

Characterizing the landscape of gene expression variance in humans

Scott Wolf^{†,1}, Diogo Melo^{†,1,2,*}, Kristina M. Garske¹, Luisa F. Pallares^{1,3},
Amanda J. Lea^{1,2,4,5}, and Julien F. Ayroles^{1,2,*}

Submission v2.0 - Jan 9th, 2023

[†] These authors contributed equally to this work.

¹ Lewis-Sigler Institute for Integrative Genomics, Princeton University

² Department of Ecology and Evolutionary Biology, Princeton University

³ Friedrich Miescher Laboratory, Max Planck Society

⁴ Department of Biological Sciences, Vanderbilt University

⁵ Child and Brain Development, Canadian Institute for Advanced Research

* Correspondence: Diogo Melo <damelo@princeton.edu>, Julien F. Ayroles <jayroles@princeton.edu>

Abstract

Gene expression variance has been linked to organismal function and fitness but remains a commonly neglected aspect of molecular research. As a result, we lack a comprehensive understanding of the patterns of transcriptional variance across genes, and how this variance is linked to context-specific gene regulation and gene function. Here, we use 57 large publicly available RNA-seq data sets to investigate the landscape of gene expression variance. These studies cover a wide range of tissues and allowed us to assess if there are consistently more or less variable genes across tissues and data sets and what mechanisms drive these patterns. We show that gene expression variance is broadly similar across tissues and studies, indicating that the pattern of transcriptional variance is consistent. We use this similarity to create both global and within-tissue rankings of variation, which we use to show that function, sequence variation, and gene regulatory signatures contribute to gene expression variance. Low-variance genes are associated with fundamental cell processes and have lower levels of genetic polymorphisms, have higher gene-gene connectivity, and tend to be associated with chromatin states associated with transcription. In contrast, high-variance genes are enriched for genes involved in immune response, environmentally responsive genes, immediate early genes, and are associated with higher levels of polymorphisms. These results show that the pattern of transcriptional variance is not noise. Instead, it is a consistent gene trait that seems to be functionally constrained in human populations. Furthermore, this commonly neglected aspect of molecular phenotypic variation harbors important information to understand complex traits and disease.

Author Summary

Gene expression variance, or the variation in the level of gene expression within a population, can have significant impacts on physiology, disease, and evolutionary adaptations. While the average level of gene expression is typically the focus of research, the variation around this average level (i.e., gene expression variance) can also be important for understanding complex traits and disease. Here, we investigate the landscape of transcriptional variance across tissues, populations, and studies. Using large publicly available RNA-seq data sets, we were able to identify the general properties associated with high- and low-variance genes, as well as factors driving variation in variance across genes. Specifically, we uncovered gene expression variance was significantly associated with gene length, nucleotide diversity, the degree of connectivity and the presence of non-coding RNA. Our results suggest that the mechanisms responsible for maintaining optimal levels of variation in high- versus low-variance differ, and that this variability is the result of different patterns of selection.

Introduction

1 Molecular phenotypes such as gene expression are powerful tools for understanding physiology, disease, and
2 evolutionary adaptations. In this context, average trait values are usually the focus of investigation, while
3 variation around the average is often considered a nuisance and treated as noise [1]. However, gene expres-
4 sion variance can be directly involved in determining fitness [2,3], can drive phenotypic variation [4], and the
5 genetic architecture of variance itself can evolve [5]. This suggests that studying gene expression variance
6 as a bona fide trait, its genetic architecture, and the evolutionary mechanisms shaping and maintaining gene-
7 specific patterns of variance has the potential to further our understanding of complex traits and disease [6-8].

8 Variability is ubiquitous in nature and is, alongside its counterpart, robustness, a fundamental feature of most
9 complex systems. But, at the same time, the degree of variability seems to differ between genes [1] suggest-
10 ing that it might be associated with biological function and therefore be shaped by selection. From a mech-
11 anistic perspective, several competing forces act to shape transcriptional variance [5,9], and the outcome of
12 the interaction between these processes is still poorly understood [10]. For example, we expect the influx of
13 new mutations to increase the variance, while the selective removal of these polymorphisms, via purifying
14 selection or selective sweeps, to decrease it [11,12]. From a quantitative trait perspective, stabilizing selection
15 should decrease variation around an optimal value, and directional selection can lead to a transient increase
16 in variance while selected alleles sweep to fixation, followed by a reduction in variance as these alleles become
17 fixed. Pleiotropic effects are also important, as they allow selection on one trait to influence the variance of
18 other traits [13,14]. Both indirect effects of directional selection on variance open the possibility that the main
19 driver of gene expression variance is not direct selection on variance but indirect effects due to selection on
20 trait means [10]. How the interaction of these processes shape gene expression variance is an open question.
21 However, some general predictions can be made. If a homogeneous pattern of stabilizing selection is the main
22 driver of gene expression variance, we would expect transcriptional variance to be consistent regardless of the
23 population, tissue, or environmental context. If idiosyncratic selection or environmental patterns are more
24 important, we could observe large differences in gene expression variance across studies.

25 A key difficulty in addressing these questions is that the constraints on gene expression variance might also
26 be dependent on the gene tissue specificity. Mean expression is known to differ across tissues [15], however,
27 to what extent differential expression (i.e., differences in mean expression level) translate into differences in
28 expression variance is not clear. Higher mean expression could lead to higher variance, but other processes
29 can also affect transcriptional variance. For example, if a gene is expressed in more than one tissue and vari-
30 ance regulation is independent across tissues, stabilizing selection on gene expression could be more intense
31 depending on the role of that gene in a particular tissue, causing a local reduction in variation that leads to
32 differences in variance across tissues (fig. 1 A). These across-tissue differences would not necessarily follow

33 mean expression. Alternatively, expression variation across tissues could be tightly coupled and, in this ex-
34 ample, selection in one tissue would lead to a reduction in variance across tissues, resulting in a consistent
35 pattern of variation (fig. 1 B). While we lack a clear picture of how tissue-specific gene expression variation is
36 regulated, Alemu et al. [16] used microarray data from several human tissues to show that epigenetic mark-
37 ers were linked to gene expression variation and that these markers were variable across tissues and between
38 high- and low-variance genes.

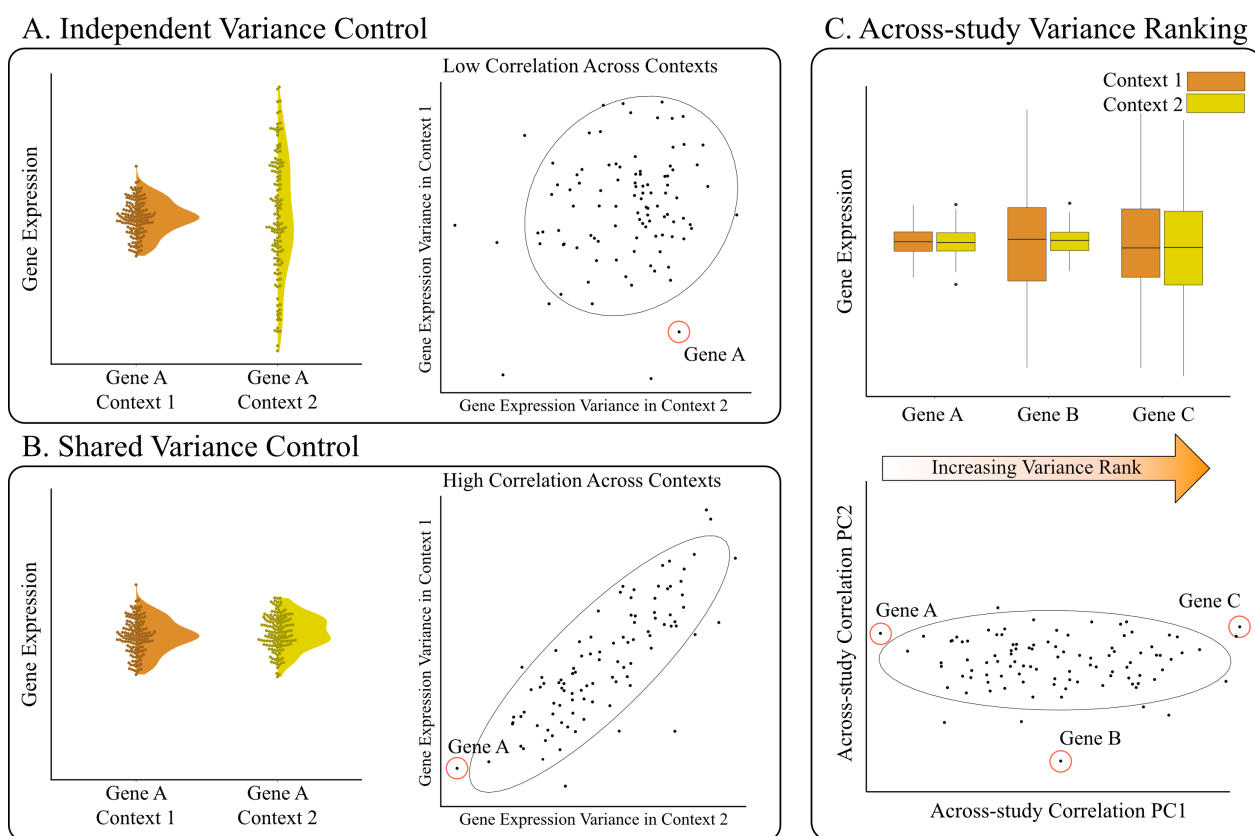


Figure 1: Example of how differences in the regulation of transcriptional variance can drive changes in the correlations between gene expression variance measures. In (A), independent regulation causes the reduction in variation to be restricted to context 1 (context here can refer to different tissues, environments, populations, studies, etc.). On the right side of panel A, independent regulation results in low correlation across contexts. In (B), a shared regulatory architecture maintains consistent variance across both conditions, leading to high similarity in transcriptional variance across contexts. In (C), we see how the similarity seen in panel B can be leveraged to create an across-context rank of gene expression variance. When transcriptional variance ranks are highly correlated, the rank of the projection onto the first principal component (PC1) allows us to summarize the across-context pattern of transcriptional variance.

39 To explore the landscape of gene expression variance and the association between transcriptional variance
40 and biological function, we use 57 publicly available human gene expression data sets spanning a wide range
41 of experimental contexts and tissues. By comparing the gene expression variance measured across such het-
42 erogeneous data sets, we show that the degree of expression variance is indeed consistent across studies and
43 tissues. We use the observed similarities to create an across-study gene expression variance ranking, which

44 orders genes from least variable to most variable. We then integrate various genomic-level functional annota-
 45 tions as well as sequence variation to probe the drivers of this variance ranking. Finally, we explore the link
 46 between gene expression variance and biological function by leveraging gene ontology and other gene annota-
 47 tions.

