

University of Groningen

## Gene expression variability - the other dimension in transcriptome analysis

de Jong, Tristan V; Moshkin, Yuri M; Guryev, Victor

*Published in:*  
Physiological Genomics

*DOI:*  
[10.1152/physiolgenomics.00128.2018](https://doi.org/10.1152/physiolgenomics.00128.2018)

**IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.**

*Document Version*  
Final author's version (accepted by publisher, after peer review)

*Publication date:*  
2019

[Link to publication in University of Groningen/UMCG research database](#)

*Citation for published version (APA):*

de Jong, T. V., Moshkin, Y. M., & Guryev, V. (2019). Gene expression variability - the other dimension in transcriptome analysis: the other dimension in transcriptome analysis. *Physiological Genomics*, 51(5), 145-158. <https://doi.org/10.1152/physiolgenomics.00128.2018>

### Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.

### Take-down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

*Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.*

1 Article type: Review

2 Title: **Gene expression variability – the other dimension in transcriptome analysis**

3

4 Authors: Tristan V. de Jong<sup>1</sup>, Yuri M. Moshkin<sup>2,3</sup>, Victor Guryev<sup>1</sup>

5

6 Affiliations:

7 <sup>1</sup> European Research Institute for the Biology of Ageing, University of Groningen,  
8 University Medical Centre Groningen, Groningen, The Netherlands.

9 <sup>2</sup> Institute of Cytology and Genetics, Siberian Branch of RAS, Novosibirsk, Russia

10 <sup>3</sup> Institute of Molecular and Cellular Biology, Siberian Branch of RAS, Novosibirsk, Russia

11

12 Corresponding authors:

13 Dr. Yuri M Moshkin; Tel: +41 76 699 42 04; E-mail: yury.moshkin@gmail.com

14 Prof. Victor Guryev; Tel: +31 6 5272 4873; E-mail: v.guryev@umcg.nl

15 **Abstract (215 words)**

16 Transcriptome sequencing is a powerful technique to study molecular changes that underlie  
17 the differences in physiological conditions and disease progression. A typical question that is  
18 posed in such studies is finding genes with significant changes between sample groups. In  
19 this respect expression variability is regarded as a nuisance factor that is primarily of  
20 technical origin and complicates the data analysis. However, it is becoming apparent that  
21 the biological variation in gene expression might be an important molecular phenotype that  
22 can affect physiological parameters.

23 In this review we explore the recent literature on technical and biological variability in gene  
24 expression, sources of expression variability, (epi-) genetic hallmarks and evolutionary  
25 constraints in genes with robust and variable gene expression. We provide an overview of  
26 recent findings on effects of external cues, such as diet and aging on expression variability  
27 and on other biological phenomena that can be linked to it. We discuss metrics and tools  
28 that were developed for quantification of expression variability and highlight the  
29 importance of future studies in this direction.

30 In order to assist the adoption of expression variability analysis, we also provide a detailed  
31 description and computer code, which can easily be utilized by other researchers. We also  
32 provide a re-analysis of recently published data to highlight the value of the analysis  
33 method.

34 **Main text**

35 Affordable sequencing has greatly advanced our understanding of changes in  
36 transcription programs and their relation to diseases. One of the sequencing-enabled  
37 technologies, transcriptome profiling by RNA-sequencing (RNA-seq) is becoming increasingly  
38 popular to study molecular phenotypes. The main advantages of this method, when  
39 compared to hybridization microarray-based approaches, include an increased sensitivity  
40 and larger dynamic range, its ability to detect unannotated transcripts and transcript  
41 isoforms and, importantly, it enables digital quantification (counting) of RNA molecules. As a  
42 result, RNA-seq has the potential to quantify lowly expressed genes, to reveal subtle  
43 changes in gene expression (115), to discover new genes, transcript isoforms and allelic  
44 variants for proteogenomics analysis (53), and, as will be discussed later, digital  
45 quantification of RNAs simplifies statistical analysis of gene expression and interpretation of  
46 its variability.

47 The typical analysis of RNA-sequencing data focuses on the finding of genes that show  
48 differential expression between groups. Such analysis can be done with tools like edgeR (58)  
49 or DEseq2 (52). The results call attention to genes which significantly change with respect to  
50 an average RNA copy number between measurable factors like age, diet, the knock-down/-  
51 out/-in of genes of interest, and so on. Unfortunately, in such analysis, variability in gene  
52 expression is often ignored as it is treated as a nuisance that only diminishes statistical  
53 power. At the same time, gene expression is naturally a stochastic process and in some  
54 cases its fluctuation, rather than the mean RNA copy number, could be significantly  
55 influenced by an experimental factor or a physiological state. Thus, while variations caused  
56 by technical factors can be considered as the true nuisance factor (80), differential  
57 variability in gene expression caused by biological factors might represent a layer of  
58 information on gene regulation just as important as changes in the mean expression levels  
59 (104). In this review we discuss recent studies exploring gene expression fluctuations, their  
60 approach to quantification of expression variability, contribution to understanding of the  
61 principles underlying physiological homeostasis and potential to uncover additional  
62 molecular phenotypes associated with disease.

63 **Sources of variability in gene expression: Poisson - “intrinsic” vs non-Poisson -**  
64 **“extrinsic” gene noise.**



65 The inter-sample differences among transcriptome profiles originate from biological  
 66 events as well as from experimental procedures. The latter represents a source of technical  
 67 noise due to the collection and storage of samples, the isolation of RNA, selection of RNA  
 68 molecules and the preparation of library (92). Library amplification and sequencing might  
 69 also introduce differences depending on instruments, read length and mode of sequencing.  
 70 All these factors have potential to complicate the analysis of biological variability in gene  
 71 expression, especially for large (inter-) national and prospective projects where data is being  
 72 produced using different versions of instruments and/or kits (58). Thus, when studying  
 73 variation in gene expression, it is important to estimate technical variability through  
 74 comparison of technical replicates prepared from the same starting material (111) and  
 75 compare it to the degree of variability seen among biologically different samples (58).

76 Putting technical variability aside, gene noise originates from the stochastic nature of  
 77 chemical reactions driving RNA synthesis (birth) and degradation (death). In a stationary  
 78 state and in the absence of upstream cellular drives, a process of RNA “birth-death” is  
 79 expected to be a stochastic Poisson process (21, 96). This process is described by just two  
 80 kinetic parameters, namely the synthesis ( $\lambda$ ) and degradation ( $\gamma$ ) rates. The expectation  
 81 (mean) and variance of RNA copy number are given by the Poisson rate ( $E[RNA] =$   
 82  $Var[RNA] = \mu$ ) represented by a constant ratio of synthesis to degradation rates:  
 83  $\mu = \frac{\lambda}{\gamma} = \hat{\lambda}$ . Gene expression noise, expressed through a squared coefficient of variation in  
 84 RNA copy number, is reciprocal to the mean of RNA copy number:  $cv^2(RNA) = \frac{Var[RNA]}{E[RNA]^2} =$   
 85  $\mu^{-1}$  (96). Here, we will refer to this as Poisson noise following (21, 66, 96). However, in  
 86 reality gene synthesis is more complex as it is regulated by so-called upstream cellular drives  
 87 (21). Because upstream cellular drives are stochastic themselves, the RNA “birth-death”  
 88 becomes a doubly stochastic (mixed) Poisson process. Consequently, this increases the gene  
 89 expression noise to the amount that is contributed by all upstream drives, which we will  
 90 refer to as non-Poisson noise following (21, 66, 96).

91 For example, promoter switching between active (ON) and inactive (OFF) states acts as  
 92 such a drive (Fig. 1). The probability of the promoter to be in ON state ( $p_{on}$ ) is a Beta-  
 93 distributed random variable, which depends on RNA degradation rate normalized  $\hat{k}_{on} = \frac{k_{on}}{\gamma}$   
 94 and  $\hat{k}_{off} = \frac{k_{off}}{\gamma}$  rates of promoter switching:  $p_{on} \sim Beta(\hat{k}_{on}, \hat{k}_{off})$ . This, in turn, defines the

95 distribution of otherwise constant Poisson rate ( $\mu = \hat{\lambda}p_{\text{on}}$ ) as Beta-Poisson (21, 72). A  
96 convenient property of mixed Poisson distributed random variables is that they allow for  
97 simple derivation of their moments (expectation and variance) just from the moments of  
98 the mixing distribution (44). That is  $E[\text{RNA}] = \hat{\lambda}E[p_{\text{on}}] = \langle \mu \rangle$  and  $Var[\text{RNA}] = \langle \mu \rangle +$   
99  $Var[\mu] = \langle \mu \rangle + \langle \mu \rangle^2 Var[p_{\text{on}}]$ , from where  $cv^2(\text{RNA}) = \langle \mu \rangle^{-1} + cv^2(\mu) = \langle \mu \rangle^{-1} +$   
100  $cv^2(p_{\text{on}})$  (Fig. 1). Thus, the total gene noise sums from Poisson noise ( $\langle \mu \rangle^{-1}$ ) and non-  
101 Poisson noise caused by upstream cellular drive, namely promoter switching ( $cv^2(\mu) =$   
102  $cv^2(p_{\text{on}})$ ).

103 Under limiting conditions of  $\hat{k}_{\text{off}} \gg \hat{k}_{\text{on}}$  and  $\hat{k}_{\text{off}} \gg 1$ , i.e. when a gene is transcribed in  
104 short bursts, the  $p_{\text{on}}$  distribution converges to Gamma ( $p_{\text{on}} \sim \text{Gamma}(\hat{k}_{\text{on}}, \hat{k}_{\text{off}})$ ) and the  
105 resulting distribution of RNA copy number is Gamma-Poisson (also known as Negative-  
106 Binomial)(72). The Gamma-Poisson representation helps understanding of how Poisson and  
107 non-Poisson noise are related to often measured in single cell experiments parameters of  
108 transcription, namely the burst size (a number of molecules synthesized in a burst) and  
109 burst frequency (93). That is because Poisson noise equals to  $\langle \mu \rangle^{-1} = (bf_b)^{-1}$  and non-  
110 Poisson noise is  $cv^2(\mu) = cv^2(p_{\text{on}}) = f_b^{-1}$ , where  $b = \lambda k_{\text{off}}^{-1}$  is a burst size and  $f_b = \hat{k}_{\text{on}}$  is a  
111 burst frequency (21, 72). Thus, non-Poisson noise is inversely related to burst frequency,  
112 which implies that changes in burst frequency are indicative of changes in non-Poisson  
113 noise. For the detailed derivations of various stochastic gene expression models under a  
114 mixed Poisson framework and further theoretical insights we refer to a compelling work by  
115 Dattani and Barahona (21).

116 In essence, the partitioning of the total gene noise into Poisson and non-Poisson,  
117 immediately corresponds to a concept of “intrinsic” and “extrinsic” gene noise (26, 94).  
118 Two-colour reporter gene assays allow for the separation of within cell variations from cell-  
119 to-cell variation in gene expression. In this assay two identical copies of a promoter drive  
120 the expression of reporters: yellow fluorescent protein (YFP) and green fluorescent protein  
121 (GFP). The single-cell fluorescence readout will show different expression levels of YFP and  
122 GFP due to the intrinsically stochastic nature of gene expression. At the same time extrinsic  
123 noise will be related to covariance between expression levels of these two reporters. Hence,  
124 the within cell gene expression fluctuations have been coined as “intrinsic” gene noise,  
125 while cell-to-cell variations were referred to as “extrinsic” gene noise. A total gene noise

126 sums, then, from “intrinsic” and “extrinsic” resulting in identical partitioning of noise as  
127 Poisson and non-Poisson.

128       However, defining gene noise through a combination of “intrinsic” and “extrinsic” noise  
129 has been subjected to criticism. First, it is not clear relative to what within biological system  
130 gene noise is “intrinsic” or “extrinsic” (68). Second, they are conditioned on each other (88).  
131 Indeed, “intrinsic” gene noise, or Poisson noise for that matter, is reciprocal to the mean  
132 gene expression. For the two-state promoter model, i.e. in the presence of upstream  
133 cellular drive caused by promoter fluctuation, the mean gene expression depends on the  
134 probability of the promoter to be in the ON state. Thus, “intrinsic” gene noise is coupled to  
135 upstream cellular drives. Likewise, “extrinsic” gene noise depends on the RNA lifetime  
136 normalized rates of promoter switching between the ON and OFF states. Thus, “extrinsic”  
137 gene noise is conditioned on the characteristic gene state variables (21, 72).

138       Having this in mind and considering that RNA “birth-death” is a doubly stochastic  
139 Poisson process, it makes more sense to stay with Poisson and non-Poisson partitioning of  
140 gene expression noise (21). Accordingly, parameters affecting the gene expression  
141 variability and thus the gene expression noise, could be classified into gene state variables  
142 (kinetic parameters of RNA synthesis/degradation), regulatory variables (concentration of  
143 transcription factors, chromatin accessibility, epigenetic state, etc.) and system state  
144 variables (aging, metabolism or other environmental factors acting on cells) (Fig. 1).

