

[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. For example, the statistical nature of human expression variation has not been reported and
7 the role of gene methylation is just beginning to be understood. Moreover, the contribution that age, sex
8 and tissue-specific factors have on expression variability are not well understood.

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

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

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 **Keywords:** Expression variability, Tissue Specificity, Essentiality, Methylation, Aging

2

3 **Background**

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

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

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 Expression variability, whatever its source, is an evolvable trait subject to natural selection, whereby each
2 genes have an optimal expression level and variance required for an organism’s fitness and selection
3 minimizes this variability [7, 10–14]. In this case, genes with low variability have been subjected to heavy
4 selection pressure to minimize population expression variance. Conversely, high variability genes have
5 been least selected for. Genes with high expression variability could be drivers of phenotypic diversity, as
6 suggested by position association between expression noise and growth [15–18]. In this interpretation,
7 genes with high variability allow for growth in fluctuating environments. Understanding the role of the
8 gene expression variability patterns across human populations will therefore provide crucial insights into
9 how genetic differences contribute to phenotypic diversity, susceptibility to disease [19, 20],
10 differentiation of disease subtypes [5], embryonic development [21, 22], and alterations in gene network
11 architecture [23].

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

20 Results

21 **Estimating expression variability.** We measured human gene expression variability (EV) [1] in post-
22 mortem non-diseased cerebellum (n=465) and frontal cortex samples (n=455) and biopsied normal breast

1
2
3
4 1 tissues (n=144). We excluded probes corresponding to non-coding transcripts as well as those with
5
6 2 missing probe coordinates, resulting in a list of 42,084 probes. To estimate a gene's EV independent of its
7
8 3 expression magnitude, we modified the method initially described by Alemu et al [1]. First, we calculated
9
10 4 the median absolute deviation (MAD) for each probe. Then we modelled the expected MAD for all genes
11
12 5 as a function of median expression using a locally weighted polynomial regression (Fig. 1A, red line). The
13
14 6 expected MAD regression curves for each tissue type exhibit a flat, negative parabolic shape where the
15
16 7 lowest and highest expression probes represent the troughs of the curve. Variability in gene expression
17
18 8 levels has previously been shown to decrease as expression approaches either extrema [7, 9, 24]. The EV
19
20 9 for each probe was calculated as the difference between its bootstrapped MAD and the expected MAD at
21
22 10 each median expression level (Fig. 1A). Positive EV values indicate that the gene has a greater expression
23
24 11 variability than genes with the same expression magnitude mean. Conversely, negative EV values imply
25
26 12 reduced population expression variability. We next plotted the probability density function of EV for each
27
28 13 tissue (Fig. 1B). The EV distributions in all three tissue types exhibit large peaks around the zero mean and
29
30 14 a long tail for positive EV probes. Breast tissue exhibited a larger shoulder of the negative EV probes
31
32 15 compared to cerebellum and frontal cortex tissues. This is likely attributable to the lower number of breast
33
34 16 samples (144 compared to 456 and 455 samples respectively).

35
36
37
38
39
40
41
42
43 17 We then confirmed the independence of EV on expression by modelling the relationship between
44
45 18 the two variables using a linear regression (Fig 1C). Based on the poor adjusted R^2 values (2×10^{-4} , 8×10^{-4} ,
46
47 19 and 5×10^{-3} for breast, cerebellum, and frontal cortex respectively) and the flat slopes, we conclude that
48
49 20 there is no substantial correlation between EV and gene expression magnitude.

50
51
52
53 21 Next, we then classified each probe into three categories based on their EV. We used the term
54
55 22 "Hyper-Variable" to describe probes whose EV was greater than $\tilde{x}_{EV} + 3 * MAD_{EV}$. Probes with an EV less
56
57 23 than $\tilde{x}_{EV} - 3 * MAD_{EV}$ were deemed "Hypo-Variable". The remaining genes that fell within the range of
58
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 $\tilde{x}_{EV} \pm 3 * MAD_{EV}$ were considered “Non-Variable”. We propose that these three distinct gene groups,
2 categorized based on EV, correspond to distinct functional and phenotypic gene characteristics.

3
4 **Figure 1. Expression variability (EV) in human breast, cerebellum, and frontal cortex tissue.** (A) Expected expression MAD for
5 curve as a function of median probe expression (solid black line). (B) Probability density function of EV. The vertical black lines
6 represent the EV classification ranges. (C) Expression variability as a function of median gene expression. Adjusted R^2 values for
7 the linear regression model shown in red were 0.0002, 0.0008, and 0.005 for breast, cerebellum, and frontal cortex tissues
8 respectively.

9
10 **Accounting for sampling error in EV classification.** We were concerned that the classification of a gene
11 into Hyper-, Hypo- and Non-Variable classes might be the result of sampling errors. To minimize this
12 possibility and to increase the accuracy of our EV classification method, we divided each of our tissue
13 samples into two equally sized probe subsets and repeated the EV analysis. This 50-50 split-retest
14 procedure was repeated 100 times. Figure 2 shows the probability distribution of a concordant EV
15 classification for each gene into Hyper-, Hypo- and Non-Variable class across the three subsets in each
16 tissue type. Fig 2 demonstrates that classification of a gene as Hyper or Hypo-variable based on a single
17 analysis of the population is problematic because of sampling bias. We see a substantial decrease in the
18 number of genes in the Hyper- and Hypo-Variable gene sets after conducting our split-retest protocol(Fig.
19 2B). Thus, our split-retest method likely increases the robustness and accuracy of EV classification

20
21 **Figure 2. Cross-Validation of EV Classifications.** (A) Probability distribution of gene EV classification accuracy between original
22 distribution and 50-50 split retest replicates ($n=100$). (B) Number of genes in each EV gene set before and after split-retest protocol.

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 **Statistical nature of Hyper-variability.** A previously unexplored aspect of expression Hyper-variability is
2 the statistical characteristics of expression amongst genes with this wide range of gene expression.
3 Specifically, high EV could be the result of a multimodal distribution of gene expression with two or more
4 distinct expression means or might simply result from a broadening of expression values around a
5 unimodal mean value. In order to distinguish between the two possibilities, we modeled each gene
6 expression as a mixture of two Gaussian distributions (Fig. 3). Next, we identified the peaks of the
7 probability density function for each Gaussian distribution and compared the distance between the peaks
8 as well as the ratio of peak heights. Probes with peaks that were greater than one median absolute
9 deviation apart and displayed a peak ratio greater than 0.1 were classified as having a bimodal expression
10 distribution. Probes that did not satisfy both criteria were considered to have a unimodal distribution.
11 Only a small minority of the Hyper-Variable (high EV) probes (15/3453 breast tissue probes, 6/2980
12 cerebellum probes, and 6/3487 frontal cortex probes) showed a bimodal distribution of gene expression.
13 The remaining majority of Hyper-Variable probes had a unimodal distribution. This indicates that high
14 expression variability is a result of a widening of possible expression values across a single mean rather
15 than the gene expression existing in two or more discrete states.