48 Results

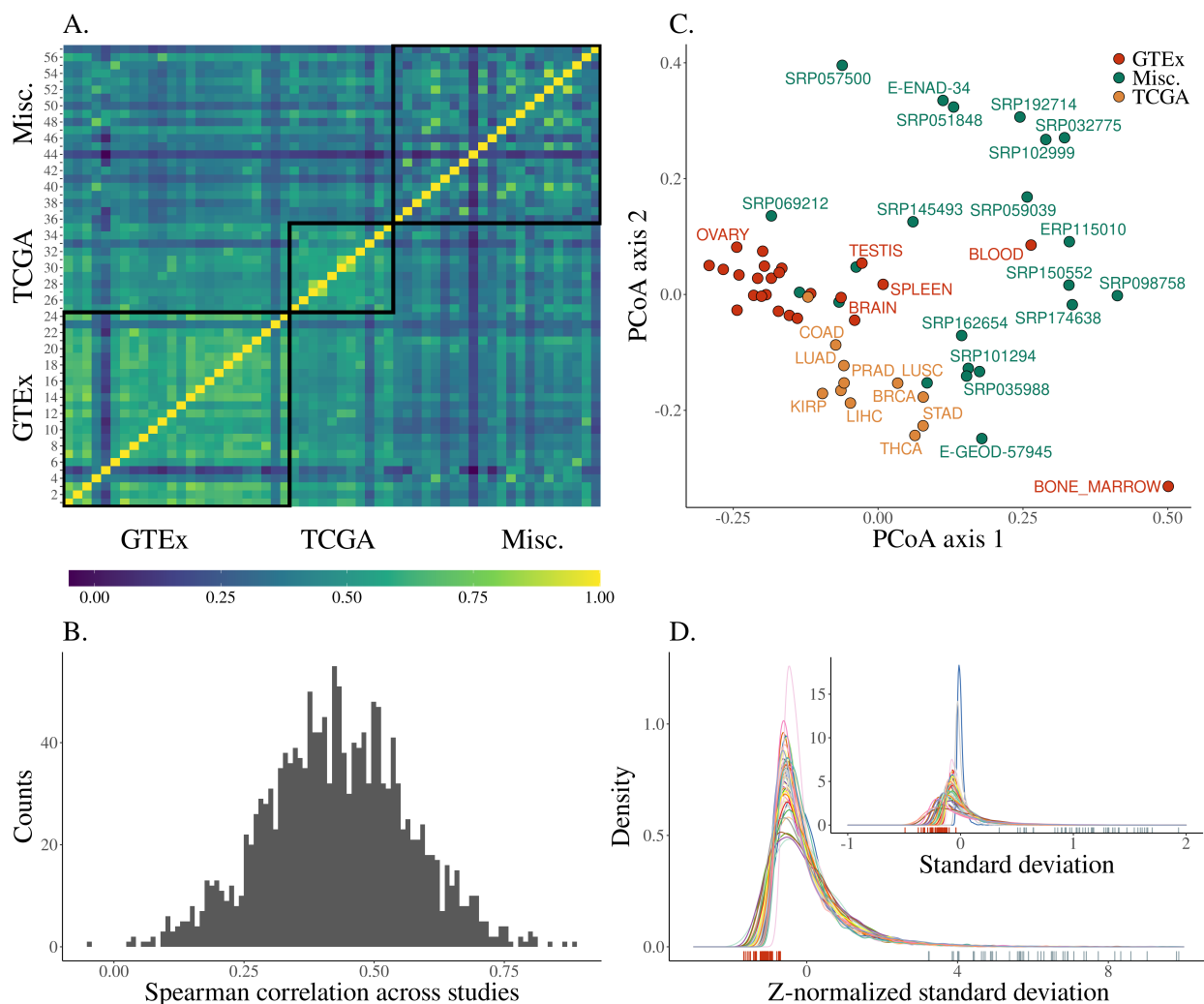


Figure 2: Overview of the distribution of transcriptional variance across studies. (A) Heatmap showing the correlation in transcriptional variance across studies (as the Spearman correlation of gene expression standard deviations). Pairs of studies with more similar patterns of gene expression variance have higher correlations. Studies are shown in the same order as in SI fig. 1, panel A. (B) Distribution of the pairwise Spearman correlations between studies shown in the previous panel. (C) PCoA using the distance between studies derived from the pairwise correlations. (D) Density plot of standard deviations after z-normalization. The inset plot shows the distribution of mean-centered standard deviations grouped by study without normalization. The corresponding rug plots show the location of the highest-ranking gene in standard deviation rank (*HBB*) (right, blue) and lowest (*WDR33*) (left, red).

49 **Data sets**

50 We use 57 publicly available human gene expression RNA-seq data sets which were derived from the publica-
51 tions listed in table 1 of the Methods section, and a complete metadata table for each study is available in the
52 supporting information (SI data 1). We only use data sets that fulfilled the following conditions: samples came
53 from bulk RNA-seq (and no single cell approaches), data sets were associated with a publication, sample-level
54 metadata was available, and the post-filtering sample size was greater than 10 (note that we did not included
55 data from non-baseline/exposure/stimulated datasets). These data sets span 13 different tissue types and the
56 post-filtering mean sample size we used for each data set was 390, with a median of 251, and ranged from 12 to
57 2931 samples. Several data sets were derived from two large consortia: GTEx [15] and TCGA [17], and we note
58 the origin of the data sets in the figures where appropriate. We refer to data sets and studies interchangeably,
59 and so each tissue in GTEx is referred to as a different study. The final list of genes used from each study can
60 be found in SI data 2.

61 **Gene expression variance**

62 For each study, transcriptional variance per gene was measured as the standard deviation (SD) of the distri-
63 bution of gene expression values for all individuals in a particular study. Mean and variance are known to
64 be correlated in RNA-seq data, both due to the nature of count data and the expectation that more highly ex-
65 pressed genes should have more variation. As our focus here is on variance, we control for both of these ex-
66 pected drivers of transcriptional variation. To achieve this, SD was calculated using a unified pipeline that
67 normalized the mean-variance relation in read-count data, controlled for batch effects, and removed outliers
68 (see Methods for details, and the calculated values for means and standard deviations are available in SI data
69 3). The observed range of gene expression SDs across genes is variable but can be normalized so that the dis-
70 tributions are comparable (fig. 2 D). This comparison reveals differences in the range of gene expression SDs
71 that can be due to any number of methodological or biological differences between the data sets. We avoid
72 having to deal with these global differences in the range of variation by using only the ranking of the genes
73 according to their gene expression SD in each study. Therefore, patterns of transcriptional variance were com-
74 pared across studies using Spearman correlations (ρ_s) between gene expression SDs. This comparison reveals
75 a broadly similar rank of gene expression variance as the correlations across studies are mostly positive and
76 high (75% of correlations are between 0.45 and 0.9, fig. 2 A and B), indicating that genes that are most vari-
77 able in one study tend to be most variable in all studies. A principal coordinate analysis [18] using $|1 - \rho_s|$ as
78 a between-study distance measure does not show clearly delineated groups, but GTEx and TCGA studies are
79 clustered among themselves and close together (fig. 2 C). This clustering indicates some effect of study source
80 on the similarity between gene expression SD across studies, which we explore in detail below.

81 To characterize what factors may explain differences in across-study similarity, we directly modeled the
82 across-study correlations using a mixed-effect linear model designed to account for the non-independence
83 in pairwise correlation data [19,20]. In this model (see Methods), we use a random effect for individual
84 study ID, a fixed effect for pairwise tissue congruence (whether a comparison is within the same tissue or
85 different tissue), and a fixed effect for pairwise study source (which pair of sources among GTEx, TCGA, and
86 miscellaneous is involved in a comparison) as predictors of the correlations (see Methods). This model (SI
87 fig. 1) shows that comparisons of studies within GTEx and TCGA have on average higher values of ρ_s , but also
88 that comparing studies across GTEx and TCGA also shows a mild increase in the average correlation (SI fig. 1
89 C). Correlations that do not involve studies from TCGA and GTEx (marked as “Misc.”) are on average lower
90 (SI fig. 1 C). While we do not have a clear explanation for this pattern, since TCGA and GTEx are independent,
91 this mild effect on the similarities could be due to the level of standardization of the data coming from these
92 two large consortia. Tissue type also affects the degree of similarity in transcriptional variance, with studies
93 using the same tissue being, on average, more similar (SI fig. 1 B). However, all these pairwise effects are mild,
94 and the largest effects on the correlations are those associated with individual studies, in particular some
95 specific tissues, i.e., comparisons involving BONE MARROW (from GTEx) and study SRPO57500 (which used
96 platelets) are on average lower (SI fig. 1 A). The only negative correlation we observe is between these two
97 studies, which also appear further away in the PCoA plot in fig. 2 C.

98 **Transcriptional variance rank**

99 The strong correlations between transcriptional variance across studies suggest that variance rank is indeed a
100 property of genes that can be robustly estimated. To estimate this gene-level rank, we devised an across-study
101 approach that allowed us to rank individual genes according to their degree of transcriptional variance by
102 averaging the ordering across all studies. We do this by calculating the score of each gene on the first principal
103 component of the across-study Spearman correlation matrix shown in fig. 2 A. This procedure is illustrated in
104 fig. 1 C. Ordering genes using these scores generate a ranked list of genes, with the most variable genes having
105 the highest rank. The position in the SD distributions shown in fig. 2 D of the most and least variable genes in
106 this rank illustrates how the extremes of the rank are indeed some of the least and most variable genes across
107 all studies. In addition, to be able to account for any residual effect of mean expression on the variance we also
108 created a similar across-study rank for mean expression. To explore tissue-specific divers or transcriptional
109 variation, we also create a set of tissue-specific SD ranks. To that end, we used the same procedure outlined
110 above but using only studies that were performed on the same tissue. Both tissue-specific and across-study
111 ranks are available in the Supporting Information (SI data 4).

112 **Biological function explains gene-level transcriptional variance**

113 As a first step toward explaining the factors that drive variation in variability between transcripts, we focused
114 on the top 5% most variable and the bottom 5% least variable genes in the across-study ranking (560 genes in
115 each group). A Gene Ontology (GO) enrichment analysis shows 59 enriched terms in the low-variance genes,
116 and 738 enriched terms in the high-variance genes (using a hypergeometric test and a conservative Benjamini-
117 Hochberg (BH) adjusted p-value threshold of 10^{-3} ; see supporting information SI data 5 for a complete listing).
118 Among the most variable genes, we observe enrichment for biological processes such as immune function,
119 response to stimulus, maintenance of homeostasis, and tissue morphogenesis (SI fig. 2 A). Furthermore, we
120 see a 7.7-fold enrichment for genes that encode secreted proteins in the most variable genes, relative to all
121 other genes (hypergeometric test, $p < 10^{-3}$).