#### 145       **Gene state determinants of expression variability.**

146       If the right conditions are met, RNA polymerase Pol II (RNAP II) binds to a promoter  
147 region and initiates transcription of the gene (81). The transcription happens in short bursts  
148 followed by a refractory period in which no transcription takes place (116). A simplified  
149 derivation of the two-state promoter model shows that non-Poisson noise depends  
150 inversely on the burst frequency, while Poisson noise is reciprocal of a product of burst size  
151 and burst frequency (21, 72). Each gene has its own bursting dynamics which, in turn,  
152 determines its noise (93). Different factors can either influence the burst frequency, a  
153 frequency of promoter activation within the mean lifetime of RNA, or the burst size, the  
154 amount of RNA produced per unit of time within a burst (82). Thus, any factor interfering

155 with promoter fluctuation, RNA synthesis or degradation kinetics is expected to modulate  
156 the within-cell and cell-to-cell variability in RNA copy number and thus gene noise.

157 In eukaryotes, the RNA “birth-death” rates are orchestrated by a complex regulatory  
158 system. With respect to the regulation of the synthesis rate, it is worth mentioning the RNA  
159 splicing process. The different proteins involved in splicing and accessibility of alternative  
160 donor/acceptor sites can modulate RNAP II elongation rate and, thus, the RNA synthesis  
161 rate. For instance, RNAP II elongation rates tend to increase throughout introns as  
162 compared to exons (42, 46). However, splice sites themselves, in mammals, but not in yeast,  
163 act as RNAP II pausing sites (19, 41). Such pausing can be bypassed by the inhibition of  
164 splicing mechanisms (65). To that, co-transcriptional checkpoints associated with splicing  
165 can further modulate the synthesis rates (3, 16). Thus RNA splicing, being intimately linked  
166 with RNA elongation, is expected to contribute to Poisson noise by modulating RNA “birth”  
167 rate.

168 The amount of RNA observed in a cell is the consequence of equilibrium between  
169 synthesis and degradation. This means not only fluctuations in the synthesis rate, but also  
170 variations in the degradation rate are likely to influence both the average expression as well  
171 as the variation in expression (57, 97). The half-life of RNA molecules depends on the length  
172 of the 3'-poly(A)-tail which is removed through deadenylation prior to degradation (67,  
173 109). As a direct consequence of the two-state promoter model, the total gene expression  
174 noise (Poisson and non-Poisson) is directly proportional to the RNA degradation rate. This  
175 implies an increased noise level for RNA species with shorter half-life and a decreased noise  
176 for the stable RNA molecules. For example, the presence of certain microRNAs have been  
177 shown to increase the rate of RNA deadenylation (107) and one can predict that such a  
178 mechanism will turn up the gene noise. Strikingly, although RNA synthesis and  
179 degradation, at first glance, are two independent processes, the RNA degradation rate was  
180 found to be regulated by transcription (13, 33). In terms of gene noise, the existence of a  
181 coupling between synthesis and degradation rates has a profound consequence as it leads  
182 to non-Poisson RNA “birth-death” process even in the absence of upstream cellular drives  
183 (96).

184 Finally, it is reasonable to assume that the kinetics of transcriptional bursts and as a  
185 result, gene noise is likely to be determined by the promoter sequence and the surrounding  
186 architecture. Indeed, the presence of a TATA-box within the promoter is known to increase  
187 not only the average expression of genes, but also its noise (11, 76, 77). TATA-box binding  
188 protein TBP associates with distinct co-activator complexes, each of which competes for the  
189 binding to the promoter, as it also mediates re-initiation of transcription by RNAP II (77, 81).  
190 Consequently, this promotes fluctuations in promoter activity, i.e. increases cell-to-cell or  
191 temporal deviations in the probability of the promoter to be in ON state. This, in turn,  
192 increases the gene expression noise, as non-Poisson noise is directly related to the  
193 fluctuations in these upstream cellular drives (21). Likewise, the complexity of the promoter  
194 can further increase the competition between distinct transcription factors and the  
195 expression noise. A simple promoter architecture, in which the promoter region is free from  
196 secondary regulation tends to generate little noise (36, 87). DNA variants in the promoter  
197 region can modulate the binding affinity of transcription factors, consequently changing  
198 both the average gene expression and expression noise (36). Besides competition for  
199 transcription factor binding within a promoter, competition between distinct promoters  
200 might also increase the gene noise by lowering the promoter burst frequencies (77). Next to  
201 that, the presence of a so called ‘speed bumps’ downstream of the transcription start site  
202 can cause RNAP II stalling (1), which might be detrimental for determination of bursting  
203 kinetics and noise. Although, further mechanistic insights into impact of gene state variables  
204 on gene noise remain to be made, the logic of doubly stochastic Poisson “birth-death”  
205 process and the two-state promoter model provide valuable tools for the dissection of gene  
206 noise determinants through the modelling of RNA “birth-death” rates.

#### 207 **Epigenetic determinants of expression variability.**

208 In eukaryotes, promoter accessibility and RNA synthesis are modulated by the  
209 epigenetic state of a gene, which sums from the DNA methylation status, nucleosome  
210 assembly and post-translational histone modifications. The epigenetic landscape is not  
211 static, as environmental cues such as diet, smoking, physical exercises and ageing can alter  
212 the epigenetic composition of the chromatin throughout organism’s lifetime (8, 29, 34, 95,  
213 102). Methylation patterns were shown to vary with circadian rhythm (5). Methylation of

214 CpG islands in promoter regions can alter transcription dynamics, resulting in the repression  
215 of transcription (10). In general, the presence of CpG islands in promoters lowers gene noise  
216 (27, 60). This might seem somewhat paradoxical, as increased CpG methylation is associated  
217 with increased nucleosome occupancy (20) and, as result, it is expected to elevate gene  
218 noise because of the lower promoter accessibility for transcription factor binding. However,  
219 occurrence of CpG islands in promoters of robustly expressed genes, i.e. in genes with low  
220 transcriptional noise, does not imply an increased methylation of their promoters. At the  
221 same time, a long-standing hypothesis suggests that DNA methylation might suppress  
222 cryptic transcription initiation from within the body of a gene, thereby reducing  
223 transcriptional noise (9, 39). Thus, it will be important to address these factors in future  
224 research on how DNA methylation partitions between promoter and gene body in genes  
225 with robust and noisy expression.

226 Assembly of eukaryotic DNA into nucleosomes adds yet another layer of complexity to  
227 gene regulation and is likely to modulate gene expression noise (17). Indeed, TATA-  
228 containing promoters favouring nucleosome assembly tend to further increase the noise  
229 due to a competition between TBP and nucleosomes (18, 83). Further, histones that  
230 constitute nucleosomes are subjected to a wide range of post-translational modifications,  
231 collectively known as a histone code (4). Transcription co-activator or co-repressor  
232 complexes recognize particular combinations of histone modifications tuning both gene  
233 expression level and expression variability (27, 108, 112). Thus, it may not be surprising that  
234 the presence of conflicting histone marks, i.e. co-occurrence of histone modifications  
235 associated with gene activation and repression, leads to an increased expression variability  
236 (27). First, bivalent histone modifications are expected to create a competitive state at the  
237 promoter and as a result, increase noise in the promoter activation. Second, bivalent  
238 histone marks might interfere with transcription initiation and elongation causing RNAP II to  
239 pause (51). In general, increased acetylation of histones and an overall “loose” chromatin  
240 structure at the promoter are associated with low expression noise, whilst a “closed”  
241 chromatin configuration and deficiency in active histone marks drive a higher noise (14, 22,  
242 63, 90, 98). In conclusion, the stochastic nature of RNA synthesis is intimately modulated by  
243 the stochastic nature of chromatin and DNA methylation states acting as upstream cellular  
244 drivers (14, 28).

245 **System state determinants of expression variability.**

246 In general, the biological processes are affected by two time-dependent factors: the  
247 circadian clock and aging. Interestingly, gene expression variability is also linked to these  
248 factors. For example, recently it has been shown that the circadian clock modulates burst  
249 frequency rather than burst size. Consequently, gene expression noise oscillates daily along  
250 with the average gene expression (63). Aging deteriorates many physiological parameters  
251 and their variability increases with time (reviewed in 15) and a clear epigenetic drift  
252 between monozygotic twins arises during aging (29). Thus, aging is expected to have a  
253 profound effect on gene expression variability (55). In accordance with this, the expression  
254 of housekeeping genes was shown to be more robust in cardiomyocytes from young mice as  
255 compared to old mice (6). To that, recent studies in mouse models provide evidence that  
256 inter-individual variability in gene expression tends to increase with age and can be reduced  
257 upon interventions aimed to slow ageing (61, 105). Further, a lower variation in gene  
258 expression was observed to correlate with the presence of H3K36me3 (27), a histone mark  
259 that was previously associated with increased longevity (86), although it is not known  
260 whether this epigenetic modification is a cause or consequence of the increased variation in  
261 gene expression. A recent study of gene expression in human skin, fat and blood samples  
262 from twin samples showed a general decrease of expression variability with age of  
263 individuals studied (101). This surprising and, perhaps, contradictory observation on linking  
264 aging and expression variability warrant further investigations of expression variability using  
265 other populations, tissue types as well as computational approaches for its quantification.

266 **Variability in gene expression might explain many biological phenomena.**

267 Variability determines plasticity, i.e. a degree to which a gene can change its expression  
268 in response to environmental fluctuations as a consequence of fluctuation-response  
269 relationship (49, 84). Plasticity of expression can serve a cell to adapt to a new  
270 environment (106). At the population level, a more varied expression of certain genes can  
271 produce individuals that are better prepared for changing conditions at the cost of reduced  
272 metabolic efficiency (12). This was shown on a microscopic scale in yeast, in which a high  
273 variability in expression of yeast plasma-membrane transporters enhanced their adaptive  
274 capabilities to a changing environment (114). Selection for the yeast TDH3 enzyme involved  
275 in the glucose metabolism was shown to have a greater impact on expression noise rather

276 than on the average level of expression, showing an example of selection for higher  
277 variability as an adaptation mechanism (59). Overall, genes involved in environmental  
278 responses show more variation in expression, which can be beneficial for non-housekeeping  
279 functions such as coping with stress or reacting to environmental queues (11, 69). Genome  
280 wide analysis of transcriptional and epigenetic variability across human immune cell types  
281 showed that neutrophils, one of the first-responder cells upon an infection, contained the  
282 largest variation in both methylation and expression and alluding that variability might be  
283 an important factor in immune response (24). Also inter-population variability has shown  
284 that genes can have similar levels of expression variability across individuals and  
285 populations, with the largest differences observed among genes associated with immune  
286 response and disease susceptibility such as chemokine receptor *CRCX4* that is important for  
287 HIV-1 infection, where variation in expression may underlie difference in disease  
288 susceptibility (50). In contrast, genes involved in growth and development (85), as well as  
289 genes which provide a universal function, such as protein synthesis or degradation generally  
290 (e.g. translation initiation and ribosomal proteins) show a relatively robust expression (62).  
291 Similarly, genes central in gene networks, like key pluripotency regulator *Pou5f1* (56) or  
292 encoding products that are critical to the survival of cells (also known as essential genes,  
293 since their deletion is lethal) and genes which code for multi-protein complexes have  
294 evolved to minimize their expression noise (30, 48, 54). Finally, a recent study in humans  
295 showed that long non-coding RNAs, such as anti-sense transcripts from the genomic loci  
296 corresponding to known protein-coding genes, display a higher inter-individual expression  
297 variability as compared to protein-coding genes (45) substantiating their role in adaptation.

298 Another biological phenomenon where the expression variability might play an  
299 important role is incomplete penetrance (71, 73). The latter study shows that in *C.elegans*  
300 mutants with more stochastic expression of *end-1* gene, a threshold for activating  
301 expression of *elt-2*, the master regulator of intestinal differentiation may or may not be  
302 reached, and hence only some of mutant embryos will develop intestine. Different levels of  
303 expression in individuals with a similar or even isogenic genetic background can explain why  
304 some individuals develop severe disease while others have a mild or even wild-type  
305 phenotype. Even individuals which are genetically identical can show phenotypic differences  
306 and even personality traits, as recently reviewed in (25). Studying transcriptomes from the



307 viewpoint of expression variability can provide new explanations for mechanisms of disease  
308 development.

309 **Prerequisites for analysis of differential variability in gene expression.**

310 Despite the high promises of differential variability analysis, several important factors  
311 should be taken into consideration when planning and performing this type of analysis.

312 *Sufficient number of samples.* While some of the studies investigating expression  
313 variability used as few as 3 samples per group (105), technical biases in library preparation  
314 and sequencing can have profound effects on the differential variability estimates. For a  
315 reproducible analysis of differential variability, a larger sample size is required in contrast to  
316 studies where a differential mean expression is tested (110). This is further exemplified  
317 below by means of power analysis in the section showcasing the differential variability  
318 analysis for mice.

319 *Avoiding batch effects.* Since technical variation can mask the effects coming from  
320 biological differences, it is important to perform all technical procedures in a single batch or,  
321 whenever that is not possible, randomly distribute samples from different groups among  
322 experiment batches.

323 *Accounting for variability in transcript structure.* While most of current studies quantify  
324 variability using number of molecules or number of sequencing reads corresponding to the  
325 gene, the structure of the transcript is rarely taken into account. Yet, variability in pre-mRNA  
326 maturation is also observed (103). At the splicing level, statistical methods were developed  
327 to identify genes with condition-specific splicing ratios (31), while variation in splicing can be  
328 defined and quantified using a recently suggested local splicing variation units (100). Future  
329 methods for differential variability analysis, therefore, should consider not only  
330 quantitative, but also structural variability of gene products.