16
17 *Figure 3. Bimodal Hyper-Variable gene expression detection. Gaussian mixture modelling method of detecting bimodal probes.*
18 *The dashed lines represent the overall gene probability density function of gene expression. The two Gaussian models are shown*
19 *in dark grey and light grey, and the dotted vertical lines represent the distribution means.*

20
21 **Tissue-specificity of EV.** Because we have calculated EV from different tissues, we were able to determine
22 the extent to which tissue-specific factors might contribute to EV. This is an important question because
23 expression variability exists not only between individuals but between different tissues in the same

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 organism. Characterizing tissue-specific EV will therefore shed light on tissue specification and identity
2 processes. As shown in Fig. 4A, there was limited overlap between the Hyper-Variable, Hypo-Variable, and
3 Non-Variable genes between the three different tissues. 13-16% of the Hyper-Variable genes were
4 classified as such in the three tissues and 18-26% of the Hypo-Variable were so classified. The Non-
5 Variable probe sets contained over 82% of genes in each tissue type, with over 71% of the measured genes
6 commonly classified as NV in all three tissue types. The poorly overlapping nature of each EV classification
7 suggests that population expression variability is largely determined by tissue-specific pathways.

8
9 **Figure 4. Tissue Specificity of EV.** (A) Venn diagrams comparing EV classifications between breast, cerebellum, and frontal cortex
10 tissues. (B) Effect of genomic position on EV. Each chromosome is divided into 100 bins (x-axis) based on the maximum gene
11 coordinate annotation, and the average EV in each bin is measured (y-axis).

12
13 **EV and gene structural characteristics.** To understand possible genomic mechanisms by which population
14 expression variability occurs, we first explored the relationship between EV and various structural features
15 of the genes. Expression variability has previously been reported to be associated with gene size, gene
16 structure, and surrounding regulatory elements [1]. However, we found no significant linear correlation
17 between EV and a gene's exon count, sequence length, transcript size, or number of isoforms (Additional
18 file 1). While certain linear models exhibited statistical significance ($p < 0.05$), the fit of the model and
19 subsequent comparison of the linear model against a local polynomial regression curve showed that the
20 correlation was either not correctly defined by a linear model or simply too small to draw a conclusion.

21 While we did not find that the physical gene characteristics were correlated to EV, previous
22 studies have shown that the position of a gene on a chromosome has considerable effects on stochastic
23 gene expression variability, independent of gene- and promoter-specific variables [25]. We next tested if

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 there is a relationship between expression variability and chromosomal position (Fig. 4B). To this end,
2 each chromosome was divided into 100 bins and the average EV all the genes within each bin was
3 determined. If there was no relationship between the two, we would expect to see EV to be uniformly
4 distributed throughout the genome. We found that EV is not uniformly distributed across the genome,
5 and individual regions of chromosomes exhibited peaks of high expression variability or troughs of low
6 expression variability. To further confirm our conclusion, we tested the cosine similarities of the
7 chromosomes within and across the tissue types (Additional file 2). The EV distributions across
8 chromosomes exhibited low similarities within each tissue type, further establishing that EV is not
9 randomly distributed throughout the genome.

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

15 *Table 1. Pearson's Chi-squared test for Essentiality in Hyper-Variable, Hypo-Variable, and Non-Variable probe sets.*

Tissue	Probe Set	Total Gene Count	Essential Gene Counts	Standardized Residuals	P-Value
Breast	Hyper	1448	180	12.22	1.50×10^{-34}
	Hypo	957	108	8.27	
	NV	33957	1657	-14.94	
Cerebellum	Hyper	1640	170	8.66	8.27×10^{-63}
	Hypo	837	83	5.54	
	NV	35257	1849	-10.42	
Frontal Cortex	Hyper	1760	196	10.6	1.52×10^{-92}
	Hypo	1254	125	7.04	
	NV	34831	1764	-12.89	

16

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 In order to examine a potential correlation between expression variability and essentiality in
2 human tissues, we first tested the independence between EV classification and annotation of human
3 essentiality (Table 2). Essentiality annotations were obtained from the CCDS [26] and MGD [27] databases.
4 Here, direct human orthologs of genes essential for prenatal, perinatal, or postnatal survival of mice were
5 classified as essential. Using the Pearson's chi-square test for the number of essential genes in each probe
6 set, we find that that the Hypo-Variable gene set in breast, cerebellum, and frontal cortex tissues were
7 significantly enriched for genes with essentiality annotation (p-value = 1.50×10^{-34} , 8.27×10^{-63} , and 1.52
8 $\times 10^{-92}$, respectively). Thus, expression variability for many essential genes is constrained in humans, likely
9 reflecting a similar biology to essential yeast genes. However, we also observe a significant enrichment of
10 essential genes within the Hyper-Variable gene sets. This was a surprise to us since essential genes are
11 thought of as being dose-sensitive and changes in the level of gene expression would be predicted to be
12 deleterious or lethal.

13 To better understand the implications of high variability in essential genes, we examined the
14 functional annotations associated with Hyper-Variable essential genes (Table 3 and Additional file 4). The
15 breast essential Hyper-Variable gene set was enriched for embryonic development, responses to growth
16 factors, cell-substrate junction assembly, regulation of epithelial cell proliferation, and positive regulation
17 of cellular component movement. The cerebellum essential Hyper-Variable gene set was enriched for cell
18 differentiation, anion transport, trans-synaptic signaling, response to growth factors, and cell projection
19 organization. Lastly, the frontal cortex essential Hyper-Variable gene set was enriched for cell
20 differentiation, secretion, tyrosine kinase signaling, cell projection organization, and heart contractions.
21 Overall, the Hyper-Variable essential gene sets tended to be enriched for morphogenic, tissue, and organ
22 system development. This suggests that tight regulation of some essential genes may only be required for
23 embryonic and morphogenic development and dose-sensitivity is lost in adults, allowing for high
24 expression variability.

1
2
3
4 **1** *Table 2. Top 5 common and unique REVIGO GO annotation subsets of Hyper-Variable and Hypo-Variable essential genes in breast,*
5 **2** *cerebellum, and frontal cortex tissues.*

	Unique Breast Annotations	Unique Cerebellum Annotations	Unique Frontal Cortex Annotations
Hyper-Variable Essential Genes	Embryo development ending in birth or egg hatching	Positive regulation of cell differentiation	Positive regulation of cell differentiation
	Cellular response to growth factor stimulus	Anion transport	Regulation of secretion
	Cell-substrate junction assembly	Trans-synaptic signalling	Transmembrane receptor protein tyrosine kinase signaling pathway
	Regulation of epithelial cell proliferation	Cellular response to growth factor stimulus	Regulation of cell projection organization
	Positive regulation of cellular component movement	Regulation of cell projection organization	Heart contraction
Hypo-Variable Essential Genes	Negative regulation of cellular component organization	DNA repair	DNA repair
	DNA repair	Negative regulation of cellular component organization	Covalent chromatic modification
	Regulation of cellular protein localization	Positive regulation of viral release from host cell	mRNA transport
	Embryo development ending in birth or egg hatching	Regulation of cellular protein localization	Progesterone receptor signaling pathway
	Apoptotic process	Regulation of cell cycle process	Regulation of cell cycle process

3
4
5 **Functional analysis of Hyper-, Hypo- and Non-Variable genes.** In order to understand the overall
6 biological significance of gene EV, we examined the functional aspects that are enriched in the Hyper-
7 Variable, Hypo-Variable, and Non-Variable probe sets by conducting a gene set enrichment analysis for
8 each probe set group. We determined the over-represented Gene Ontology (GO) terms that were unique
9 in each tissue type, as well as GO terms that were common in all three tissue types. The resulting GO
10 annotations were simplified and visualized using a REVIGO treemap. The top five terms for each tissue
11 type can be found in Table 1, while the complete list of GO term treemaps can be found in Additional file
12 3.