122 Among the least variable genes, we see enrichment for housekeeping functions such as mRNA processing, cell
123 cycle regulation, methylation, histone modification, translation, transcription, and DNA repair (SI fig. 2 B);
124 accordingly, we also find a 2.0-fold enrichment in previously characterized human housekeeping genes [21]
125 (hypergeometric test, $p < 10^{-3}$). The genes exhibiting the lowest variance are also enriched for genes that have
126 been previously shown to have a high probability of being loss-of-function intolerant (pLI) [22] (1.2-fold enrich-
127 ment, hypergeometric test, $p < 10^{-3}$). Genes with a high pLI have been shown to be important in housekeeping
128 functions and have higher mean expression values across a broad set of tissues and cell types [22]. The ob-
129 servation that genes with low variance are enriched for both housekeeping genes and genes with high pLI is
130 consistent with this previous report; and we further see that the mean expression of genes positively corre-
131 lates with pLI (partial Spearman correlation $\rho_s = 0.32$, $p < 10^{-3}$), showing the opposite relationship between
132 variance and mean expression when considering pLI.

133 In the previous analysis, we explored the relationship between transcriptional variance and function by start-
134 ing from the extremes of the variance distribution and searching for GO enrichment among these high- and
135 low-variance genes. We also approach the problem from the opposite direction, starting from the genes as-
136 sociated with each GO term and searching for enrichment for high- or low-variance genes among them. To
137 this end, we gathered all biological process GO terms in level 3 (i.e., terms that are at a distance of 3 from the
138 top of the GO hierarchy). Using level-3 terms gives us a good balance between number of terms and genes
139 per term. We separated the genes associated with at least one of these level-3 terms into expression variance
140 deciles, with the first decile having the lowest variance. We then counted how many genes in each decile have
141 been associated with each term. If variance rank is not associated with the GO annotations, terms should have
142 an equal proportion of genes in each decile. We measured how far from this uniform allocation each term is by
143 measuring the Shannon entropy of the proportion of genes in each decile. Higher entropy is associated with a

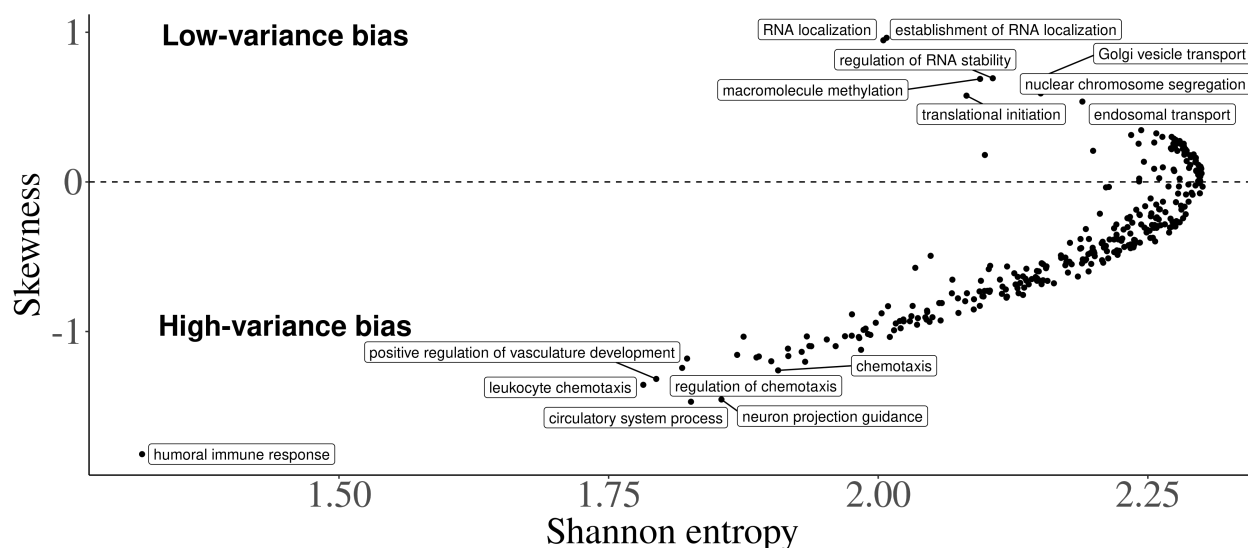


Figure 3: Relationship between skew and entropy of rank decile distributions for each GO term. High entropy terms, to the right of the plot, are associated with a more egalitarian proportion of genes in each of the SD rank deciles. The terms on the left of the plot are associated with more genes in some particular decile. The skewness in the y-axis measures if the high- or low-variance deciles are more represented for a particular term. Terms on the positive side of the y-axis are associated with low-variance genes, and terms on the negative side of the y-axis are associated with high-variance genes. The GO terms are filtered for gene counts greater than 100, as in fig. 4. Some of the top high- and low-skewness terms are labeled for illustration.

144 more uniform distribution of genes across deciles. GO terms with low entropy indicate some deciles are over-
 145 represented in the genes associated with that term. We also measured skewness for each term, which should
 146 be zero if no decile is over-represented, negative if high-variance terms are over-represented, and positive
 147 if low-variance deciles are over-represented. The relation between skewness and entropy for each GO term
 148 can be seen in fig. 3. Positive-skew low-entropy terms, those enriched with low-variance genes, are associated
 149 with housekeeping functions, like RNA localization, translation initiation, methylation, and chromosome seg-
 150 regation (fig. 4 A). Likewise, terms with negative skew and low entropy, enriched for high-variance genes, are
 151 related to immune response, tissue morphogenesis, chemotaxis—all dynamic biological functions related to
 152 interacting with the environment (fig. 4 B).

153 Both GO analyses suggest a strong association between biological function and the degree of transcriptional
 154 variance. Genes associated with baseline fundamental functions, expected to be under strong stabilizing selec-
 155 tion, are also low-variance; high-variance genes are associated with responding to external stimuli (i.e., tissue
 156 reorganization and immune response).

157 Environmental sensitivity predicts transcriptional variance

158 As suggested by our GO enrichment analyses, one mechanism that may generate consistent variability in gene
 159 expression is the response to environmental inputs. In other words, high-variance genes may be those that

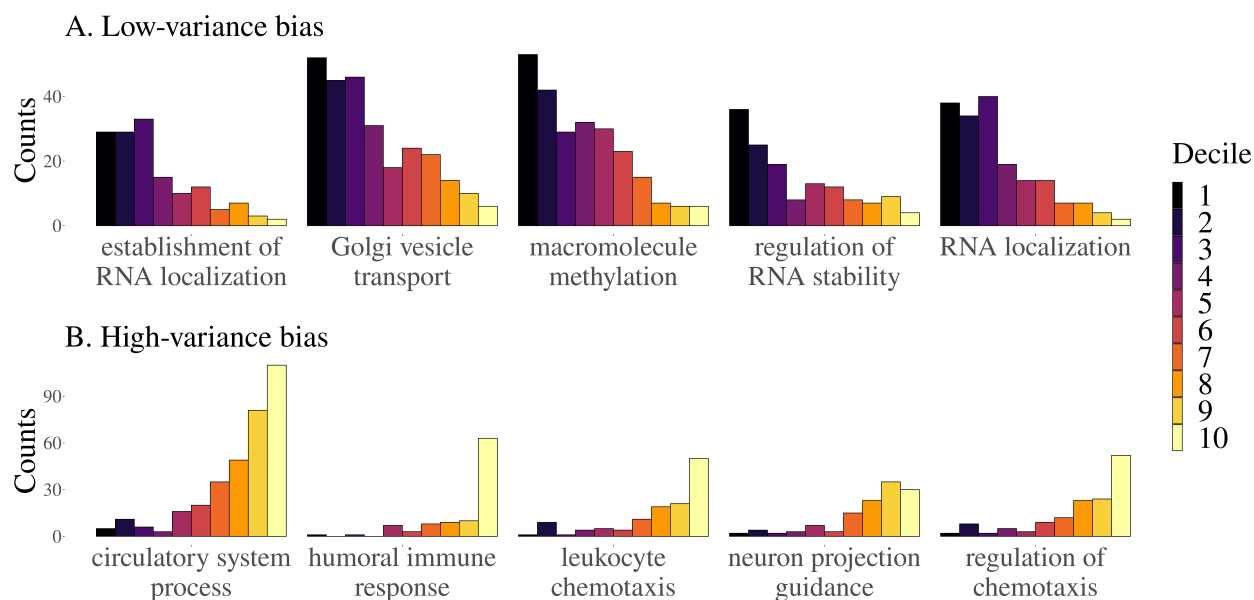


Figure 4: Distributions of decile ranks of level-3 GO terms. Each plot shows the count of genes in each decile of the rank. Only GO terms that are associated with at least 100 genes are used. We sort these terms by the skewness of the distribution. The top panel (A) shows the 5 most positively skewed terms, and the bottom panel (B) shows the 5 most negatively skewed terms.

