1  **Widespread population variability of intron size in evolutionary old genes:**

2  **implications for gene expression variability**

3  **Intronic CNVs and gene regulation**

4

5  Maria Rigau[1], David Juan[2], Alfonso Valencia[1,3,¶,*] and Daniel Rico[4,¶,*]

6

7  [1]Barcelona Supercomputing Centre (BSC), Barcelona, 08034, Spain.

8  [2]Institut de Biologia Evolutiva, Consejo Superior de Investigaciones Científicas–Universitat

9  Pompeu Fabra, Parc de Recerca Biomèdica de Barcelona, Barcelona, 08003, Spain

10  [3]Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, 08010, Spain.

11  [4]Institute of Cellular Medicine, Newcastle University, Newcastle upon Tyne, NE2 4HH, United

12  Kingdom.

13

14  [¶] These authors contributed equally to this work.

15  [*] Corresponding authors

16  E-mails: alfonso.valencia@bsc.es, daniel.rico@newcastle.ac.uk

## Abstract

Introns were originally thought to be "junk DNA" without function but accumulating evidence has shown that they can have important functions in the regulation of gene expression. In humans and other mammals, introns can be extraordinarily large and together they account for the majority of the sequence in human protein-coding loci. However, little is known about their structural variation in human populations and the potential functional impact of this genomic variation. To address this, we have studied how copy number variants (CNVs) differentially affect exonic and intronic sequences of protein-coding genes. Using five different CNV maps, we found that CNV gains and losses are consistently underrepresented in coding regions. However, we found purely intronic losses in protein-coding genes more frequently than expected by chance, even in essential genes. Following a phylogenetic approach, we dissected how CNV losses differentially affect genes depending on their evolutionary age. Evolutionarily young genes frequently overlap with deletions that partially or entirely eliminate their coding sequence, while in evolutionary ancient genes the losses of intronic DNA are the most frequent CNV type. A detailed characterisation of these events showed that the loss of intronic sequence can be associated with significant differences in gene length and expression levels in the population. In summary, we show that genomic variation is shaping gene evolution in different ways depending on the age and function of genes. CNVs affecting introns can exert an important role in maintaining the variability of gene expression in human populations, a variability that could be related with human adaptation.

## Author summary

Most human genes have introns that have to be removed after a gene is transcribed from DNA to RNA because they not encode information to translate RNA into proteins. As mutations in introns do not affect protein sequences, they are usually ignored when looking for normal or pathogenic genomic variation. However, introns comprise about half of the human genome and they can have important regulatory roles. We show that deletions of intronic regions appear more frequent than previously expected in the healthy population, with a significant proportion of genes with evolutionary ancient and essential functions carrying them. This finding was very surprising, as ancient genes tend to have high conservation of their coding sequence. However, we show that deletions of their non-coding intronic sequence can produce considerable changes in their length, significant drops of GC content that could affect splicing or occur in introns harboring regulatory elements. Finally, we found that a significant number of these intronic deletions are associated with under- or over-expression of the affected genes, showing that intronic deletions can be responsible for gene expression variability in ancient genes with highly conserved protein sequences. Our data suggests that the frequent gene length variation in ancient genes resulting from intronic CNVs might have an important role in the fine-tuning of their regulation in different individuals.

## Introduction

Most eukaryotic protein coding genes contain introns that are removed from the messenger RNA during the process of splicing. Although the potential functions of introns remain elusive, a number of cases support the idea that the potential energetic disadvantage that they represent for the cell might be compensated by a number of acquired functionalities [1–3]. For example, introns make possible the expression of multiple transcription products from a single gene by alternative splicing and facilitate the formation of new genes by exon shuffling [3,4].

Human introns are longer than those of other vertebrates and invertebrates [5,6] and their lengths are very variable, contrarily to exon lengths. This variability in intron length leads to substantial differences in size among human genes, which cause differences in the time taken to transcribe a gene from seconds to over 24 hours [7]. Introns can influence several steps of gene transcription [8,9] and it has been seen that a considerable amount of intronic sequence is important in regulating gene expression [10].