1
2
3
4 **1** *Table 3. Top 5 common and unique REVIGO GO annotations in the Hyper-Variable and Hypo-Variable gene sets of breast,*
5
6 **2** *cerebellum, and frontal cortex tissues.*

	Common Annotations	Unique Breast Annotations	Unique Cerebellum Annotations	Unique Frontal Cortex Annotations
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	

37 **3**
38
39 **4** The breast Hyper-Variable gene set was uniquely enriched for epithelial cell differentiation,
40
41 **5** primary alcohol metabolism, and positive regulation of cellular component movement. The cerebellum
42
43 **6** Hyper-Variable gene set was uniquely enriched for regulation of nervous system development,
44
45 **7** transmembrane transport, and neuron death. The frontal cortex Hyper-Variable gene set was enriched
46
47 **8** for histamine secretion, regulation of cell morphogenesis, and trans-synaptic signalling. The breast,
48
49 **9** cerebellum, and frontal cortex Hyper-Variable gene sets were commonly enriched for regulation of tissue
50
51 **10** remodeling, inflammatory responses, and responses to inorganic substances. Of note, many of the
52
53 **11** enriched GO annotations of the Hyper-Variable genes are involved in signalling pathways. These pathways
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 1 require dynamic control based on internal and external stimuli and their high EV likely represents differing
5
6 2 environmental or hormonal conditions amongst the individuals.
7
8

9
10 3 In the case of the Hypo-Variable gene sets, all three tissue types were enriched for protein
11
12 4 catabolism and metabolism, ribonucleoprotein complexes, and negative regulation of autophagy. The
13
14 5 breast Hypo-Variable gene set was enriched for Golgi vesicle transport, nucleoside metabolism, and
15
16 6 protein catabolism. The cerebellum Hypo-Variable gene set was enriched for DNA conformation change,
17
18 7 modification-dependent macromolecule catabolism, and retrograde transport.
19
20

21
22 8 Genes with higher expression variability have previously been shown to be functionally and
23
24 9 physically involved with the physical cell periphery, localizing in the membrane, transmembrane, or
25
26 10 extracellular matrix regions [23]. Our results corroborated these findings as the Hyper-Variable gene sets
27
28 11 from breast, cerebellum, and frontal cortex were enriched for GO annotations associated with cell surface
29
30 12 signalling pathways, as well as cellular component ontologies enriched at the plasma membrane. In
31
32 13 contrast, genes with low variability genes tended to regulate nucleic acid and metabolic pathways,
33
34 14 localizing in the cell interior. These Hypo-Variable genes are likely involved in complex, dose-sensitive
35
36 15 gene networks and require tight regulation of their expression to function correctly.
37
38
39
40

41
42 16 **DNA methylation and expression variability.** One factor that has been postulated to regulate EV is DNA
43
44 17 methylation. While the relationship between methylation and gene expression is complex, low promoter
45
46 18 methylation is associated with high levels of gene expression [28–31]. Like gene expression, DNA
47
48 19 methylation is highly variable at the cell, tissue, and individual level [32], suggesting that EV could result
49
50 20 from variations in gene methylation. To explore this idea, we used DNA methylation annotations that
51
52 21 were available in 724 out of 911 brain tissue samples.
53
54
55

56
57 22 DNA methylation in CpG sites is thought to be bimodal, meaning that the gene is either
58
59 23 hypomethylated or hypermethylated [31]. In order to differentiate between low, medium, and high
60
61
62
63
64
65

1
2
3
4 1 methylation states in our samples, we modelled gene methylation using Gaussian mixture models for the
5
6 2 mean methylation for each gene. The distribution of gene methylation in both cerebellum and frontal
7
8 3 cortex tissue was best modelled as a three-component system. The first component was a sub-population
9
10 4 Gaussian mixture while the second and third components were modelled as single Gaussian distributions.
11
12 5 Genes whose methylation fell within the first component were classified as Non-Methylated genes. Genes
13
14 6 were classified as Medium Methylated for those in the second component and Highly Methylated if they
15
16 7 were in third. The distribution of methylation amongst the genes is predominantly bimodal with only a
17
18 8 minority of genes being Medium Methylated (Fig. 4A). In contrast, over 62% of cerebellum genes are non-
19
20 9 methylated and 23% highly methylated. Similarly, 58% of frontal cortex genes are non-methylated and
21
22 10 22% are highly methylated).

23
24
25
26
27
28
29 11 Next, we explored the correlation between methylation and expression based on the EV. When
30
31 12 we subset the methylation distribution by EV classification (Fig. 4B), we observe that Hypo-Variable genes
32
33 13 have a visibly different methylation pattern than Hyper- or Non-Variable genes insofar as Hypo-Variable
34
35 14 genes are visibly overrepresented in the Non-Methylated gene group compared to both the Hyper-
36
37 15 Variable and Non-Variable genes.

38
39
40
41 16
42
43
44 17 **Figure 5. Methylation in human cerebellum and frontal cortex tissue.** (A) Probability density function of average gene
45
46 18 methylation. Gaussian mixture models were used to classify the genes into Non-, Medium- and Highly- methylated clusters. (B)
47
48 19 Probability density function of average gene methylation by EV classification. The dashed vertical lines represent the methylation
49
50 20 state cluster cut-offs generated by the Gaussian mixture modelling. The y-axis is scaled by the square root of the methylation
51
52 21 density.

53
54
55 22

56
57
58
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 To further quantify the overrepresentation of Hypo-Variable genes in the Non-Methylated gene
2 group, we conducted a chi-squared test of independence between the methylation state clusters and the
3 EV classifications (Table 4). Both the cerebellum and frontal cortex tissues exhibited a significant
4 relationship between the methylations clusters and EV classifications ($p = 7.57 \times 10^{-36}$ and $p = 1.58 \times 10^{-59}$,
5 respectively). By examining the standardized residuals of the chi-square test of independence, we
6 quantitatively confirmed the enrichment of Non-Methylated genes within the Hypo-Variable gene set. We
7 also observe a significant enrichment of Highly Methylated genes in the Non-Variable gene set as well as
8 an enrichment of Medium Methylated genes in the Hyper-Variable gene set. The high significance and
9 non-overlapping enrichments across each of the three groups suggests that there is a strong relationship
10 between methylation and EV classifications: low gene methylation is important for the tight expression
11 constraint in Hypo-Variable genes while high gene methylation contributes to Non-variable expression.