160 are environmentally sensitive, while low-variance genes may be robust to environmental stimuli or pertur-
 161 bations (or alternatively, responsive to all stimuli, such that they are always highly expressed across individ-
 162 uals). To understand the relationship between environmental sensitivity and variance, we drew on gene ex-
 163 pression data from a recently generated catalog of environmentally responsive genes in lymphoblastoid cell
 164 lines (LCLs). This catalog was generated by exposing 544 LCLs derived from individuals included in the 1000
 165 Genomes Project to each of 11 in vitro exposures (including immune signaling molecules, hormones, and man-
 166 made chemicals), as well as a control; these manipulations of the cellular environment were followed by mRNA-
 167 seq and differential expression analyses comparing each treatment to its control [23]. Using lists of environ-
 168 mentally responsive genes derived from this study, we found that high-variance genes were more likely to
 169 respond to at least one in vitro exposure, relative to genes not classified as high- or low-variance (Fisher's ex-
 170 act test: $p < 0.05$, odds = 1.524); as predicted, the same is not true for low-variance genes (Fisher's exact test: p
 171 = 0.993, odds = 0.797). When we analyzed each exposure separately, we found that high-variance genes were
 172 more likely to be responsive to 4 out of the 10 environments we explored (1 environment was dropped due
 173 to a lack of differentially expressed genes in the original experiment; Fisher's exact test, FDR < 10%; SI table
 174 2). These exposures included key immune stimuli and hormones such as interferon gamma and dexametha-
 175 sone (a synthetic glucocorticoid). In contrast, we found that low-variance genes showed the opposite pattern:
 176 they are significantly underrepresented among environmentally responsive genes across 4/10 environments
 177 (Fisher's exact test, FDR < 10%; SI table 2). Though not all of our environment-specific tests reached statistical

178 significance, it is also worth noting that almost all 10 environments showed concordance in effect size direc-
179 tion (i.e., high-variance genes tended to be overrepresented among environmentally responsive genes and
180 low-variance genes tended to be underrepresented). While the above analyses show that high-variance genes
181 tend to overlap with genes induced by a given exposure, we hypothesized that genes that are similarly induced
182 by many different exposures may in fact exhibit moderate or low variance. In other words, genes induced by
183 many stimuli may always be highly expressed across individuals, and thus low variance, while genes induced
184 by select stimuli may only be upregulated in a subset of the population, and thus exhibit high variance. In
185 support of this idea, we found that, among genes that responded to at least one environments in the LCL exper-
186 iment, high-variance genes responded to a median of only one environment, while both low-variance genes
187 and the background set responded to a median of 4/10 environments (generalized linear model comparing
188 high-variance to background and low-variance, $p < 10^{-7}$ and $p < 10^{-11}$, respectively). Thus, high-variance genes
189 are indeed more likely to be environmentally sensitive, but in a highly select and stimulus-specific manner,
190 which we hypothesize drives their between-individual heterogeneity. We note that all analyses presented in
191 this section focused on the composite set of high- and low-variance genes defined across tissues, but we obtain
192 similar results when focusing on blood, the tissue in our dataset most similar to LCLs (SI table 2).

193 **Evolutionary forces at play in shaping transcriptional variance**

194 We use three gene-level summary statistics, nucleotide diversity (π), gene expression connectivity, and the
195 rate of adaptive substitutions (α), as a proxy to assess whether selection might be involved in shaping gene
196 expression variance. For all the correlations in this section, we use partial Spearman correlations that include
197 the mean gene expression rank as a covariate, which accounts for any residual mean-variance correlation.
198 Nucleotide diversity in the gene region is used as a proxy for the impact of cis-regulatory genetic variation on
199 transcriptional variance. As expected, low-variance genes tend to have lower levels of polymorphisms (partial
200 Spearman correlation, $\rho_s = 0.184$, $p < 10^{-10}$). Gene-gene connectivity, a proxy for gene regulatory interactions
201 and selective constraints [24], is, in turn, negatively correlated with the expression variance (partial Spearman
202 correlation, $\rho_s = -0.024$, $p < 10^{-2}$), supporting the expectation that highly connected genes are more constrained
203 in their variation. Finally, we also find that low-variance genes tend to have fewer substitutions by comparing
204 the across-study rank with α (partial Spearman correlation, $\rho_s = -0.044$, $p < 10^{-2}$), in line with the expectation
205 that genes under stronger selection should be less variable. Despite all associations being significant and in
206 the expected direction, their effect sizes are very small, suggesting a weak link between these broad measures
207 and transcriptional variance.

208 **Specific gene regulatory signatures are associated with transcriptional variance**

209 To assess how local epigenetic features relate to gene expression variance we calculate the proportion of the
210 gene (± 10 kb) that corresponds to epigenetic signatures of gene regulation defined through ChromHMM [25]
211 chromatin states. Chromatin states associated with distal (i.e., non-promoter) gene regulation are positively
212 correlated with the across-study variance rank, regardless of whether the regulatory effect on gene expres-
213 sion is positive or negative (fig. 5; see across-study correlations in SI fig. 3A). For example, both the proportion
214 of gene regions made up of enhancers and repressed genomic states are positively correlated with gene ex-
215 pression variance (BH adjusted Spearman correlation, $p < 0.05$). In contrast, histone modifications associated
216 with active promoters, as well as transcribed states, are inversely correlated with gene expression variance (SI
217 fig. 3A), whereas they are positively correlated with the mean rank (SI fig. 3B). Taken together, these results
218 are compatible with gene expression variance being regulated through distal (i.e., non-promoter) gene regula-
219 tory mechanisms, rather than the overall active transcriptional state of a gene region, as is the case with mean
220 gene expression.

221 Given that ChromHMM chromatin states are available for specific tissues, we asked whether the regulatory
222 signatures associated with the across-study variance rank are recapitulated at the tissue level. Many of the
223 across-study correlations are recapitulated at the tissue-specific level (with two exceptions noted below), in-
224 cluding a strong and highly consistent positive correlation between the proportion of gene regions made up
225 of enhancer states and that gene's expression variance, and an inverse relationship between gene expression
226 variance and histone marks associated with gene transcription (SI fig. 3A). Two blood associations stand out
227 as being different from the consistent effects across the other tissue-level and across-study associations. First,
228 the weak (i.e., histone marks associated with both activating and repressive functions) promoter state is posi-
229 tively correlated with transcriptional variance in all comparisons except blood. Second, the consistent inverse
230 correlation of gene expression variance with weak transcription is reversed in blood, such that there is a pos-
231 itive correlation between histone marks associated with weak transcription and blood gene expression vari-
232 ance (SI fig. 3A). Taken together, these results suggest that, rather than genes with a bivalent promoter state
233 (i.e., poised genes) exhibiting more expression variance, blood high-variance genes are more likely already
234 expressed at basal levels (i.e., weakly transcribed), as discussed previously [26].

235 Immediate early genes (IEGs) respond quickly to external signals without requiring *de novo* protein synthe-
236 sis, and a bivalent state has been reported to be associated with IEG promoters [reviewed in 27]. Given our
237 results that genes with high expression variance are enriched for cellular signaling and response mechanisms
238 (SI fig. 2 A), and bivalent promoter states are correlated with the gene expression variance rank (SI fig. 3A),
239 we hypothesized that IEGs would be enriched within genes in the top expression variance ranks. This was

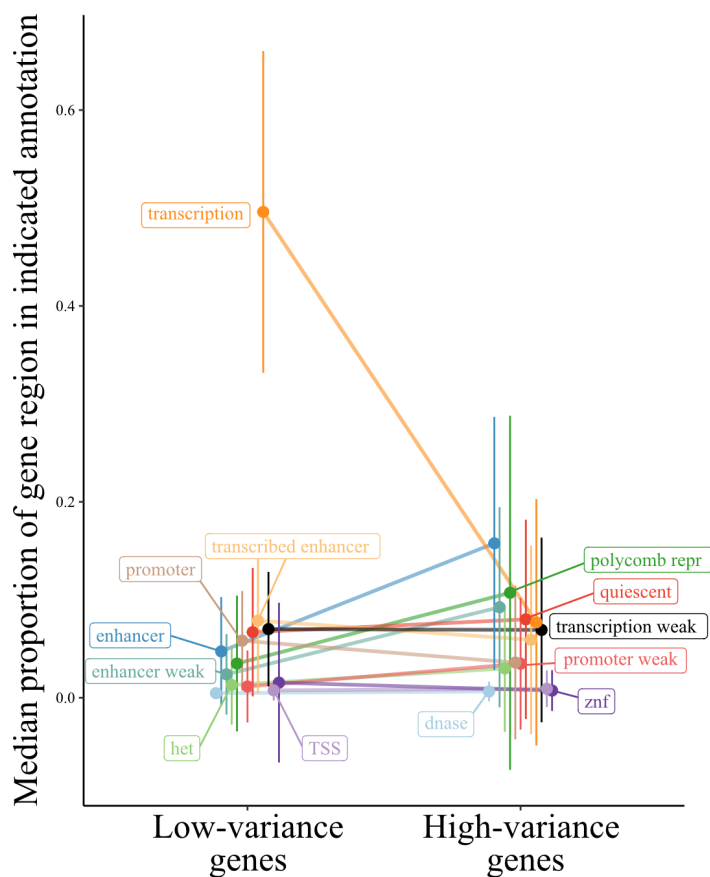


Figure 5: Proportion of gene regions made up of ChromHMM chromatin states for low- and high-variance genes. The line plot contrasts the proportion of gene regions made up of the indicated chromatin states for genes in the top and bottom 5% of the across-study variance rank metric. Ends denote the median proportion of gene regions made up of the chromatin state, and error bars represent the standard error of the mean. States colored black are not significant, all others exhibit significant differences between low- and high-variance genes (BH adjusted Wilcoxon signed-rank test, $p < 0.05$). Het indicates heterochromatin; TSS, transcription start sites; znf, zinc finger genes. The mean rank version of this analysis is shown in SI fig. 4.

240 the case for all tissue-level gene expression variance ranks (enrichment ratios range from 3.3-8.8, Bonferroni-
241 adjusted hypergeometric test, $p < 0.05$), except for blood (enrichment ratio = 1.2, hypergeometric test, $p = 0.3$).
242 Thus, once again blood stands out when attempting to understand genomic regulatory drivers of expression
243 variance. In all, while high-variance genes are generally shared across tissues and enriched for immune and
244 environmental signaling pathways, it seems that the gene regulatory mechanisms governing their expression
245 are distinct between immune cell types and other tissues studied here.

246 **Linking expression variance and disease**

247 To explore the link between transcription variance and genes known to be associated with human diseases, we
248 used a data set designed to provide causal relationships between gene expressions and complex traits/diseases
249 (based on a probabilistic transcriptome-wide association study (PTWAS) [28]). Using the list of significant
250 gene-disease pairs at 5% FDR provided by Zhang et al. [28], we performed a hypergeometric enrichment test for
251 the top 5% high- and low-variance genes in our across-study rank and in all tissue-specific gene variance ranks.
252 We use both across-study and tissue-specific ranks because some genes only appear in the tissue-specific rank
253 due to their limited tissue-specific gene expression. In the high-variance group, we find no enrichment in the
254 across-study rank, but we do find enrichment of genes annotated for allergy, immune disease, and endocrine
255 system disease among the high-variance genes in several tissue-specific variance ranks. For example, among
256 high-variance genes in the colon, we see enrichment for endocrine system disease (1.77-fold, hypergeomet-
257 ric test, $p < 10^{-4}$). Among high-variance genes in the immune cells, we see enrichment for endocrine system
258 disease (1.67-fold, hypergeometric test, $p < 10^{-3}$), allergy (1.7-fold, hypergeometric test, $p < 10^{-3}$), and immune
259 disease (1.32-fold, hypergeometric test, $p < 10^{-2}$). Among high-variance genes in the thyroid, we see enrichment
260 for endocrine system disease (1.9-fold, hypergeometric test, $p < 10^{-5}$), allergy (1.85-fold, hypergeometric test, p
261 $< 10^{-4}$), and immune disease (1.45-fold, hypergeometric test, $p < 10^{-4}$). These are all rather similar and suggest a
262 stable pattern of high-variance gene expression across these tissues, with enrichment for these three classes
263 of diseases. The link with immune diseases is expected given the high enrichment for immune-related genes
264 in the high-variance group [8]. As for the low-variance group, we found strong enrichment for genes associ-
265 ated with psychiatric and neurological disorders in the across-study rank and in some tissue-specific ranks
266 (breast, liver, and stomach; ~1.2-fold enrichment, hypergeometric test, $p < 0.05$, for all cases). The psychiatric
267 disease link is consistent with previous work [7] and is discussed below; however, the enrichment among the
268 low-variance genes is weaker.