Introns contribute to the control of gene expression by their inclusion of regulatory regions and non-coding functional RNA genes or directly by their length [3,11,12]. Indeed, intron size is highly conserved in genes associated with developmental patterning [13], suggesting that genes that require a precise time coordination of their transcription are reliant on a consistent transcript length. Highly expressed genes are enriched in housekeeping essential functions [14] and tend to have shorter introns [15]. It has been suggested that selection could be acting to reduce the costs of transcription by keeping short introns in highly expressed genes [15]. Genes transcribed early in development [16–18] and genes involved in rapid biological responses [19] also conserve intron-

4

73    poor structures. Shorter introns would allow these genes to be transcribed faster and thus they may

74    be particularly sensitive to changes in the time to be transcribed. Interestingly, Keane and Seoighe

75    [20] recently found that intron lengths of coexpressed genes or genes participating in the same

76    protein complexes tend to coevolve (their intron sizes show a significant correlation across species)

77    possibly because a precise temporal regulation of the expression of these genes is required.

78    Despite their potential importance in regulating transcription levels, transcription timing and

79    splicing, little attention has been payed to the potential role of introns in human population

80    variability studies. Given that direct associations between intronic mutations and certain diseases

81    have been reported [24–27], we need to characterise the normal genetic variability in introns so we

82    can better distinguish normal from pathogenic variations.

83    It is increasingly apparent that one of the most important sources of variability is the presence of

84    copy number variants (CNVs). CNVs are defined as imbalanced structural variants that result in the

85    gain or loss of >50 bp of genomic sequence [28] and appear in more than one individual of the

86    population. CNVs can be classified in gains (regions that are found duplicated when compared with

87    expected number from the reference genome, which is 2 for autosomes), losses (homozygously or

88    heterozygously deleted regions) or gain/loss CNVs (regions that are found duplicated in some

89    individuals - or alleles - and deleted in others). Microarray and next generation sequencing

90    approaches have shown that CNVs are more important and frequent than originally thought. CNVs

91    may have neutral, advantageous or pathological consequences [29]. Initially, CNVs thought to

92    account for about 1% of the entire human genome [30], current estimates range from 4.8 to 9.5%

93    [31].

94    Here, we have studied the effect of intronic variants using CNV maps of high resolution. Most of

95    these CNVs have been detected using whole genome sequencing (WGS) data, which allows to

96    determine the exact genomic boundaries of these variants and thus their overlap with exons and/or

97    introns. Despite significant advances in the detection of CNVs, discovery and genotyping of these

98    variants remain challenging [32]. To gain consistency, we have analyzed in parallel 5 recent CNV

99    maps obtained by different groups with different experimental and analysis systems [31,33–36]. We

100    have been able to compare the effect of the CNVs overlapping totally or partially with the genes of

101    different evolutionary ages, studying in depth the effect of the intronic variants. We show how

102    intronic variation results in widespread gene length variability in human populations and the

103    potential impact of this variability in splicing and gene expression.

## Results

### Most genic CNVs fall within introns

CNVs can affect genes in different ways depending on the degree of overlap with them. Some CNVs cover entire genes (from now on *whole gene CNVs*), other CNVs overlap with part of the coding sequence but not the whole gene (*exonic CNVs*) and other CNVs are found within intronic regions (*intronic CNVs*, **Fig 1A**). As defined, intronic CNVs do not overlap with exons from any annotated transcript isoforms or with exons from overlapping genes.

To analyze the impact of CNVs on protein coding genes in healthy humans, we used five recently published, high resolution CNV maps [31,33–36]. Each of the maps has been derived from a different number of individuals, from different populations and using different techniques and algorithms for CNV detection (**S1 Fig** and **S1 Table)**. Due to these differences, each dataset provides us with a different set of CNVs (**S1 Fig**), which we analysed independently. We only considered the variants present in at least 2 individuals in a dataset, filtering out the variants mapped in sex chromosomes and the private variants within each map.

The total number of autosomal protein coding genes overlapping with common CNVs varies depending on the filtered map, ranging from 1,694 (according to Handsaker's map [34]) to 5,610 (Sudmant-Nature's map [36]), with a total of 7,267 genes (out of 19,430 autosomal protein-coding genes) affected by CNVs when aggregating all 5 maps. Remarkably, only 402 (5.5%) of all genes affected by CNVs coincide in the 5 maps (**S2 Fig**). However, this overlap is non-random ($P < 2.2e$-16).