331 The first two points are rather general experimental design considerations, while the  
332 latter is more specific for RNA-sequencing based profiling of gene expression.

333 **Statistical inference of gene expression variability**

334 Several metrics have been proposed to measure the variability of gene expression, such  
335 as variance ( $\sigma^2$ ), the (squared) coefficient of variation ( $cv$ , also known as signal to noise  
336 ratio), Fano factor (also known as noise strength), and their robust counterparts median  
337 absolute deviation from the median (MAD), (quartile) coefficient of dispersion (COD or  
338 QCOD), etc. (74, 83, 99) (Table 1).

339 Applicability and interpretation of these metrics depend on how gene expression data  
 340 was obtained and processed. For example, variance stabilizing transformations (VST,  $f(x)$ )  
 341 of microarray hybridization intensities or normalized RNA counts (such as CPM - counts per  
 342 million or FPKM – fragments per kilobase of transcript per million) transform mean and  
 343 variance as  $E[f(X)] \approx f(\mu_X)$  and  $Var[f(X)] \approx (f'(\mu_X))^2 \sigma_X^2$  respectively, following the 1<sup>st</sup>-  
 344 order Taylor expansion, where  $\mu_X$  and  $\sigma_X^2$  are original mean and variance respectively.  
 345 Among commonly used VSTs are the logarithm ( $\log_2(X)$ ) and generalized logarithm  
 346 ( $g\log_2(X) = \log_2(X + \sqrt{X^2 + 1})$ ) functions (38). This implies that the variance of  $\log_2$  or  
 347  $g\log_2$  transformed variables corresponds to the squared coefficient of variation of the  
 348 original variable ( $cv_X^2$ ) as  $Var[\log_2(X)] \approx \log(2)^{-2} \frac{\sigma_X^2}{\mu_X^2} = \log(2)^{-2} cv_X^2$  and  
 349  $Var[g\log_2(X)] \approx \log(2)^{-2} \frac{\sigma_X^2}{\mu_X^2 + 1} \approx \log(2)^{-2} cv_X^2$  (for  $\mu_X^2 \gg 1$ ). Thus, it makes no sense to  
 350 estimate neither  $cv$  nor Fano factor for VST transformed variables as their variance is  
 351 already proportional to  $cv_X^2$ . Similar logic applies to robust measures of variability as  
 352  $MAD[\log_2(X)] \approx \text{median}(|\log_2(X_i/\bar{X})|)$  and  $MAD[g\log_2(X)] \approx \text{median}(|\log_2(X_i/\bar{X})|)$   
 353 (for  $X_i \gg 1$ ), and additional normalization of  $MAD$  to the median of VST transformed  
 354 variable is unnecessary.

355 In contrast, when dealing with untransformed variables emitted by Poisson or mixed-  
 356 Poisson processes (such as RNA-sequencing counts), normalization to the mean is required  
 357 due to the presence of the mean-variance relationships.  $Var[X] = \sigma_X^2 = \mu_X$  for Poisson and  
 358  $Var[X] = \sigma_X^2 = \mu_X + \alpha_X \mu_X^2$  for mixed-Poisson distributed RNA counts, where  $\alpha_X > 0$  is the  
 359 overdispersion parameter (44). Then, Fano factor turns to be handy for estimation of  
 360 deviation from Poisson process, as  $F = \sigma_X^2 / \mu_X > 1$  indicates overdispersion, while  
 361  $cv_X^2 = \mu_X^{-1} + \alpha_X$  partitions noise into two asymptotically orthogonal parameters of mixed-  
 362 Poisson distributions, which we refer to as Poisson and non-Poisson noise. In the section  
 363 showcasing the differential variability analysis for mice we demonstrate statistical inference  
 364 of both  $\mu_X$  and  $\alpha_X$  parameters for genes' RNA counts.

365 So far, statistical inference of expression variability is limited to only a few tools. For  
 366 instance, tools, such as AEGS and pathVar aim to discover biological pathways, for which the  
 367 expression variability changes. AEGS is a webserver that uses case-control data and allows  
 368 to identify which of pre-defined gene sets (e.g. genes belonging to the same gene ontology

369 category) are more variable expressed and ranks variability of individual genes within each  
370 set (32). The tool is also available as standalone program and can, in principle, be easily  
371 integrated into RNA-Seq analysis pipelines. PathVar enables case-control pathway-based  
372 interpretation of gene expression variability, but can also compare a single group of samples  
373 against a background distribution (99). This tool is available from Bioconductor collection of  
374 packages, provides a broad choice of variability measures and can also become part of  
375 routine transcriptome analysis.

376 Another tool, MDseq employs a generalized linear model (GLM) to estimate statistically  
377 significant changes in both expression mean and variability in response to experimental  
378 factors (74). Although MDseq considerably expands the standard GLM approach employed  
379 in many tools for differential gene expression analysis, its current implementation seems to  
380 be limited to a fixed effect negative binomial (NB) model (74). To that, MDseq  
381 parametrization of the NB implies a linear mean-variance relationship for RNA counts:  
382  $Var(X) = \mu\phi$ , while many RNA-seq studies suggest a quadratic relationship (58). In fact, a  
383 quadratic mean-variance relationship originates from the mixed-Poisson nature of a  
384 stochastic process driving the RNA synthesis and degradation (21, 40, 66, 72).

385 In brief, for a mixed-Poisson processes, the Poisson rate ( $\mu$ ), represented by a ratio of  
386 RNA synthesis to degradation rates, is assumed to be a random variable with expectation  
387  $E(\mu) = \mu$  and the variance defined by an underlying mixing distribution  $g_\mu(\mu)$ . The mixed  
388 Poisson distribution of RNA counts takes the following general form:  $P(X = x) =$   
389  $\int_0^\infty \frac{e^{-\mu}\mu^x}{x!} g_\mu(\mu) d\mu$ , where mixing distribution  $g_\mu(\mu)$  can take on any parametric form  
390 depending on upstream cellular drives (21). For example, promoter switching between  
391 active and inactive states (bursts) leads under limiting conditions to a gamma distribution of  
392 the Poisson rate ( $\mu$ ). As a result, the cell-to-cell distribution of the RNA copy number follows  
393 a gamma-Poisson distribution (also known as a negative binomial, NB) (21, 72). Likewise, the  
394 NB distribution empirically fits well to RNA sequencing counts from both tissues and cell  
395 populations (58).

396 For any mixed-Poisson process, i.e. independent of a specific form of the  $g_\mu(\mu)$ , a total  
397 variance and noise (a squared coefficient of variation of RNA counts) sums from the Poisson  
398 (1<sup>st</sup> summand) and non-Poisson (2<sup>nd</sup> summand) parts as:  $Var[X] = \mu + \alpha\mu^2$ ,  $cv^2(X) =$   
399  $\mu^{-1} + \alpha$  respectively (44, 79). Non-Poisson variation –  $\alpha$  is often referred to as the

400 overdispersion parameter or the biological coefficient of variation ( $\alpha = bcv^2$ ) (58). The  
401 Poisson and non-Poisson variation are also assigned as “intrinsic” and “extrinsic”  
402 respectively (68). Thus, the goal of differential gene expression analysis is to estimate the  
403 average RNA copy number -  $\mu$ , while that of differential gene noise analysis is to estimate  
404 overdispersion -  $\alpha$  from a distribution of RNA counts.

#### 405 **A showcase for differential gene expression variability analysis using GAMLSS**

406 Here we propose to utilize GAMLSS to assess the effects of biological factors on a  
407 gene’s Poisson ( $\mu^{-1}$ ) and non-Poisson ( $\alpha$ ) variation. GAMLSS stands for generalized additive  
408 model for location, scale and shape and offers immense flexibility for semi-parametric  
409 mixed effect modelling of up to four distribution parameters (78, 91).

410 The suggested analysis strategy has several advantages. First, GAMLSS comes with an  
411 extensive list of mixed-Poisson distributions along with their zero inflated/adjusted variants  
412 (79). Second, GAMLSS allows for the fitting of mixed effect models to RNA counts. And third,  
413 smoothing terms (splines) can also be used to model non-linear relations of mixed-Poisson  
414 distribution parameters with continuous experimental variables such as age. These factors  
415 combined give it a much better control in the modelling of differential gene expression and  
416 variability.

417 To demonstrate GAMLSS at work, we provide a brief re-analysis of age-dependent  
418 changes in the overdispersion (non-Poisson variation) for genes expressed in liver samples  
419 taken from young and old C57BL/6J mice (61). All computer programs used here and  
420 description of the analysis are available as GitHub repository  
421 (<https://github.com/Vityay/ExpVarQuant>).

422 We modeled genes’ RNA counts using the  $NB(\mu, \alpha)$  distribution parametrized with  
423 respect to the mean ( $\mu$ ) and non-Poisson variation ( $\alpha$ ) in such a way that the quadratic  
424 mean-variance relationship holds. The probability mass function for independent and  
425 identically distributed RNA counts ( $X$ ) for a given gene:  $X \stackrel{\text{ind}}{\sim} NB(\mu, \alpha)$  is defined as:

$$426 \quad P(X = x) = \frac{\Gamma(\frac{1}{\alpha} + x)}{\Gamma(\frac{1}{\alpha})\Gamma(x+1)} \left(\frac{1}{1+\alpha\mu}\right)^{\frac{1}{\alpha}} \left(\frac{\alpha\mu}{1+\alpha\mu}\right)^x,$$

427 with expectation (mean) and variance of RNA counts:  $E[X] = \mu$ ,  $Var[X] = \mu + \alpha\mu^2$ ,  
428 and  $cv^2(X) = \mu^{-1} + \alpha$ .

429 Then, we specified a GAMLSS model to account for the age (young – 5 months and old –  
430 20 months old mice) effect on both the mean RNA counts and the overdispersion:

431  $\log(X_i) \sim age_j \beta_{\mu_j} + \log(N_i),$

432  $\log(\alpha) \sim age_j \beta_{\alpha_j},$

433 where  $i = 1, \dots, n$  is  $i^{th}$  observation of gene's mRNA counts ( $X_i$ );  $j = 1, \dots, p$  is  $j^{th}$  factor  
 434 level (young – 5 weeks, old – 20 weeks); and  $\log(N_i)$  is offset vector represented by library  
 435 sizes. The first equation of GAMLSS specifies a model of a factor effect, namely  $age_j$ , on  
 436 library size ( $N_i$ ) normalized mean mRNA counts ( $\mu_j = e^{\beta_{\mu_j}}$ ,  $cpm_j = 10^6 \mu_j$ ). Essentially, this  
 437 part of the model corresponds to a GLM model of differential gene expression (58),  
 438 however, GAMLSS allows for more flexibility as random effects and smoothing terms can  
 439 also be included (91). The second equation of GAMLSS models a factor effect on non-  
 440 Poisson noise ( $\alpha$ ), where  $\beta_{\alpha_j}$  is a maximum-likelihood estimation of overdispersion  
 441 parameter ( $\alpha_j = e^{\beta_{\alpha_j}}$ ).

442 Significance values of age-mediated changes in  $\mu$  and  $\alpha$  parameters of the  $NB(\mu, \alpha)$   
 443 were assessed for each gene with likelihood ratio tests (LR). For a given gene, LR test  
 444 statistic for changes in mean RNA counts between old and young mice was calculated as  
 445 following:

446 
$$D_{\mu} = -2 \log \frac{\text{likelihood for reduced model}}{\text{likelihood for GAMLSS model}} = -2 \log \frac{\mathcal{L}(\mu_0, \alpha_j | X_i)}{\mathcal{L}(\mu_j, \alpha_j | X_i)},$$

447 Where the reduced model omits factor effect (age) from the model of  $\mu$ :  $\log(X_i) \sim \beta_{\mu_0} +$   
 448  $\log(N_i)$ , while the age effect on non-Poisson noise was still accounted for. It can be readily  
 449 noted that the estimation of differential gene expression by GAMLSS differs from that by  
 450 classical GLM as the latter estimates only the shared overdispersion (58). In brief, the GLM  
 451 model is specified as:

452  $\log(X_i) \sim age_j \beta_{\mu_j} + \log(N_i),$

453  $\log(\alpha) \sim \beta_{\alpha_0}$

454 in GAMLSS notation and the LR test statistic is calculated as:

455 
$$D_{\mu_{GLM}} = -2 \log \frac{\text{likelihood for null model}}{\text{likelihood for GLM model}} = -2 \log \frac{\mathcal{L}(\mu_0, \alpha_0 | X_i)}{\mathcal{L}(\mu_j, \alpha_0 | X_i)},$$

456 where null model omits factor effect on both  $\mu$  and  $\alpha$ . Finally, LR test statistic for changes in  
 457 non-Poisson noise was calculated by comparing GLM model (as reduced model for  $\alpha$ ) with  
 458 full GAMLSS model:

459 
$$D_{\alpha} = -2 \log \frac{\text{likelihood for GLM model}}{\text{likelihood for GAMLSS model}} = -2 \log \frac{\mathcal{L}(\mu_j, \alpha_0 | X_i)}{\mathcal{L}(\mu_j, \alpha_j | X_i)}.$$

460  $D_\mu$ ,  $D_{\mu_{GLM}}$  and  $D_\alpha$  are asymptotically  $\chi^2$ -distributed with degrees of freedom equal to a  
461 difference between the number of compared models' parameters. Thus, from this example  
462 it is clear that GAMLSS is an extension of a GLM model allowing for the estimation of factor  
463 effects on both parameters of the distribution of RNA counts, namely mean and  
464 overdispersion (non-Poisson noise).