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

	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

15
16
17 **Effects of age, sex, and PMI on variability.** To further understand the biological relevance of EV, we
18 focused on the Hyper-Variable genes to identify potential mechanisms of decreased constraint on gene
19 expression across the samples. We systematically analyzed EV as a function of sex, age, and post-mortem
20 interval (PMI). The breast tissue lacked these clinical annotations and were excluded from this analysis.
21 We employed a probe-wise linear regression analysis to model the relationship between Hyper-Variable

1 gene expression and age, sex, and PMI (Table 5). The resulting p-values were adjusted for multiple
2 comparisons using the Benjamini-Hochberg procedure and considered significant when the adjusted p-
3 value was less than 0.01.

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

	PMI			Sex			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

7
8 PMI might be a source of apparent expression variability because an extended PMI might
9 compromise sample RNA integrity and lead to degradation of labile RNA [33]. Brain samples had PMI times
10 ranging from 1 hour to 94 hours (mean = 36.14 hr), but we observe a negligible number of probes that are
11 correlated with PMI. This suggests that sample integrity is unlikely to be a source of EV changes. Somewhat
12 more surprisingly, however, is the low number of probes that are correlated with sex. Only 2 out of 1640
13 Hyper-Variable cerebellum genes and 16 out of 1760 Hyper-Variable frontal cortex genes show sex-
14 dependent differences in EV. While other studies have shown widespread sex differences in post-mortem
15 adult brain gene expression [34], we conclude that in our analysis, EV is overwhelmingly sex-independent.

16 However, we observe that age has a substantial effect on gene expression variability. Age is
17 correlated with over 31% of Hyper-Variable cerebellum probes and over 41% of Hyper-Variable frontal
18 cortex probes. This means that the expression of these genes becomes either more or less constrained
19 during aging. In the cerebellum, there were 247 Hyper-Variable genes whose expression increased as a
20 function of age and 267 genes with decreased expression. Similarly, the frontal cortex contained 373
21 genes with increased expression and 354 genes with reduced expression. Given that age is correlated with

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 a considerable number of Hyper-Variable genes, we classified the age of the samples in the cerebellum
2 and frontal cortex tissues into three age clusters according to BIC for expectation-maximization (EM)
3 initialized by hierarchical clustering for parameterized Gaussian mixture models. The oldest cluster
4 contained samples whose ages were between 58 and 98 ($\bar{x}_1 = 79$). The second cluster ranged between
5 32 and 57 years ($\bar{x}_2 = 45$), while the youngest age cluster contained samples aged 1 through 31 ($\bar{x}_3 =$
6 17).

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

14 While the cerebellum is generally considered a regulator of motor processes, it is also implicated
15 in cognitive and non-motor functions [35]. Many of these age-dependent upregulated Hyper-Variable
16 genes corroborate previous studies exploring the relationship between brain aging and changes in gene
17 expression, including cellular responses to chemical stimuli (gold cluster). In particular, reactive oxygen
18 and nitrogen species have been shown to change ion transport channel activity, and serve as an important
19 mechanism in brain aging [36]. While all the genes selected were age-regulated, some genes exhibit
20 outlier samples whose expression remains high across all genes in the dark orange cluster, regardless of
21 age. These genes are more likely to be overexpressed in the samples as age increases and are enriched
22 for peripheral nervous system neuron development and neuron apoptotic pathways. Similar enrichments
23 of neurogenic and chemical stimuli response pathways are seen in the upregulated frontal cortex genes
24 (gold cluster). The dark orange cluster in the upregulated frontal cortex age-dependent genes exhibits an

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 entirely sample-specific over- or under-expression of genes. These bimodally expressed genes are
2 enriched for glial cell differentiation, adenosine receptor signaling pathways, and antigen processing.
3 Lastly, we see a random scattering of expression in the yellow cluster of the frontal cortex heatmap that
4 steadily increases with age. These genes are enriched for glial cell differentiation, cellular response to
5 alcohol, and defense responses to fungus.

6

7 **Figure 6. Hierarchical clustering of Hyper-Variable genes by age in (A) cerebellum tissue, and (B) frontal cortex tissue.** The
8 vertical axis represents the age-regulated Hyper-Variable genes while the samples were clustered by age and plotted on the
9 horizontal axis. The top heatmaps represent the positively correlated age-regulated genes while the bottom heatmaps represent
10 the negatively correlated age-regulated genes. The age clusters decrease in age from left to right in both heatmaps and
11 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].

12

13 Most of the downregulated age-dependent Hyper-variable genes in the cerebellum fall into the
14 green cluster where expression of the genes in the cluster increases with age. These genes are involved
15 in leukocyte-mediated immunity and defense responses to other organisms, which is supported by
16 previous studies [37]. Interestingly, the yellow cluster exhibits U-shaped expression levels, whereby the
17 lowest expression is seen in the middle age cluster. These genes are enriched for optic nerve
18 development, response to interferon-gamma, and synaptic signalling. In the frontal cortex, the majority
19 of downregulated age-dependent genes fall in the red cluster, and are enriched for ion transport, cell
20 morphogenesis, and trans-synaptic signalling. Overall, the functional annotations of the age-regulated

1
2
3
4 1 Hyper-Variable gene clusters suggest that population EV is one outcome of age-dependent gene
5
6 2 expression changes.

7
8
9 3 We next investigated a possible impact of methylation changes on gene expression in the Up- and
10
11 4 Down-regulated Hyper-Variable genes. Fig. 7 shows the histogram distribution of correlation between the
12
13 5 sample-specific gene expression and gene methylation. We observe no strong correlation between
14
15 6 expression and methylation, suggesting age-dependent changes in expression of the age-regulated Hyper-
16
17 7 Variable genes are not the result of methylation changes.

18
19
20
21
22 8
23
24
25 9 *Figure 7. Expression and methylation correlation plot. Histogram of Pearson correlation coefficient between paired gene*
26
27 10 *expression and gene methylation levels in Hyper-Variable and Hypo-Variable genes.*
28
29

30 11

31 32 33 12 **Discussion**

34
35
36
37 13 Gene expression variability in a population is the cumulative result of intrinsic genetic factors,
38
39 14 extrinsic environmental factors, and stochastic noise. A fundamental issue in biology is understanding the
40
41 15 cause of expression variability within an individual organism and between isogenic and genetically
42
43 16 dissimilar individuals of a population. In this report, we study population gene expression variability in
44
45 17 human breast, cerebellum, and frontal cortex tissues.