269 Discussion

270 Using large publicly available data sets allowed us to probe the landscape of transcriptional variance in humans.
271 We find a broadly similar pattern of transcriptional variance, evidenced by the high correlations between gene
272 expression variance across most studies, consistent with measurements of expression variance in single cells
273 and in populations of cells for various tissues [6,16,29]. Leveraging this similarity between gene expression
274 variance across tissues and contexts, we developed a multivariate strategy to create a single rank of expression
275 variance, which allowed us to order almost 13k genes (~65% of the genes expressed in humans) according to
276 their transcriptional variance. Using this rank, we were able to study the general properties associated with
277 high- and low-variance genes as well as factors driving variation in variance across genes.

278 Some differences in gene expression variance were driven by technical aspects of gene expression measure-
279 ment (with data derived from large consortia showing more similar patterns of variance across genes), and by
280 tissue (with studies using the same tissues also showing higher similarities). This suggests that careful consid-
281 eration of sample sizes and experimental design are fundamental to the study of gene expression variance, and
282 the usual small samples of RNA-seq studies might be underpowered for the study of this particular aspect of
283 gene expression. However, both the effects of study origin and tissue were small, and the largest drivers of dif-
284 ferences across studies were idiosyncratic differences related to single data sets, with tissues known to have
285 divergent gene expression patterns (i.e., bone marrow, blood, testis, and platelets) also showing the largest
286 differences in gene expression variance. Understanding the consequences of these differences in variance for
287 specific tissues is still an open field. It is clear, however, that differences in variance are informative beyond
288 the differences in mean expression. Even after we account for differences in mean expression, differences in
289 gene expression variance carry information about tissue origin and function.

290 Functional analyses using GO enrichment indicated a clear link between function and gene expression vari-
291 ance. On the one hand, genes with high transcriptional variance were enriched for biological functions related
292 to response to environmental stimuli, such as immune function and tissue reconstruction. On the other hand,
293 low-variance genes were enriched for basic cell functions, (e.g., RNA processing, translation, DNA methyl-
294 ation, and cell duplication). These results are consistent with previous analyses of gene expression variance on
295 a tissue-by-tissue basis [16]. This pattern of enrichment is also observed when we look at enrichment for high-
296 or low-variance genes within the genes associated with each term in the GO hierarchy. Basic cell function
297 terms are enriched for low-variance genes, and terms involved in response to external stimulus are enriched
298 for high-variance genes.

299 While indirect, all these patterns point to a selective structuring of gene expression variance. Stabilizing and
300 purifying selection are consistent: genes expected to be under strong stabilizing selection, those linked with
301 fundamental baseline biological processes, are indeed overrepresented in the least variable genes. These same

302 genes are also expected to be under strong purifying selection and to show low levels of polymorphisms, which
303 we observe. Likewise, genes whose function is constrained by myriad interactions with several other genes,
304 those with high connectivity, are less variable. Furthermore, genes involved with direct interaction with the
305 environment, which must change their pattern of expression depending on external conditions, are expected
306 to be more variable, and again we see a strong enrichment of environmentally responsive genes among the
307 most variable. Given this strong link between function and variance, it is not surprising that the gene variance
308 ranking is similar across data sets.

309 One interesting aspect of the GO term analysis shown in fig. 3 and fig. 4 is that there is no GO biological process
310 term associated with enrichment for intermediate variance genes: the low-entropy terms have either positive
311 or negative skew, never zero skew. In other words, there is no annotated biological process for which the as-
312 sociated genes are kept at some intermediary level of variation. For the GO terms we used, either there is no
313 relation between the transcriptional variance and the biological process, or there is a strong bias toward high
314 or low-variance genes. This suggests that selective shaping of gene expression has two modes, corresponding
315 with (1) biological processes under strong stabilizing selection (i.e., variance-reducing selection) or (2) biolog-
316 ical processes under disruptive selection (i.e., variance-increasing selection). In short, we find strong support
317 for the idea that there are genes with consistently more (or less) variable expression levels, and that these
318 differences in variance are the result of different patterns of selection.

319 Following Alemu et al. [16], we observe that epigenetic signatures of gene regulation, such as enhancer histone
320 marks, make up a higher proportion of the surrounding genomic regions of genes that exhibit higher variance
321 in expression. In contrast, an accumulation of strong promoter elements and overall transcriptional activation
322 is associated with genes with lower expression variance. These results suggest the presence of distinct modes
323 of regulation for genes with high vs. low variance. Combined, the differences in the types of genomic regula-
324 tory features surrounding the high- and low-variance genes and their distinct functional annotations suggest
325 different mechanisms of regulation of their gene expression variance [16]. This heterogeneity could lead to
326 detectable differences in selection signatures between distal regulatory elements and promoters depending
327 on the transcriptional variance. This heterogeneity in regulation for high and low-variance genes suggests
328 that important biological information has been overlooked given the focus that the field has placed on under-
329 standing gene expression robustness, in the sense of reducing variation [30–33]. For example, Siegal and Leu
330 [30] provide several examples of known regulatory mechanisms for reducing gene expression variance, but
331 no examples for the maintenance of high gene expression variance. We posit that it should be possible to go
332 beyond the usual characterization of mechanisms of gene expression robustness, in the sense of reducing vari-
333 ation, and to explore mechanisms for the *robustness of plasticity*, that is, the maintenance of high levels of gene
334 expression variation given environmental cues.

335 Given the broad consistency of gene expression variance in healthy tissues, a natural question is how do these
336 well-regulated levels of variation behave in disease conditions. We find some suggestive links between tissue-
337 specific variance ranks and disease, but these links need to be further explored using more specific methods.
338 Comparing two HapMap populations, Li et al. [6] showed that gene expression variance was similar in both
339 populations and that high-variance genes were enriched for genes related to HIV susceptibility, consistent
340 with our observation of enrichment for immune-related genes among those with more variable expression.
341 In a case-control experiment, Mar et al. [7] showed that expression variance was related to disease status in
342 Schizophrenia and Parkinson's disease patients, with altered genes being non-randomly distributed across sig-
343 naling networks. These authors also find a link between gene network connectivity and expression variance, in
344 agreement with the effect we find using the gene expression variance rank. The pattern of variance alteration
345 differed across diseases, with Parkinson's patients showing increased expression variance, and Schizophrenia
346 patients showing more constrained patterns of expression. The authors hypothesize that the reduced variance
347 in Schizophrenia patients reduces the robustness of their gene expression networks, what we refer to as a loss
348 of plasticity. This suggests that several types of shifts in gene expression variation are possible, each with dif-
349 ferent outcomes. We highlight three distinct possibilities: First, low-variance genes, under strong stabilizing
350 selection, could become more variable under stress, indicating a reduced capacity for maintaining homeosta-
351 sis. Second, high-variance genes, expected to be reactive to changes in the environment, could become less
352 variable, indicating a reduced capacity to respond to external stimuli. Third, the covariance between different
353 genes could be altered, leading to decoherence between interdependent genes [34]. Any one of these changes
354 in expression variance patterns could have physiological consequences and exploring these differences should
355 be a major part of linking gene expression to cell phenotypes and function (see Hagai et al. [8] for example).
356 Genes are also expected to differ in their capacity to maintain an optimal level of gene expression variance
357 [32]. Variation in robustness is linked to gene regulatory networks and epigenetic gene expression regulation
358 [31,35] and, therefore, should differ across high- and low-variance genes. Our results suggest that the mecha-
359 nisms responsible for maintaining optimal levels of variation in high- and low-variance could differ and that
360 this variability is the result of different patterns of selection.

361 **Methods**

362 **Data sources**

363 We selected 57 human RNA-seq data sets from the public gene expression repositories recount3 [36] and Expression Atlas
364 [37]. We only used data sets with an associated publication, for which raw read count and sample-level metadata were
365 available. Because we are interested in individual-level variation of gene expression, we exclude single-cell studies. Meta-
366 data and details on the included data sets can be found in the supporting information. We use the word "studies" to refer

367 to independent data sets, which could have been generated by the same consortium. For example, the GTEx data are sepa-
 368 rated by tissue, and we refer to each tissue as a separate study. We divide our data sets into three categories depending on
 369 their origin: GTEx, TCGA, and Miscellaneous.

Table 1: Data set source references. Columns show the study ID, with the corresponding tissue in parenthesis, and the source publication.

Study ID	Citation
ADIPOSE_TISSUE (Fat), ADRENAL_GLAND (Adrenal), BLOOD (Blood), BLOOD_VESSEL (Blood_vessel), BONE_MARROW (Marrow), BRAIN (Neuron), HEART (Heart), BREAST (Breast), SALIVARY_GLAND (Salivary), COLON (Colon), LIVER (Liver), NERVE (Neuron), LUNG (Lung), PANCREAS (Pancreas), MUSCLE (Muscle), THYROID (Thyroid), OVARY (Ovary), STOMACH (Stomach), ESOPHAGUS (Esophagus), SPLEEN (Spleen), PROSTATE (Prostate), SKIN (Skin), PITUITARY (Pituitary), TESTIS (Testis)	The GTEx Consortium, 2020 - [38]
LUSC (Lung), STAD (Stomach), COAD (Colon), LUAD (Lung), BRCA (Breast), KIRC (Kidney), KIRP (Kidney), LIHC (Liver), THCA (Thyroid), PRAD (Prostate), UCEC (Uterus)	The Cancer Genome Atlas Research Network et al., 2013 - [17]
SRP150552 (Blood)	Altman et al., 2019 - [39]
SRP101294 (Fat)	Armenise et al., 2017 - [40]
SRP057500 (Platelets)	Best et al., 2015 - [41]
SRP051848 (Immune)	Breen et al., 2015 - [42]
SRP187978 (Liver)	Çalışkan et al., 2019 - [43]
E-ENAD-34 (Immune)	Chen et al., 2016 - [44]
SRP059039 (Blood)	DeBerg et al., 2018 - [45]
SRP174638 (Immune)	Dufort et al., 2019 - [46]
E-GEOD-57945 (Colon)	Haberman et al., 2014 - [47]
SRP162654 (Blood)	Harrison et al., 2019 - [48]
SRP095272 (Blood)	Jadhav et al., 2019 - [49]
SRP102999 (Blood)	Kuan et al., 2017 - [50]
SRP145493 (Immune)	Kuan et al., 2019 - [51]
E-GEUV-1 (Immune)	Lappalainen et al., 2013 - [52]
SRP035988 (Skin)	Li et al., 2014 - [53]
SRP192714 (Blood)	Michlmayr et al., 2020 - [54]
ERP115010 (Blood)	Roe et al., 2020 - [55]
E-ENAD-33 (Neuron)	Schwartzentruber et al., 2018 - [56]