465 We excluded genes with zero counts in any of the samples from the analysis as this  
466 might bias the estimation of non-Poisson variation. In fact, an excess of zeros in RNA-seq  
467 data imposes a certain problem for statistical inference of the distribution parameters for  
468 RNA counts. Indeed, in many cases it is impossible to discriminate whether observing a zero  
469 is the result of a gene being silenced or whether it is observed due to an insufficient  
470 sequencing depth causing dropouts of the lowly expressed genes. In principle, the former  
471 case corresponds to a zero-adjusted model, while the latter – to a zero-inflated model, and  
472 both could be fitted by GAMLSS. However, neither of these assumptions alone resolves the  
473 uncertainty that zero values introduce to transcriptome analysis.

474 Having estimated the parameters  $\mu$  and  $\alpha$  for liver genes expressed in young and old  
475 mice, we noted that their absolute values were practically uncorrelated ( $\rho(\mu, \alpha) \rightarrow 0$ ). This  
476 could be attributed directly to the given parametrization of the  $NB(\mu, \alpha)$ , which implies an  
477 asymptotic independence of the estimated parameters. It follows from the Fisher  
478 information matrix as its element  $I_{\mu\alpha} = -E \frac{\partial^2}{\partial\mu\partial\alpha} \log(P(X|\mu, \alpha)) = 0$ . To that, changes in  
479 the mean gene expression and the non-Poisson variation occurring with age were also  
480 almost uncorrelated ( $\rho(\Delta\mu, \Delta\alpha) \rightarrow 0$ ). Testing under the assumption that the cellular RNA  
481 concentration (total number of RNA molecules per cell) is the same for the samples taken  
482 from young and old mice, we scored about a comparable number of genes for which the  
483 mean RNA counts either increased or decreased significantly with age (Fig. 2A, 3A).  
484 Estimation of the mean also yielded the estimation of the Poisson variation as they are  
485 reciprocal to each other (Poisson variation =  $\mu^{-1}$ ). In contrast to the Poisson variation, non-  
486 Poisson variation increased with age (Fig. 2B). Importantly, applying the GAMLSS model  
487 enabled for the identification of genes for which the non-Poisson variation, but not the  
488 mean, changed significantly with age (Fig. 2B, 3B).

489 However, it must be noted that the relative standard errors of overdispersion estimates  
490 tend to be larger than that of mean estimates. As a result, this lowers the statistical power

491 of likelihood ratio test for factor effects on non-Poisson variation. This is evident from the  
492 power analysis of LR tests for fold changes in mean and overdispersion (Fig. 2C, D). Although  
493 a derivation of the analytical form for the power of LR tests for complex models deems  
494 impossible, this can be circumvented by a simulation method. To this end, a thousand pairs  
495 of samples of NB distributed random variables were generated with the given parameters  
496  $\mu_0$  (counts) and  $\alpha_0$  (non-Poisson noise) for reference samples and fold changes ( $\delta$ ) in one of  
497 the NB parameters for test samples. Then, LR tests were applied comparing simulated  
498 reference samples  $NB(\mu_0, \alpha_0)$  with test samples  $NB(\delta\mu_0, \alpha_0)$  and  $NB(\mu_0, \delta\alpha_0)$ . The power  
499 of LR tests for  $\mu_0 \neq \delta\mu_0$  (Fig. 2C) and  $\alpha_0 \neq \delta\alpha_0$  (Fig. 2D) was then estimated as proportion  
500 of true positives at significance level of  $< 0.05$ . Obviously for all tested configurations of NB  
501 ( $\mu_0$ : {10,100,1000} and  $\alpha_0$ : {0.1,0.25,0.5}) the power of LR tests for mean and  
502 overdispersion increased with an increasing sample size. To that, the power of LR tests for  
503 fold changes in mean counts (Fig. 2C) is higher than that of non-Poisson noise (Fig. 2D).  
504 Unexpectedly though, the power of LR tests tends to increase, especially for the tests  
505 comparing overdispersion, with increasing  $\mu_0$  irrespectively of the presence or absence of  
506 an offset parameter, which simulates library size. This suggests that an increase of sample  
507 size and sequencing depth (library size) will eventually increase the statistical power of tests  
508 aimed at comparing changes in mean expression and non-Poisson noise.

#### 509 **The expression variability analysis provides additional insights into dataset**

510 To identify biological pathways associated with the age-mediated increase in non-  
511 Poisson variation, we fitted a ridge regression model to the  $\log_2$  fold change in  
512 overdispersion using KEGG annotations of genes as a model matrix (Fig. 4A) (35, 43). Such  
513 an approach circumvents the problem of pathways overrepresentation analysis associated  
514 with the necessity to select a threshold for statistical significance. It is also well suited for  
515 the analysis of non-Poisson variation when a common trend for genes is to increase in  
516 variability with age. As a result, the KEGG-pathway ridge regression model revealed several  
517 pathways, such as the complement and coagulation cascades, amino acid (Val, Leu, Ile)  
518 degradation, chemokine signaling and others for which non-Poisson variation increased in  
519 aged mice (Fig. 3B, 4B).

#### 520 **Fluctuation-response relationship for RNA counts**

521 Gene expression noise is thought to drive gene expression plasticity due to a  
522 fluctuation-response relationship (49, 84). This implies that an absolute change in the

523 expectation ( $\mu$ ) of some measurable quantity ( $X$ ) in response to an influence is proportional  
524 to its initial variance:  $|\mu_1 - \mu_0| \sim \text{Var}(X)$ . However, this relationship holds only true for  
525 Gaussian-like distributed quantities under the assumption of a fixed variance:  
526  $\text{Var}(X_1) \sim \text{Var}(X_0)$ . Nonetheless, if log transformed RNA counts approximate a Gaussian-  
527 like distribution, then the fluctuation-response relationship takes on the following form:  
528  $|\log(\mu_1/\mu_0)| \sim \alpha = bcv^2$ , as a result of the Taylor expansions for the moments for genes  
529 expressed at large copy number ( $\mu \gg 1$ ). We noted a modest, but significant, positive  
530 correlation between absolute  $\log_2$  fold changes in the mean gene expression for old and  
531 young mice with non-Poisson variation for young mice (Fig. 5A). A lack of a stronger  
532 correlation could be due to the violation of the fluctuation-response assumption of a fixed  
533 variance or overdispersion for log-transformed variables. In general, this substantiates the  
534 fluctuation-response relationship for the RNA copy number.

535

536 **Estimates of gene variation from tissues retain information on gene state**  
537 **determinants of non-Poisson noise.**

538 Finally, we wondered if the estimate of non-Poisson variation from RNA-sequencing  
539 data of cell populations contain information on gene state determinants. To this end, we  
540 compared the genes' non-Poisson variation estimates with their promoter DNA-sequence  
541 composition. First, we noted that on average, that the non-Poisson variation was higher for  
542 genes that were regulated by TATA-containing promoters (Fig. 5B). Second, in accordance  
543 with the fluctuation-response relationship (Fig. 5A), aging induced more pronounced  
544 changes in the mean expression of genes with TATA-containing promoters (Fig. 5B, C).  
545 Overall, this result is in agreement with the TATA-mediated promoter fluctuation caused by  
546 a competition between distinct TBP-co-activator complexes (77, 82, 87) and it substantiates  
547 that gene state signals are retained in cell population estimates of non-Poisson variation.

548 To conclude this brief showcase of GAMLSS, we advocate for the use of this framework  
549 to dissect the determinants of both the mean RNA counts and the non-Poisson variation as  
550 two independent parameters of gene expression network.

551 **Combining other -omics data with RNA-seq can lead to new discoveries.**

552 A connection between the gene expression variability measured on different levels:  
553 cell-to-cell, inter-individual and inter-population has been suggested previously (23, 25). The  
554 rapid development of accessible and cost-efficient methods for single-cell RNA-seq (scRNA-



555 seq) will provide us with improved estimates of cell-to-cell variability in gene expression  
556 (70). Flow cytometry techniques can help in the further separation into (so called / the  
557 suggested) macro-heterogeneity, which is the variability that encompasses both the on- and  
558 off- state of genes, as well as the micro-heterogeneity, which represents the variability in  
559 gene expression of genes in different (37). Further, recently generated large transcriptome  
560 datasets for hundreds of individuals (2, 47) should increase our understanding of  
561 transcriptome variability at population level.

562       Apart from transcriptomics data, large sets of epigenetics data will be of great value.  
563 For example, the changing landscape of histone modifications with age has been established  
564 (89), as has the property of histone modifications to be associated with the average gene  
565 expression and variation in gene expression (108). Similarly, the beneficial effects of  
566 alterations in diet have been shown to extend the lifespan of mice (7), as has the  
567 methylation of genes and the consequent variation in expression been shown to contribute  
568 to the pathophysiology of mice on a high fat diet (113). In line with these two observations,  
569 it has been shown that the suppression of inter-individual variation has positive effects on  
570 the lifespans of *C. elegans* (75).

571       Finally, when speaking of gene expression variability, it is important to consider how the  
572 variability in RNA copy number translates to variability at a protein level. Often there seems  
573 to be a discrepancy between the amount of RNA transcribed and the amount of the  
574 matching protein being produced within samples (64). Yet, many principles of gene noise  
575 have been derived by quantifying reporter genes expression on protein level, such as two-  
576 color reporter assay (26, 94). To that, derivations of protein fluctuations from theoretical  
577 models of stochastic gene expression highlight the contribution of RNA-level noise to  
578 protein-level noise (68). Thus, it makes it reasonable to propose that gene expression  
579 variability might propagate from RNA to protein, from protein to cell, from cell to tissue and  
580 from tissue to organism.

581       To conclude, the analysis of differential transcriptome variability complements the  
582 standard analysis of differential gene expression and reveals another dimension of  
583 expression analysis. With the further development of tools and with a wider acceptance of  
584 these methods, we will advance our understanding of the mechanisms underlying the  
585 regulation of transcription, common physiological traits and disease predispositions.

586 **Table 1.** Commonly used measures of variability

Coefficient of variation (signal to noise ratio)	$cv = \sigma/\mu$
Fano factor (noise strength)	$F = \sigma^2/\mu$
Median Absolute Deviation from the Median	$MAD = \text{median}( X_i - \tilde{X} )$
Coefficient of Dispersion	$COD = MAD/\tilde{X}$
Quartile Coefficient of Dispersion	$QCOD = (Q_3 - Q_1)/(Q_3 + Q_1)$

587  $\tilde{X}$  - median;  $Q_1$  and  $Q_3$  are the 1<sup>st</sup> and 3<sup>rd</sup> quartiles respectively.

588

589 **Figure legends**

590 **Figure 1.** A model depicting factors influencing the gene expression variability/noise. Key  
 591 equations depicting the partitioning of variance and squared coefficient of variations into  
 592 Poisson (blue, Pois.) and non-Poisson (red, non-Pois.) variability/noise are shown. Such  
 593 partitioning holds true for any mixed-Poisson distribution, where the Poisson rate  $\mu$  is a  
 594 random variable distributed with expectation  $\langle \mu \rangle$  and variance  $Var[\mu]$ . Key equations for  
 595 the expectation ( $E[RNA]$ ), variance ( $Var[RNA]$ ) and noise ( $cv^2(RNA)$ ) for two-state  
 596 promoter model are expressed in terms of burst size ( $b$ ) and burst frequency ( $f_b$ ). See text  
 597 for further details.

598

599 **Figure 2.** A GAMLSS analysis of age-mediated changes in gene expression and non-Poisson  
 600 noise.

601 **A)** Boxplots of a GAMLSS estimations of the mean mRNAs copy numbers (counts per million  
 602 mapped reads, cpm) for genes expressed in the liver of young (5 months, n = 6) and old (20  
 603 months, n = 6) C57BL/6J mice (left panel). Scatter plot of genes' mean mRNA copy number  
 604 in young and old mice (middle panel) and a boxplot of  $\log_2$  fold changes in expression  
 605 between old and young mice (right panel). Significantly up- and down-regulated genes (false  
 606 discovery rate,  $FDR \leq 0.05$ ) are indicated in red and blue respectively. In boxplots, the box  
 607 spans the interquartile range (IQR) from 25% ( $Q_1$ ) to 75% ( $Q_3$ ) and the middle line indicates  
 608 50% (median). Whiskers span to 1.5 IQR from the lower ( $Q_1$ ) and upper ( $Q_3$ ) quartiles or are  
 609 truncated to the min or max values, if those are within 1.5 IQR.

610 **B)** GAMLSS estimation of non-Poisson variability in mRNAs copy numbers (left panel). A  
 611 scatter plot of genes' estimates of non-Poisson variability in young and old mice (middle

612 panel) and a boxplot of  $\log_2$  fold changes in non-Poisson variability with age (right panel).  
613 Genes for which the non-Poisson noise increased or decreased significantly with age are  
614 marked in red or blue respectively.