46
47
48
49 18 Our investigation into human gene expression variability yielded several main findings. First, we
50
51 19 find that Hyper-Variability in gene expression is fundamentally unimodal and does not represent
52
53 20 population switching between two or more discrete expression stages. In addition, both Hypo-Variable
54
55 21 (highly constrained expression) and Hyper-Variable (lowly constrained expression) gene sets are enriched
56
57 22 for essential genes and that both Hyper- and Hypo-variability are largely regulated by tissue specific
58
59
60
61
62
63
64
65

1
2
3
4 1 factors. We also find that gene methylation has an important role in expression variability. Lastly, we find
5
6 2 that Hypervariability is primarily associated with age and not sex.
7
8

9
10 3 Our observation that Hyper-Variable genes have a unimodal expression pattern is significant
11
12 4 because population gene expression diversity could manifest itself in multimodal or unimodal
13
14 5 distributions. Unimodal distributions represent a continuum of expression levels in a population while
15
16 6 multimodal distributions represent two or more distinct expression states between which cells might
17
18 7 switch between [39]. Since Hyper-Variability is almost exclusively unimodal, high EV in both brain and
19
20 8 breast is the result of a wide range of permissible expression levels in a population rather than the
21
22 9 existence of discrete expression subpopulations. In the case of organismal phenotypes, the benefits of a
23
24 10 unimodal phenotype distribution is highest in environments with very rapid or very noisy changes [40].
25
26 11 The unimodal distribution of Hyper-Variability indicates that similar evolutionary constraints apply to gene
27
28 12 expression.
29
30
31
32
33

34 13 The enrichment of essential genes in the Hypo-Variable gene sets is in agreement with previous
35
36 14 findings in yeast showing that essential yeast genes are likely to have low expression variability. However,
37
38 15 we detected a significant number of essential genes amongst the Hyper-Variable genes in breast,
39
40 16 cerebellum, and frontal cortex tissue. Inactivation of these essential genes leads to pre- or neonatal
41
42 17 fatality in mice and humans [41] and functional enrichment analysis indicates that the essential genes are
43
44 18 indeed involved in developmental pathways. This was a surprise to us since we expected that expression
45
46 19 of developmental genes should be tightly regulated, yet we observe highly variable expression being
47
48 20 tolerated in obligate developmental pathways. It is possible that these “essential” genes are required for
49
50 21 embryonic development but have different post-embryonic roles and may not be essential post-natally.
51
52 22 Alternatively, it is possible that these essential genes are not dose-sensitive in humans, meaning that only
53
54 23 a certain level of baseline expression is required and expression above this baseline might be well
55
56 24 tolerated. One additional possibility is that their protein abundance could be regulated translationally
57
58
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 rather than transcriptionally. Inefficient translation of certain genes may have been selected for during
2 evolution to prevent fluctuations in protein concentrations³². Perhaps a combination of these factors is at
3 play. Regardless, tight regulation of some embryonically essential pathways is not required in fully
4 developed adults, allowing for a wider range of expression values within the human population, resulting
5 in high expression variability.

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

17 We find that there is limited overlap in gene identity between Hyper- and Hypo-Variable genes in
18 breast and brain tissue. This suggests that expression variability is controlled by tissue-specific factors.
19 We propose that tissue identity is created and preserved, at least in part, by changes in gene expression
20 control pathways. Thus, genes that are Hypo-Variable in any given tissue have a constrained expression
21 pattern because they are likely to be important in the tissue-specific function and physiology of that
22 organelle. While there is limited overlap of genes within the corresponding EV gene sets of different
23 tissues, the Hyper-Variable gene sets of the different tissues have similar functional enrichments and
24 cellular protein localizations. Specifically, proteins encoded by Hyper-Variable genes tend to localize at

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 cell periphery and are enriched for cell surface signalling pathways and tissue development, including
2 tissue remodeling and ion transport. We therefore propose that high expression variability in genes
3 associated with tissue development pathways is an important component of tissue identity.

4 We did not observe any substantial sex dependent effects in expression variability. However, an
5 important conclusion of our study is that many Hyper-Variable genes have age-dependent expression
6 variability: that is, their expression increases or decreases during aging. One main cause of accelerated
7 brain aging and a causal factor of neurodegeneration is a reduction in immunological functions [42, 43].
8 We see evidence of downregulated immune responses in the cerebellum, specifically leukocyte mediated
9 immunity, defense responses to other organisms, and interferon-gamma response pathways. Many
10 studies also suggest that aging is associated with the upregulation of inflammatory responses [44], which
11 is a pathogenic mechanism implicated in many age-related diseases, including cardiovascular disease,
12 Alzheimer’s disease, and Parkinson’s disease [45]. Consistent with this idea, we see an enrichment of
13 acute inflammatory response in the cerebellum gold cluster. Another mechanism that has been implicated
14 with age-related diseases, such as Alzheimer’s disease and Parkinson’s disease, is synaptic dysfunction
15 that can affect neuroendocrine signaling [46–48]. We see a downregulation of ion transport and trans-
16 synaptic signaling in the frontal cortex, which are key components of neurotransmission and membrane
17 excitability, and whose downregulation likely causes deficiencies in these complex processes.
18 Furthermore, we see an upregulation of genes associated with glial cell differentiation in the frontal cortex
19 across multiple gene clusters. Initially thought of as cells that merely support neurons, emerging research
20 shows that neuron-astrocyte-microglia interactions are crucial for the functional organization of the brain
21 [49]. In addition, genes specific to astrocytes and oligodendrocytes, two different types of glial cells, have
22 been shown to shift regional expression patterns upon aging, and are better predictors of biological age
23 than neuronal-specific genes [50]. This suggests that the Hyper-Variability and age-dependent

1
2
3
4 1 upregulation of genes associated with glial cell differentiation is caused by differences in normal brain
5
6 2 aging between the samples.
7
8

9
10 3 Without examining the mechanistic control of individuals genes, it is difficult to determine if
11
12 4 changes in gene expression result in repression or activation of their associated pathways. For example,
13
14 5 we see an upregulation in neurogenesis associated genes during aging in both the cerebellum and the
15
16 6 frontal cortex, despite the common theory that neurodegeneration is a ubiquitous effect of normal brain
17
18 7 aging. An emerging concept in neuroscience is that homeostatic plasticity of neurons is maintained
19
20 8 through local adjustments of neural activities [51]. This overexpression of genes in pathways whose
21
22 9 function is known to decline over time may be a compensatory mechanism for an inefficient, aging system.
23
24 10 Within the cerebellum, a decline in neuronal function that occurs with aging may cause an upregulation
25
26 11 of genes associated with neurogenesis pathways. In addition to mitigating neuronal dysfunction, localized
27
28 12 increases in neurogenesis may be induced in response to cerebral diseases or acute injuries for self-repair
29
30 13 [52]. Lastly, chronic antidepressant usage has also been shown to result in an increase in neurogenesis
31
32 14 [53], suggesting that psychopharmaceuticals can alter neurochemistry and mimic compensatory anti-
33
34 15 aging responses. Overall, EV plays an important role in aging, specifically in immune responses and
35
36 16 inflammation, neurotransmission, and neurogenesis. Age-dependent gene expression could reflect a loss
37
38 17 of regulatory control or be a part of a regulated pathway of development.
39
40
41
42
43
44
45

46 18 In summary, our work shows that gene expression variability in the human population is likely to
47
48 19 be important in development, methylation, and in aging. As such, the EV of a gene is an important
49
50 20 functional characteristic of the gene itself. Therefore, the classification of a gene as one with
51
52 21 Hypervariability or Hypovariability in a human population or in a specific tissue should be useful in the
53
54 22 identification of important genes that functionally regulate development or disease.
55
56
57
58
59 23

1 **Methods**

2
3
4
5
6
7
8 **2 Illumina gene expression and methylation microarray data.** The analysis was conducted on two separate
9
10 3 datasets, both utilizing the Illumina HumanHT-12 V3.0 expression BeadChip. The first dataset provides
11
12 4 high quality RNA-derived transcriptional profiling of breast-adjacent tissue from 144 samples***. The
13
14
15 5 associated genotype and expression data have been deposited at the European Genome-Phenome
16
17 6 Archive (EGA, <http://www.ebi.ac.uk/ega/>), which is hosted by the European Bioinformatics Institute,
18
19 7 under accession number EGAS00000000083. The microarray readings were preprocessed using the
20
21
22 8 author's own custom script based on existing functionality within the beadarray package in R and were
23
24
25 9 reported as a log2 intensity. This dataset is referred to as breast tissue.

