

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.

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

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

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

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

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

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

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

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

*To whom correspondence should be addressed.

Abstract

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.

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

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

Introduction

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

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

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

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

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

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

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

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

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

The null hypothesis that *pathVar* evaluates is:

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

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

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

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

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

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

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 counts $E_i = n * p_i$ are the expected number of genes in each level of expression variability within a specific pathway. The Chi-squared test achieves its computational efficiency because it is based on an asymptotic approximation, where the test becomes more accurate as n , the total number of genes in a pathway increases.

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

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

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

Results

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

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

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

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

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

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

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

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

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

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

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

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

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

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

Discussion

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

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

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

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

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

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

Conclusion

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

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

Acknowledgements

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

References

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

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

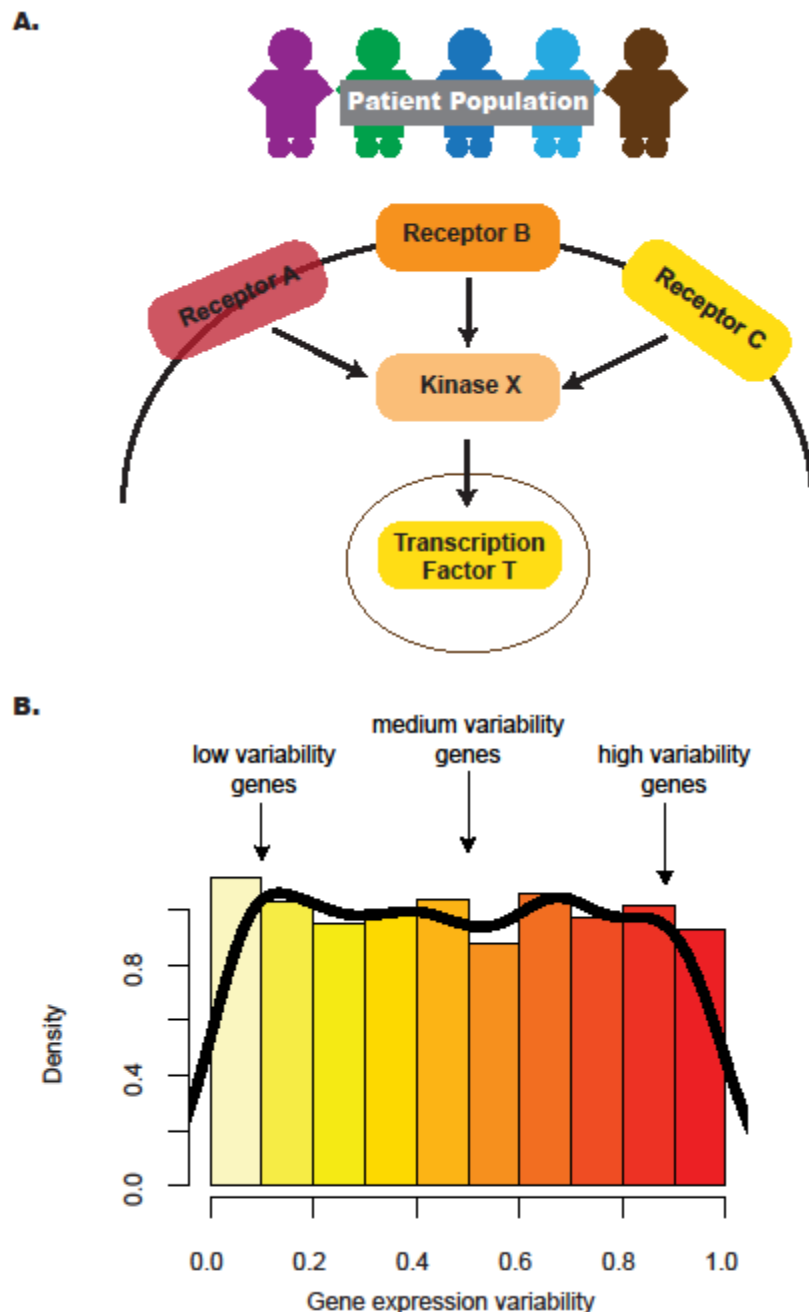


Fig 1. The distribution of gene expression variability highlights the regulatory control that different genes in the pathway are subjected to. **A.** Absolute gene expression is a proxy for how genes are transcriptionally regulated between samples. Studying the consistency of how genes are expressed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

