

[Click here to view linked References](#)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Human Gene Expression Variability and its Dependence on Methylation and Aging

Nasser Bashkeel¹, Theodore J. Perkins², Mads Kærn³, Jonathan M. Lee^{4,*}

¹ Department of Biochemistry, Microbiology, and Immunology, University of Ottawa, 451 Smyth Rd, Ottawa, Ontario, K1H 8L1, Canada (nbash058@uottawa.ca)

² Ottawa Hospital Research Institute, 501 Smyth Rd, Ottawa, Ontario, K1H 8L6 Canada (theodore.j.perkins@gmail.com)

³ Department of Cellular and Molecular Medicine, University of Ottawa, 451 Smyth Rd, Ottawa, Ontario, K1H 8L1, Canada (Mads.Kaern@uottawa.ca)

⁴ Department of Biochemistry, Microbiology, and Immunology, University of Ottawa, 451 Smyth Rd, Ottawa, Ontario, K1H 8L1, Canada (jlee@uottawa.ca)

* Corresponding Author

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 **Abstract**

2 **Background:** Phenotypic variability of human populations is partly the result of gene polymorphism and
3 differential gene expression. As such, understanding the molecular basis for diversity requires identifying
4 genes with both high and low population expression variance and identifying the mechanisms underlying
5 their expression control. Key issues remain unanswered with respect to expression variability in human
6 populations. The role of gene methylation as well as the contribution that age, sex and tissue-specific
7 factors have on expression variability are not well understood.

8 **Results:** Here we used a novel method that accounts for sampling error to classify human genes based on
9 their expression variability in normal human breast and brain tissues. We find that high expression
10 variability is almost exclusively unimodal, indicating that variance is not the result of segregation into
11 distinct expression states. Genes with high expression variability differ markedly between tissues and we
12 find that genes with high population expression variability are likely to have age-, but not sex-dependent
13 expression. Lastly, we find that methylation likely has a key role in controlling expression variability insofar
14 as genes with low expression variability are likely to be non-methylated.

15 **Conclusions:** We conclude that gene expression variability in the human population is likely to be
16 important in tissue development and identity, methylation, and in natural biological aging. The expression
17 variability of a gene is an important functional characteristic of the gene itself and the classification of a
18 gene as one with Hyper-Variability or Hypo-Variability in a human population or in a specific tissue should
19 be useful in the identification of important genes that functionally regulate development or disease.

20 **Keywords:** Expression Variability, Tissue Specificity, Essentiality, Methylation, Aging

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 Background

2 Within the last decade, many studies have established that gene expression patterns vary
3 between individuals, across tissue types[1], and within isogenic cells in a homogenous environment[2].
4 These differences in gene expression lead to phenotypic variability across a population. Differential gene
5 expression gene expression is typically detected by analyzing expression data from a population of
6 samples in two or more genetic or phenotypic states, for example a cancerous and non-cancerous sample
7 or between two different individuals. Various differential gene expression algorithms, such as edgeR and
8 DESeq, are then used to identify genes whose expression mean differs significantly between the states.
9 While differential co-expression analyses have successfully been used to identify novel disease-related
10 genes[3], the statistical methods used in these analyses consider gene expression variance within the
11 sample population as a component of the statistical significance estimate. However, expression variability
12 within populations has been emerging as an informative metric of cell state an informative metric of a
13 phenotypic state, particularly as it relates to human disease[4, 5].

14 There are several sources of expression variability in a population. The first are polymorphisms
15 that contribute, both genetically and epigenetically, to promoter activity, message stability and
16 transcriptional control. Another source of gene expression variability is plasticity, whereby an organism
17 adjusts gene expression to alter its phenotype in response to a changing environment[6]. However, gene
18 expression patterns can also vary among genetically identical cells in a constant environment[7–10]. This
19 is commonly described as “noise”.

20 Expression variability, whatever its source, is an evolvable trait subject to natural selection,
21 whereby each genes have an optimal expression level and variance required for an organism’s fitness and
22 selection minimizes this variability[7, 10–14]. In this case, genes with low variability have been subjected
23 to heavy selection pressure to minimize population expression variance. Conversely, high variability genes

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 have been selected for high variance. Genes with high expression variability could be drivers of phenotypic
2 diversity, as suggested by position association between expression noise and growth[15–18]. In this
3 interpretation, genes with high variability allow for growth in fluctuating environments. Understanding
4 the role of the gene expression variability patterns across human populations and in isogenic mice will
5 therefore provide crucial insights into how genetic differences contribute to phenotypic diversity,
6 susceptibility to disease[19, 20], differentiation of disease subtypes[5], development[21–24], and
7 alterations in gene network architecture[25].

8 In this analysis, we used a novel method to analyze global gene expression variability in non-
9 diseased human breast, cerebellum, and frontal cortex tissues. Our method differs from other protocols
10 in that we account for sampling error in our analysis as well as estimate expression variability independent
11 of expression magnitude. In addition, we analyzed gene methylation in conjunction with expression
12 variability. Our work suggests that expression variability is an important part of the development and
13 aging process and that identifying genes with very high or very low expression variability is one way to
14 identify physiologically and important genes.

16 Results

17 **Estimating expression variability.** We measured human gene expression variability (EV)[1] in post-
18 mortem non-diseased cerebellum (n = 465) and frontal cortex samples (n = 455) and biopsied normal
19 breast tissues (n = 144). Gene expression was measured using the Illumina HumanHT-12 V3.0 expression
20 BeadChip. We excluded probes corresponding to non-coding transcripts as well as those with missing
21 probe coordinates, resulting in a list of 42,084 probes. We chose to estimate EV of a microarray probe
22 independent of its expression magnitude. In this respect, neither the coefficient of variation nor variance

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 are suitable. The former has a bias for genes with low mean expression and the latter has a bias for high
2 mean expression genes. We modified the method initially described by Alemu et al[1]. First, we calculated
3 the median absolute deviation (MAD) for each probe. Then we modelled the expected MAD for all probes
4 as a function of median expression using a locally weighted polynomial regression (Fig. 1A, red line). The
5 expected MAD regression curves for each tissue type exhibit a flat, negative parabolic shape where the
6 lowest and highest expression probes represent the troughs of the curve. Variability in gene expression
7 levels has previously been shown to decrease as expression approaches either extrema[7, 9, 26]. The EV
8 for each probe was calculated as the difference between its bootstrapped MAD and the expected MAD at
9 each median expression level (Fig. 1A). Positive EV values indicate that the probe has a greater expression
10 variability than probes with the same expression magnitude mean. Conversely, negative EV values imply
11 reduced population expression variability. We next plotted the kernel density estimation function of EV
12 for each tissue (Fig. 1B). The EV distribution in all three tissue types exhibit large peaks around the zero
13 mean and a long tail for positive EV probes. Breast tissue exhibited a larger shoulder of the negative EV
14 probes compared to cerebellum and frontal cortex tissues. This is likely attributable to the lower number
15 of breast samples (144 compared to 456 and 455 samples respectively).

16 We then confirmed the independence of EV on expression by modelling the relationship between
17 the two variables using a linear regression (Fig 1C) and calculating the Kendall rank correlation coefficient
18 for each tissue type (Table 1). Based on the poor adjusted R² values and Kendall rank correlation
19 coefficients, we conclude that there is no substantial correlation between probe EV and expression
20 magnitude.

21 Next, we then classified each probe into three categories based on their EV. We used the term
22 “Hyper-Variable” to describe probes whose EV was greater than $\tilde{x}_{EV} + 3 * MAD_{EV}$. Probes with an EV less
23 than $\tilde{x}_{EV} - 3 * MAD_{EV}$ were deemed “Hypo-Variable”. The remaining probes that fell within the range of
24 $\tilde{x}_{EV} \pm 3 * MAD_{EV}$ were considered “Non-Variable”. A probe classified with a “Non-Variable” EV means

1
2
3
4 1 that its bootstrapped MAD is similar to the MAD of all genes with similar expression magnitude. It is
5
6 2 important to note that these probes still have expression variability across the population. We propose
7
8
9 3 that these three distinct groups, categorized based on EV, correspond to distinct functional and
10
11 4 phenotypic gene characteristics.
12
13
14
15
16

17 6 *Table 1. Correlation analysis of EV and probe expression. Adjusted R² values were calculated using a linear regression model.*

| | Breast | Cerebellum | Frontal Cortex |
|---|--------------------|--------------------|--------------------|
| Kendall Rank Correlation Coefficient | -0.208 | -0.201 | -0.213 |
| Linear Regression Adjusted R ² Value | 2x10 ⁻⁴ | 8x10 ⁻⁴ | 5x10 ⁻³ |

23
24 7
25 8 **Statistical nature of Hyper-variability.** A previously unexplored aspect of expression Hyper-variability is
26
27 9 the statistical characteristics of expression amongst genes with this wide range of gene expression.
28
29
30 10 Specifically, high EV could be the result of a multimodal distribution of gene expression with two or more
31
32 11 distinct expression means or might simply result from a broadening of expression values around a
33
34
35 12 unimodal mean value. In order to distinguish between the two possibilities, we modeled each probe
36
37 13 expression as a mixture of two Gaussian distributions prior to estimating probe EV (Fig. 2). Next, we
38
39 14 identified the peaks of the kernel density estimation function for each Gaussian distribution and
40
41 15 compared the distance between the peaks as well as the ratio of peak heights. Probes with peaks that
42
43
44 16 were greater than one median absolute deviation apart and displayed a peak ratio greater than 0.1 were
45
46
47 17 classified as having a bimodal expression distribution. Probes that did not satisfy both criteria were
48
49 18 considered to have a unimodal distribution. Only a small minority of the probes (16/41,968 breast tissue
50
51 19 probes, 6/41,968 cerebellum probes, and 6/41,968 frontal cortex probes) showed a bimodal distribution
52
53
54 20 of gene expression. The remaining majority of Hyper-Variation probes had a unimodal distribution. This
55
56 21 indicates that high expression variability is a result of a widening of possible expression values across a
57
58
59 22 single mean rather than the gene expression existing in two or more discrete states.
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2 **Accounting for sampling error in EV classification.** We were concerned that the classification of a probe
3 into Hyper-, Hypo- and Non-Variable classes might be the result of sampling errors. To minimize this
4 possibility and to increase the accuracy of our EV classification method, we divided each of our tissue
5 samples into two equally sized sample probe subsets and repeated the EV analysis. This 50-50 split-retest
6 procedure was repeated 100 times with each iterative retest using a random split of the probes. Fig 1B
7 shows the kernel density estimation function of a concordant EV classification for each probe into Hyper-
8 , Hypo- and Non-Variable class across the three subsets in each tissue type. Fig 3A demonstrates that
9 classification of a probe as Hyper or Hypo-Variable based on a single analysis of the population is
10 problematic due to sampling bias. We see a substantial decrease in the number of probes in the Hyper-
11 and Hypo-Variable probe sets after conducting our split-retest protocol (Fig. 3B and Table 2). Thus, our
12 split-retest method likely increases the robustness and accuracy of EV classification.