26
27
28 10 The second gene expression and the methylation datasets were catalogued by the North
29
30 11 American Brain Expression Consortium and UK Human Brain Expression Database (UKBEC) [34, 54]. The
31
32 12 expression data was obtained from the Gene Expression Omnibus (GEO) database [55] under accession
33
34
35 13 number GSE36192. A total of 911 tissue samples were analyzed from frozen brain tissue from the
36
37 14 cerebellum and frontal cortex from 396 subjects. The microarray readings were processed using a cubic
38
39
40 15 spline normalization method in Illumina Genome Studio Gene Expression Module v3.2.7. The expression
41
42 16 levels were log2 transformed before any analysis. The methylation data was also obtained from GEO
43
44
45 17 under accession number GSE36194. A total of 724 tissue samples were analyzed from frozen brain tissue
46
47 18 from the cerebellum and frontal cortex from 318 subjects. The methylation microarray readings were
48
49 19 processed using BeadStudio Methylation Module v3.2.0 with no normalization.

50
51
52 20
53
54
55 21 **Preprocessing the datasets.** Since the brain expression and methylation datasets were individually
56
57 22 processed by different tissue banks and in several batches, we corrected for the batch effect using the
58
59
60 23 limma package in R. The breast tissue dataset was previously batch corrected by the authors. Next, we

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 subset the data into 16 groups based on the clinical annotations provided by the UKBEC database. These
2 annotations included tissue type (Cerebellum and Frontal Cortex), sex (Male and Female), and age group
3 (0-25 years, 25-50 years, 50-75 years, or 75+ years). For each of the 16 groups, we calculated the median
4 expression for each probe and performed a hierarchical clustering via multiscale bootstrap resampling
5 using the pvclust package in R. Using a p-value threshold of 0.01, we see that the ideal clustering method
6 is to subset the data by tissue type, dividing the expression and methylation brain datasets into
7 cerebellum and frontal cortex tissue datasets.

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

15
16 **Identification and removal of bimodal expression probes.** Probes expressions that exhibited a bimodal
17 distribution were thought of as having two exclusive phenotypic states. However, our focus in this analysis
18 was to examine the factors affecting the tightly regulated expression of Hypo-Variable probes or the highly
19 variable gene expression of Hyper-Variable probes. In order to identify if a gene's expression was
20 unimodal or bimodal, we modeled each gene expression as a mixture of two gaussian distributions using
21 the mixtools package in R. Next, we identified the peaks of the probability density functions for each
22 gaussian distribution and compared the distance between the peaks as well as the ratio of peak heights.
23 Probes with peaks that were greater than one MAD apart and displayed a peak ratio greater than 0.1 were

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 treated as having a bimodal expression and subsequently removed from the analysis. Probes that did not
2 satisfy these criteria were considered to have a unimodal distribution and were kept for further analysis.

3
4 **EV gene set classification.** We classified the probes into three distinct probes sets based on their
5 expression variability:

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

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

11
12 **Bootstrapping EV gene set classifications.** To statistically validate our EV classifications, we split our data
13 into two equally sized subsets and repeated the previously explained EV method. This 50-50 split-retest
14 procedure was repeated 100 times per tissue. Next, we determined the accuracy our of original
15 classifications by comparing original classification of each gene with the 50-50 split classifications using a
16 binomial test with a probability of success greater than 0.5. In this hypothesis, a “success” is defined as
17 consistent EV classification across all three subsets, and gene classifications were considered significant
18 with a p-value < 0.05.

19
20 **Structural analysis of EV genes.** Data regarding the structural features of the genes was obtained from
21 the GRCh38/hg38 assembly of UCSC Table Browser [56]. Linear regression analyses were conducted to
22 find any correlation between gene EV and their structural features. For the linear regression analysis of

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 transcript size, we individually examined the largest and smallest transcripts separately. The sequence
2 lengths excluded introns, 3' and 5' UTR exons, and any upstream or downstream regions.

3
4 **Gene cluster analysis.** The GO term enrichment analyses were conducted using ConsensusPathDB gene
5 set over-representation analysis [26]. The complete list of unique Illumina HumanHT-12 V3.0 expression
6 BeadChip genes was used as a background list of genes. The resulting GO terms were then filtered
7 manually using a q-value cutoff of 0.05. Common and unique GO terms were summarized using REVIGO
8 [57] and visualized through treemaps by the provided R scripts. The parameters used were a medium
9 allowed similarity (0.7) using Homo sapiens database of GO terms.

10
11 **Enrichment analyses.** Using the Pearson's chi-square test, we tested for enrichment of essential genes in
12 each gene set relative to the total number of essential genes in the Illumina HumanHT-12 V3.0 expression
13 BeadChip. A list of 20,029 protein coding genes from the CCDS database was used to test for essentiality
14 enrichment [26]. Only genes that are solely classified as essential are considered in the analysis, resulting
15 in a list of 2377 essential genes present in the dataset. Once the number of annotated genes and gene
16 sets were deemed dependent variables, we determine the enrichment of annotated genes using the
17 Pearson residuals.

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

20
21 **Hierarchical clustering of age-dependent Hyper-Variable genes.** With the exception of a few groups, the
22 hierarchical clustering groups with the opposite sex and the same age groups tended to cluster together.

1
2
3
4 1 While the p-values of the sex and age groupings during the hierarchical clustering were too high to warrant
5
6 2 further subsetting of the brain dataset samples into distinct groups, they were significant enough to
7
8
9 3 inspect on a gene-by-gene basis.

10
11
12 4 We used a multiple linear regression model to measure the changes in expression of the Hyper-
13
14 5 Variable probes as a function of age, sex, and post-mortem interval (PMI):

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

17
18
19
20 7 where Y_i is the expression level of a probe and β_n is the coefficient for each term. The p-values were
21
22
23 8 calculated using a type III sum of squares regression and adjusted for multiple comparisons using the
24
25 9 Benjamini-Hochberg method. Probes that exhibit an FDR < 0.01 were considered significant for the
26
27
28 10 specific coefficient, and the sign of the coefficient determines if the probe is positively or negatively
29
30 11 correlated with the factor.

31
32
33 12 The choice to use three age clusters as the optimal number of clusters to examine changes of EV
34
35 13 across age samples was determined using an expectation-maximization (EM) algorithm initialized by
36
37
38 14 hierarchical clustering for parameterized Gaussian mixture models in the mclust package of R. The
39
40 15 Bayesian information criterion for each hierarchical clustering model was determined, and both the
41
42 16 cerebellum and frontal cortex displayed identical optimal numbers of age clusters. Once the samples were
43
44
45 17 correctly clustered by age, the gene clusters were selected by cutting the gene dendrograms manually.
46
47 18 The gene expressions were then visualized as heatmaps using the gplots package in R.