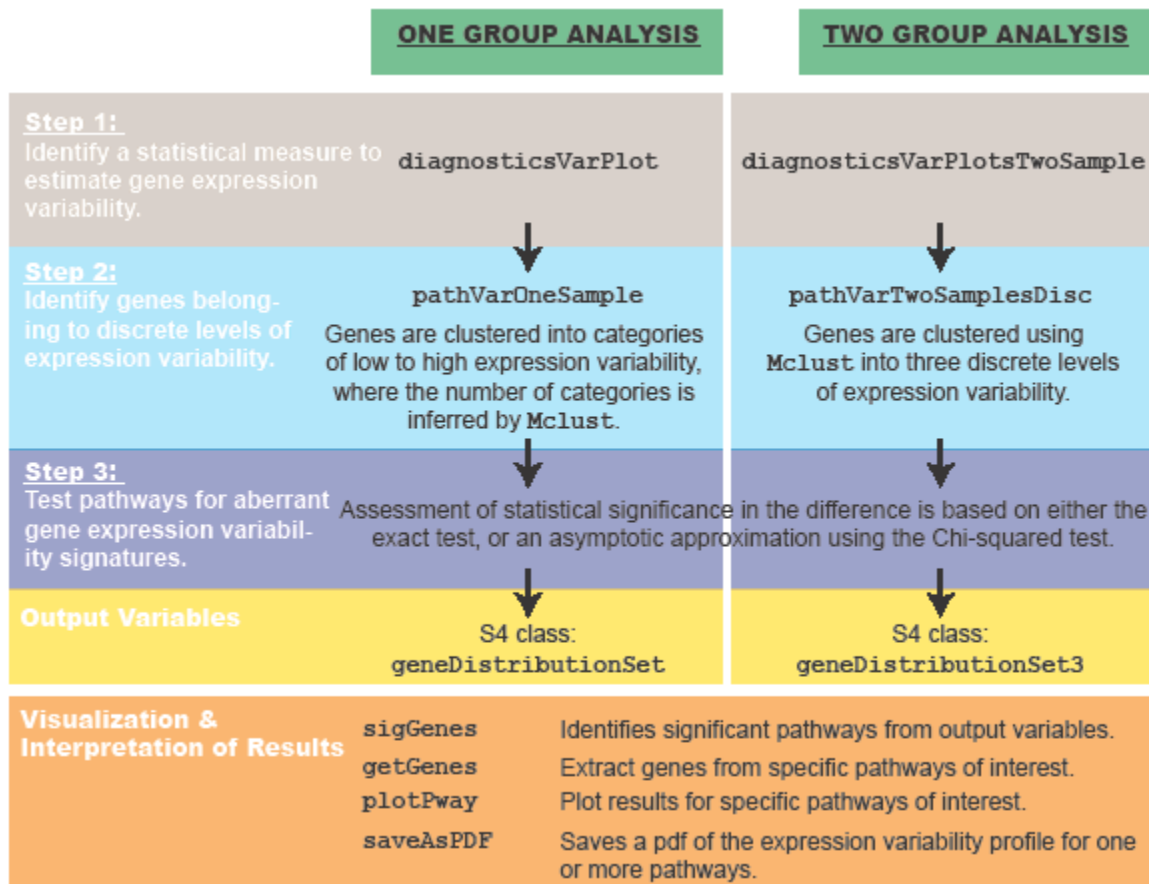


Fig. 2. Overview of *pathVar*, including the main functions in the R package.