615 **C)** Heatmap depicting a power analysis of the likelihood ratio (LR) test for fold changes ( $\delta$ ) in  
616  $\mu_0$  (mean counts). For each power analysis (1000) pairs of samples from reference  
617  $NB(\mu_0, \alpha_0)$  and test  $NB(\delta\mu_0, \alpha_0)$  distributions were simulated with  $\mu_0 \in \{10, 100, 1000\}$ ,  
618  $\alpha_0 \in \{0.1, 0.25, 0.5\}$  and  $\delta \in \{\frac{1}{4}, \frac{1}{3}, \dots, 3, 4\}$ . Sample sizes were  $\{5, 10, \dots, 100, 1000\}$ . Null  
619 hypothesis:  $H_0: \mu_0 = \delta\mu_0$  were rejected at significance level of 0.05 and power was  
620 calculated as the probability of rejecting  $H_0$ . Red indicates high power, white -low.

621 **D)** Heatmap depicting a power analysis of the likelihood ratio (LR) test for fold change ( $\delta$ ) in  
622  $\alpha_0$  (non-Poisson noise).

623

624 **Figure 3.** Examples of differentially expressed genes (**A**) and genes showing increase in non-  
625 Poisson variability with age (**B**). Upper panel, boxplots of selected liver genes' mRNAs copy  
626 numbers (expressed as  $\log_2(\text{cpm})$ ) for young (green,  $n = 6$ ) and old (red,  $n = 6$ ) C57BL/6J  
627 mice. Whiskers extend to minimum and maximum values. Middle panel, boxplots of  
628  $\log_2(\text{cpm})$  residual values corrected for genes' grand mean expression for young and old  
629 mice ( $\sim \text{gene}$ ). Lower panel, boxplots of  $\log_2(\text{cpm})$  residuals corrected for genes' group-wise  
630 mean expression in young and old mice ( $\sim \text{gene}:\text{age}$ ). The middle panel serves to illustrate  
631 differential gene expression, while the lower panel shows whether the gene expression  
632 variability is affected by age. Genes were selected based on significance of the age-  
633 mediated changes in mean mRNA counts (**A**,  $\text{FDR}_{\text{cpm}} \leq 0.05$ ) or changes in non-Poisson  
634 variability (**B**,  $\text{FDR}_{\text{non-Pois. variability}} \leq 0.05$ ). For (**B**), note an increase in  $\log_2(\text{cpm})$  variability for  
635 selected genes in population of 20 weeks old mice due to an increase in non-Poisson  
636 variability with age as compared to 5 weeks mice. Left panel in (**B**) shows genes associated  
637 with complement and coagulation cascades according to KEGG annotation, the right panel  
638 shows a selection of 30 genes with the highest statistically significant gain in non-Poisson  
639 variability.

640

641 **Figure 4.** Pathway analysis of age-mediated changes in non-Poisson variability.

642 **A)** Ridge regression model predicting age-mediated changes in non-Poisson variability based  
643 on the genes' KEGG pathway annotations.

644 **B)** Top 20 KEGG pathways associated with age-mediated increase in non-Poisson variability.  
645 Pathways were selected based on the ranking of model coefficients.

646

647 **Figure 5. A)** Relationships between the initial non-Poisson variability in young mice and the  
648 age-mediated responses in the mean mRNAs counts. Gene expression responses are  
649 represented as absolute  $\log_2$  ratios (top panel) of mean mRNA counts in old and young  
650 mice. GAMLSS estimates of genes' non-Poisson variability in young mice are given as ranked  
651 values ranging from lowest (1) to highest (10). Spearman correlation coefficients are shown.  
652 Trend lines were generated by LOESS local regression.

653 **B-C)** TATA-box associated with increased non-Poisson variability and age-mediated response  
654 in mean expression levels. **(B)** Boxplots show the initial non-Poisson variability in 5 months  
655 old mice (young, upper panel) and absolute changes in the mean gene expression (lower  
656 panel) for mouse genes classified according to all possible combinations of four promoter  
657 motifs: the TATA-box, Initiator (Inr), CCAAT-box and GC-box. A group of genes lacking any of  
658 those is labelled as "none". **(C)** Scatterplot of genes' group-wise medians in the initial non-  
659 Poisson variability at age of 5 months and the absolute changes in mean gene expression  
660 levels between old and young mice. Genes containing a TATA-box in any of these  
661 combinations in their promoters tend to have a higher non-Poisson variability and respond  
662 stronger to age with respect to the changes in mean expression levels. The Pearson  
663 correlation coefficient and significance are indicated.

664 **References**

665

- 666 1. **Adelman K, Henriques T.** Transcriptional speed bumps revealed in high  
667 resolution. *Nature* 560: 560–561, 2018.
- 668 2. **Aguet F, Brown AA, Castel SE, Davis JR, He Y, Jo B, Mohammadi P, Park Y,**  
669 **Parsana P, Segrè A V., Strober BJ, Zappala Z, Cummings BB, Gelfand ET,**  
670 **Hadley K, Huang KH, Lek M, Li X, Nedzel JL, Nguyen DY, Noble MS, Sullivan TJ,**  
671 **Tukiainen T, MacArthur DG, Getz G, Addington A, Guan P, Koester S, Little**  
672 **AR, Lockhart NC, Moore HM, Rao A, Struewing JP, Volpi S, Brigham LE, Hasz**  
673 **R, Hunter M, Johns C, Johnson M, Kopen G, Leinweber WF, Lonsdale JT,**  
674 **McDonald A, Mestichelli B, Myer K, Roe B, Salvatore M, Shad S, Thomas JA,**  
675 **Walters G, Washington M, Wheeler J, Bridge J, Foster BA, Gillard BM, Karasik**  
676 **E, Kumar R, Miklos M, Moser MT, Jewell SD, Montroy RG, Rohrer DC, Valley D,**  
677 **Mash DC, Davis DA, Sobin L, Barcus ME, Branton PA, Abell NS, Balliu B,**  
678 **Delaneau O, Frésard L, Gamazon ER, Garrido-Martín D, Gewirtz ADH, Gliner**  
679 **G, Gloude-mans MJ, Han B, He AZ, Hormozdiari F, Li X, Liu B, Kang EY,**  
680 **McDowell IC, Ongen H, Palowitch JJ, Peterson CB, Quon G, Ripke S, Saha A,**  
681 **Shabalin AA, Shimko TC, Sul JH, Teran NA, Tsang EK, Zhang H, Zhou Y-H,**  
682 **Bustamante CD, Cox NJ, Guigó R, Kellis M, McCarthy MI, Conrad DF, Eskin E,**  
683 **Li G, Nobel AB, Sabatti C, Stranger BE, Wen X, Wright FA, Ardlie KG,**  
684 **Dermitzakis ET, Lappalainen T, Aguet F, Ardlie KG, Cummings BB, Gelfand**  
685 **ET, Getz G, Hadley K, Handsaker RE, Huang KH, Kashin S, Karczewski KJ, Lek**  
686 **M, Li X, MacArthur DG, Nedzel JL, Nguyen DT, Noble MS, Segrè A V.,**  
687 **Trowbridge CA, Tukiainen T, Abell NS, Balliu B, Barshir R, Basha O, Battle A,**  
688 **Bogu GK, Brown A, Brown CD, Castel SE, Chen LS, Chiang C, Conrad DF, Cox**  
689 **NJ, Damani FN, Davis JR, Delaneau O, Dermitzakis ET, Engelhardt BE, Eskin E,**  
690 **Ferreira PG, Frésard L, Gamazon ER, Garrido-Martín D, Gewirtz ADH, Gliner**  
691 **G, Gloude-mans MJ, Guigo R, Hall IM, Han B, He Y, Hormozdiari F, Howald C,**  
692 **Kyung Im H, Jo B, Yong Kang E, Kim Y, Kim-Hellmuth S, Lappalainen T, Li G, Li**  
693 **X, Liu B, Mangul S, McCarthy MI, McDowell IC, Mohammadi P, Monlong J,**  
694 **Montgomery SB, Muñoz-Aguirre M, Ndungu AW, Nicolae DL, Nobel AB, Oliva**  
695 **M, Ongen H, Palowitch JJ, Panousis N, Papasaikas P, Park Y, Parsana P, Payne**  
696 **AJ, Peterson CB, Quan J, Reverter F, Sabatti C, Saha A, Sammeth M, Scott AJ,**  
697 **Shabalin AA, Sodaei R, Stephens M, Stranger BE, Strober BJ, Sul JH, Tsang EK,**  
698 **Urbut S, van de Bunt M, Wang G, Wen X, Wright FA, Xi HS, Yeger-Lotem E,**  
699 **Zappala Z, Zaugg JB, Zhou Y-H, Akey JM, Bates D, Chan J, Chen LS,**  
700 **Claussnitzer M, Demanelis K, Diegel M, Doherty JA, Feinberg AP, Fernando**  
701 **MS, Halow J, Hansen KD, Haugen E, Hickey PF, Hou L, Jasmine F, Jian R, Jiang**  
702 **L, Johnson A, Kaul R, Kellis M, Kibriya MG, Lee K, Billy Li J, Li Q, Li X, Lin J, Lin**  
703 **S, Linder S, Linke C, Liu Y, Maurano MT, Molinie B, Montgomery SB, Nelson J,**  
704 **Neri FJ, Oliva M, Park Y, Pierce BL, Rinaldi NJ, Rizzardi LF, Sandstrom R, Skol**  
705 **A, Smith KS, Snyder MP, Stamatoyannopoulos J, Stranger BE, Tang H, Tsang**  
706 **EK, Wang L, Wang M, Van Wittenberghe N, Wu F, Zhang R, Nierras CR,**  
707 **Branton PA, Carithers LJ, Guan P, Moore HM, Rao A, Vaught JB, Gould SE,**  
708 **Lockart NC, Martin C, Struewing JP, Volpi S, Addington AM, Koester SE, Little**  
709 **AR, Brigham LE, Hasz R, Hunter M, Johns C, Johnson M, Kopen G, Leinweber**  
710 **WF, Lonsdale JT, McDonald A, Mestichelli B, Myer K, Roe B, Salvatore M, Shad**  
711 **S, Thomas JA, Walters G, Washington M, Wheeler J, Bridge J, Foster BA,**  
712 **Gillard BM, Karasik E, Kumar R, Miklos M, Moser MT, Jewell SD, Montroy RG,**

- 713 Rohrer DC, Valley DR, Davis DA, Mash DC, Undale AH, Smith AM, Tabor DE,  
714 Roche N V, McLean JA, Vatanian N, Robinson KL, Sobin L, Barcus ME,  
715 Valentino KM, Qi L, Hunter S, Hariharan P, Singh S, Um KS, Matose T,  
716 Tomaszewski MM, Barker LK, Mosavel M, Siminoff LA, Traino HM, Flicek P,  
717 Juettemann T, Ruffier M, Sheppard D, Taylor K, Trevanion SJ, Zerbino DR,  
718 Craft B, Goldman M, Haeussler M, Kent WJ, Lee CM, Paten B, Rosenbloom KR,  
719 Vivian J, Zhu J. Genetic effects on gene expression across human tissues. *Nature*  
720 550: 204–213, 2017.
- 721 3. Alexander RD, Innocente SA, David Barrass J, Beggs JD. Splicing-Dependent  
722 RNA Polymerase Pausing in Yeast. *Mol Cell* 40: 582–593, 2010.
- 723 4. Allis CD, Jenuwein T. The molecular hallmarks of epigenetic control. *Nat Rev*  
724 *Genet* 17: 487, 2016.
- 725 5. Azzi A, Dallmann R, Casserly A, Rehrauer H, Patrignani A, Maier B, Kramer A,  
726 Brown SA. Circadian behavior is light-reprogrammed by plastic DNA  
727 methylation. *Nat Neurosci* 17: 377–82, 2014.
- 728 6. Bahar R, Hartmann CH, Rodriguez KA, Denny AD, Busuttill RA, Dollé MET,  
729 Calder RB, Chisholm GB, Pollock BH, Klein CA, Vijg J. Increased cell-to-cell  
730 variation in gene expression in ageing mouse heart. *Nature* 441: 1011–1014,  
731 2006.
- 732 7. Barrington WT, Wulfridge P, Wells AE, Rojas CM, Howe SYF, Perry A, Hua K,  
733 Pellizzon MA, Hansen KD, Voy BH, Bennett BJ, Pomp D, Feinberg AP,  
734 Threadgill DW. Improving Metabolic Health Through Precision Dietetics in Mice.  
735 *Genetics* 208: 399–417, 2018.
- 736 8. Bauer T, Trump S, Ishaque N, Thürmann L, Gu L, Bauer M, Bieg M, Gu Z,  
737 Weichenhan D, Mallm J-P, Röder S, Herberth G, Takada E, Mücke O, Winter  
738 M, Junge KM, Grützmann K, Rolle-Kampczyk U, Wang Q, Lawerenz C, Borte M,  
739 Polte T, Schlesner M, Schanne M, Wiemann S, Georg C, Stunnenberg HG, Plass  
740 C, Rippe K, Mizuguchi J, Herrmann C, Eils R, Lehmann I. Environment-induced  
741 epigenetic reprogramming in genomic regulatory elements in smoking mothers  
742 and their children. *Mol Syst Biol* 12: 861, 2016.
- 743 9. Bird AP. Gene number, noise reduction and biological complexity. *Trends Genet*  
744 11: 94–100, 1995.
- 745 10. Blackledge NP, Klose RJ. CpG island chromatin: A platform for gene regulation.  
746 *Epigenetics* 6: 147–152, 2011.
- 747 11. Blake WJ, Balázs G, Kohanski MA, Isaacs FJ, Murphy KF, Kuang Y, Cantor CR,  
748 Walt DR, Collins JJ. Phenotypic Consequences of Promoter-Mediated  
749 Transcriptional Noise. *Mol Cell* 24: 853–865, 2006.
- 750 12. Bódi Z, Farkas Z, Nevozhay D, Kalapis D, Lázár V, Csörgő B, Nyerges Á,  
751 Szamecz B, Fekete G, Papp B, Araújo H, Oliveira JL, Moura G, Santos MAS,  
752 Székely Jr T, Balázs G, Pál C. Phenotypic heterogeneity promotes adaptive  
753 evolution. *PLoS Biol* 15: e2000644, 2017.
- 754 13. Braun KA, Young ET. Coupling mRNA Synthesis and Decay Downloaded from.  
755 *Mol Cell Biol* 34: 4078–4087, 2014.
- 756 14. Brown CR, Mao C, Falkovskaia E, Jurica MS, Boeger H. Linking Stochastic  
757 Fluctuations in Chromatin Structure and Gene Expression. *PLoS Biol* 11:  
758 e1001621, 2013.
- 759 15. Cellerino A, Ori A. What have we learned on aging from omics studies? *Semin Cell*  
760 *Dev Biol* 70: 177–189, 2017.
- 761 16. Chathoth KT, David Barrass J, Webb S, Beggs JD. A Splicing-Dependent

