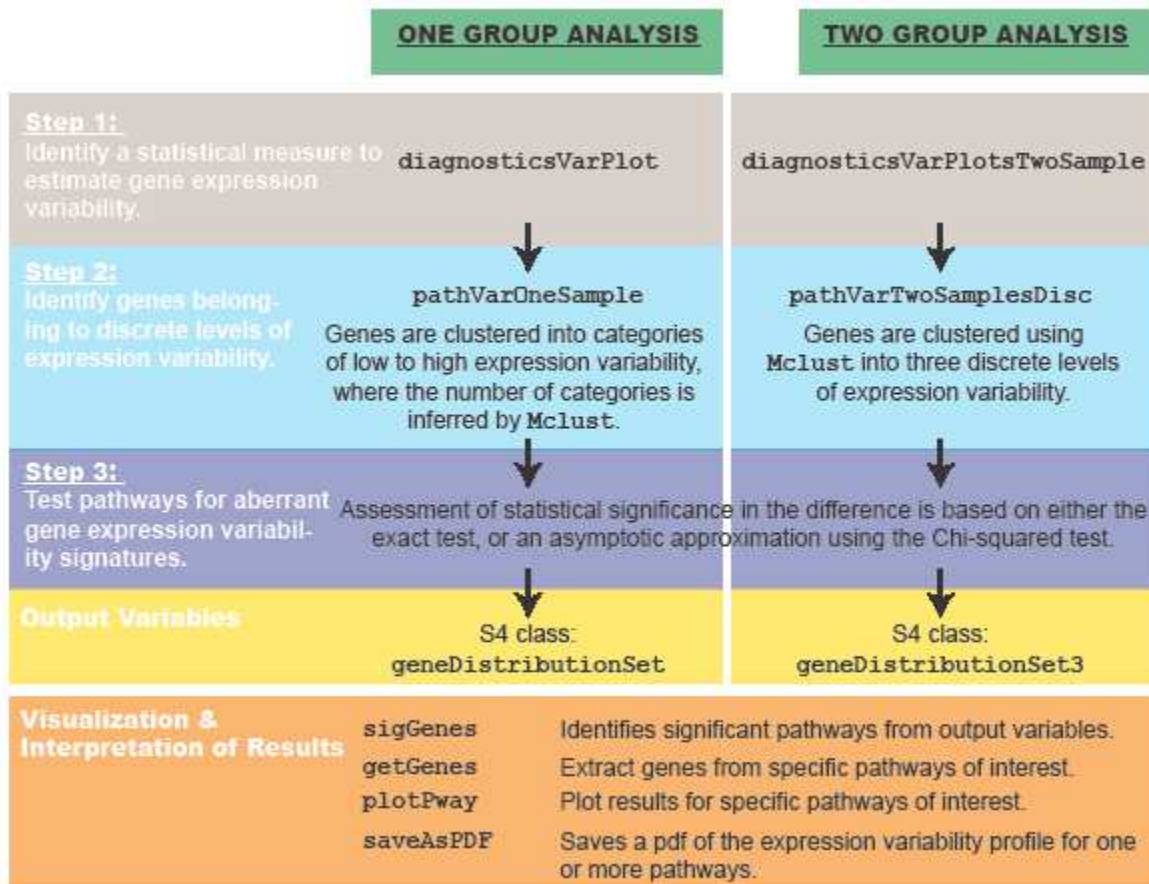
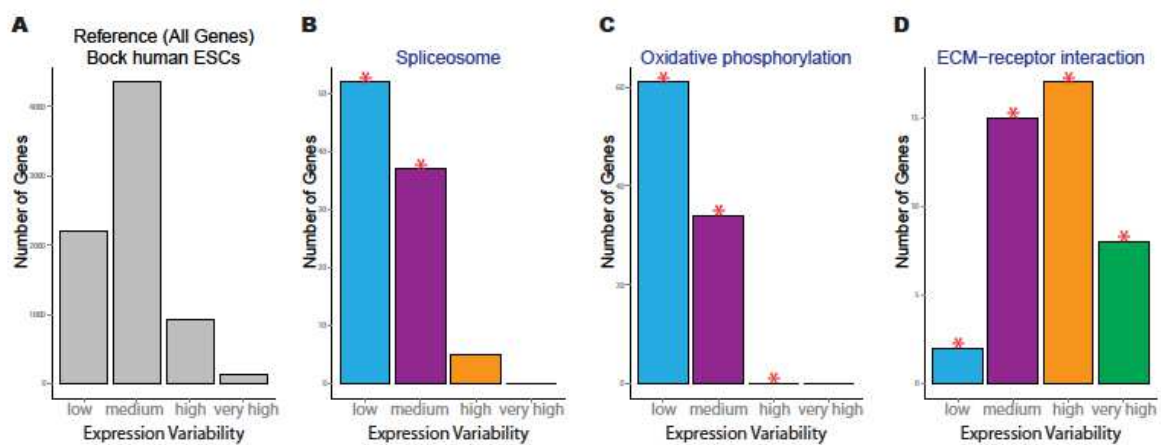