Table 2. Count summary of probes before and after 50-50 split-retest procedure. Hypervariable and Hypovariable probes that were not retained after the split-retest were relabeled as "Non-Variable".

| Probe Set | Tissue | Number of Probes Before Retesting | Number of Probes After Retesting | % of Probes After Retesting |
|---------------|------------|-----------------------------------|----------------------------------|-----------------------------|
| Hypervariable | Breast | 3125 | 1448 | 46.34 |
| | Cerebellum | 2987 | 1640 | 54.90 |
| | Frontal | 2949 | 1760 | 59.68 |
| Hypovariable | Breast | 4371 | 957 | 21.89 |
| | Cerebellum | 2619 | 837 | 31.96 |
| | Frontal | 3019 | 1254 | 41.54 |
| Non-Variable | Breast | 34456 | 39547 | 114.78 |
| | Cerebellum | 36356 | 39485 | 108.61 |
| | Frontal | 35994 | 38948 | 108.21 |

17 **Tissue-specificity of EV..** We next mapped Hyper-, Hypo and Non-Variable probes onto their respective
18 genes. Individual genes can have multiple probes attached to them and we refer to the identified genes

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 as being “probe-mapped”. A probe-mapped gene is assigned to a Variability group if one or more of its
2 probes have that characteristic Variability. Thus, the possibility exists that an individual gene could be
3 placed in one or more Variability groups based on differential behavior of probes mapped to that gene.
4 However, the number of genes that have are classified in one or more Variability groups involved is small
5 (Breast: 2.22%, Cerebellum: 2.76%, Frontal Cortex: 3.18%).

6 Because we have calculated EV from different tissues, we were able to determine the extent to
7 which tissue-specific factors might contribute to EV. This is an important question because expression
8 variability exists not only between individuals but between different tissues in the same organism.. As
9 shown in Fig. 4A, only a small minority of Hyper-Variable and Hypo-Variable probe-mapped gene sets are
10 shared between the three tissues. 16% of the Hyper-Variable probe-mapped genes were classified as such
11 in the three tissues and 18-26% of the Hypo-Variable were so classified. The Non-Variable probe-mapped
12 gene sets contained over 82% of genes in each tissue type, with over 71% of the measured genes
13 commonly classified as NV in all three tissue types.

14
15 **EV and gene structural characteristics.** To understand possible genomic mechanisms by which population
16 expression variability occurs, we first explored the relationship between EV and various structural features
17 of the genes. Expression variability has previously been reported to be associated with gene size, gene
18 structure, and surrounding regulatory elements[1]. However, we found no significant linear correlation
19 between EV and a gene’s exon count, sequence length, transcript size, or number of isoforms (Additional
20 file 1). While certain linear models exhibited statistical significance ($p < 0.05$), the fit of the model and
21 subsequent comparison of the linear model against a local polynomial regression curve showed that the
22 correlation was either too small to draw a conclusion or not correctly defined by a linear model.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 While we did not find that the physical gene characteristics were correlated to EV, previous
2 studies have shown that the position of a gene on a chromosome has considerable effects on stochastic
3 gene expression variability [27]. We next tested if there is a relationship between expression variability
4 and chromosomal position (Fig. 4B). To this end, each chromosome was divided into 100 bins and the
5 mean EV all the genes within each bin determined. We display mean EV so that the graphed value does
6 not depend on the probe density. However, bins that have a small number of probes may skew positional
7 values. We therefore introduced a minimal threshold for number of probes in each bin. Any bin with less
8 than 3 probes would be considered to have a zero EV value. We found that EV is not uniformly distributed
9 across the genome, and individual regions of chromosomes exhibited peaks of high expression variability
10 or troughs of low expression variability. To further confirm our conclusion, we tested the cosine
11 similarities of the chromosomes within and across the tissue types (Additional file 2). This similarity
12 analysis is consistent with the idea that EV is not randomly distributed throughout the genome.
13 Furthermore, chromosomal EV distributions across chromosomes exhibited low similarities with each
14 other. Because the probes used for the three different tissues are identical, this conclusion is not affected
15 by probe density.

16
17 **Functional analysis of Hyper-, Hypo- and Non-Variable genes.** In order to understand the overall
18 biological significance of EV, we examined the functional aspects that are enriched in the Hyper-Variable,
19 Hypo-Variable, and Non-Variable probe-mapped gene sets by conducting a gene set enrichment analysis
20 in each category. We conducted a functional enrichment analyses of the gene symbols corresponding to
21 the probes in each probe-mapped gene set. We determined the over-represented Gene Ontology (GO)
22 terms that were unique in each tissue type, as well as GO terms that were common in all three tissue
23 types. The resulting GO annotations were simplified and visualized using a REVIGO treemap. The top five
24 terms for each tissue type can be found in Table 3, while the complete list of GO term treemaps can be

1 found in Additional file 3. It should be noted that the GO term “Proteolysis involved in cellular catabolism”
 2 appears both in the “Common Probe-Mapped Genes” and “Breast-Specific Probe Mapped Genes” for the
 3 Hypo-Variable set. The genes involved in both cases are unique but they are members of the same GO
 4 pathway.

5
 6 *Table 3. Top 5 common and tissue-specific REVIGO GO annotations in the Hyper-Variable and Hypo-Variable probe mapped gene*
 7 *sets of breast, cerebellum, and frontal cortex tissues.*

| | Common Probe-Mapped Genes | Breast-Specific Probe-Mapped Genes | Cerebellum-Specific Probe-Mapped Genes | Frontal Cortex-Specific Probe-Mapped Genes |
|-----------------------|---|--|--|---|
| Hyper-Variable | Regulation of bone remodeling | Epithelial cell differentiation | Regulation of nervous system development | Histamine secretion |
| | Regulation of inflammatory response | Primary alcohol metabolism | Regulation of transmembrane transport | Regulation of cell morphogenesis |
| | Response to zinc ion | Positive regulation of cellular component movement | Regulation of neuron death | Trans-synaptic signaling |
| | Carboxylic acid biosynthesis | Response to corticosteroid | Negative regulation of response to external stimulus | Regulation of neurological system process |
| | Regulation of ion transport | Transmembrane receptor protein tyrosine kinase signaling pathway | Response to calcium ion | Dephosphorylation |
| Hypo-Variable | Proteolysis involved in cellular protein catabolism | Golgi vesicle transport | DNA conformation change | ncRNA metabolism |
| | Ribonucleoprotein complex assembly | Nucleoside monophosphate metabolism | Modification-dependent macromolecule catabolism | Response to interleukin-1 |
| | Regulation of cellular amino acid metabolism | Proteolysis involved in cellular protein catabolism | Response to camptothecin | Regulation of enter of bacterium into host cell |
| | Innate immune response activating cell surface receptor signaling pathway | Cellular response to nitrogen starvation | Retrograde transport, endosome to Golgi | |
| | Negative regulation of autophagy | Mitochondrial respiratory chain complex I assembly | Regulation of ubiquitin-protein transferase activity | |

8
 9
 10 The breast Hyper-Variable probe-mapped gene set was uniquely enriched for epithelial cell
 11 differentiation, primary alcohol metabolism, and positive regulation of cellular component movement.
 12 The cerebellum Hyper-Variable probe-mapped gene set was uniquely enriched for regulation of nervous

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 system development, transmembrane transport, and neuron death. The frontal cortex Hyper-Variable
2 probe-mapped gene set was enriched for histamine secretion, regulation of cell morphogenesis, and
3 trans-synaptic signalling. The breast, cerebellum, and frontal cortex Hyper-Variable probe-mapped gene
4 sets were commonly enriched for regulation of tissue remodeling, inflammatory responses, and responses
5 to inorganic substances. Of note, many of the enriched GO annotations of the Hyper-Variable genes are
6 involved in signalling pathways.

7 In the case of the Hypo-Variable probe-mapped gene sets, all three tissue types were enriched for
8 protein catabolism and metabolism, ribonucleoprotein complexes, and negative regulation of autophagy.

9 In this respect, many of the shared Hypo-Variable genes could be considered housekeeping genes. The
10 breast Hypo-Variable probe-mapped gene set was enriched for Golgi vesicle transport, nucleoside
11 metabolism, and protein catabolism. The cerebellum Hypo-Variable probe-mapped gene set was enriched
12 for DNA conformation change, modification-dependent macromolecule catabolism, and retrograde
13 transport.

14
15 **Essentiality enrichment in variable genes.** Previous studies in yeast have shown that gene expression
16 variability is reduced in genes that are essential for survival. It is believed that evolution has selected for
17 transcriptional networks that limit stochastic expression variation of essential genes[13]. If this were true
18 for humans, we would expect a significant number of essential genes to exhibit Hypo-Variable expression
19 and a depletion of essential genes within the Hyper-Variable probe sets.