48
49
50 19

51 52 53 54 20 **List of Abbreviations**

55
56
57 21
58
59
60 22 EGA European Genome-Phenome Archive

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

10 **Declarations**

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

12 **Consent for publication.** Not applicable

13 **Availability of data and material.**

14 The datasets analyzed in this study are available in the GEO repository under the accession number
15 GSE36192 (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE36192>) and GSE36194
16 (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE36194>). The remaining dataset analyzed this
17 study is available from European Genome-phenome Archive but restrictions apply to the availability of
18 these data and are not publicly available. Data are however available from the authors upon reasonable
19 request and with permission of EGA.

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

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 **Funding.** This work was supported by an operating grant from the Canadian Breast Cancer Foundation

2 (JML).

3 **Authors' contributions.** All authors contributed to the analysis of the results. NB performed the

4 computational data analysis, prepared figures, and was a major contributor in writing the manuscript. JL

5 conceived, coordinated, and supervised the work. All authors read and approved the final manuscript.

6 **Acknowledgements.** Not applicable

7 **Additional Files**

9 Additional file 1: Structural analysis of genes as a function of EV

10 Additional file 2: EV correlation between different tissue types

11 Additional file 3: Complete list of GO term treemaps for all genes

12 Additional file 4: Complete list of GO term treemaps for essential genes

13 Additional file 5: Complete list of GO term treemaps for age-regulated Hyper-Variable genes

15 **References**

16 1. Alemu EY, Carl JW, Corrada Bravo H, Hannenhalli S. Determinants of expression variability. *Nucleic*
17 *Acids Res.* 2014;42:3503–14.

18 2. Roberfroid S, Vanderleyden J, Steenackers H. Gene expression variability in clonal populations: Causes
19 and consequences. *Crit Rev Microbiol.* 2016;42:969–84.

20 3. Choi JK, Yu U, Yoo OJ, Kim S. Differential coexpression analysis using microarray data and its
21 application to human cancer. *Bioinformatics.* 2005;21:4348–55.

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 4. Ho JWK, Stefani M, dos Remedios CG, Charleston MA. Differential variability analysis of gene
2 expression and its application to human diseases. *Bioinformatics*. 2008;24:i390–8.

3 5. Ecker S, Pancaldi V, Rico D, Valencia A. Higher gene expression variability in the more aggressive
4 subtype of chronic lymphocytic leukemia. *Genome Med*. 2015;7:8.

5 6. Chen E-H, Hou Q-L, Wei D-D, Jiang H-B, Wang J-J. Phenotypic plasticity, trade-offs and gene expression
6 changes accompanying dietary restriction and switches in *Bactrocera dorsalis* (Hendel) (Diptera:
7 Tephritidae). *Sci Rep*. 2017;7. doi:10.1038/s41598-017-02106-3.

8 7. Newman JRS, Ghaemmaghami S, Ihmels J, Breslow DK, Noble M, DeRisi JL, et al. Single-cell proteomic
9 analysis of *S. cerevisiae* reveals the architecture of biological noise. *Nature*. 2006;441:840–6.

10 8. Singh GP. Coupling Between Noise and Plasticity in *E. coli*. *G3 Genes Genomes Genet*. 2013;3:2115–
11 20.

12 9. Taniguchi Y, Choi PJ, Li G-W, Chen H, Babu M, Hearn J, et al. Quantifying *E. coli* Proteome and
13 Transcriptome with Single-Molecule Sensitivity in Single Cells. *Science*. 2010;329:533–8.

14 10. Silander OK, Nikolic N, Zaslaver A, Bren A, Kikoin I, Alon U, et al. A Genome-Wide Analysis of
15 Promoter-Mediated Phenotypic Noise in *Escherichia coli*. *PLoS Genet*. 2012;8.
16 doi:10.1371/journal.pgen.1002443.

17 11. Wolf L, Silander OK, van Nimwegen E. Expression noise facilitates the evolution of gene regulation.
18 *eLife*. 4. doi:10.7554/eLife.05856.

19 12. Barkai N, Shilo B-Z. Variability and Robustness in Biomolecular Systems. *Mol Cell*. 2007;28:755–60.

20 13. Lehner B. Selection to minimise noise in living systems and its implications for the evolution of gene
21 expression. *Mol Syst Biol*. 2008;4:170.

22 14. Lehner B. Conflict between Noise and Plasticity in Yeast. *PLoS Genet*. 2010;6.
23 doi:10.1371/journal.pgen.1001185.

24 15. Blake WJ, Balázsi G, Kohanski MA, Isaacs FJ, Murphy KF, Kuang Y, et al. Phenotypic Consequences of
25 Promoter-Mediated Transcriptional Noise. *Mol Cell*. 2006;24:853–65.

26 16. Bishop AL, Rab FA, Sumner ER, Avery SV. Phenotypic heterogeneity can enhance rare-cell survival in
27 ‘stress-sensitive’ yeast populations. *Mol Microbiol*. 2007;63:507–20.

28 17. Ackermann M, Stecher B, Freed NE, Songhet P, Hardt W-D, Doebeli M. Self-destructive cooperation
29 mediated by phenotypic noise. *Nature*. 2008;454:987–90.

30 18. Zhang Z, Qian W, Zhang J. Positive selection for elevated gene expression noise in yeast. *Mol Syst*
31 *Biol*. 2009;5:299.

32 19. Ward MC, Gilad Y. Human genomics: Cracking the regulatory code. *Nature*. 2017;550:190–1.

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 20. Li J, Liu Y, Kim T, Min R, Zhang Z. Gene Expression Variability within and between Human Populations
2 and Implications toward Disease Susceptibility. *PLOS Comput Biol.* 2010;6:e1000910.

3 21. Hough SR, Laslett AL, Grimmond SB, Kolle G, Pera MF. A Continuum of Cell States Spans Pluripotency
4 and Lineage Commitment in Human Embryonic Stem Cells. *PLOS ONE.* 2009;4:e7708.

5 22. Kalmar T, Lim C, Hayward P, Muñoz-Descalzo S, Nichols J, Garcia-Ojalvo J, et al. Regulated
6 Fluctuations in Nanog Expression Mediate Cell Fate Decisions in Embryonic Stem Cells. *PLOS Biol.*
7 2009;7:e1000149.

8 23. Mar JC, Matigian NA, Mackay-Sim A, Mellick GD, Sue CM, Silburn PA, et al. Variance of Gene
9 Expression Identifies Altered Network Constraints in Neurological Disease. *PLOS Genet.*
10 2011;7:e1002207.

11 24. Carey LB, van Dijk D, Sloot PMA, Kaandorp JA, Segal E. Promoter Sequence Determines the
12 Relationship between Expression Level and Noise. *PLoS Biol.* 2013;11.
13 doi:10.1371/journal.pbio.1001528.

14 25. Batenchuk C, St-Pierre S, Tepliakova L, Adiga S, Szuto A, Kabbani N, et al. Chromosomal Position
15 Effects Are Linked to Sir2-Mediated Variation in Transcriptional Burst Size. *Biophys J.* 2011;100:L56–8.