Fig. 2. Overview of *pathVar*, including the main functions in the R package.

455
456
457
458

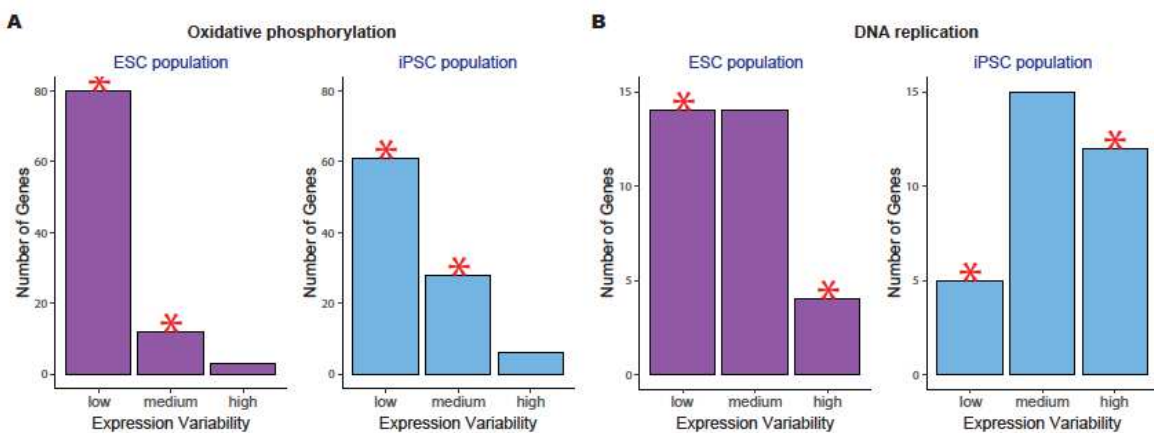


459

460

461 **Fig. 3. Example of four significant KEGG pathways for one-group pathVar analysis of the Bock**
 462 **embryonic stem cell data. A.** Variability count distribution for the reference. **B.** Spliceosome pathway
 463 (hsa03040), **C.** oxidative phosphorylation (hsa00190), **D.** ECM-receptor interaction (hsa04512). The red
 464 stars indicate a significant difference between the pathway and reference distribution for a specific level of
 465 expression variability.

466



467

468

Fig 4. Example of two significant KEGG pathways when comparing human embryonic stem cells (ESC) and induced pluripotent stem cell (iPSC) data using the two-group *pathVar* analysis. A. Oxidative phosphorylation (hsa00190), **B.** DNA replication (hsa03030). In both pathways, a higher number of genes with lower variability are present in ESCs versus iPSCs. The red stars indicate a significant difference between the two groups for a specific level of expression variability.

469