21 *Table 4. Pearson's Chi-squared test for Essentiality in Hyper-Variable, Hypo-Variable, and Non-Variable probe mapped gene sets.*

| Tissue | Probe Set | Total Gene Count | Essential Gene Counts | Standardized Residuals | P-Value |
|--------|-----------|------------------|-----------------------|------------------------|---------|
|--------|-----------|------------------|-----------------------|------------------------|---------|

| | | | | | | |
|--|----------------|-------|-------|------|-------|--------------------------|
| | Breast | Hyper | 1448 | 165 | 8.65 | 1.48 x 10 ⁻²² |
| | | Hypo | 957 | 103 | 4.94 | |
| | | NV | 39547 | 2095 | -9.87 | |
| | Cerebellum | Hyper | 1640 | 160 | 5.88 | 4.85 x 10 ⁻¹⁰ |
| | | Hypo | 837 | 76 | 2.69 | |
| | | NV | 39485 | 2128 | -6.42 | |
| | Frontal Cortex | Hyper | 1760 | 181 | 7.28 | 1.43 x 10 ⁻¹⁶ |
| | | Hypo | 1254 | 121 | 4.15 | |
| | | NV | 38948 | 2062 | -8.38 | |

In order to examine a potential correlation between expression variability and essentiality in human tissues, we first tested the independence between EV classification and annotation of human essentiality (Table 4). Essentiality annotations were obtained from the CCDS[28] and MGD[29] databases. Here, direct human orthologs of genes essential for prenatal, perinatal, or postnatal survival of mice were classified as essential. Using the Pearson's chi-square test using the `chisq.test` function[30] in R for the number of essential genes in each probe set (Additional File 4), we find that that the Hypo-Variable probe-mapped gene set in breast, cerebellum, and frontal cortex tissues were significantly enriched for genes with essentiality annotation. Thus, expression variability for many essential genes is constrained in humans, likely reflecting a similar biology to essential yeast genes. However, we surprisingly observe a significant enrichment of essential genes within the Hyper-Variable probe-mapped gene sets.

To better understand the implications of high variability in essential genes, we examined the functional annotations associated with Hyper-Variable essential genes (Table 5 and Additional file 5). The breast essential Hyper-Variable probe-mapped gene set was enriched for chordate embryonic development, cellular response to growth factor stimulus, mesenchymal cell apoptotic process, carboxylic acid biosynthesis, and cell-substrate junction assembly. The cerebellum essential Hyper-Variable probe-mapped gene set was enriched for regulation of cell development, epithelial cell migration, positive regulation of cell proliferation, cellular response to growth factor stimulus, and anterograde synaptic signalling. Lastly, the frontal cortex essential Hyper-Variable probe-mapped gene set was

1 enriched for positive regulation of cell differentiation, transmembrane receptor protein tyrosine kinase
 2 signalling pathway, epithelial cell migration, regulation of actin cytoskeleton organization, and regulation
 3 of lipase activity. Overall, the Hyper-Variable essential probe-mapped gene sets tended to be enriched for
 4 morphogenic, tissue, and organ system development.

5
 6 *Table 5. Top 5 common and unique REVIGO GO annotation subsets of Hyper-Variable and Hypo-Variable essential genes in breast,*
 7 *cerebellum, and frontal cortex tissues.*

| | Breast-Specific Probe-Mapped Genes | Cerebellum-Specific Probe-Mapped Genes | Frontal Cortex-Specific Probe-Mapped Genes |
|---------------------------------------|---|---|---|
| | Chordate embryonic development | Regulation of cell development | Positive regulation of cell differentiation |
| Hyper-Variable Essential Genes | Cellular response to growth factor stimulus | Epithelial cell migration | Transmembrane receptor protein tyrosine kinase signalling pathway |
| | Mesenchymal cell apoptotic process | Positive regulation of cell proliferation | Epithelial cell migration |
| | Carboxylic acid biosynthesis | Cellular response to growth factor stimulus | Regulation of actin cytoskeleton organization |
| | Cell-substrate junction assembly | Anterograde trans-synaptic signalling | Regulation of lipase activity |
| | DNA repair | DNA repair | DNA repair |
| Hypo-Variable Essential Genes | Regulation of cellular protein localization | Protein oligomerization | Peptide transport |
| | Mitochondrial genome maintenance | Positive regulation of viral process | Regulation of type I interferon production |
| | Chordate embryonic development | Negative regulation of cell cycle | Response to UV |
| | Protein modification by small protein removal | Lysosomal transport | Phosphorylation |

8
 9 **DNA methylation and expression variability.** One factor that has been postulated to regulate EV is DNA
 10 methylation. While the relationship between methylation and gene expression is complex, low promoter
 11 methylation is associated with high levels of gene expression[31–34]. Like gene expression, DNA
 12 methylation is highly variable at the cell, tissue, and individual level[35], suggesting that EV could result
 13 from variations in gene methylation. To explore this idea, we used DNA methylation annotations that
 14 were available in 724 out of 911 brain tissue samples.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 DNA methylation in CpG sites is thought to be bimodal, meaning that the gene is either
2 hypomethylated or hypermethylated[34]. In order to differentiate between low, medium, and high
3 methylation states in our samples, we modelled gene methylation using Gaussian mixture models for the
4 mean methylation for each gene. The distribution of gene methylation in both cerebellum and frontal
5 cortex tissue was best modelled as a three-component system. The first component was a sub-population
6 Gaussian mixture while the second and third components were modelled as single Gaussian distributions.
7 Genes whose methylation fell within the first component were classified as Non-Methylated genes. Genes
8 were classified as Medium Methylated for those in the second component and Highly Methylated if they
9 were in third. The distribution of methylation amongst the genes is predominantly bimodal with only a
10 minority of genes being Medium Methylated (Fig. 5A). In contrast, over 62% of cerebellum genes are non-
11 methylated and 23% highly methylated. Similarly, 58% of frontal cortex genes are non-methylated and
12 22% are highly methylated).

13 Next, we explored the correlation between methylation and expression based on the EV. When
14 we subset the methylation distribution by EV classification (Fig. 5B), we observe that Hypo-Variable genes
15 have a visibly different methylation pattern than Hyper- or Non-Variable genes insofar as Hypo-Variable
16 genes are visibly overrepresented in the Non-Methylated gene group compared to both the Hyper-
17 Variable and Non-Variable genes.

18 To further quantify the overrepresentation of Hypo-Variable genes in the Non-Methylated gene
19 group, we conducted a chi-squared test of independence between the methylation state clusters and the
20 EV classifications (Table 6 and Additional file 4). Both the cerebellum and frontal cortex tissues exhibited
21 a significant relationship between the methylation clusters and EV classifications ($p = 7.57 \times 10^{-36}$ and $p =$
22 1.58×10^{-59} , respectively). By examining the standardized residuals of the chi-square test of independence,
23 we quantitatively confirmed the enrichment of Non-Methylated genes within the Hypo-Variable probe-
24 mapped gene set. We also observe a significant enrichment of Highly Methylated genes in the Non-

1
2
3
4 1 Variable gene set as well as an enrichment of Medium Methylated genes in the Hyper-Variable probe-
5
6 2 mapped gene set. This indicates that methylation and EV classification are correlated.
7
8
9

10 3
11
12 4 *Table 6. Pearson's Chi-Squared Test Standardized Residuals. We tested the independence between the methylation state clusters*
13
14 5 *and the EV classifications in cerebellum and frontal cortex tissues and found a significant relationship between the two variables*
15
16 6 *($p = 7.57 \times 10^{-36}$ and $p = 1.58 \times 10^{-59}$, respectively).*
17
18

| | Cerebellum Tissue | | | Frontal Cortex Tissue | | |
|----------------|-------------------|-------------------|-------------------|-----------------------|-------------------|-------------------|
| | Non-Methylated | Medium Methylated | Highly Methylated | Non-Methylated | Medium Methylated | Highly Methylated |
| Hypo-Variable | 11.98 | -5.69 | -9.04 | 14.84 | -7.11 | -10.79 |
| Non-Variable | -7.52 | 0.06 | 8.59 | -10.00 | -0.04 | 11.73 |
| Hyper-Variable | 0.07 | 4.21 | -3.58 | -0.23 | 6.23 | -5.47 |

19
20
21
22
23
24
25
26
27 7
28
29 8
30
31
32 9 **Effects of age, sex, and PMI on variability.** To further understand the biological relevance of EV, we
33
34 10 focused on the Hyper-Variable genes to identify potential mechanisms of decreased constraint on gene
35
36 11 expression across the samples. We systematically analyzed expression as a function of sex, age, and post-
37
38 12 mortem interval (PMI). The breast tissue dataset lacked these clinical annotations and was excluded from
39
40 13 this analysis. We employed a probe-wise linear regression analysis to model the relationship between
41
42 14 Hyper-Variable probe expression and age, sex, and PMI. The resulting p-values were adjusted for multiple
43
44 15 comparisons using the Benjamini-Hochberg procedure and considered significant when the adjusted p-
45
46 16 value was less than 0.01. The total number of Hyper-Variable probes with sex, PMI or age as co- are shown
47
48 17 in Table 7.
49
50
51
52
53
54 18
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 1 *Table 7. Probe-Wise Multiple Linear Regression of Sex, PMI, and Age. Probes that exhibit an FDR < 0.01 are considered significant*
5
6 2 *for the specific coefficient.*

| | Sex | | | PMI | | | Age | | |
|------------|-----|------|-------|-----|------|-------|-----|------|-------|
| | Up | Down | Total | Up | Down | Total | Up | Down | Total |
| Cerebellum | 12 | 10 | 22 | 2 | 0 | 2 | 247 | 267 | 514 |
| Frontal | 8 | 15 | 23 | 7 | 9 | 16 | 373 | 354 | 727 |

3
4 PMI might be a source of apparent expression variability because an extended PMI might
5 compromise sample RNA integrity and lead to degradation of labile RNA[36]. Brain samples had PMI times
6 ranging from 1 hour to 94 hours (mean = 36.14 hr), but we observe a negligible number of probes that are
7 correlated with PMI (2 out of 1640 and 16 out of 1760 probes for cerebellum and frontal cortex,
8 respectively). This suggests that sample integrity is unlikely to be a source of EV changes. Somewhat more
9 surprisingly, however, is the low number of probes that are correlated with sex. Only 22 out of 1640
10 Hyper-Variable cerebellum probes and 23 out of 1760 Hyper-Variable frontal cortex probes show sex-
11 dependent differences in EV. While other studies have shown widespread sex differences in post-mortem
12 adult brain gene expression[37], EV is not substantially dependent on sex in our analysis.

13 However, we observe that age has a substantial effect on expression variability. Age is correlated
14 with over 31% of Hyper-Variable cerebellum probes and over 41% of Hyper-Variable frontal cortex probes.
15 This means that the expression of these probes becomes either more or less constrained during aging. In
16 the cerebellum, there were 247 Hyper-Variable probes whose expression increased as a function of age
17 and 267 genes with decreased expression. Similarly, the frontal cortex contained 373 probes with
18 increased expression and 354 probes with reduced expression. Given that age is correlated with a
19 considerable number of Hyper-Variable probes, we classified the age of the samples in the cerebellum
20 and frontal cortex tissues into three age clusters according to BIC for expectation-maximization (EM)
21 initialized by hierarchical clustering for parameterized Gaussian mixture models. The oldest cluster
22 contained samples whose ages were between 58 and 98 ($\bar{x}_1 = 79$). The second cluster ranged between

1
2
3
4 1 32 and 57 years ($\bar{x}_2 = 45$), while the youngest age cluster contained samples aged 1 through 31 ($\bar{x}_3 =$
5
6 2 17).