7

124    Most of the CNVs overlapping with genes fall within intronic regions (~63% of all CNVs) without

125    any overlap with exons. More surprisingly, of the purely intronic CNVs detected, over the 94% are

126    losses or gain/loss CNVs. This is in stark contrast with whole gene CNVs, which tend to be

127    exclusively gains (55% of the cases) or gain/loss CNVs (25% of the cases) (**Fig 1B-F**). There is a

128    significant enrichment of purely intronic losses ($P < 0.001$; permutation testing) in 4 out of 5 maps,

129    with 6 to 15% more deletions falling in introns than expected by chance, depending on the CNV

130    map. (No significant differences with the expected values were found with Sudmant-Science's map,

131    $P = 0.6683$). In contrast, in protein-coding genes, there were 13-70% fewer CNV deletions

132    overlapping with exons than would be expected by chance, depending on the map ($P < 0.05$ in all

133    maps).

134    Given the potential regulatory role of introns and the high frequency of purely intronic deletions in

135    the population, we focused on the impact of CNVs on introns. For all the subsequent analyses we

136    restricted our set of CNVs to loss and gain/loss CNVs, as they together represent sequences that are

137    lost in some individuals. It is important to note that three maps (Sudmant-Nature's [36], Zarrei's

138    [31] and Abyzov's [33]) together represent the 86% of all intronic deletions from our datasets. The

139    methods used to generate the other two maps (Handsaker's [34] and Sudmant-Science's [35]) tend

140    to detect less losses and larger CNV regions, that result in maps with fewer purely intronic deletions

141    (**S1 Fig**).

142    For the above reasons, we focused our analyses of intronic deletions on the three maps with more

143    intronic deletions (Sudmant-Nature's [36], Zarrei's [31] and Abyzov's [33]). We checked that this

144    enrichment of intronic losses was still significant when controlling for different intron sizes (**S3**

145    **Fig**) and DNA replication timing ($P < 5e{-}04$; permutation testing).

146    Finally, we tested whether intronic losses were distributed equally between essential and non-

147    essential genes. We separated the protein-coding genes into two groups: those that have been

148    reported to be essential after CRISPR-based genomic targeting [37,38] or gene-trap insertional

149    mutagenesis methodology [39], and those which were not found to be essential. Strikingly, we

150    observed that the proportion of intronic deletions is higher than expected by chance in both

151    essential and non-essential genes **(S2 Table)**. The fact that these intronic deletions can appear in

152    essential genes suggests that they might be an unexpected source of genetic variation that could

153    potentially influence the regulation of functionally relevant genes in human populations.

154

155    **Intronic losses accumulate in evolutionarily old genes, while losses in coding regions are more**

156    **frequent in young genes**

157    Intrigued by the overrepresentation of intronic deletions in human protein-coding genes, we next

158    investigated in more detail the quantitative and functional impact of these deletions. We have

159    previously reported that the evolutionarily younger a gene is, the more likely it is to carry whole

160    gene CNVs in the population [40]. Here we confirm that most ancient genes are depleted of CNVs

161    that affect their coding regions, while primate-specific genes are enriched in CNVs **(S4 Fig)**. This

162    pattern was also observed when CNV gains were excluded (**Fig 2A**). The generation of random

163    background models revealed that ancient genes were significantly depleted of coding region losses

164    (both exonic and whole gene) ($P < 0.05$), while these were enriched in young genes ($P < 0.05$) (**Fig**

165    **2A**).

9

166   Surprisingly, we observed an opposite trend for purely intronic deletions: the proportion of ancient

167   genes with intronic deletions was higher than that of young genes, and also higher than expected by

168   chance (**Fig 2B**). This finding was confirmed with additional analyses considering only genes with

169   introns and adjusting by the different size distributions of introns (**S5 Fig**). The introns of essential

170   genes tend to be shorter [21,22] and essential genes also tend to be ancient [41]. Therefore, we

171   compared the intron size of genes within the same age groups and found that the introns of essential

172   genes are shorter than those of non-essential genes of the same age (**S6A and S6B Fig**).