Study ID	Citation
SRP181886 (Neuron)	Srinivasan et al., 2020 - [57]
SRP098758 (Blood)	Suliman et al., 2018 - [58]
SRP032775 (Blood)	Tran et al., 2016 - [59]
SRP069212 (Liver)	Yang et al., 2017 - [60]

370 Processing pipeline

371 We use a standardized pipeline to measure gene expression variance while removing extraneous sources of variation. Be-
372 cause we are interested in variation under non-perturbed conditions, data from case-control studies were filtered to keep
373 only control samples. Technical replicates were summed. For each study, we filtered genes that did not achieve a minimum
374 of 1 count per million (cpm) reads in all samples and a mean of 5 cpm reads across samples. To account for library size and
375 the mean-variance relation in RNA-seq count data, we applied a variance stabilizing transformation implemented in the
376 function `vst` from the DESeq2 R package [61] to the genes passing the read-count filters. This mean-variance correction
377 was verified by plotting mean-variance relations before and after correction, and these plots can be seen in the support-
378 ing information (SI appendix 1). Various technical covariates (like experimental batch, sex, etc.) were manually curated
379 from the metadata associated with each study and accounted for using an independent linear fixed-effects model for each
380 study. A list of covariates used for each study is available in the supporting information (SI data 1). Outlier individuals in
381 the residual distribution were removed using a robust Principal Component Analysis (PCA) approach of automatic outlier
382 detection described in [62]. This procedure first estimates robust Principal Components for each study and then measures
383 the Mahalanobis distance between each sample and the robust mean. Samples that are above the 0.99 percentile in Maha-
384 lanobis distance to the mean are marked as outliers and removed. We verify that the batch effect correction and outlier
385 removal are reasonable by using PCA scatter plots after each step of the pipeline to check the result for residual problems
386 like groupings or other artifacts. These PCA plots before and after batch correction and outlier removal are also included
387 in SI appendix 1. After all sample filtering, the mean sample size we used for each data set was 390, with a median of 251,
388 and ranged from 12 to 2931 samples. Gene expression standard deviations (SDs) are measured as the residual standard
389 deviations after fixed effect correction and outlier removal. We choose standard deviation as a measure of variation to
390 have a statistic on a linear scale, and we do not use the coefficient of variation because we have already corrected for mean
391 differences and for the mean-variance relation inherent to RNA-seq count data [1]. The full annotated pipeline is available
392 on GitHub at github.com/ayroles-lab/expressionVariance-code.

393 Correlations in transcriptional variance

394 We assessed the similarity in gene expression variance across studies by using a across-study Spearman correlation matrix
395 of the measured SDs. Only genes present in all studies were used to calculate the Spearman correlation matrix, ~4200 genes
396 in total. Using Spearman correlations avoids problems related to overall scaling or coverage differences, and allows us to
397 assess if the same genes are usually more or less variable across studies. To investigate the factors involved in determining

398 correlations between studies, we used a Bayesian varying effects model to investigate the effect of study origin and tissue
399 on the correlations across studies. This model is designed to take the non-independent nature of a set of correlations into
400 account when modeling the correlation between gene expression SDs. This is accomplished by adding a per-study random
401 effect, see [20] for details. The Fisher z-transformed Spearman correlations across studies ($z(\rho_{ij})$) are modeled as:

$$\begin{aligned}z(\rho_{ij}) &\sim N(\mu_{ij}, \sigma) \\ \mu_{ij} &= \mu_0 + \alpha_i + \alpha_j + \beta X \\ \alpha_i &\sim N(0, \sigma_\alpha)\end{aligned}$$

402 The α_i terms account for the non-independence between the pairs of correlations and estimate the idiosyncratic contri-
403 bution of each study to all the correlations it is involved in. The fixed effects encoded in the design matrix X measure
404 the effects of tissue congruence and study-origin congruence. We also explored a version of this model that included the
405 effect of sample size on the pairwise correlations, but sample size did not have a relevant effect and so was dropped in the
406 final model. All fixed effect parameters (β) and per-study parameters (α_i) receive weakly informative normal priors with
407 a standard deviation of one quarter. For the overall variance (σ) we use a unit exponential prior, and for the intercept
408 (μ_0) a unit normal prior. This model was fit in Stan [63] via the *rethinking* R package [64], using eight chains, with 4000
409 warm-up iterations and 2000 sampling iterations per chain. Convergence was assessed using R-hat diagnostics [65], and
410 we observed no warnings or divergent transitions.

411 **Gene expression SD rank:** Given that most of the variation in the Spearman correlation across studies is explained by a
412 single principal component (PC1 accounts for 62% of the variation in the across-study Spearman correlation matrix, while
413 PC2 accounts for only 5%; see SI fig. 5), we use the ranked projections of gene expression SDs in this principal component
414 (PC1) to create an across-study rank of gene variation. The higher the rank, the higher the expression SD of a given gene.
415 Genes that were expressed in at least 50% of the studies were included in the rank. To project a particular gene onto the
416 PC1 of the across-study correlation matrix, we impute missing values using a PCA-based imputation [66]. The imputation
417 procedure has minimal effect on the ranking and imputing missing SD ranks at the beginning or at the end of the ranks
418 produces similar results. We also create a tissue-specific variance ranking, using the same ranking procedure but joining
419 studies done in the same tissue type. For this tissue-level ranking, we only use genes that are expressed in all studies of a
420 given tissue, and in this case, no imputation is required. For tissues that are represented by a single study, we use the SD
421 ranking for that study as the tissue rank.

422 **Gene expression mean rank:** We also use the same strategy to create a mean gene expression rank, repeating the process
423 but using mean expression instead of standard deviation. All ranks are available in the supporting information.

424 Gene level statistics

425 **Genetic variation:** Genetic variation measures were obtained from the PopHuman project, which provides a comprehen-
426 sive set of genomic information for human populations derived from the 1000 Genomes Project. Gene-level metrics were
427 used when available. If only window-based metrics are available, we assembled gene-level information from 10 kb window

428 tracks where each window that overlaps with a given gene was assigned to the gene and the mean metric value is reported.
429 In parallel, we use the PopHumanScan data set, which expands PopHuman by compiling and annotating regions under
430 selection. Similarly, we used gene-level information when possible, and for tracks with only window-based metrics, gene-
431 level information was assembled from the 10 kb windows using the same assignment method described above. Nucleotide
432 diversity (π), the average pairwise number of differences per site among the chromosomes in a population [67], provides
433 insight into the genetic diversity within a population, in this case, the CEU population within 1000 genomes.

434 **Gene connectivity:** For each data set, we calculated the average weighted connectivity for all genes by creating a fully
435 connected gene-by-gene graph in which each edge is weighted by the Spearman correlation between gene expression levels
436 across samples. We then trimmed this graph by keeping only edges for which the Spearman correlation is significant at
437 a BH false discovery rate of 1%. In this trimmed network, we then took the average of the Spearman correlation of all
438 remaining edges for each gene. So, for each study, we have a measure of the average correlation of each gene with every
439 other gene. The average connectivity for each gene is the average across all studies in which that gene is expressed.

440 **Cross-tissue vs. tissue-level chromatin states:** We use the universal [68] and tissue-specific [69] ChromHMM [25] chro-
441 matin states to compare the non-overlapping genome segmentation to cross-tissue and tissue-level gene expression vari-
442 ance metrics. We use the proportion of the gene regions (gene \pm 10 kb) made up of each of the ChromHMM chromatin
443 states.

444 **Correlations:** We use the ppcor R package v1.1 [70] to run the pairwise partial Spearman correlations between gene-level
445 statistics and the gene expression variance rank while controlling for the mean expression rank. P-values are corrected
446 using the Benjamini-Hochberg procedure and comparisons with an adjusted $p < 0.05$ are considered significant.

447 **Gene function assessment**

448 **GO term enrichment:** All gene ontology (GO) analyses were done using the clusterProfiler R package v4.2.2 [71] and the
449 Org.Hs.eg.db database package v3.14.0 [72]. GO and all further enrichment analyses used the hypergeometric test to assess
450 the significance of the enrichment.

451 **Environmentally responsive genes:** We used the list of environmentally responsive genes available in the supporting
452 information from Lea et al. [23]. When we overlapped the list of LCL-expressed genes with our list of ranked genes, we
453 retained 9282 genes in the cross-tissue analyses presented in the main text, and 5574 genes in the blood-specific analyses
454 presented in SI table 1. We used Fisher's exact tests to ask whether high-variance genes were more likely to be respon-
455 sive to >0 environments relative to all genes not included in the low-variance category (and vice versa for low-variance
456 genes). We also used Fisher's exact tests followed by Benjamini-Hochberg false discovery rate correction to ask whether
457 high-variance genes were more likely to be responsive to each individual environment relative to all genes not included
458 in the low-variance category (and vice versa for low-variance genes). Finally, we used generalized linear models with a
459 Poisson error structure to test for an effect of gene category (high-variance, low-variance, or neither) on the number of
460 environments that a gene responded to (ranging from 0-11).

461 **Housekeeping genes:** Human housekeeping genes were identified as genes that are expressed with low variance in all 52

462 human cell and tissue types, assessed in over 10,000 samples [21]. We test for enrichment of housekeeping genes in the
463 genes within the highest and lowest 5% of gene expression variance rank.

464 **Probability of being loss-of-function intolerant (pLI):** Genes that are likely haploinsufficient (i.e., intolerant of heterozy-
465 gous loss-of-function variants) were detected as those with fewer than expected protein-truncating variants (PTVs) in
466 ExAC [73]. We use genes with a pLI > 0.9 to test for the enrichment of loss-of-function intolerant genes in the genes exhibit-
467 ing the highest and lowest 5% gene expression variance estimates.

468 **Secreted genes:** We use The Protein Atlas [74] to extract information on which proteins are secreted [75] and test for en-
469 richment of genes with secreted products in the genes within the highest and lowest 5% of gene expression variance rank.