16 26. Pruitt KD, Harrow J, Harte RA, Wallin C, Diekhans M, Maglott DR, et al. The consensus coding
17 sequence (CCDS) project: Identifying a common protein-coding gene set for the human and mouse
18 genomes. *Genome Res.* 2009;19:1316–23.

19 27. Blake JA, Eppig JT, Kadin JA, Richardson JE, Smith CL, Bult CJ. Mouse Genome Database (MGD)-2017:
20 community knowledge resource for the laboratory mouse. *Nucleic Acids Res.* 2017;45 Database
21 issue:D723–9.

22 28. Cedar H. DNA methylation and gene activity. *Cell.* 1988;53:3–4.

23 29. Moore LD, Le T, Fan G. DNA Methylation and Its Basic Function. *Neuropsychopharmacology.*
24 2013;38:23–38.

25 30. Irvine RA, Lin IG, Hsieh C-L. DNA Methylation Has a Local Effect on Transcription and Histone
26 Acetylation. *Mol Cell Biol.* 2002;22:6689–96.

27 31. Wagner JR, Busche S, Ge B, Kwan T, Pastinen T, Blanchette M. The relationship between DNA
28 methylation, genetic and expression inter-individual variation in untransformed human fibroblasts.
29 *Genome Biol.* 2014;15:R37.

30 32. Zhang B, Zhou Y, Lin N, Lowdon RF, Hong C, Nagarajan RP, et al. Functional DNA methylation
31 differences between tissues, cell types, and across individuals discovered using the M&M algorithm.
32 *Genome Res.* 2013;23:1522–40.

33 33. Birdsill AC, Walker DG, Lue L, Sue LI, Beach TG. POSTMORTEM INTERVAL EFFECT ON RNA AND GENE
34 EXPRESSION IN HUMAN BRAIN TISSUE. *Cell Tissue Bank.* 2011;12:311–8.

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 34. Trabzuni D, Ramasamy A, Imran S, Walker R, Smith C, Weale ME, et al. Widespread sex differences in
2 gene expression and splicing in the adult human brain. *Nat Commun.* 2013;4:ncomms3771.

3 35. Harada CN, Natelson Love MC, Triebel K. Normal Cognitive Aging. *Clin Geriatr Med.* 2013;29:737–52.

4 36. Annunziato L, Pannaccione A, Cataldi M, Secondo A, Castaldo P, Di Renzo G, et al. Modulation of ion
5 channels by reactive oxygen and nitrogen species: a pathophysiological role in brain aging? *Neurobiol*
6 *Aging.* 2002;23:819–34.

7 37. Montecino-Rodriguez E, Berent-Maoz B, Dorshkind K. Causes, consequences, and reversal of
8 immune system aging. *J Clin Invest.* 2013;123:958–65.

9 38. Kærn M, Elston TC, Blake WJ, Collins JJ. Stochasticity in gene expression: from theories to
10 phenotypes. *Nat Rev Genet.* 2005;6:nrg1615.

11 39. Balaban NQ, Merrin J, Chait R, Kowalik L, Leibler S. Bacterial Persistence as a Phenotypic Switch.
12 *Science.* 2004;305:1622–5.

13 40. Garcia-Bernardo J, Dunlop MJ. Phenotypic Diversity Using Bimodal and Unimodal Expression of
14 Stress Response Proteins. *Biophys J.* 2016;110:2278–87.

15 41. Georgi B, Voight BF, Bućan M. From Mouse to Human: Evolutionary Genomics Analysis of Human
16 Orthologs of Essential Genes. *PLoS Genet.* 2013;9. doi:10.1371/journal.pgen.1003484.

17 42. Streit WJ, Xue Q-S. The Brain's Aging Immune System. *Aging Dis.* 2010;1:254–61.

18 43. Lucin KM, Wyss-Coray T. Immune activation in brain aging and neurodegeneration: too much or too
19 little? *Neuron.* 2009;64:110–22.

20 44. Singh P, Goode T, Dean A, Awad SS, Darlington GJ. Elevated Interferon Gamma Signaling Contributes
21 to Impaired Regeneration in the Aged Liver. *J Gerontol A Biol Sci Med Sci.* 2011;66A:944–56.

22 45. Wu D, Meydani SN. Age-associated changes in immune and inflammatory responses: impact of
23 vitamin E intervention. *J Leukoc Biol.* 2008;84:900–14.

24 46. Azpurua J, Eaton BA. Neuronal epigenetics and the aging synapse. *Front Cell Neurosci.* 2015;9.
25 doi:10.3389/fncel.2015.00208.

26 47. Hebert LE, Beckett LA, Scherr PA, Evans DA. Annual Incidence of Alzheimer Disease in the United
27 States Projected to the Years 2000 Through 2050. *Alzheimer Dis Assoc Disord.* 2001;15:169–73.

28 48. Levy G, Schupf N, Tang M-X, Cote LJ, Louis ED, Mejia H, et al. Combined effect of age and severity on
29 the risk of dementia in Parkinson's disease. *Ann Neurol.* 2002;51:722–9.

30 49. Cerbai F, Lana D, Nosi D, Petkova-Kirova P, Zecchi S, Brothers HM, et al. The Neuron-Astrocyte-
31 Microglia Triad in Normal Brain Ageing and in a Model of Neuroinflammation in the Rat Hippocampus.
32 *PLOS ONE.* 2012;7:e45250.

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 50. Soreq L, Rose J, Soreq E, Hardy J, Tratzuni D, Cookson MR, et al. Major Shifts in Glial Regional
2 Identity Are a Transcriptional Hallmark of Human Brain Aging. *Cell Rep.* 2017;18:557–70.

3 51. Braegelmann K, Streeter K, Fields D, Baker T. Plasticity in respiratory motor neurons in response to
4 reduced synaptic inputs: a form of homeostatic plasticity in respiratory control? *Exp Neurol.* 2017;287 Pt
5 2:225–34.

6 52. Galvan V, Jin K. Neurogenesis in the aging brain. *Clin Interv Aging.* 2007;2:605–10.

7 53. Malberg JE, Eisch AJ, Nestler EJ, Duman RS. Chronic Antidepressant Treatment Increases
8 Neurogenesis in Adult Rat Hippocampus. *J Neurosci.* 2000;20:9104–10.

9 54. Hernandez DG, Nalls MA, Moore M, Chong S, Dillman A, Tratzuni D, et al. Integration of GWAS SNPs
10 and tissue specific expression profiling reveal discrete eQTLs for human traits in blood and brain.
11 *Neurobiol Dis.* 2012;47:20–8.

12 55. Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, et al. NCBI GEO: archive for
13 functional genomics data sets—update. *Nucleic Acids Res.* 2013;41:D991–5.

14 56. Karolchik D, Hinrichs AS, Furey TS, Roskin KM, Sugnet CW, Haussler D, et al. The UCSC Table Browser
15 data retrieval tool. *Nucleic Acids Res.* 2004;32 suppl_1:D493–6.

16 57. Supek F, Bošnjak M, Škunca N, Šmuc T. REVIGO Summarizes and Visualizes Long Lists of Gene
17 Ontology Terms. *PLOS ONE.* 2011;6:e21800.