173   Even if introns are shorter in essential genes, we found a significant proportion of them present

174   intronic deletions in the Sudmant-Nature's map [36] (**S6C Fig**) suggesting that intronic size

175   variation in the population might be more important than we originally thought.

176

177   **Intronic deletions result in population variation of gene lengths**

178   The percentage of each intron that can be lost due to CNV losses is highly variable, from 0.03% to

179   96.8%, representing a loss of the 0.01% to 77.5% of the total genic size. (**Fig 3A-C**). Some

180   examples of genes with a notable change in size after a single intronic deletion are the neuronal

181   glutamate transporter SLC1A1, with a loss of the 37% of its genic size, TCTN3 (tectonic family

182   member 3), which loses 45% of its gene size and the LINGO2 (Leucine Rich Repeat And Ig

183   Domain Containing 2, alias LERN3 or LRRN6C) gene with a loss of the 34% of its size (**Fig 3D**).

184   Remarkably, these genes are highly conserved at the protein level and are amongst the 20% of

185   genes most intolerant to functional variation according to the ranking of the RVIS (Residual

186   Variation Intolerance Score) gene scores, which is based on the amount of genetic variation of each

10

187    gene at an exome level [42]. This result shows that genes with a very conserved coding sequence

188    with a general depletion of deletions [36] can have important losses of intronic regions, which

189    might affect their regulation without affecting their protein structure.

190    Intronic deletions can impact regulatory regions such as enhancers or CTCF binding sites, which

191    are enriched (P < 1e-04) and impoverished (P < 1e-04) in introns, respectively. Indeed, we find very

192    frequently deletions in introns with these regulatory features (P < 2.2e-16). However, the direct

193    overlap of the deletions with both enhancers and CTCF binding sites is significantly lower than

194    expected by chance (P < 1e-04, **S3 Table**). This suggests intronic losses can occur close to

195    regulatory features within introns but that deleting part of a regulatory feature might often have a

196    deleterious impact.

197    Besides altering the size of the introns or disrupting regulatory regions, intronic deletions could also

198    affect splicing, which is required to avoid the translation of introns. In genes with long introns, the

199    recognition of introns and exons by splicing machinery is based on their differential GC content

200    (Amit et al. 2012, Gelfman et al 2013) as the lower GC content in introns facilitates their

201    recognition. Presumably, this recognition mechanism has contributed to the expansion of introns in

202    higher eukaryotes (Hollander et al 2016). We analyzed the GC content of introns with deletions and

203    found that the deleted sequences had a significantly higher GC content to that of the introns where

204    they are located ($P = 1.8e-28$). Moreover, we observed that the loss of these fragments decreased

205    significantly the overall GC content of the remaining introns ($P = 2.23e-16$). Our results suggest

206    that the deletion of GC rich regions within introns could lower the overall GC content of the intron,

207    increasing the difference of GC content between introns and their flanking exons, what could

208     facilitate exon definition during splicing (**S8** and **S9 Fig**).


209     We have shown that deletions within introns are widespread in introns of varying sizes and can

210     produce important changes on the sequence composition and the regulatory architecture of many

211     protein-coding genes, which might be relevant for transcription and splicing of those genes through

212     different mechanisms. Intronic deletions constitute a previously unexpected source of variation in

213     gene and transcript length across individuals and can subtly affect ancient genes with important

214     functions that don't tolerate more drastic alterations.


215


216     **Effect of intronic deletions on gene expression**


217     Multiallelic CNVs affecting whole genes have been shown to correlate with gene expression:

218     generally, the higher the number of copies of the gene, the higher its expression levels [34,36]. Our

219     data suggests that the intronic size variation could also impact the expression of the affected genes.

220     Therefore, we looked into the possible effect of intronic hemizygous deletions on gene expression

221     variation at the population level, comparing the effects with hemizygous deletions in coding (whole

222     gene and exonic) and intergenic non-coding deletions. We used available RNA-seq data from

223     Geuvadis [43] that was derived from lymphoblastoid cell lines for 445 individuals for whom we

224     have the matching CNV data (Sudmant Nature's map [36]). In order to look for differences in gene

225     expression we selected variants for which we had at least 2 hemizygous individuals (individuals

226     with copy number = 1) and at least 2 wild-type individuals (copy number = 2) and we compared the

227     expression levels among these two groups (**Fig 4A** and **S10 Fig**).