7
8
9 3 To further explore this effect, we examined the age-dependent changes in expression of the
10
11 4 Hyper-Variable probes across the three clusters. In each tissue type, we labeled probes whose expression
12
13 5 was positively correlated with age as “Upregulated”, while the negatively correlated probes were termed
14
15 6 “Downregulated”. Then, we used a hierarchical clustering method with an expression heatmap to visualize
16
17 7 how these upregulated and downregulated probes are expressed throughout the age clusters (Fig. 6). The
18
19 8 resulting probe hierarchical trees were clustered into groups via manual tree cutting. The complete list of
20
21 9 GO term treemaps for significant gene clusters can be found in Additional file 6.
22
23
24
25
26

27 10 While the cerebellum is generally considered a regulator of motor processes, it is also implicated
28
29 11 in cognitive and non-motor functions[38]. Many of these age-dependent upregulated Hyper-Variable
30
31 12 genes corroborate previous studies exploring the relationship between brain aging and changes in gene
32
33 13 expression, including cellular responses to chemical stimuli (gold cluster). In particular, reactive oxygen
34
35 14 and nitrogen species have been shown to change ion transport channel activity, and serve as an important
36
37 15 mechanism in brain aging[39]. While all the genes selected were age-regulated, some genes exhibit outlier
38
39 16 samples whose expression remains high across all genes in the dark orange cluster, regardless of age.
40
41 17 These genes are more likely to be overexpressed in the samples as age increases and are enriched for
42
43 18 peripheral nervous system neuron development and neuron apoptotic pathways. Similar enrichments of
44
45 19 neurogenic and chemical stimuli response pathways are seen in the upregulated frontal cortex genes (gold
46
47 20 cluster). The dark orange cluster in the upregulated frontal cortex age-dependent genes exhibits a sample-
48
49 21 specific over- or under-expression of genes. These bimodally expressed genes are enriched for glial cell
50
51 22 differentiation, adenosine receptor signaling pathways, and antigen processing. Lastly, we see a random
52
53 23 scattering of expression in the yellow cluster of the frontal cortex heatmap that steadily increases with
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 1 age. These genes are enriched for glial cell differentiation, cellular response to alcohol, and defense
5
6 2 responses to fungus.
7
8

9
10 3 Most of the downregulated age-dependent Hyper-Variable genes in the cerebellum fall into the
11
12 4 green cluster where expression of the genes in the cluster increases with age. These genes are involved
13
14 5 in leukocyte-mediated immunity and defense responses to other organisms, which is supported by
15
16 6 previous studies[40]. Interestingly, the yellow cluster exhibits U-shaped expression levels, whereby the
17
18 7 lowest expression is seen in the middle age cluster. These genes are enriched for optic nerve
19
20 8 development, response to interferon-gamma, and synaptic signalling. In the frontal cortex, the majority
21
22 9 of downregulated age-dependent genes fall in the red cluster, and are enriched for ion transport, cell
23
24 10 morphogenesis, and trans-synaptic signalling. Overall, the functional annotations of the age-regulated
25
26 11 Hyper-Variable gene clusters suggest that population EV is one outcome of age-dependent gene
27
28 12 expression changes.
29
30
31
32
33

34 13 We next investigated a possible impact of methylation status on gene expression in the Up- and
35
36 14 Down-regulated Hyper-Variable genes. Fig. 7 shows the histogram distribution of correlation between
37
38 15 paired gene expression and gene methylation for each gene. We observe no strong correlation between
39
40 16 expression and methylation, suggesting age-dependent changes in expression of the age-regulated Hyper-
41
42 17 Variable genes are not the result of methylation changes.
43
44
45
46
47
48
49

50 19 **Discussion**

51
52

53 20 Gene expression variability in a population is the cumulative result of intrinsic genetic factors,
54
55 21 extrinsic environmental factors, and stochastic noise. A fundamental issue in biology is understanding the
56
57 22 cause of expression variability within an individual organism and between isogenic and genetically
58
59
60
61
62
63
64
65

1
2
3
4 1 dissimilar individuals of a population [42]. Expression variability has been postulated to be part of
5
6 2 evolution, differentiation and organ homeostasis [43,44]. In this report, we study population gene
7
8
9 3 expression variability in human breast, cerebellum, and frontal cortex tissues.

10
11 4 Our investigation into human gene expression variability yielded several main findings. First, we
12
13 5 find that Hyper-Variability in population gene expression is fundamentally unimodal and does not
14
15 6 represent population switching between two or more discrete expression stages. In addition, both Hypo-
16
17 7 Variable (highly constrained expression) and Hyper-Variable (lowly constrained expression) probe-
18
19 8 mapped gene sets are enriched for essential genes. We observe only a small (16-26%, Figure 4A) overlap
20
21 9 in Hyper- and Hypo-Variable probe-mapped gene sets between the three tissues, consistent with the idea
22
23 10 that EV could be controlled by tissue-specific factors. We also find that gene methylation could have a
24
25 11 role in expression variability. Lastly, we find that only a small number of Hyper-Variable probe-mapped
26
27 12 genes exhibit co-variability with sex (22/1640 cerebellum probes, and 23/1760 frontal cortex probes). On
28
29 13 the other hand, substantially more Hyper-Variable probes exhibit a strong linear association with age
30
31 14 (514/1640 in Cerebellum and 727/1760 in Frontal Cortex).

32
33
34
35
36
37
38 15 A confounding issue with our study is the bulk nature of the tissue samples used. It is likely that
39
40 16 multiple cell types are found in each tissue sample and that the magnitude of this heterogeneity varies
41
42 17 between samples. This issue is not unique to our study and is common to all non-single cell sequencing
43
44 18 studies. With respect to expression variance, cell type heterogeneity is likely to manifest itself in the
45
46 19 identification of a gene as Hyper-Variable based on the fluctuating presence of a cell type with a unique
47
48 20 gene expression profile. This could be one explanation for the presence of cell-type specific process in the
49
50 21 Hyper-Variable genes associated with aging (e.g. Glial Cell Differentiation) or in the Frontal cortex-specific
51
52 22 Hyper-Variable genes (e.g. Histamine Secretion).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 However, tissue heterogeneity is only one possible explanation for Hyper-Variability. We have
2 several reasons to suspect that tissue heterogeneity and concomitant sampling heterogeneity does not
3 fatally compromise our analysis. First, we used large samples sizes ($n > 400$) which would help mitigate (but
4 not completely eliminate) heterogeneity issues. Secondly, we identified common Hyper-Variable genes
5 between the breast, cerebellum and frontal cortex. Because of the drastic tissue type differences between
6 these three tissues, we propose that tissue composition heterogeneity is a poor explanation for high
7 variance gene expression common across these three tissue types. Rather, we propose that this common
8 high variability reflects an important functional descriptor of the genes involved. Lastly, we observed that
9 Hyper-Variable probes have an almost exclusively unimodal expression pattern (41,956/41,968 breast
10 tissue probes, 41,962/41,968 cerebellum probes, and 41,962/41,968 frontal cortex probes). This is
11 significant because it suggests that high EV is not the result of a chance observation of rare cell types with
12 an unusual gene expression pattern. Nonetheless, we acknowledge that this study has not taken tissue
13 heterogeneity into account and is a caveat to our interpretations of Hyper-Variability. Ideally, single cell
14 analysis or sorting of the cell samples will clarify the issue. In one single cell study, Osorio et al [45] used
15 single cell RNA-Seq to estimate gene expression variability in genetically identical human cells of three
16 different types. Their analysis revealed that within these lines, subsets of genes with high and low
17 expression variability could be found. They also found a positive correlation between a gene's expression
18 variability within a specific cell group to its variability between individuals in a population. Some genes,
19 notably those with GO annotations for B cell activation involved in the immune response, cytokine
20 receptor activity, cellular response to drug, and regulation of tyrosine phosphorylation of STAT protein,
21 have a strong correlation between expression variability in single cells and in that in the population. Thus,
22 it is likely that some of the HyperVariable genes we identified from our individuals will be genes with
23 highly variable expression amongst cells of the same type.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 On the other hand, our identification of Hypo-Variable probe-mapped genes is not affected by
2 any potential tissue and sampling heterogeneity. These Hypo-Variable probes exhibit a restricted range
3 of expression values in each of the samples, independent of sample heterogeneity. Shared GO annotations
4 provided by the functional enrichment analysis of the Hypo-Variable probe-mapped genes in breast and
5 brain tissues (Table 3) indicate that many of these genes are likely to have housekeeping functions. The
6 definition of what constitutes a housekeeping gene is arbitrary but, in a traditional sense, it implies a
7 strong requirement in all cell types of an organism and a limited tolerance for variations in gene
8 expression. Some common Hypo-variable genes that would typically be considered housekeeping ones
9 include genes for Ribonucleoprotein Complex Assembly and Regulation of Cellular Amino Acid
10 Metabolism and Proteolysis. However, we were surprised to find a broad range of functional annotations
11 amongst the Hypo-Variable genes. Amongst these are Negative Regulation of Autophagy, Cellular
12 Response to Nitrogen Starvation, and Response to Interleukin-1, which would be typically be thought of
13 as induced processes rather than obligate ones. Thus, tissues tightly regulate the expression of genes in a
14 wide variety of processes and Hypo-Variability, similar to Hyper-Variability, is likely to be an important
15 physiological characteristic of a gene.