- 762 Transcriptional Checkpoint Associated with Prespliceosome Formation. *Mol Cell*  
763 53: 779–790, 2014.
- 764 17. **Chereji R V, Kan T-W, Grudniewska MK, Romashchenko A V, Berezikov E,**  
765 **Zhimulev IF, Guryev V, Morozov A V, Moshkin YM.** Genome-wide profiling of  
766 nucleosome sensitivity and chromatin accessibility in *Drosophila melanogaster*.  
767 *Nucleic Acids Res* 44: 1036–51, 2016.
- 768 18. **Choi JK, Kim Y-J.** Intrinsic variability of gene expression encoded in nucleosome  
769 positioning sequences. *Nat Genet* 41: 498–503, 2009.
- 770 19. **Churchman LS, Weissman JS.** Nascent transcript sequencing visualizes  
771 transcription at nucleotide resolution. *Nature* 469: 368–73, 2011.
- 772 20. **Collings CK, Anderson JN.** Links between DNA methylation and nucleosome  
773 occupancy in the human genome. *Epigenetics Chromatin* 10: 18, 2017.
- 774 21. **Dattani J, Barahona M.** Stochastic models of gene transcription with upstream  
775 drives: exact solution and sample path characterization. *J R Soc Interface* 14:  
776 20160833, 2017.
- 777 22. **Dey SS, Foley JE, Limsirichai P, Schaffer D V, Arkin AP.** Orthogonal control of  
778 expression mean and variance by epigenetic features at different genomic loci.  
779 *Mol Syst Biol* 11: 806, 2015.
- 780 23. **Dong D, Shao X, Deng N, Zhang Z.** Gene expression variations are predictive for  
781 stochastic noise. *Nucleic Acids Res* 39: 403–413, 2011.
- 782 24. **Ecker S, Chen L, Pancaldi V, Bagger FO, Fernández JM, Carrillo de Santa Pau**  
783 **E, Juan D, Mann AL, Watt S, Casale FP, Sidiropoulos N, Rapin N, Merkel A,**  
784 **BLUEPRINT Consortium HG, Stunnenberg HG, Stegle O, Frontini M, Downes**  
785 **K, Pastinen T, Kuijpers TW, Rico D, Valencia A, Beck S, Soranzo N, Paul DS.**  
786 Genome-wide analysis of differential transcriptional and epigenetic variability  
787 across human immune cell types. *Genome Biol* 18: 18, 2017.
- 788 25. **Ecker S, Pancaldi V, Valencia A, Beck S, Paul DS.** Epigenetic and Transcriptional  
789 Variability Shape Phenotypic Plasticity. *BioEssays* 40: 1700148, 2018.
- 790 26. **Elowitz MB, Levine AJ, Siggia ED, Swain PS.** Stochastic Gene Expression in a  
791 Single Cell. *Science* 297: 1183–1186, 2002.
- 792 27. **Faure AJ, Schmiedel JM, Lehner B.** Systematic Analysis of the Determinants of  
793 Gene Expression Noise in Embryonic Stem Cells. *Cell Syst* 5: 471–484.e4, 2017.
- 794 28. **Feinberg AP, Irizarry RA.** Evolution in health and medicine Sackler colloquium:  
795 Stochastic epigenetic variation as a driving force of development, evolutionary  
796 adaptation, and disease. *Proc Natl Acad Sci U S A* 107 Suppl 1: 1757–64, 2010.
- 797 29. **Fraga MF, Ballestar E, Paz MF, Ropero S, Setien F, Ballestar ML, Heine-Suñer**  
798 **D, Cigudosa JC, Urioste M, Benitez J, Boix-Chornet M, Sanchez-Aguilera A,**  
799 **Ling C, Carlsson E, Poulsen P, Vaag A, Stephan Z, Spector TD, Wu Y-Z, Plass C,**  
800 **Esteller M.** Epigenetic differences arise during the lifetime of monozygotic twins.  
801 *Proc Natl Acad Sci U S A* 102: 10604–9, 2005.
- 802 30. **Fraser HB, Hirsh AE, Giaever G, Kumm J, Eisen MB.** Noise Minimization in  
803 Eukaryotic Gene Expression. *PLoS Biol* 2: e137, 2004.
- 804 31. **Gonzalez-Porta M, Calvo M, Sammeth M, Guigo R.** Estimation of alternative  
805 splicing variability in human populations. *Genome Res* 22: 528–538, 2012.
- 806 32. **Guan J, Chen M, Ye C, Cai JJ, Ji G.** AEGS: identifying aberrantly expressed gene  
807 sets for differential variability analysis. *Bioinformatics* 34: 881–883, 2018.
- 808 33. **Haimovich G, Choder M, Singer RH, Trcek T.** The fate of the messenger is pre-  
809 determined: a new model for regulation of gene expression. *Biochim Biophys Acta*  
810 1829: 643–53, 2013.

- 811 34. **Hannum G, Guinney J, Zhao L, Zhang L, Hughes G, Satta S, Klotzle B, Bibikova**  
812 **M, Fan J-B, Gao Y, Deconde R, Chen M, Rajapakse I, Friend S, Ideker T, Zhang**  
813 **K.** Genome-wide methylation profiles reveal quantitative views of human aging  
814 rates. *Mol Cell* 49: 359–367, 2013.
- 815 35. **Hastie T, Tibshirani R, Friedman J.** Basis Expansions and Regularization. In: *The*  
816 *Elements of Statistical Learning*. Springer, p. 139–189, 2009.
- 817 36. **Hornung G, Bar-Ziv R, Rosin D, Tokuriki N, Tawfik DS, Oren M, Barkai N.**  
818 Noise-mean relationship in mutated promoters. *Genome Res* 22: 2409–2417,  
819 2012.
- 820 37. **Huang S.** Non-genetic heterogeneity of cells in development: more than just  
821 noise. *Development* 136: 3853–3862, 2009.
- 822 38. **Huber W, von Heydebreck A, Sultmann H, Poustka A, Vingron M.** Variance  
823 stabilization applied to microarray data calibration and to the quantification of  
824 differential expression. *Bioinformatics* 18: S96–S104, 2002.
- 825 39. **Huh I, Zeng J, Park T, Yi S V.** DNA methylation and transcriptional noise.  
826 *Epigenetics Chromatin* 6: 9, 2013.
- 827 40. **Iyer-Biswas S, Jayaprakash C.** Mixed Poisson distributions in exact solutions of  
828 stochastic autoregulation models. *Phys Rev E* 90: 052712, 2014.
- 829 41. **Johnson TL, Ares M.** SMITten by the Speed of Splicing. *Nat Rev Endocrinol* 106:  
830 219–246, 2016.
- 831 42. **Jonkers I, Lis JT.** Getting up to speed with transcription elongation by RNA  
832 polymerase II HHS Public Access. *Nat Rev Mol Cell Biol* 16: 167–177, 2015.
- 833 43. **Kanehisa M, Sato Y, Kawashima M, Furumichi M, Tanabe M.** KEGG as a  
834 reference resource for gene and protein annotation. *Nucleic Acids Res* 44: D457–  
835 D462, 2016.
- 836 44. **Karlis D, Xekalaki E.** Mixed Poisson Distributions. *Internat Stat Rev* 73: 35–58,  
837 2005.
- 838 45. **Kornienko AE, Dotter CP, Guenzl PM, Gisslinger H, Gisslinger B, Cleary C,**  
839 **Kralovics R, Pauler FM, Barlow DP.** Long non-coding RNAs display higher  
840 natural expression variation than protein-coding genes in healthy humans.  
841 *Genome Biol* 17: 14, 2016.
- 842 46. **Kwak H, Lis JT.** Control of Transcriptional Elongation. *Annu Rev Genet* 47: 483–  
843 508, 2013.
- 844 47. **Lappalainen T, Sammeth M, Friedländer MR, 't Hoen PAC, Monlong J, Rivas**  
845 **MA, González-Porta M, Kurbatova N, Griebel T, Ferreira PG, Barann M,**  
846 **Wieland T, Greger L, van Iterson M, Almlöf J, Ribeca P, Pulyakhina I, Esser D,**  
847 **Giger T, Tikhonov A, Sultan M, Bertier G, MacArthur DG, Lek M, Lizano E,**  
848 **Buermans HPJ, Padioleau I, Schwarzmayer T, Karlberg O, Ongen H, Kilpinen**  
849 **H, Beltran S, Gut M, Kahlem K, Amstislavskiy V, Stegle O, Pirinen M,**  
850 **Montgomery SB, Donnelly P, McCarthy MI, Flicek P, Strom TM, Lehrach H,**  
851 **Schreiber S, Sudbrak R, Carracedo Á, Antonarakis SE, Häsler R, Syvänen A-C,**  
852 **van Ommen G-J, Brazma A, Meitinger T, Rosenstiel P, Guigó R, Gut IG, Estivill**  
853 **X, Dermitzakis ET, Dermitzakis ET.** Transcriptome and genome sequencing  
854 uncovers functional variation in humans. *Nature* 501: 506–511, 2013.
- 855 48. **Lehner B.** Selection to minimise noise in living systems and its implications for  
856 the evolution of gene expression. *Mol Syst Biol* 4: 170, 2008.
- 857 49. **Lehner B, Kaneko K.** Fluctuation and response in biology. *Cell Mol Life Sci* 68:  
858 1005–1010, 2011.
- 859 50. **Li J, Liu Y, Kim T, Min R, Zhang Z.** Gene Expression Variability within and