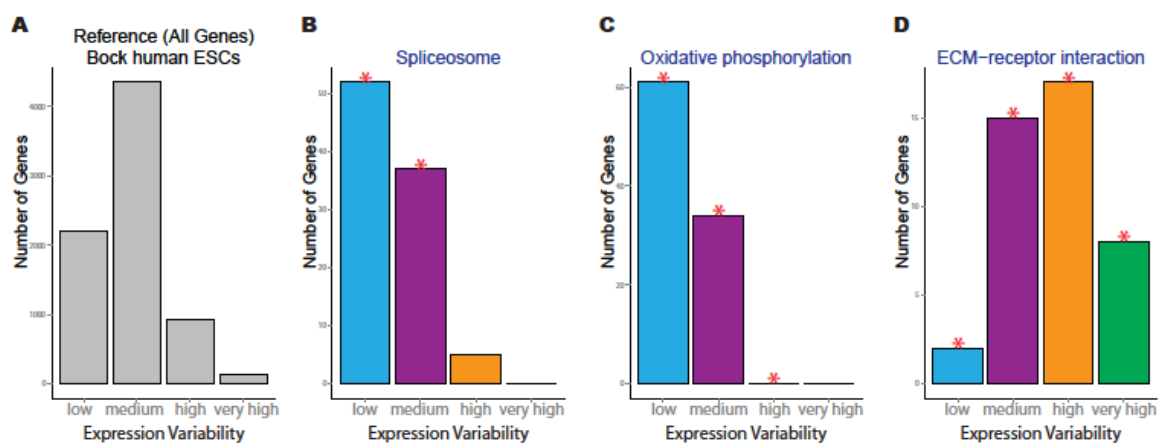


Fig. 3. Example of four significant KEGG pathways for one-group pathVar analysis of the Bock embryonic stem cell data. A. Variability count distribution for the reference. B. Spliceosome pathway (hsa03040), C. oxidative phosphorylation (hsa00190), D. ECM-receptor interaction (hsa04512). The red stars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

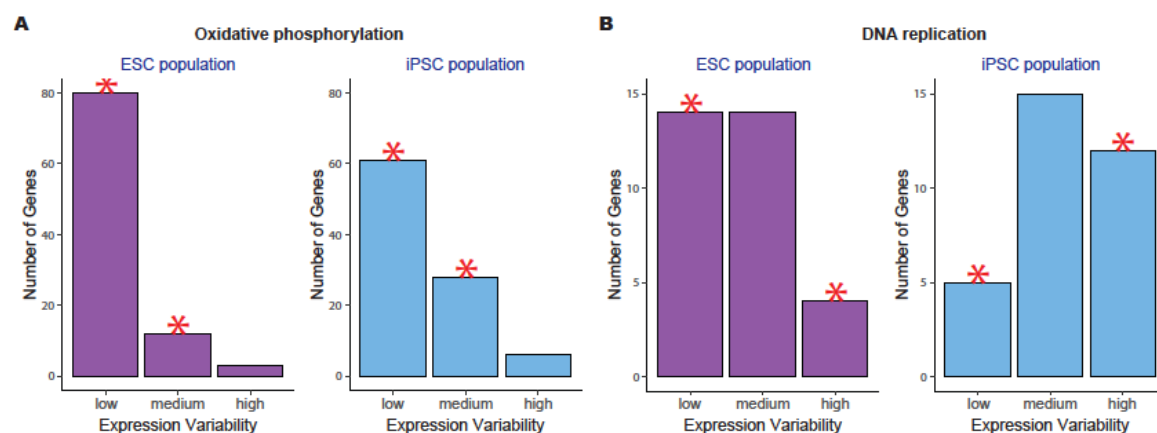


Fig 4. Example of two significant KEGG pathways when comparing human embryonic stem cells (ESC) and induced pluripotent stem cell (iPSC) data using the two-group *pathVar* analysis. A. Oxidative phosphorylation (hsa00190), **B.** DNA replication (hsa03030). In both pathways, a higher number of genes 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.