16 The enrichment of essential genes in the Hypo-Variable probe-mapped gene sets is in agreement
17 with previous findings in yeast showing that essential yeast genes are likely to have low expression
18 variability. However, we detected a significant number of essential genes amongst the Hyper-Variable
19 probe-mapped gene sets in breast, cerebellum, and frontal cortex tissue. Inactivation of these essential
20 genes leads to pre- or neonatal fatality in mice and humans[46]. This was a surprise to us since we
21 expected that expression of developmental genes should be tightly regulated. Our functional enrichment
22 analysis indicates that these Hyper-Variable genes are enriched for morphogenic, tissue, and organ system
23 development, consistent with an “essential” function yet we observe highly variable expression being
24 tolerated. One possible explanation would be tissue heterogeneity in the samples (see above). Another

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 possibility is these “essential” genes are required for embryonic development but have different post-
2 embryonic roles and may not be essential postnatally. Alternatively, it is possible that these essential
3 genes are not dose-sensitive in humans, meaning that only a certain level of baseline expression is
4 required and expression above this baseline might be well tolerated. One additional possibility is that their
5 protein abundance could be regulated translationally rather than transcriptionally. Inefficient translation
6 of certain genes may have been selected for during evolution to prevent fluctuations in protein
7 concentrations[32]. Perhaps a combination of these factors is at play.

8 The non-random distribution of Hyper-Variable and Hypo-Variable genes across the genome
9 suggests that EV is dependent on epigenetic factors. Examining the methylation status of the genes
10 allowed us to determine the relationship between gene methylation and expression variability. Firstly, we
11 find that Non-Variable genes in the cerebellum and frontal cortex are likely to have high gene methylation.
12 Secondly, we find that Hypo-Variable genes are likely to be non-methylated. We propose a model for
13 methylation-dependent expression variability where the highly constrained levels of Hypo-Variable gene
14 expression require non-methylated genes. We speculate that the lack of methylation allows
15 transcriptional regulators requiring non-methylated DNA for binding to tightly control gene expression.
16 On the other hand, high gene methylation reduces transcription noise and epigenetically inhibits
17 promoter variability in human populations. Future studies should investigate the role that these putative
18 regulators of expression play on EV, including cis-regulatory elements and transcription factors.

19 We find that there is limited (<26%) overlap in gene identity between Hyper- and Hypo-Variable
20 probes in breast and brain tissue. Indeed, the chromosomal pattern of EV differs between tissue types.
21 Our favored explanation for this is that tissue identity is created and preserved, at least in part, by changes
22 in gene expression control pathways. Thus, genes mapped by Hypo-Variable probes in any given tissue
23 have a constrained expression pattern because they are likely to be important in the tissue-specific
24 function and physiology of that organelle. While there is limited overlap of genes within the corresponding

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 EV probe-mapped gene sets of different tissues, the Hyper-Variable probe-mapped gene sets of the
2 different tissues have similar functional enrichments and cellular protein localizations. Specifically,
3 proteins encoded by genes mapped by Hyper-Variable probes tend to localize at the cell periphery and
4 are enriched for cell surface signalling pathways and tissue development, including tissue remodeling and
5 ion transport. In this respect, our work is broadly consistent with previous findings on transcript
6 abundance in mice[23, 24]. We therefore propose that tissue identity involves high expression variability
7 in specific tissue development pathways.

8 We did not observe any substantial sex dependent effects in expression variability. However, an
9 important conclusion of our study is that many Hyper-Variable probes have age-dependent expression
10 variability: that is, their expression significantly increases or decreases during aging. One main cause of
11 accelerated brain aging and a causal factor of neurodegeneration is a reduction in immunological
12 functions[47, 48]. We see evidence of downregulated immune responses in the cerebellum, specifically
13 Leukocyte Mediated Immunity, Defense Responses to Other Organisms, and Interferon-Gamma Response
14 pathways. Many studies also suggest that aging is associated with the upregulation of inflammatory
15 responses[49], which is a pathogenic mechanism implicated in many age-related diseases, including
16 cardiovascular disease, Alzheimer’s disease, and Parkinson’s disease[50]. Consistent with this idea, we see
17 an enrichment of acute inflammatory response in the cerebellum gold cluster. Another mechanism that
18 has been implicated with age-related diseases, such as Alzheimer’s disease and Parkinson’s disease, is
19 synaptic dysfunction that can affect neuroendocrine signaling[51–53]. We see a downregulation of ion
20 transport and trans-synaptic signaling in the frontal cortex, which are key components of
21 neurotransmission and membrane excitability, and whose downregulation likely causes deficiencies in
22 these complex processes. Furthermore, we see an upregulation of genes associated with glial cell
23 differentiation in the frontal cortex across multiple gene clusters. Initially thought of as cells that merely
24 support neurons, emerging research shows that neuron-astrocyte-microglia interactions are crucial for

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 the functional organization of the brain[54]. In addition, genes specific to astrocytes and
2 oligodendrocytes, two different types of glial cells, have been shown to shift regional expression patterns
3 upon aging, and are better predictors of biological age than neuronal-specific genes[55]. This suggests
4 that the Hyper-Variability and age-dependent upregulation of genes associated with glial cell
5 differentiation or an increase in the number of glial cells in the samples.

6 Without examining the mechanistic control of individuals genes, it is difficult to determine if
7 changes in gene expression result in repression or activation of their associated pathways. For example,
8 we see an upregulation in neurogenesis-associated genes during aging in both the cerebellum and the
9 frontal cortex, despite the common theory that neurodegeneration is a ubiquitous effect of normal brain
10 aging. An emerging concept in neuroscience is that homeostatic plasticity of neurons is maintained
11 through local adjustments of neural activities[56]. This overexpression of genes in pathways whose
12 function is known to decline over time may be a compensatory mechanism for an inefficient, aging system.
13 Within the cerebellum, a decline in neuronal function that occurs with aging may cause an upregulation
14 of genes associated with neurogenesis pathways. In addition to mitigating neuronal dysfunction, localized
15 increases in neurogenesis may be induced in response to cerebral diseases or acute injuries for self-
16 repair[57]. Lastly, chronic antidepressant usage has also been shown to result in an increase in
17 neurogenesis[58], suggesting that psychopharmaceuticals can alter neurochemistry and mimic
18 compensatory anti-aging responses. Overall, EV plays an important role in aging, specifically in immune
19 responses and inflammation, neurotransmission, and neurogenesis. Age-dependent gene expression
20 could reflect a loss of regulatory control or be a part of a regulated pathway of development.

21
22 **Conclusion.**

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 Our work shows that gene expression variability in the human population is likely to be important in
2 development, tissue-specific identity, methylation, and in aging. As such, the EV of a gene is an important
3 feature of the gene itself. Therefore, the classification of a gene as one with Hypervariability or
4 Hypovariability in a human population or in a specific tissue should be useful in the identification of
5 important genes that functionally regulate development or disease. In addition, we propose that the split-
6 retest procedure describer here is a useful technique for quantifying gene expression differences in a
7 sample population.

8

9 **Methods**

10 **Illumina gene expression and methylation microarray data.** The analysis was conducted on two separate
11 datasets, both utilizing the Illumina HumanHT-12 V3.0 expression BeadChip. The first dataset provides
12 high quality RNA-derived transcriptional profiling of breast-adjacent tissue from 144 samples. The
13 associated genotype and expression data have been deposited at the European Genome-Phenome
14 Archive (EGA, <http://www.ebi.ac.uk/ega/>), which is hosted by the European Bioinformatics Institute,
15 under accession number EGAS00000000083. The microarray readings were preprocessed using the
16 author's own custom script based on existing functionality within the beadarray package[59] in R and
17 were reported as a log2 intensity. This dataset is referred to as breast tissue.

18 The second gene expression and the methylation datasets were catalogued by the North
19 American Brain Expression Consortium and UK Human Brain Expression Database (UKBEC)[37, 60]. The
20 expression data was obtained from the Gene Expression Omnibus (GEO) database[61] under accession
21 number GSE36192. A total of 911 tissue samples were analyzed from frozen brain tissue from the
22 cerebellum and frontal cortex from 396 subjects (Table 8). The microarray readings were processed using

1 a cubic spline normalization method in Illumina Genome Studio Gene Expression Module v3.2.7. The
 2 expression levels were log2 transformed before any analysis. The methylation data was also obtained
 3 from GEO under accession number GSE36194. A total of 724 tissue samples were analyzed from frozen
 4 brain tissue from the cerebellum and frontal cortex from 318 subjects. The methylation microarray
 5 readings were processed using BeadStudio Methylation Module v3.2.0 with no normalization.

6 *Table 8. Description of brain sample dataset cohorts. Clinical annotations were not available for breast tissue samples.*

| Clinical Annotation | Dataset | Min | Q1 | Median | Mean | Q3 | Max |
|---------------------|----------------|--------------------|------------------|--------------------|------------------|----|-----|
| Age | Expression | 1 | 24 | 46 | 47.79 | 71 | 98 |
| | Methylation | 1 | 21 | 44 | 47.48 | 74 | 96 |
| PMI | Expression | 1 | 14 | 25 | 36.14 | 61 | 94 |
| | Methylation | 1 | 14 | 21 | 26.65 | 36 | 62 |
| | Dataset | Females (n) | Males (n) | Females (%) | Males (%) | | |
| Sex | Expression | 289 | 622 | 31.72% | 68.28% | | |
| | Methylation | 243 | 481 | 33.56% | 66.44% | | |

7
 8 **Preprocessing the datasets.** Since the brain expression and methylation datasets were individually
 9 processed by different tissue banks and in several batches, we corrected for the batch effect using the
 10 limma package[62] in R. The breast tissue dataset was previously batch corrected by the authors. Next,
 11 we subset the data into groups based on the available clinical annotations provided by the NABEC/UKBEC
 12 database. These annotations included tissue type (Cerebellum and Frontal Cortex), sex (Male and Female),
 13 and age (ranging from 0 to 98 years old). We clustered the age annotations into groups using a K-Means
 14 clustering algorithm (Additional File 7), whereby the optimal number of clusters was determined using
 15 the elbow method. After four clusters, the change in total within-clusters sum of squares did not explain
 16 a significant amount of additional variance, therefore k=3 was chosen as the optimal number of clusters
 17 for the age annotation. We then converted the continuous, numeric age annotation into three categorical
 18 age groups (0-21 years, 22-73 years, 74+ years).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 We then compared the 12 possible clinical annotation permutations to determine the optimal method to
2 subset the brain samples. For each of the 12 groups, we calculated the median expression for each probe
3 and performed a hierarchical clustering via multiscale bootstrap resampling using the pvclust package[63]
4 in R (Additional File 7). Using an approximately unbiased (AU) p-value of 0.99, analogous to a p-value
5 significance level of 0.01, the ideal clustering method was to subset the data solely by tissue type. Thus,
6 we divided the brain dataset into the cerebellum tissue and frontal cortex tissue datasets. Due to the
7 paired nature of the methylation and expression data, the methylation brain dataset was also subset into
8 cerebellum and frontal cortex tissue subsets.