470 **Immediate early genes (IEGs):** Human IEGs were curated from the literature in [76] as genes that respond to experimental
471 stimulation through up-regulation within the first 60 minutes of the experiment. We use the hypergeometric test to assess
472 the significance of the enrichment. Immediate early genes (IEGs): Human IEGs were curated from the literature in [76] as
473 genes that respond to experimental stimulation through up-regulation within the first 60 minutes of the experiment.

474 **Disease annotations:** We use the gene annotations for involvement with diseases provided by the supporting information
475 Table S2 from Zhang et al. [28] and test for enrichment for disease annotations in the genes within the highest and lowest
476 5% of gene expression variance rank.

477 **Code availability**

478 Code for reproducing all analyses and figures, along with a walk-through, is available at [github.com/ayroles-](https://github.com/ayroles-lab/expressionVariance-code)
479 [lab/expressionVariance-code](https://github.com/ayroles-lab/expressionVariance-code).

480 **Supporting information**

481 Supporting information is available at github.com/ayroles-lab/expressionVariance-manuscript.

- 482 1. SI figure 1 - Modeling the correlations between transcriptional variance across studies.
- 483 2. SI figure 2 - GO enrichment analysis of the most and least variable genes.
- 484 3. SI figure 3 - Across-study and tissue-specific gene expression variance and mean correlations with non-overlapping
485 chromatin states through ChromHMM.
- 486 4. SI figure 4 - Proportion of gene regions made up of ChromHMM chromatin states for genes in the top and bottom 5%
487 of the across-study mean rank metric.
- 488 5. SI figure 5 - Scree plot showing variance explained by each PC of the across-study Spearman correlation matrix of
489 gene expression standard deviations.
- 490 6. SI table 1 - Variance and mean rank metrics and the corresponding ChromHMM annotations used.
- 491 7. SI table 2 - Enrichment analysis of environmentally responsive genes in LCLs.
- 492 8. SI appendix 1 - Diagnostics plots for processing pipeline.

- 493 9. SI data 1 - Study metadata - Metadata file describing the data used in the study as well as some intermediate process-
494 ing information.
- 495 10. SI data 2 - Study gene lists - List of genes included in each study after filtering.
- 496 11. SI data 3 - Gene expression means and standard deviations - Tables with final calculated means and standard devia-
497 tions.
- 498 12. SI data 4 - Gene ranks - Gene expression mean and variance ranks, across-study and tissue-specific.
- 499 13. SI data 5 - GO enrichment - Combined table describing gene ontology enrichment in the top 5% and bottom 5% of
500 genes as ranked by variance.

501 Author Contributions

502 **Conceptualization:** S.W., D.M., L.P., and J.A. **Analysis:** S.W., D.M., K.G., and A.L. **Draft:** D.M. **Review and Editing:** S.W.,
503 D.M., K.G., L.P., A.L., and J.A. **Funding Acquisition:** S.W., D.M., K.G., L.P., A.L., and J.A.

504 Acknowledgments

505 We thank all members of the Ayroles lab for their support. We thank Noah Rose and Cara Weisman for their thoughtful com-
506 ments. We thank Pedro Madrigal for help with the Expression Atlas interface. S.W. is supported by the National Science
507 Foundation Graduate Research Fellowship Program (DGE-2039656). D.M. is funded by a fellowship from the Princeton Pres-
508 idential Postdoctoral Research Fellows Program. K.G. is funded by National Institutes of Health (NIH) grant F32ES034668.
509 L.P. was funded by a Long-Term Postdoctoral Fellowship from the Human Frontiers Science Program and is funded by
510 the Max Planck Society. J.A. is funded by grants from the NIH: National Institute of Environmental Health Sciences (R01-
511 ES029929) and National Institute of General Medical Sciences (NIGMS) (R35GM124881). This study was supported in part by
512 the Lewis-Sigler Institute for Integrative Genomics at Princeton University. We also acknowledge that the work reported in
513 this paper was substantially performed using the Princeton Research Computing resources at Princeton University which
514 is a consortium of groups led by the Princeton Institute for Computational Science and Engineering (PICSciE) and Office of
515 Information Technology's Research Computing. A.L. is supported by the Canadian Institute for Advanced Research Global
516 Scholars Program, the Searle Scholars Program, and through the NIH/NIGMS (R35GM147267).

517 References

- 518 1. Jong TV de, Moshkin YM, Guryev V. Gene expression variability: The other dimension in transcriptome analysis.
519 *Physiol Genomics*. 2019 May;51(5):145–58.
- 520 2. Fraser HB, Hirsh AE, Giaever G, Kumm J, Eisen MB. Noise minimization in eukaryotic gene expression. *PLoS Biol*.
521 2004 Jun;2(6):e137.
- 522 3. Wang Z, Zhang J. Impact of gene expression noise on organismal fitness and the efficacy of natural selection. *Proc*
523 *Natl Acad Sci U S A*. 2011 Apr;108(16):E67–76.

- 524 4. Hansen TF, Pélabon C. Evolvability: A Quantitative-Genetics perspective. *Annu Rev Ecol Evol Syst.* 2021
525 Nov;52(1):153–75.
- 526 5. Bruijning M, Metcalf CJE, Jongejans E, Ayroles JF. The evolution of variance control. *Trends Ecol Evol.* 2020
527 Jan;35(1):22–33.
- 528 6. Li J, Liu Y, Kim T, Min R, Zhang Z. Gene expression variability within and between human populations and impli-
529 cations toward disease susceptibility. *PLoS Comput Biol.* 2010 Aug;6(8).
- 530 7. Mar JC, Matigian NA, Mackay-Sim A, Mellick GD, Sue CM, Silburn PA, et al. Variance of gene expression identifies
531 altered network constraints in neurological disease. *PLoS Genet.* 2011 Aug;7(8):e1002207.
- 532 8. Hagai T, Chen X, Miragaia RJ, Rostom R, Gomes T, Kunowska N, et al. Gene expression variability across cells and
533 species shapes innate immunity. *Nature.* 2018 Nov;563(7730):197–202.
- 534 9. Houle D. How should we explain variation in the genetic variance of traits? *Genetica.* 1998;102-103(1-6):241–53.
535
- 536 10. Hansen TF. Epigenetics: Adaptation or contingency. In: Benedikt Hallgrímsson BKH, editor. *Epigenetics: Linking
537 genotype and phenotype in development and evolution.* University of California press Berkeley, CA; 2011. p. 357–76.
- 538 11. Schmutzer M, Wagner A. Gene expression noise can promote the fixation of beneficial mutations in fluctuating
539 environments. *PLoS Comput Biol.* 2020 Oct;16(10):e1007727.
- 540 12. Pettersson ME, Nelson RM, Carlborg O. Selection on variance-controlling genes: Adaptability or stability. *Evolu-
541 tion.* 2012 Dec;66(12):3945–9.
- 542 13. Wagner GP, Booth G, Bagheri-Chaichian H. A POPULATION GENETIC THEORY OF CANALIZATION. *Evolution.* 1997
543 Apr;51(2):329–47.
- 544 14. Pavlicev M, Hansen TF. Genotype-Phenotype Maps Maximizing Evolvability: Modularity Revisited. *Evol Biol.* 2011
545 Dec;38(4):371–89.
- 546 15. GTEx Consortium. Genetic effects on gene expression across human tissues. *Nature.* 2017 Oct;550(7675):204–13.
547
- 548 16. Alemu EY, Carl JW Jr, Corrada Bravo H, Hannenhalli S. Determinants of expression variability. *Nucleic Acids Res.*
549 2014 Apr;42(6):3503–14.
- 550 17. Cancer Genome Atlas Research Network, Weinstein JN, Collisson EA, Mills GB, Shaw KRM, Ozenberger BA, et al.
551 The cancer genome atlas Pan-Cancer analysis project. *Nat Genet.* 2013 Oct;45(10):1113–20.
- 552 18. Gower JC. Some distance properties of latent root and vector methods used in multivariate analysis. *Biometrika.*
553 1966 Dec;53(3-4):325–38.
- 554 19. Dias FS, Betancourt M, Rodríguez-González PM, Borda-de-Água L. Analysing the distance decay of community sim-
555 ilarity in river networks using bayesian methods. *Sci Rep.* 2021 Nov;11(1):21660.
- 556 20. Dias FS, Betancourt M, Rodríguez-González PM, Borda-de-Água L. BetaBayes—A bayesian approach for comparing
557 ecological communities. *Diversity.* 2022 Oct;14(10):858.

- 558 21. Hounkpe BW, Chenou F, Lima F de, De Paula EV. HRT atlas v1.0 database: Redefining human and mouse house-
559 keeping genes and candidate reference transcripts by mining massive RNA-seq datasets. *Nucleic Acids Res.* 2020
Jul;49(D1):D947-55.
- 560 22. Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, et al. Analysis of protein-coding genetic vari-
561 ation in 60,706 humans. *Nature.* 2016;536(7616):285-91.
- 562 23. Lea AJ, Peng J, Ayroles JF. Diverse environmental perturbations reveal the evolution and context-dependency of
563 genetic effects on gene expression levels. *Genome Res.* 2022 Oct;32(10):1826-39.
- 564 24. Mähler N, Wang J, Terebieniec BK, Ingvarsson PK, Street NR, Hvidsten TR. Gene co-expression network connectiv-
565 ity is an important determinant of selective constraint. *PLoS Genet.* 2017 Apr;13(4):e1006402.
- 566 25. Ernst J, Kellis M. ChromHMM: Automating chromatin-state discovery and characterization. *Nature methods.*
567 2012;9(3):215-6.
- 568 26. Rogatsky I, Adelman K. Preparing the first responders: Building the inflammatory transcriptome from the ground
569 up. *Mol Cell.* 2014 Apr;54(2):245-54.
- 570 27. Bahrami S, Drabløs F. Gene regulation in the immediate-early response process. *Adv Biol Regul.* 2016 Sep;62:37-49.
571
- 572 28. Zhang Y, Quick C, Yu K, Barbeira A, GTEx Consortium, Luca F, et al. PTWAS: Investigating tissue-relevant causal
573 molecular mechanisms of complex traits using probabilistic TWAS analysis. *Genome Biol.* 2020 Sep;21(1):232.
- 574 29. Dong D, Shao X, Deng N, Zhang Z. Gene expression variations are predictive for stochastic noise. *Nucleic Acids Res.*
575 2011 Jan;39(2):403-13.
- 576 30. Siegal ML, Leu JY. On the nature and evolutionary impact of phenotypic robustness mechanisms. *Annu Rev Ecol
577 Evol Syst.* 2014 Nov;45:496-517.
- 578 31. Payne JL, Wagner A. Mechanisms of mutational robustness in transcriptional regulation. *Front Genet.* 2015
579 Oct;6(October):1-0.
- 580 32. Macneil LT, Walhout AJM. Gene regulatory networks and the role of robustness and stochasticity in the control of
581 gene expression. *Genome Res.* 2011 May;21(5):645-57.
- 582 33. Denby CM, Im JH, Yu RC, Pesce CG, Brem RB. Negative feedback confers mutational robustness in yeast transcrip-
583 tion factor regulation. *Proc Natl Acad Sci U S A.* 2012 Mar;109(10):3874-8.
- 584 34. Lea A, Subramaniam M, Ko A, Lehtimäki T, Raitoharju E, Kähönen M, et al. Genetic and environmental perturba-
585 tions lead to regulatory decoherence. *Elife.* 2019 Mar;8.
- 586 35. Chalancon G, Ravarani CNJ, Balaji S, Martinez-Arias A, Aravind L, Jothi R, et al. Interplay between gene expression
587 noise and regulatory network architecture. *Trends Genet.* 2012 May;28(5):221-32.
- 588 36. Wilks C, Zheng SC, Chen FY, Charles R, Solomon B, Ling JP, et al. recount3: Summaries and queries for large-scale
589 RNA-seq expression and splicing. *Genome Biol.* 2021 Nov;22(1):323.
- 590 37. Papatheodorou I, Moreno P, Manning J, Fuentes AMP, George N, Fexova S, et al. Expression atlas update: From
tissues to single cells. *Nucleic Acids Res.* 2020 Jan;48(D1):D77-83.