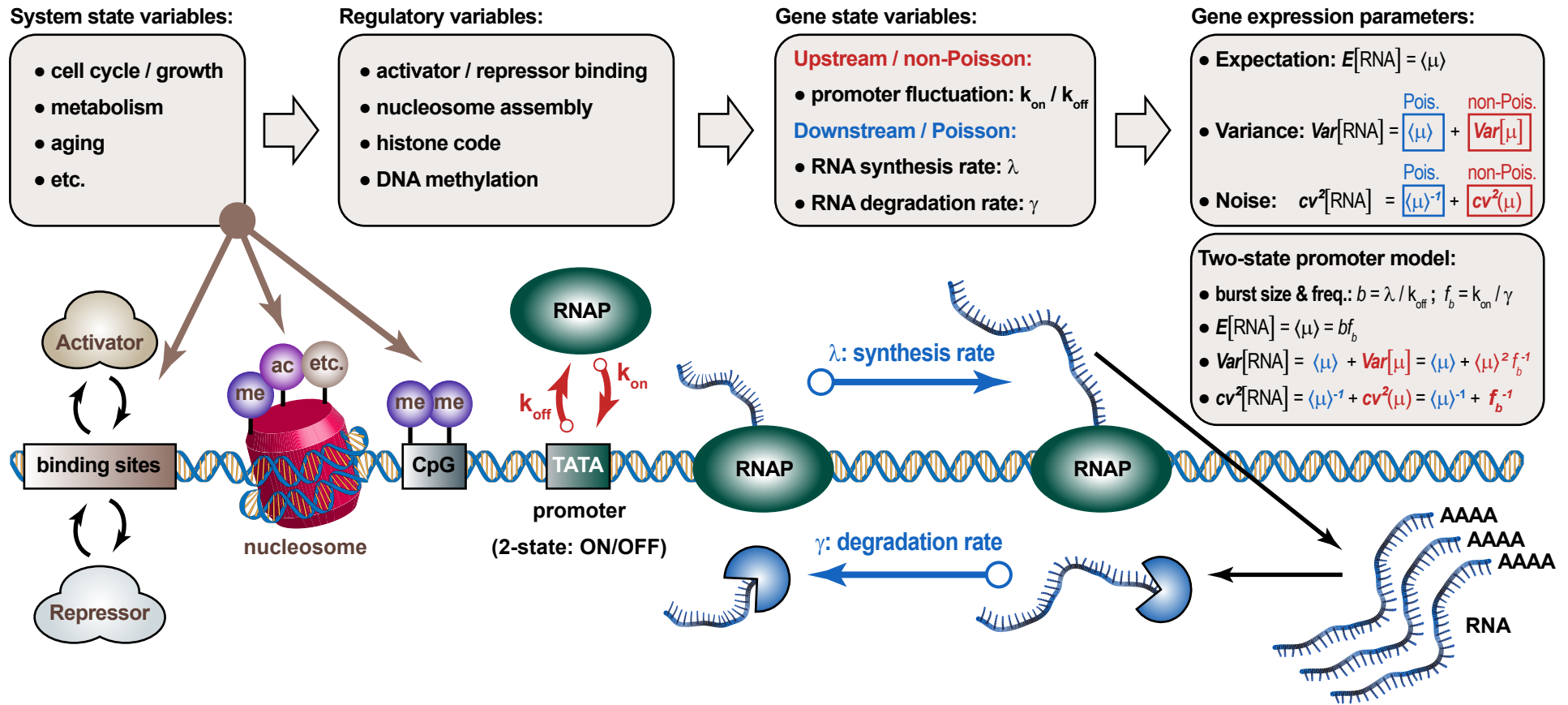


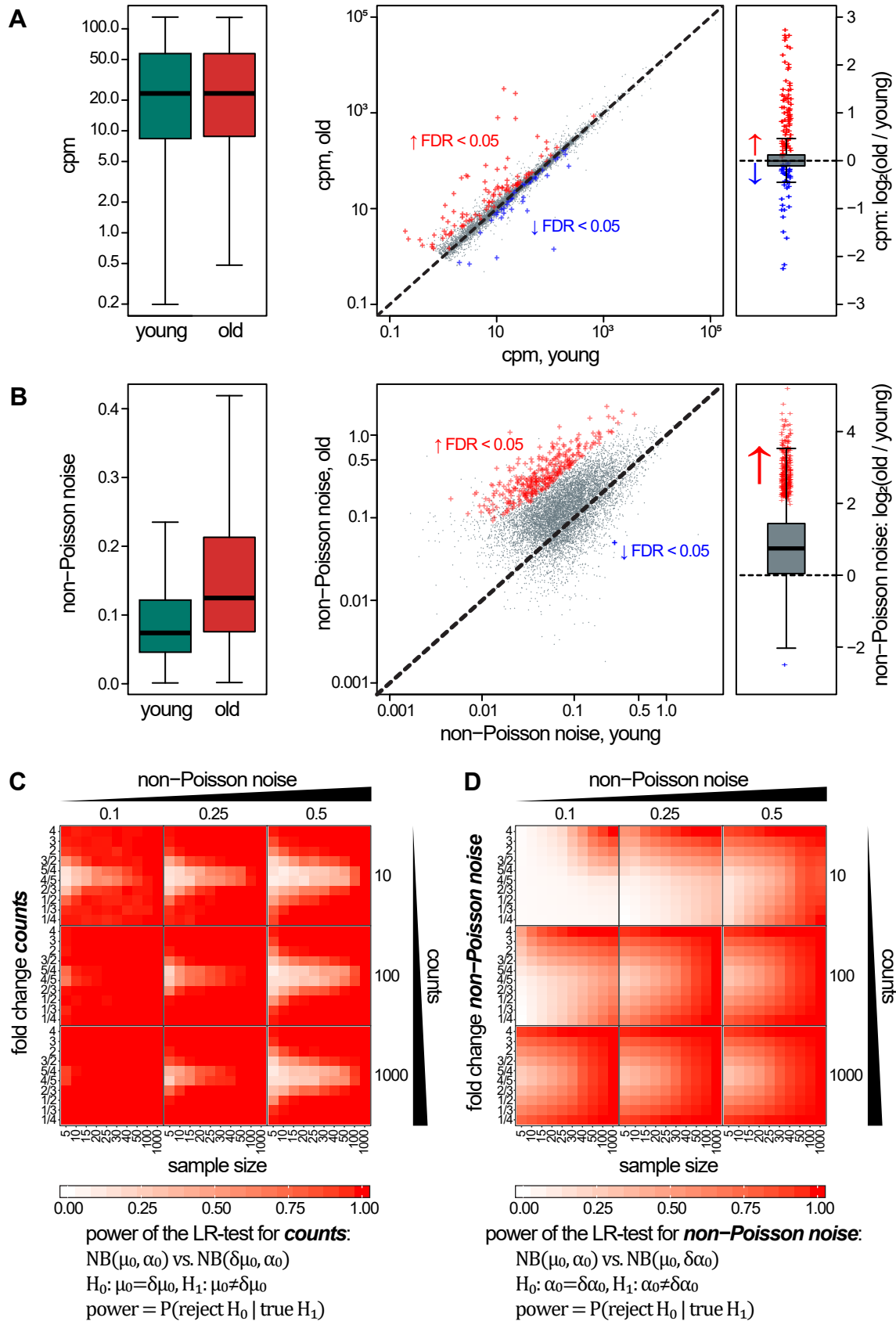
- 860 between Human Populations and Implications toward Disease Susceptibility.  
861 *PLoS Comput Biol* 6: e1000910, 2010.
- 862 51. **Liu J, Wu X, Zhang H, Pfeifer GP, Lu Q.** Dynamics of RNA Polymerase II Pausing  
863 and Bivalent Histone H3 Methylation during Neuronal Differentiation in Brain  
864 Development. *Cell Rep* 20: 1307–1318, 2017.
- 865 52. **Love MI, Huber W, Anders S.** Moderated estimation of fold change and  
866 dispersion for RNA-seq data with DESeq2. *Genome Biol* 15: 550, 2014.
- 867 53. **Low TY, van Heesch S, van den Toorn H, Giansanti P, Cristobal A, Toonen P,  
868 Schafer S, Hübner N, van Breukelen B, Mohammed S, Cuppen E, Heck AJR,  
869 Guryev V.** Quantitative and qualitative proteome characteristics extracted from  
870 in-depth integrated genomics and proteomics analysis. *Cell Rep* 5: 1469–78, 2013.
- 871 54. **MacNeil LT, Walhout AJM.** Gene regulatory networks and the role of robustness  
872 and stochasticity in the control of gene expression. *Genome Res* 21: 645–657,  
873 2011.
- 874 55. **Martinez-Jimenez CP, Eling N, Chen H-C, Vallejos CA, Kolodziejczyk AA,  
875 Connor F, Stojic L, Rayner TF, Stubbington MJT, Teichmann SA, de la Roche  
876 M, Marioni JC, Odom DT.** Aging increases cell-to-cell transcriptional variability  
877 upon immune stimulation. *Science* 355: 1433–1436, 2017.
- 878 56. **Mason EA, Mar JC, Laslett AL, Pera MF, Quackenbush J, Wolvetang E, Wells  
879 CA.** Gene expression variability as a unifying element of the pluripotency  
880 network. *Stem cell reports* 3: 365–77, 2014.
- 881 57. **McAdams HH, Arkin A.** Stochastic mechanisms in gene expression. *Proc Natl  
882 Acad Sci U S A* 94: 814–9, 1997.
- 883 58. **McCarthy DJ, Chen Y, Smyth GK.** Differential expression analysis of multifactor  
884 RNA-Seq experiments with respect to biological variation. *Nucleic Acids Res* 40:  
885 4288–4297, 2012.
- 886 59. **Metzger BPH, Yuan DC, Gruber JD, Duveau F, Wittkopp PJ.** Selection on noise  
887 constrains variation in a eukaryotic promoter. *Nature* 521: 344–7, 2015.
- 888 60. **Morgan MD, Marioni JC.** CpG island composition differences are a source of gene  
889 expression noise indicative of promoter responsiveness. *Genome Biol* 19: 1–13,  
890 2018.
- 891 61. **Müller C, Zidek LM, Ackermann T, de Jong T, Liu P, Kliche V, Zaini MA,  
892 Kortman G, Harkema L, Verbeek DS, Tuckermann JP, von Maltzahn J, de  
893 Bruin A, Guryev V, Wang Z-Q, Calkhoven CF.** Reduced expression of C/EBPβ-  
894 LIP extends health and lifespan in mice. *Elife* 7: e34985, 2018.
- 895 62. **Newman JRS, Ghaemmaghami S, Ihmels J, Breslow DK, Noble M, DeRisi JL,  
896 Weissman JS.** Single-cell proteomic analysis of *S. cerevisiae* reveals the  
897 architecture of biological noise. *Nature* 441: 840–846, 2006.
- 898 63. **Nicolas D, Zoller B, Suter DM, Naef F.** Modulation of transcriptional burst  
899 frequency by histone acetylation. *Proc Natl Acad Sci U S A* 115: 7153–7158, 2018.
- 900 64. **Nie L, Wu G, Culley DE, Scholten JCM, Zhang W.** Integrative Analysis of  
901 Transcriptomic and Proteomic Data: Challenges, Solutions and Applications. *Crit  
902 Rev Biotechnol* 27: 63–75, 2007.
- 903 65. **Nojima T, Gomes T, Grosso ARF, Kimura H, Dye MJ, Dhir S, Carmo-Fonseca M,  
904 Proudfoot NJ.** Mammalian NET-Seq Reveals Genome-wide Nascent Transcription  
905 Coupled to RNA Processing. *Cell* 161: 526–540, 2015.
- 906 66. **Park SJ, Song S, Yang G-S, Kim PM, Yoon S, Kim J-H, Sung J.** The Chemical  
907 Fluctuation Theorem governing gene expression. *Nat Commun* 9: 297, 2018.
- 908 67. **Parker R, Song H.** The enzymes and control of eukaryotic mRNA turnover. *Nat*