9
10 **Estimating expression variability.** To calculate a magnitude-independent measure of variability for
11 expression and methylation, we used a modified method described in Alemu et al[1]. Briefly, we first
12 calculated a bootstrapped estimate of the median absolute deviation of each gene using 1000 bootstrap
13 replicates. Next, a local polynomial regression curve (loess function with default parameters on R version
14 3.4.2) was used to determine the expected gene expression MAD as a function of the median value. No
15 additional smoothing was used for the regression curve. We calculated gene EV as the difference between
16 the bootstrapped MAD and the expected MAD at each gene's median expression level.

17
18 **Identification and removal of bimodal expression probes.** Probes expressions that exhibited a bimodal
19 distribution were thought of as having two exclusive phenotypic states. However, our focus in this analysis
20 was to examine the factors affecting the tightly regulated expression of Hypo-Variable probes or the highly
21 variable gene expression of Hyper-Variable probes. In order to identify if a gene's expression was
22 unimodal or bimodal, we modeled each gene expression as a mixture of two gaussian distributions using
23 the mixtools package[64] in R. Next, we identified the peaks of the kernel density estimation functions for

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 each gaussian distribution and compared the distance between the peaks as well as the ratio of peak
2 heights. Probes with peaks that were greater than one MAD apart and displayed a peak ratio greater than
3 0.1 were treated as having a bimodal expression and subsequently removed from the analysis. Probes
4 that did not satisfy these criteria were considered to have a unimodal distribution and were kept for
5 further analysis.

6
7 **EV gene set classification.** We classified the probes into three distinct probes sets based on their
8 expression variability:

$$\tilde{x}_{EV} \pm 3 * MAD_{EV} \quad (1)$$

10 where \tilde{x}_{EV} is the EV median for each dataset, and MAD_{EV} is the bootstrapped estimate of the median
11 absolute deviation of EV using 1000 replicates. Probes whose EV fell within the range were considered
12 Non-Variable, those above this range termed Hyper-Variable, and the remaining were considered Hypo-
13 Variable.

14 For the subsequent analyses, we used the probe sets for initial classifications then proceeded with the list
15 of corresponding gene symbols. As such, there is a small subset of duplicate gene symbols in different EV
16 classifications. However, the small number of duplicate genes does not significantly affect the results of
17 the analyses.

18
19 **Bootstrapping EV gene set classifications.** To statistically validate our EV classifications, we split our data
20 into two equally sized subsets and repeated the previously explained EV method. This 50-50 split-retest
21 procedure was repeated 100 times per tissue. Next, we determined the accuracy our of original
22 classifications by comparing original classification of each gene with the 50-50 split classifications using a

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 binomial test with a probability of success greater than 0.5. In this hypothesis, a “success” is defined as
2 consistent EV classification across all three subsets, and gene classifications were considered significant
3 with a p-value < 0.05. We also calculated the methylation variability (MV) using an identical method to
4 EV, but did not find significant correlations between any MV classes and EV classes based on Spearman's
5 rank-order correlation (Additional File 8).

6
7 **Structural analysis of EV genes.** Data regarding the structural features of the genes was obtained from
8 the GRCh38/hg38 assembly of UCSC Table Browser[65]. Linear regression analyses were conducted to find
9 any correlation between gene EV and their structural features. For the linear regression analysis of
10 transcript size, we individually examined the largest and smallest transcripts separately. The sequence
11 lengths excluded introns, 3' and 5' UTR exons, and any upstream or downstream regions.

12
13 **Gene cluster analysis.** The GO term enrichment analyses were conducted using ConsensusPathDB gene
14 set over-representation analysis[28]. The complete list of unique Illumina HumanHT-12 V3.0 expression
15 BeadChip genes was used as a background list of genes. The resulting GO terms were then filtered
16 manually using a q-value cutoff of 0.05. Common and unique GO terms were summarized using
17 REVIGO[66] and visualized through treemaps by the provided R scripts. The parameters used were a
18 medium allowed similarity (0.7) using Homo sapiens database of GO terms.

19
20 **Enrichment analyses.** Using the Pearson's chi-square test, we tested for enrichment of essential genes in
21 each probe-mapped gene set relative to the total number of essential genes in the Illumina HumanHT-12
22 V3.0 expression BeadChip. A list of 20,029 protein coding genes from the CCDS database was used to test

1
2
3
4 1 for essentiality enrichment[28]. Only genes that are solely classified as essential are considered in the
5
6 2 analysis, resulting in a list of 2377 essential genes present in the dataset. Once the number of annotated
7
8 3 genes and gene sets were deemed dependent variables, we determine the enrichment of annotated
9
10 4 genes using the Pearson residuals.
11
12
13

14 5 The Pearson's chi-square test was also used to test the enrichment of methylation clusters across
15
16 6 the Hyper-Variable, Hypo-Variable, and Non-Variable probe sets.
17
18
19

20 7
21
22
23 8 **Classification of methylation status.** In order to merge the brain expression dataset with the brain
24
25 9 methylation dataset, we first identified the corresponding ID_REF to match the samples from each
26
27 10 dataset. Since we could not match specific expression probe mappings to specific methylation probe
28
29 11 mappings of CpG islands, we calculated the median probe values with a single gene mapping for both
30
31 12 expression and methylation for each sample. This resulted in a list of median expression and median
32
33 13 methylation of each gene for each sample. Next, we calculated the correlation between paired expression
34
35 14 and methylation values for each gene. Lastly, we classified the genes into one of three methylation
36
37 15 clusters based on their median methylation using Gaussian mixture models for each tissue type. In both
38
39 16 the cerebellum and frontal cortex tissue, the distribution of median gene methylation was best modelled
40
41 17 as a three-component system. The first component was a sub-population Gaussian mixture while the
42
43 18 second and third components were modelled as single Gaussian distributions. Genes whose methylation
44
45 19 fell within the first component were classified as Non-Methylated genes. Genes were classified as Medium
46
47 20 Methylated for those in the second component and Highly Methylated if they were in third.
48
49
50
51
52

53
54 21
55
56
57 22 **Hierarchical clustering of age-dependent Hyper-Variable genes.** With the exception of a few groups, the
58
59 23 hierarchical clustering groups with the opposite sex and the same age groups tended to cluster together.
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 While the p-values of the sex and age groupings during the hierarchical clustering were too high to warrant
2 further subsetting of the brain dataset samples into distinct groups, they were significant enough to
3 inspect on a gene-by-gene basis.

4 We used a multiple linear regression model to measure the changes in expression of the Hyper-
5 Variable probes as a function of age, sex, and post-mortem interval (PMI):

$$6 \quad Y_i = \beta + \beta_1 Age + \beta_2 Sex + \beta_3 PMI \quad (2)$$

7 where Y_i is the expression level of a probe and β_n is the coefficient for each term. The p-values were
8 calculated using a type III sum of squares regression and adjusted for multiple comparisons using the
9 Benjamini-Hochberg method. Probes that exhibit an FDR < 0.01 were considered significant for the
10 specific coefficient, and the sign of the coefficient determines if the probe is positively or negatively
11 correlated with the factor.

12 The choice to use three age clusters as the optimal number of clusters to examine changes of EV
13 across age samples was determined using an expectation-maximization (EM) algorithm initialized by
14 hierarchical clustering for parameterized Gaussian mixture models in the mclust package of R. The
15 Bayesian information criterion for each hierarchical clustering model was determined, and both the
16 cerebellum and frontal cortex displayed identical optimal numbers of age clusters. Once the samples were
17 correctly clustered by age, the gene clusters were selected by cutting the gene dendrograms manually.
18 The gene expressions were then visualized as heatmaps using the gplots package[67] in R.

19

20 **List of Abbreviations**

21
22 AU Approximately Unbiased

| | | | |
|----|----|-------|------------------------------------|
| 1 | | | |
| 2 | | | |
| 3 | | | |
| 4 | 1 | EGA | European Genome-Phenome Archive |
| 5 | | | |
| 6 | | | |
| 7 | 2 | EM | Expectation-Maximization |
| 8 | | | |
| 9 | | | |
| 10 | 3 | EV | Expression Variability |
| 11 | | | |
| 12 | | | |
| 13 | 4 | GEO | Gene Expression Omnibus |
| 14 | | | |
| 15 | | | |
| 16 | 5 | GO | Gene Ontology |
| 17 | | | |
| 18 | | | |
| 19 | 6 | MAD | Median Absolute Deviation |
| 20 | | | |
| 21 | | | |
| 22 | 7 | MV | Methylation Variability |
| 23 | | | |
| 24 | | | |
| 25 | 8 | PMI | Post-Mortem Interval |
| 26 | | | |
| 27 | | | |
| 28 | | | |
| 29 | 9 | UKBEC | UK Human Brain Expression Database |
| 30 | | | |
| 31 | | | |
| 32 | 10 | | |
| 33 | | | |
| 34 | | | |

11 **Declarations**

12 **Ethics approval and consent to participate.** Not applicable

13 **Consent for publication.** Not applicable

14 **Availability of data and material.**

15 The datasets analyzed in this study are available in the GEO repository under the accession number

16 GSE36192 (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE36192>) and GSE36194

17 (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE36194>). The remaining dataset analyzed this

18 study is available from European Genome-phenome Archive but restrictions apply to the availability of

19 these data and are not publicly available. Data are however available from the authors upon reasonable

20 request and with permission of EGA. All code has been deposited at <https://github.com/nbashkeel/EV>.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 **Competing interests.** The authors declare that they have no competing interests.

2 **Funding.** This work was supported by funding from the Canadian Breast Cancer Foundation (JML) and
3 the Natural Sciences and Engineering Research Council (NSERC of Canada). The funding bodies had no
4 role in the design of the study, the collection, analysis, and interpretation of data nor in writing the
5 manuscript.

6
7 **Authors' contributions.** NB performed the computational data analysis, prepared figures, and wrote the
8 manuscript. JL conceived, coordinated, supervised the work and helped write the manuscript. TP and
9 MK assisted in the computational analysis and helped write the manuscript. All authors read and
10 approved the final manuscript.