591

592 38. GTEx Consortium. The GTEx consortium atlas of genetic regulatory effects across human tissues. *Science*. 2020
593 Sep;369(6509):1318–30.

594 39. Altman MC, Gill MA, Whalen E, Babineau DC, Shao B, Liu AH, et al. Transcriptome networks identify mechanisms
595 of viral and nonviral asthma exacerbations in children. *Nat Immunol*. 2019 May;20(5):637–51.

596 40. Armenise C, Lefebvre G, Carayol J, Bonnel S, Bolton J, Di Cara A, et al. Transcriptome profiling from adipose tissue
during a low-calorie diet reveals predictors of weight and glycemic outcomes in obese, nondiabetic subjects. *Am J
597 Clin Nutr*. 2017 Sep;106(3):736–46.

598 41. Best MG, Sol N, Kooi I, Tannous J, Westerman BA, Rustenburg F, et al. RNA-Seq of Tumor-Educated platelets enables
599 Blood-Based Pan-Cancer, multiclass, and molecular pathway cancer diagnostics. *Cancer Cell*. 2015 Nov;28(5):666–
76.

600 42. Breen MS, Maihofer AX, Glatt SJ, Tylee DS, Chandler SD, Tsuang MT, et al. Gene networks specific for innate im-
601 munity define post-traumatic stress disorder. *Mol Psychiatry*. 2015 Dec;20(12):1538–45.

602 43. Çalışkan M, Manduchi E, Rao HS, Segert JA, Beltrame MH, Trizzino M, et al. Genetic and epigenetic fine mapping
603 of complex trait associated loci in the human liver. *Am J Hum Genet*. 2019 Jul;105(1):89–107.

604 44. Chen L, Ge B, Casale FP, Vasquez L, Kwan T, Garrido-Martín D, et al. Genetic drivers of epigenetic and transcrip-
605 tional variation in human immune cells. *Cell*. 2016 Nov;167(5):1398–1414.e24.

606 45. DeBerg HA, Zaidi MB, Altman MC, Khaenam P, Gersuk VH, Campos FD, et al. Shared and organism-specific
607 host responses to childhood diarrheal diseases revealed by whole blood transcript profiling. *PLoS One*. 2018
Jan;13(1):e0192082.

608 46. Dufort MJ, Greenbaum CJ, Speake C, Linsley PS. Cell type-specific immune phenotypes predict loss of insulin se-
609 cretion in new-onset type 1 diabetes. *JCI Insight*. 2019 Feb;4(4).

610 47. Haberman Y, Tickle TL, Dexheimer PJ, Kim MO, Tang D, Karns R, et al. Pediatric crohn disease patients exhibit
611 specific ileal transcriptome and microbiome signature. *J Clin Invest*. 2014 Aug;124(8):3617–33.

612 48. Harrison GF, Sanz J, Boulais J, Mina MJ, Grenier JC, Leng Y, et al. Natural selection contributed to immunological
613 differences between hunter-gatherers and agriculturalists. *Nat Ecol Evol*. 2019 Aug;3(8):1253–64.

614 49. Jadhav B, Monajemi R, Gagalova KK, Ho D, Draisma HHM, Wiel MA van de, et al. RNA-Seq in 296 phased trios
615 provides a high-resolution map of genomic imprinting. *BMC Biol*. 2019 Jun;17(1):50.

616 50. Kuan PF, Waszczuk MA, Kotov R, Clouston S, Yang X, Singh PK, et al. Gene expression associated with PTSD in
617 world trade center responders: An RNA sequencing study. *Transl Psychiatry*. 2017 Dec;7(12):1297.

618 51. Kuan PF, Yang X, Clouston S, Ren X, Kotov R, Waszczuk M, et al. Cell type-specific gene expression patterns asso-
619 ciated with posttraumatic stress disorder in world trade center responders. *Transl Psychiatry*. 2019 Jan;9(1):1.

620 52. Lappalainen T, Sammeth M, Friedländer MR, Hoen PAC 't, Monlong J, Rivas MA, et al. Transcriptome and genome
621 sequencing uncovers functional variation in humans. *Nature*. 2013 Sep;501(7468):506–11.

- 622 53. Li B, Tsoi LC, Swindell WR, Gudjonsson JE, Tejasvi T, Johnston A, et al. Transcriptome analysis of psoriasis in a large
623 case-control sample: RNA-seq provides insights into disease mechanisms. *J Invest Dermatol.* 2014 Jul;134(7):1828-
624 38.
- 624 54. Michlmayr D, Kim EY, Rahman AH, Raghunathan R, Kim-Schulze S, Che Y, et al. Comprehensive immunoprofiling
625 of pediatric zika reveals key role for monocytes in the acute phase and no effect of prior dengue virus infection.
626 *Cell Rep.* 2020 Apr;31(4):107569.
- 626 55. Roe J, Venturini C, Gupta RK, Gurry C, Chain BM, Sun Y, et al. Blood transcriptomic stratification of short-term risk
627 in contacts of tuberculosis. *Clin Infect Dis.* 2020 Feb;70(5):731-7.
- 628 56. Schwartzentruber J, Foskolou S, Kilpinen H, Rodrigues J, Alasoo K, Knights AJ, et al. Molecular and functional vari-
629 ation in iPSC-derived sensory neurons. *Nat Genet.* 2018 Jan;50(1):54-61.
- 630 57. Srinivasan K, Friedman BA, Etxeberria A, Huntley MA, Brug MP van der, Foreman O, et al. Alzheimer's patient
631 microglia exhibit enhanced aging and unique transcriptional activation. *Cell Rep.* 2020 Jun;31(13).
- 632 58. Suliman S, Thompson EG, Sutherland J, Weiner J 3rd, Ota MOC, Shankar S, et al. Four-Gene Pan-African blood
633 signature predicts progression to tuberculosis. *Am J Respir Crit Care Med.* 2018 May;197(9):1198-208.
- 634 59. Tran TM, Jones MB, Ongoiba A, Bijker EM, Schats R, Venepally P, et al. Transcriptomic evidence for modulation of
635 host inflammatory responses during febrile plasmodium falciparum malaria. *Sci Rep.* 2016 Aug;6:31291.
- 636 60. Yang Y, Chen L, Gu J, Zhang H, Yuan J, Lian Q, et al. Recurrently deregulated lncRNAs in hepatocellular carcinoma.
637 *Nat Commun.* 2017 Feb;8:14421.
- 638 61. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2.
639 *Genome Biol.* 2014;15(12):550.
- 640 62. Chen X, Zhang B, Wang T, Bonni A, Zhao G. Robust principal component analysis for accurate outlier sample de-
641 tection in RNA-Seq data. *BMC Bioinformatics.* 2020 Jun;21(1):269.
- 642 63. Carpenter B, Gelman A, Hoffman MD, Lee D, Goodrich B, Betancourt M, et al. Stan: A probabilistic programming
643 language. *J Stat Softw.* 2017;76(1).
- 644 64. McElreath R. *Statistical rethinking: A bayesian course with examples in r and stan.* Chapman; Hall/CRC; 2020.
645
- 646 65. Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A, Rubin DB. *Bayesian data analysis, third edition.* CRC Press;
647 2013.
- 648 66. Husson F, Josse J, Narasimhan B, Robin G. Imputation of mixed data with multilevel singular value decomposition.
649 *J Comput Graph Stat.* 2019 Jul;28(3):552-66.
- 650 67. Nei M, Li WH. Mathematical model for studying genetic variation in terms of restriction endonucleases. *Proc Natl*
651 *Acad Sci U S A.* 1979 Oct;76(10):5269-73.
- 652 68. Vu H, Ernst J. Universal annotation of the human genome through integration of over a thousand epigenomic
653 datasets. *Genome biology.* 2022;23(1):1-37.

- 654 69. Ernst J, Kellis M. Large-scale imputation of epigenomic datasets for systematic annotation of diverse human tissues.
655 Nat Biotechnol. 2015 Apr;33(4):364–76.
- 656 70. Kim S. Ppcor: An r package for a fast calculation to semi-partial correlation coefficients. Communications for sta-
657 tistical applications and methods. 2015;22(6):665.
- 658 71. Wu T, Hu E, Xu S, Chen M, Guo P, Dai Z, et al. clusterProfiler 4.0: A universal enrichment tool for interpreting omics
659 data. Innovation (N Y). 2021 Aug;2(3):100141.
- 660 72. Carlson M. Org.hs.eg.db: Genome wide annotation for human. R package version 3.14.0. 2021.
661
- 662 73. Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, et al. Analysis of protein-coding genetic vari-
663 ation in 60,706 humans. Nature. 2016 Aug;536(7616):285–91.
- 664 74. Uhlén M, Fagerberg L, Hallström BM, Lindskog C, Oksvold P, Mardinoglu A, et al. Tissue-based map of the human
665 proteome. Science. 2015;347(6220):1260419.
- 666 75. Uhlén M, Karlsson MJ, Hober A, Svensson AS, Scheffel J, Kotel D, et al. The human secretome. Science signaling.
667 2019;12(609):eaazo274.
- 668 76. Arner E, Daub CO, Vitting-Seerup K, Andersson R, Lilje B, Drabløs F, et al. Transcribed enhancers lead waves of
669 coordinated transcription in transitioning mammalian cells. Science. 2015 Feb;347(6225):1010–4.