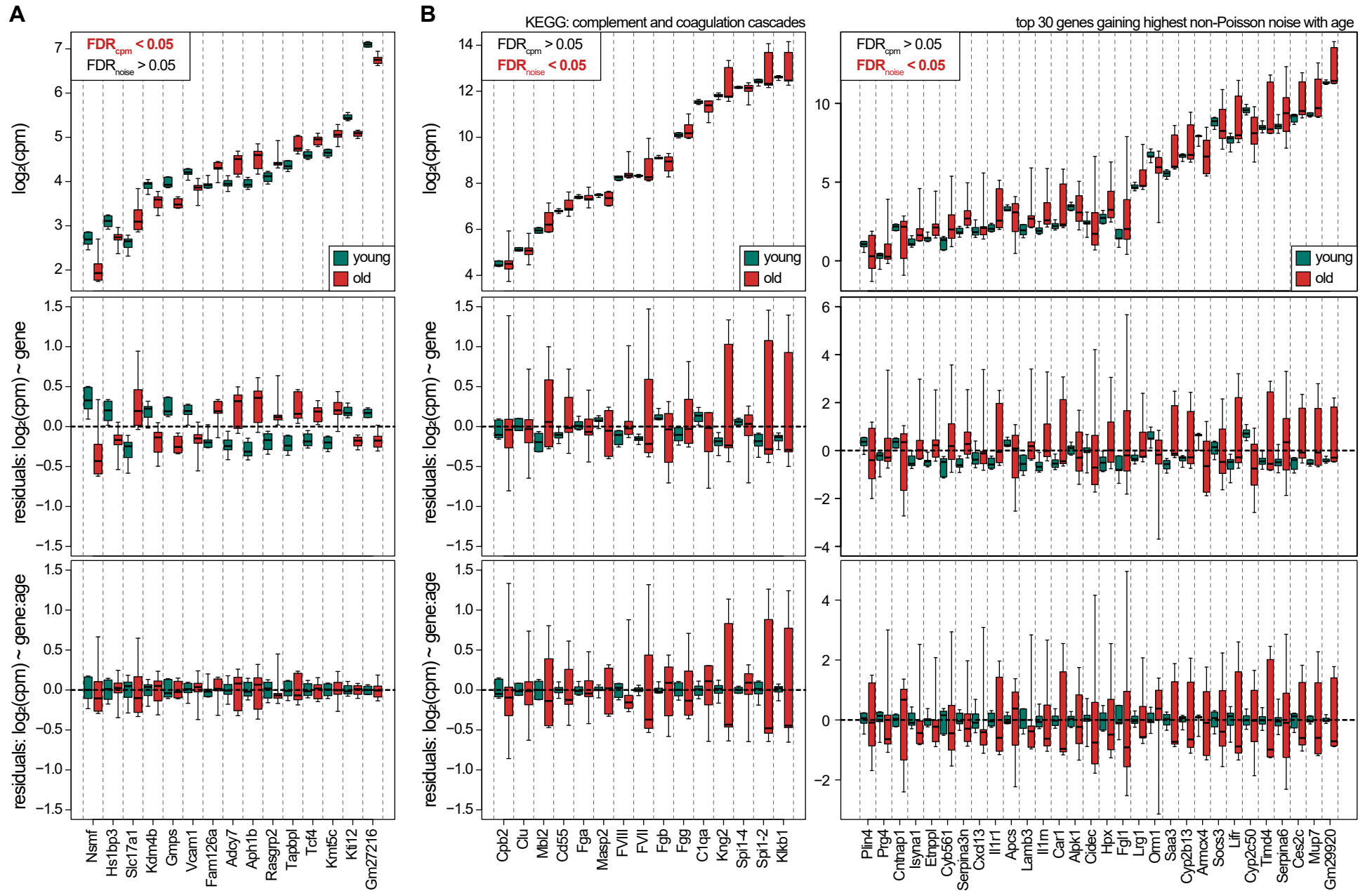
- 909 *Struct Mol Biol* 11: 121–127, 2004.
- 910 68. **Paulsson J.** Models of stochastic gene expression. *Phys Life Rev* 2: 157–175, 2005.
- 911 69. **Pedraza J, van Oudenaarden A.** Noise propagation in gene networks. *Science*
- 912 307: 1965–1974, 2005.
- 913 70. **Potter SS.** Single-cell RNA sequencing for the study of development, physiology
- 914 and disease. *Nat Rev Nephrol* 14: 479–492, 2018.
- 915 71. **Raj A, van Oudenaarden A.** Nature, Nurture, or Chance: Stochastic Gene
- 916 Expression and Its Consequences. *Cell* 135: 216–226, 2008.
- 917 72. **Raj A, Peskin CS, Tranchina D, Vargas DY, Tyagi S.** Stochastic mRNA synthesis
- 918 in mammalian cells. *PLoS Biol* 4: e309, 2006.
- 919 73. **Raj A, Rifkin SA, Andersen E, van Oudenaarden A.** Variability in gene
- 920 expression underlies incomplete penetrance. *Nature* 463: 913–918, 2010.
- 921 74. **Ran D, Daye ZJ.** Gene expression variability and the analysis of large-scale RNA-
- 922 seq studies with the MDSeq. *Nucleic Acids Res* 45: e127–e127, 2017.
- 923 75. **Rangaraju S, Solis GM, Thompson RC, Gomez-Amaro RL, Kurian L, Encalada**
- 924 **SE, Niculescu AB, Salomon DR, Petrascheck M.** Suppression of transcriptional
- 925 drift extends *C. elegans* lifespan by postponing the onset of mortality. *Elife* 4:
- 926 e08833, 2015.
- 927 76. **Raser JM, O’Shea EK.** Control of Stochasticity in Eukaryotic Gene Expression.
- 928 *Science* 304: 1811–1814, 2004.
- 929 77. **Ravarani CNJ, Chalancon G, Breker M, de Groot NS, Babu MM.** Affinity and
- 930 competition for TBP are molecular determinants of gene expression noise. *Nat*
- 931 *Commun* 7: 10417, 2016.
- 932 78. **Rigby RA, Stasinopoulos DM.** Generalized additive models for location, scale and
- 933 shape. *J R Stat Soc Ser C* 54: 507–554, 2005.
- 934 79. **Rigby RA, Stasinopoulos DM, Akantziliotou C.** A framework for modelling
- 935 overdispersed count data, including the Poisson-shifted generalized inverse
- 936 Gaussian distribution. *Comput Stat Data Anal* 53: 381–393, 2008.
- 937 80. **Risso D, Ngai J, Speed TP, Dudoit S.** Normalization of RNA-seq data using factor
- 938 analysis of control genes or samples. *Nat Biotechnol* 32: 896–902, 2014.
- 939 81. **Sainsbury S, Bernecky C, Cramer P.** Structural basis of transcription initiation
- 940 by RNA polymerase II. *Nat Rev Mol Cell Biol* 16: 129–143, 2015.
- 941 82. **Sanchez A, Choubey S, Kondev J.** Regulation of Noise in Gene Expression. *Annu*
- 942 *Rev Biophys* 42: 469–491, 2013.
- 943 83. **Sanchez A, Golding I.** Genetic determinants and cellular constraints in noisy gene
- 944 expression. *Science* 342: 1188–93, 2013.
- 945 84. **Sato K, Ito Y, Yomo T, Kaneko K.** On the relation between fluctuation and
- 946 response in biological systems. *Proc Natl Acad Sci* 100: 14086–14090, 2003.
- 947 85. **Sears KE, Maier JA, Rivas-Astroza M, Poe R, Zhong S, Kosog K, Marcot JD,**
- 948 **Behringer RR, Cretekos CJ, Rasweiler JJ, Rapti Z.** The Relationship between
- 949 Gene Network Structure and Expression Variation among Individuals and Species.
- 950 *PLoS Genet* 11: e1005398, 2015.
- 951 86. **Sen P, Dang W, Donahue G, Dai J, Dorsey J, Cao X, Liu W, Cao K, Perry R, Lee**
- 952 **JY, Wasko BM, Carr DT, He C, Robison B, Wagner J, Gregory BD, Kaerberlein M,**
- 953 **Kennedy BK, Boeke JD, Berger SL.** H3K36 methylation promotes longevity by
- 954 enhancing transcriptional fidelity. *Genes Dev* 29: 1362–76, 2015.
- 955 87. **Sharon E, van Dijk D, Kalma Y, Keren L, Manor O, Yakhini Z, Segal E.** Probing
- 956 the effect of promoters on noise in gene expression using thousands of designed
- 957 sequences. *Genome Res* 24: 1698–706, 2014.

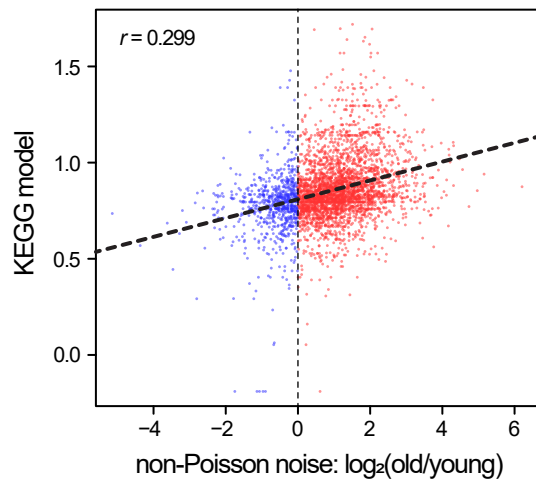
- 958 88. **Sherman MS, Lorenz K, Lanier MH, Cohen BA.** Cell-to-Cell Variability in the  
959 Propensity to Transcribe Explains Correlated Fluctuations in Gene Expression.  
960 *Cell Syst* 1: 315–325, 2015.
- 961 89. **Slieker RC, van Iterson M, Luijk R, Beekman M, Zhernakova D V., Moed MH,**  
962 **Mei H, van Galen M, Deelen P, Bonder MJ, Zhernakova A, Uitterlinden AG,**  
963 **Tigchelaar EF, Stehouwer CDA, Schalkwijk CG, van der Kallen CJH, Hofman A,**  
964 **van Heemst D, de Geus EJ, van Dongen J, Deelen J, van den Berg LH, van**  
965 **Meurs J, Jansen R, 't Hoen PAC, Franke L, Wijmenga C, Veldink JH, Swertz MA,**  
966 **van Greevenbroek MMJ, van Duijn CM, Boomsma DI, Slagboom PE, Heijmans**  
967 **BT.** Age-related accrual of methylomic variability is linked to fundamental ageing  
968 mechanisms. *Genome Biol* 17: 1–13, 2016.
- 969 90. **Small EC, Xi L, Wang J-P, Widom J, Licht JD.** Single-cell nucleosome mapping  
970 reveals the molecular basis of gene expression heterogeneity. *Proc Natl Acad Sci U*  
971 *SA* 111: E2462-71, 2014.
- 972 91. **Stasinopoulos MD, Rigby RA, Heller GZ, Voudouris V, Bastiani F De, Rigby RA,**  
973 **Heller GZ, Voudouris V, Bastiani F De.** Flexible Regression and Smoothing. CRC  
974 Press, 2017.
- 975 92. **Sultan M, Amstislavskiy V, Risch T, Schuette M, Dökel S, Ralser M, Balzereit**  
976 **D, Lehrach H, Yaspo M-L.** Influence of RNA extraction methods and library  
977 selection schemes on RNA-seq data. *BMC Genomics* 15: 675, 2014.
- 978 93. **Suter DM, Molina N, Gatfield D, Schneider K, Schibler U, Naef F.** Mammalian  
979 genes are transcribed with widely different bursting kinetics. *Science* 332: 472–4,  
980 2011.
- 981 94. **Swain PS, Elowitz MB, Siggia ED.** Intrinsic and extrinsic contributions to  
982 stochasticity in gene expression. *Proc Natl Acad Sci* 99: 12795–12800, 2002.
- 983 95. **Tan Q, Heijmans BT, Hjelmberg JVB, Soerensen M, Christensen K,**  
984 **Christiansen L.** Handling blood cell composition in epigenetic studies on ageing.  
985 *Int J Epidemiol* 46: 1717–1718, 2017.
- 986 96. **Thattai M.** Universal Poisson Statistics of mRNAs with Complex Decay Pathways.  
987 *Biophys J* 110: 301–305, 2016.
- 988 97. **Thattai M, van Oudenaarden A.** Intrinsic noise in gene regulatory networks.  
989 *Proc Natl Acad Sci* 98: 8614–8619, 2001.
- 990 98. **Tirosh I, Barkai N.** Two strategies for gene regulation by promoter nucleosomes.  
991 *Genome Res* 18: 1084–91, 2008.
- 992 99. **de Torrente L, Zimmerman S, Taylor D, Hasegawa Y, Wells CA, Mar JC.**  
993 *pathVar*: a new method for pathway-based interpretation of gene expression  
994 variability. *PeerJ* 5: e3334, 2017.
- 995 100. **Vaquero-Garcia J, Barrera A, Gazzara MR, González-Vallinas J, Lahens NF,**  
996 **Hogenesch JB, Lynch KW, Barash Y.** A new view of transcriptome complexity  
997 and regulation through the lens of local splicing variations. *Elife* 5: e11752, 2016.
- 998 101. **Viñuela A, Brown AA, Buil A, Tsai P-C, Davies MN, Bell JT, Dermitzakis ET,**  
999 **Spector TD, Small KS.** Age-dependent changes in mean and variance of gene  
1000 expression across tissues in a twin cohort. *Hum Mol Genet* 27: 732–741, 2018.
- 1001 102. **Voisin S, Eynon N, Yan X, Bishop DJ.** Exercise training and DNA methylation in  
1002 humans. *Acta Physiol (Oxf)* 213: 39–59, 2015.
- 1003 103. **Wan Y, Larson DR.** Splicing heterogeneity: separating signal from noise. *Genome*  
1004 *Biol* 19: 86, 2018.
- 1005 104. **Wang Z, Zhang J.** Impact of gene expression noise on organismal fitness and the  
1006 efficacy of natural selection. *Proc Natl Acad Sci* 108: E67–E76, 2011.

- 1007 105. **White RR, Milholland B, MacRae SL, Lin M, Zheng D, Vijg J.** Comprehensive  
1008 transcriptional landscape of aging mouse liver. *BMC Genomics* 16: 899, 2015.
- 1009 106. **Wolf L, Silander OK, van Nimwegen E.** Expression noise facilitates the evolution  
1010 of gene regulation. *Elife* 4: e05856, 2015.
- 1011 107. **Wu L, Fan J, Belasco JG.** MicroRNAs direct rapid deadenylation of mRNA. *Proc*  
1012 *Natl Acad Sci* 103: 4034–4039, 2006.
- 1013 108. **Wu S, Li K, Li Y, Zhao T, Li T, Yang Y-F, Qian W.** Independent regulation of gene  
1014 expression level and noise by histone modifications. *PLOS Comput Biol* 13:  
1015 e1005585, 2017.
- 1016 109. **Yamashita A, Chang T-C, Yamashita Y, Zhu W, Zhong Z, Chen C-YA, Shyu A-B.**  
1017 Concerted action of poly(A) nucleases and decapping enzyme in mammalian  
1018 mRNA turnover. *Nat Struct Mol Biol* 12: 1054–63, 2005.
- 1019 110. **Yip SH, Sham PC, Wang J.** Evaluation of tools for highly variable gene discovery  
1020 from single-cell RNA-seq data. *Brief. Bioinform.* bby011, 2018.
- 1021 111. **Yu Y, Fuscoe JC, Zhao C, Guo C, Jia M, Qing T, Bannon DI, Lancashire L, Bao W,**  
1022 **Du T, Luo H, Su Z, Jones WD, Moland CL, Branham WS, Qian F, Ning B, Li Y,**  
1023 **Hong H, Guo L, Mei N, Shi T, Wang KY, Wolfinger RD, Nikolsky Y, Walker SJ,**  
1024 **Duerksen-Hughes P, Mason CE, Tong W, Thierry-Mieg J, Thierry-Mieg D, Shi**  
1025 **L, Wang C.** A rat RNA-Seq transcriptomic BodyMap across 11 organs and 4  
1026 developmental stages. *Nat Commun* 5: 3230, 2014.
- 1027 112. **Zaugg JB, Luscombe NM.** A genomic model of condition-specific nucleosome  
1028 behavior explains transcriptional activity in yeast. *Genome Res* 22: 84–94, 2012.
- 1029 113. **Zhang H-M, Diaz V, Walsh ME, Zhang Y.** Moderate lifelong overexpression of  
1030 tuberous sclerosis complex 1 (TSC1) improves health and survival in mice. *Sci*  
1031 *Rep* 7: 834, 2017.
- 1032 114. **Zhang Z, Qian W, Zhang J.** Positive selection for elevated gene expression noise  
1033 in yeast. *Mol Syst Biol* 5: 299, 2009.
- 1034 115. **Zhao S, Fung-Leung W-P, Bittner A, Ngo K, Liu X.** Comparison of RNA-Seq and  
1035 microarray in transcriptome profiling of activated T cells. *PLoS One* 9: e78644,  
1036 2014.
- 1037 116. **Zoller B, Nicolas D, Molina N, Naef F.** Structure of silent transcription intervals  
1038 and noise characteristics of mammalian genes. *Mol Syst Biol* 11: 823, 2015.
- 1039







**A****B**