11 **Acknowledgements.** We thank Stephane Aris-Brosou, Mathieu Lavalley, Martin Pelchat and Redaet
12 Daniel for helpful discussion.

13

1 References

1. Alemu EY, Carl JW, Corrada Bravo H, Hannenhalli S. Determinants of expression variability. *Nucleic Acids Res.* 2014;42:3503–14.
2. Roberfroid S, Vanderleyden J, Steenackers H. Gene expression variability in clonal populations: Causes and consequences. *Crit Rev Microbiol.* 2016;42:969–84.
3. Choi JK, Yu U, Yoo OJ, Kim S. Differential coexpression analysis using microarray data and its application to human cancer. *Bioinformatics.* 2005;21:4348–55.
4. Ho JWK, Stefani M, dos Remedios CG, Charleston MA. Differential variability analysis of gene expression and its application to human diseases. *Bioinformatics.* 2008;24:i390–8.
5. Ecker S, Pancaldi V, Rico D, Valencia A. Higher gene expression variability in the more aggressive subtype of chronic lymphocytic leukemia. *Genome Med.* 2015;7:8.
6. Chen E-H, Hou Q-L, Wei D-D, Jiang H-B, Wang J-J. Phenotypic plasticity, trade-offs and gene expression changes accompanying dietary restriction and switches in *Bactrocera dorsalis* (Hendel) (Diptera: Tephritidae). *Sci Rep.* 2017;7. doi:10.1038/s41598-017-02106-3.
7. Newman JRS, Ghaemmaghami S, Ihmels J, Breslow DK, Noble M, DeRisi JL, et al. Single-cell proteomic analysis of *S. cerevisiae* reveals the architecture of biological noise. *Nature.* 2006;441:840–6.
8. Singh GP. Coupling Between Noise and Plasticity in *E. coli*. *G3 Genes Genomes Genet.* 2013;3:2115–20.
9. Taniguchi Y, Choi PJ, Li G-W, Chen H, Babu M, Hearn J, et al. Quantifying *E. coli* Proteome and Transcriptome with Single-Molecule Sensitivity in Single Cells. *Science.* 2010;329:533–8.
10. Silander OK, Nikolic N, Zaslaver A, Bren A, Kikoin I, Alon U, et al. A Genome-Wide Analysis of Promoter-Mediated Phenotypic Noise in *Escherichia coli*. *PLoS Genet.* 2012;8. doi:10.1371/journal.pgen.1002443.
11. Wolf L, Silander OK, van Nimwegen E. Expression noise facilitates the evolution of gene regulation. *eLife.* 4. doi:10.7554/eLife.05856.
12. Barkai N, Shilo B-Z. Variability and Robustness in Biomolecular Systems. *Mol Cell.* 2007;28:755–60.
13. Lehner B. Selection to minimise noise in living systems and its implications for the evolution of gene expression. *Mol Syst Biol.* 2008;4:170.
14. Lehner B. Conflict between Noise and Plasticity in Yeast. *PLoS Genet.* 2010;6. doi:10.1371/journal.pgen.1001185.
15. Blake WJ, Balázsi G, Kohanski MA, Isaacs FJ, Murphy KF, Kuang Y, et al. Phenotypic Consequences of Promoter-Mediated Transcriptional Noise. *Mol Cell.* 2006;24:853–65.

- 1
- 2
- 3
- 4 1 16. Bishop AL, Rab FA, Sumner ER, Avery SV. Phenotypic heterogeneity can enhance rare-cell survival in
- 5 2 'stress-sensitive' yeast populations. *Mol Microbiol.* 2007;63:507–20.
- 6
- 7 3 17. Ackermann M, Stecher B, Freed NE, Songhet P, Hardt W-D, Doebeli M. Self-destructive cooperation
- 8 4 mediated by phenotypic noise. *Nature.* 2008;454:987–90.
- 9
- 10 5 18. Zhang Z, Qian W, Zhang J. Positive selection for elevated gene expression noise in yeast. *Mol Syst*
- 11 6 *Biol.* 2009;5:299.
- 12
- 13 7 19. Ward MC, Gilad Y. Human genomics: Cracking the regulatory code. *Nature.* 2017;550:190–1.
- 14
- 15 8 20. Li J, Liu Y, Kim T, Min R, Zhang Z. Gene Expression Variability within and between Human Populations
- 16 9 and Implications toward Disease Susceptibility. *PLOS Comput Biol.* 2010;6:e1000910.
- 17
- 18 10 21. Hough SR, Laslett AL, Grimmond SB, Kolle G, Pera MF. A Continuum of Cell States Spans Pluripotency
- 19 11 and Lineage Commitment in Human Embryonic Stem Cells. *PLOS ONE.* 2009;4:e7708.
- 20
- 21 12 22. Kalmar T, Lim C, Hayward P, Muñoz-Descalzo S, Nichols J, Garcia-Ojalvo J, et al. Regulated
- 22 13 Fluctuations in Nanog Expression Mediate Cell Fate Decisions in Embryonic Stem Cells. *PLOS Biol.*
- 23 14 2009;7:e1000149.
- 24
- 25 15 23. Pritchard CC, Hsu L, Delrow J, Nelson PS. Project normal: Defining normal variance in mouse gene
- 26 16 expression. *Proc Natl Acad Sci U S A.* 2001;98:13266–71.
- 27
- 28 17 24. Vedell PT, Svenson KL, Churchill GA. Stochastic variation of transcript abundance in C57BL/6J mice.
- 29 18 *BMC Genomics.* 2011;12:167.
- 30
- 31 19 25. Mar JC, Matigian NA, Mackay-Sim A, Mellick GD, Sue CM, Silburn PA, et al. Variance of Gene
- 32 20 Expression Identifies Altered Network Constraints in Neurological Disease. *PLOS Genet.*
- 33 21 2011;7:e1002207.
- 34
- 35 22 26. Carey LB, van Dijk D, Sloom PMA, Kaandorp JA, Segal E. Promoter Sequence Determines the
- 36 23 Relationship between Expression Level and Noise. *PLoS Biol.* 2013;11.
- 37 24 doi:10.1371/journal.pbio.1001528.
- 38
- 39 25 27. Batenchuk C, St-Pierre S, Tepliakova L, Adiga S, Szuto A, Kabbani N, et al. Chromosomal Position
- 40 26 Effects Are Linked to Sir2-Mediated Variation in Transcriptional Burst Size. *Biophys J.* 2011;100:L56–8.
- 41
- 42 27 28. Pruitt KD, Harrow J, Harte RA, Wallin C, Diekhans M, Maglott DR, et al. The consensus coding
- 43 28 sequence (CCDS) project: Identifying a common protein-coding gene set for the human and mouse
- 44 29 genomes. *Genome Res.* 2009;19:1316–23.
- 45
- 46 30 29. Blake JA, Eppig JT, Kadin JA, Richardson JE, Smith CL, Bult CJ. Mouse Genome Database (MGD)-2017:
- 47 31 community knowledge resource for the laboratory mouse. *Nucleic Acids Res.* 2017;45 Database
- 48 32 issue:D723–9.
- 49
- 50 33 30. R Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria: R
- 51 34 Foundation for Statistical Computing; 2019. <https://www.R-project.org/>.
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60
- 61
- 62
- 63
- 64
- 65

- 1
- 2
- 3
- 4 1 31. Cedar H. DNA methylation and gene activity. *Cell*. 1988;53:3–4.
- 5
- 6 2 32. Moore LD, Le T, Fan G. DNA Methylation and Its Basic Function. *Neuropsychopharmacology*.
- 7 3 2013;38:23–38.
- 8
- 9
- 10 4 33. Irvine RA, Lin IG, Hsieh C-L. DNA Methylation Has a Local Effect on Transcription and Histone
- 11 5 Acetylation. *Mol Cell Biol*. 2002;22:6689–96.
- 12
- 13 6 34. Wagner JR, Busche S, Ge B, Kwan T, Pastinen T, Blanchette M. The relationship between DNA
- 14 7 methylation, genetic and expression inter-individual variation in untransformed human fibroblasts.
- 15 8 *Genome Biol*. 2014;15:R37.
- 16
- 17
- 18 9 35. Zhang B, Zhou Y, Lin N, Lowdon RF, Hong C, Nagarajan RP, et al. Functional DNA methylation
- 19 10 differences between tissues, cell types, and across individuals discovered using the M&M algorithm.
- 20 11 *Genome Res*. 2013;23:1522–40.
- 21
- 22
- 23 12 36. Birdsill AC, Walker DG, Lue L, Sue LI, Beach TG. POSTMORTEM INTERVAL EFFECT ON RNA AND GENE
- 24 13 EXPRESSION IN HUMAN BRAIN TISSUE. *Cell Tissue Bank*. 2011;12:311–8.
- 25
- 26 14 37. Trabzuni D, Ramasamy A, Imran S, Walker R, Smith C, Weale ME, et al. Widespread sex differences in
- 27 15 gene expression and splicing in the adult human brain. *Nat Commun*. 2013;4:ncomms3771.
- 28
- 29
- 30 16 38. Harada CN, Natelson Love MC, Triebel K. Normal Cognitive Aging. *Clin Geriatr Med*. 2013;29:737–52.
- 31
- 32 17 39. Annunziato L, Pannaccione A, Cataldi M, Secondo A, Castaldo P, Di Renzo G, et al. Modulation of ion
- 33 18 channels by reactive oxygen and nitrogen species: a pathophysiological role in brain aging? *Neurobiol*
- 34 19 *Aging*. 2002;23:819–34.
- 35
- 36 20 40. Montecino-Rodriguez E, Berent-Maoz B, Dorshkind K. Causes, consequences, and reversal of
- 37 21 immune system aging. *J Clin Invest*. 2013;123:958–65.
- 38
- 39
- 40 22 41. Kærn M, Elston TC, Blake WJ, Collins JJ. Stochasticity in gene expression: from theories to
- 41 23 phenotypes. *Nat Rev Genet*. 2005;6:nrg1615.
- 42
- 43 24 42. Pelkmans L: Cell Biology. Using cell-to-cell variability--a new era in molecular biology. *Science* 2012,
- 44 25 336(6080):425-426.
- 45
- 46 26 43. Eldar A, Elowitz MB: Functional roles for noise in genetic circuits. *Nature* 2010, 467(7312):167-173.
- 47
- 48 27 44. Dueck H, Khaladkar M, Kim TK, Spaethling JM, Francis C, Suresh S, Fisher SA, Seale P, Beck SG, Bartfai
- 49 28 T *et al*: Deep sequencing reveals cell-type-specific patterns of single-cell transcriptome variation.
- 50 29 *Genome Biol* 2015, 16:122.
- 51
- 52
- 53 30 45. Osorio D, Yu X, Zhong Y, Li G, Yu P, Serpedin E, Huang J, Cai JJ: Extent, heritability, and functional
- 54 31 relevance of single cell expression variability in highly homogeneous populations of human cells. *bioRxiv*
- 55 32 2019:574426.
- 56
- 57 33 46. Georgi B, Voight BF, Bućan M. From Mouse to Human: Evolutionary Genomics Analysis of Human
- 58 34 Orthologs of Essential Genes. *PLoS Genet*. 2013;9. doi:10.1371/journal.pgen.1003484.
- 59
- 60
- 61
- 62
- 63
- 64
- 65

- 1
- 2
- 3
- 4 1 47. Streit WJ, Xue Q-S. The Brain's Aging Immune System. *Aging Dis.* 2010;1:254–61.
- 5
- 6 2 48. Lucin KM, Wyss-Coray T. Immune activation in brain aging and neurodegeneration: too much or too
- 7 3 little? *Neuron.* 2009;64:110–22.
- 8
- 9
- 10 4 49. Singh P, Goode T, Dean A, Awad SS, Darlington GJ. Elevated Interferon Gamma Signaling Contributes
- 11 5 to Impaired Regeneration in the Aged Liver. *J Gerontol A Biol Sci Med Sci.* 2011;66A:944–56.
- 12
- 13 6 50. Wu D, Meydani SN. Age-associated changes in immune and inflammatory responses: impact of
- 14 7 vitamin E intervention. *J Leukoc Biol.* 2008;84:900–14.
- 15
- 16
- 17 8 51. Azpurua J, Eaton BA. Neuronal epigenetics and the aging synapse. *Front Cell Neurosci.* 2015;9.
- 18 9 doi:10.3389/fncel.2015.00208.
- 19
- 20 10 52. Hebert LE, Beckett LA, Scherr PA, Evans DA. Annual Incidence of Alzheimer Disease in the United
- 21 11 States Projected to the Years 2000 Through 2050. *Alzheimer Dis Assoc Disord.* 2001;15:169–73.
- 22
- 23 12 53. Levy G, Schupf N, Tang M-X, Cote LJ, Louis ED, Mejia H, et al. Combined effect of age and severity on
- 24 13 the risk of dementia in Parkinson's disease. *Ann Neurol.* 2002;51:722–9.
- 25
- 26
- 27 14 54. Cerbai F, Lana D, Nosi D, Petkova-Kirova P, Zecchi S, Brothers HM, et al. The Neuron-Astrocyte-
- 28 15 Microglia Triad in Normal Brain Ageing and in a Model of Neuroinflammation in the Rat Hippocampus.
- 29 16 *PLOS ONE.* 2012;7:e45250.
- 30
- 31
- 32 17 55. Soreq L, Rose J, Soreq E, Hardy J, Trabzuni D, Cookson MR, et al. Major Shifts in Glial Regional
- 33 18 Identity Are a Transcriptional Hallmark of Human Brain Aging. *Cell Rep.* 2017;18:557–70.
- 34
- 35 19 56. Braegelmann K, Streeter K, Fields D, Baker T. Plasticity in respiratory motor neurons in response to
- 36 20 reduced synaptic inputs: a form of homeostatic plasticity in respiratory control? *Exp Neurol.* 2017;287 Pt
- 37 21 2:225–34.
- 38
- 39
- 40 22 57. Galvan V, Jin K. Neurogenesis in the aging brain. *Clin Interv Aging.* 2007;2:605–10.
- 41
- 42 23 58. Malberg JE, Eisch AJ, Nestler EJ, Duman RS. Chronic Antidepressant Treatment Increases
- 43 24 Neurogenesis in Adult Rat Hippocampus. *J Neurosci.* 2000;20:9104–10.
- 44
- 45 25 59. Dunning MJ, Smith ML, Ritchie ME, Tavaré S. beadarray: R classes and methods for Illumina bead-
- 46 26 based data. *Bioinformatics.* 2007;23:2183–4.
- 47
- 48
- 49 27 60. Hernandez DG, Nalls MA, Moore M, Chong S, Dillman A, Trabzuni D, et al. Integration of GWAS SNPs
- 50 28 and tissue specific expression profiling reveal discrete eQTLs for human traits in blood and brain.
- 51 29 *Neurobiol Dis.* 2012;47:20–8.
- 52
- 53 30 61. Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, et al. NCBI GEO: archive for
- 54 31 functional genomics data sets—update. *Nucleic Acids Res.* 2013;41:D991–5.
- 55
- 56
- 57 32 62. Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, et al. limma powers differential expression
- 58 33 analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.* 2015;43:e47.
- 59
- 60
- 61
- 62
- 63
- 64
- 65

- 1
- 2
- 3
- 4 1 63. Suzuki R, Shimodaira H. Pvcust: an R package for assessing the uncertainty in hierarchical clustering.
- 5 2 Bioinformatics. 2006;22:1540–2.
- 6
- 7
- 8 3 64. Benaglia T, Chauveau D, Hunter DR, Young DS. mixtools: An R Package for Analyzing Mixture Models.
- 9 4 J Stat Softw. 2009;32. <https://doaj.org>.
- 10
- 11 5 65. Karolchik D, Hinrichs AS, Furey TS, Roskin KM, Sugnet CW, Haussler D, et al. The UCSC Table Browser
- 12 6 data retrieval tool. Nucleic Acids Res. 2004;32 suppl_1:D493–6.
- 13
- 14 7 66. Supek F, Bošnjak M, Škunca N, Šmuc T. REVIGO Summarizes and Visualizes Long Lists of Gene
- 15 8 Ontology Terms. PLOS ONE. 2011;6:e21800.
- 16
- 17
- 18 9 67. Warnes GR, Bolker B, Bonebakker L, Gentleman R, Liaw WHA, Lumley T, et al. gplots: Various R
- 19 10 Programming Tools for Plotting Data. 2015. [https://www.scienceopen.com/document?vid=0e5d8e31-](https://www.scienceopen.com/document?vid=0e5d8e31-1fe4-492f-a3d8-8cd71b2b8ad9)
- 20 11 [1fe4-492f-a3d8-8cd71b2b8ad9](https://www.scienceopen.com/document?vid=0e5d8e31-1fe4-492f-a3d8-8cd71b2b8ad9). Accessed 11 Feb 2019.
- 21
- 22

23 12

24

25

26 13 **Additional Files**

27

28

29

30 14 Additional File 1: Structural analysis of genes as a function of EV

31

32

33 15 Additional File 2: EV Correlation between different tissue types

34

35

36 16 Additional File 3: Complete list of GO term treemaps for all genes

37

38

39 17 Additional File 4: Chi-Squared enrichment analysis methodology

40

41

42 18 Additional File 5: Complete list of GO term treemaps for essential genes

43

44

45 19 Additional File 6: Complete list of GO term treemaps for age-regulated Hyper-Variable genes

46

47

48 20 Additional File 7: Preprocessing of Brain Samples

49

50

51 21 Additional File 8: Methylation Variability

52

53

54

55 22

56

57

58 23

59

60

61

62

63

64

65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 **Figure Legends**

2 **Figure 1. Expression variability (EV) in human breast, cerebellum, and frontal cortex tissue. (A)**

3 Expected expression MAD for curve as a function of median probe expression (solid black line). (B)

4 Kernel density estimation function of EV. The vertical black lines represent the EV classification ranges.

5 (C) Expression variability as a function of median gene expression. Adjusted R^2 values for the linear

6 regression model shown in red were 0.0002, 0.0008, and 0.005 and the associated Kendall rank

7 correlation coefficients were -0.208, -0.201, -0.213 for breast, cerebellum, and frontal cortex tissues

8 respectively.

10 **Figure 2. Bimodal Hyper-Variable gene expression detection.** Gaussian mixture modelling method of

11 detecting bimodal probes. The dashed lines represent the overall gene kernel density estimation

12 function of gene expression. The two Gaussian models are shown in dark grey and light grey, and the

13 dotted vertical lines represent the distribution means.

15 **Figure 3. Cross-Validation of EV Classifications.** (A) Relative frequency of EV classification accuracy

16 between original distribution and 50-50 split retest replicates (n=100). (B) Number of probes in each EV

17 probe set before and after split-retest protocol.

19 **Figure 4. Tissue Specificity of EV.** (A) Venn diagrams comparing EV classifications of probe mapped

20 genes sets between breast, cerebellum, and frontal cortex tissues. (B) Effect of genomic position on EV.

21 Each chromosome is divided into 100 bins (x-axis) based on the maximum gene coordinate annotation,

22 and the average EV in each bin is measured (y-axis). Bins with an average EV greater than 0 are

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 represented in green, while those with a negative EV are represented in red. Bins with less than three
2 probes were assigned an average EV of zero.

3
4 **Figure 5. Methylation in human cerebellum and frontal cortex tissue.** (A) Kernel density estimation
5 function of average gene methylation. Gaussian mixture models were used to classify the genes into
6 Non-, Medium- and Highly- methylated clusters. (B) Kernel density estimation function of average gene
7 methylation by EV classification. The dashed vertical lines represent the methylation state cluster cut-
8 offs generated by the Gaussian mixture modelling.

9
10 **Figure 6. Hierarchical clustering of Hyper-Variable genes by age in (A) cerebellum tissue, and (B)**
11 **frontal cortex tissue.** The vertical axis represents the age-regulated Hyper-Variable genes while the
12 samples were clustered by age and plotted on the horizontal axis. The top heatmaps represent the
13 positively correlated age-regulated genes while the bottom heatmaps represent the negatively
14 correlated age-regulated genes. The age clusters decrease in age from left to right in both heatmaps and
15 correspond to the following age ranges: $\bar{x}_1 = 79$ [58,98], $\bar{x}_2 = 45$ [32,57], and $\bar{x}_3 = 17$ [1,31].

16
17 **Figure 7. Expression and methylation correlation.** Histogram of Pearson correlation coefficient between
18 paired gene expression and gene methylation levels in the Hyper-Variable and Hypo-Variable probe sets.