228    We first studied the effect of intronic deletions on gene expression and we observed significant

229    differences in gene expression in 52 out of the 1,474 genes with intronic deletions (3.5%) in

230    lymphoblastoid cell lines. This percentage is higher than expected by chance (P = 1e-4) **(Fig 4)**,

231    being the expected values the total of differentially expressed genes (DEGs) when randomizing the

232    individuals carrying the mutation. Of the DEGs, 62% were downregulated and the other 38%

233    upregulated, suggesting that intronic deletions might result both in enhancing or repressing gene

234    expression.

235    We investigated if deletions in introns of genes showing differential expression tended to overlap

236    with regulatory features, but we did not observe any significant enrichment (P = 1). Even though

237    first introns are known to be particularly important for gene regulation [3,44], there was no

238    significant enrichment of DEGs with their first intron affected (P = 0.86). These results suggest that

239    other mechanisms independent of intronic regulatory regions might be responsible for these

240    changes in gene expression. It is also possible that a combination of multiple different mechanisms

241    may be necessary to explain the observed effects. In addition, we cannot rule out that the lack of

242    association between intronic regulatory features and gene expression changes is due to the small

243    number of DEGs in this cell type and/or lack of detailed enough epigenomic annotations.

244    We wondered how the impact of non-coding intronic deletions in gene expression compared to

245    those of non-coding intergenic deletions. We focused on intergenic regions that show long-range

246    interactions with promoters of protein-coding genes - what it is generally assumed to reflect a

247    regulatory function for these intergenic regions [45]. The impact of noncoding intronic versus

248    intergenic deletions on gene expression was therefore studied. We used promoter-capture Hi-C

13

249    published data for B-lymphocytes [46] to link deletions in intergenic regions with interacting genes.

250    Significant changes in gene expression were seen in 11 out of 872 (1.26%) genes identified to have

251    a deletion in an intergenic contacting region. Contrary to the effect of intronic deletions within the

252    same gene, this percentage of DEGs was not different to that expected by chance (P = 0.08).

253    Therefore, our data suggests that variation within intronic regions may have a more significant

254    impact on gene regulation than intergenic regions.

255    The effect on gene expression appears to be greater when coding regions were affected, compared

256    to purely intronic sequence losses: 15 out of 51 (29.4%, P < 1e-4) whole gene deletion CNVs

257    resulted in significant downregulation of gene expression and 30 out of 239 genes with partial

258    exonic deletions that were differentially expressed (12.6%, 28 down- and 2 up-regulated (P < 1e-4).

259    However, given the higher frequency of intronic deletions in the population, the absolute number of

260    DEGs with intronic deletions (52 genes) was similar to the total of DEGs with coding deletions (45

261    genes, **Fig 4B** and **S4 Table**). Moreover, while coding losses mostly associate to gene down-

262    regulation, intronic losses are frequently associated to gene up-regulation. This shows the potential

263    global relevance of intronic deletions on gene expression, especially considering their frequency in

264    ancient genes (27.9 % in genes older than Sarcopterygii) is almost the double than the one for

265    coding deletions (14.6%, **Fig 2**). In summary, these data suggest that intronic variants could have an

266    important regulatory impact on ancient genes.

14

267 **Discussion**

268 Several studies have explored copy number variation in healthy humans [31,33–36,47–50], many of

269 which have reported common CNVs overlapping with protein coding regions less frequently than

270 expected by chance [29,31,35,36,47]. However, little attention has been paid to introns and, to the

271 best of our knowledge, no previous study has dissected the differential impact of CNVs on exonic

272 and intronic regions. Taking five recently published CNV maps [31,33–36], we observed an

273 enrichment of deletions in introns resulting in gene length differences among individuals.

274 The different CNV maps used were built using different datasets and CNV calling algorithms,

275 resulting in very different numbers of CNV sizes and types. Still, we think that each of these studies

276 has their own limitations and probably none of them actually reflects all the variability in the

277 genome. Therefore, instead of merging them into one map, we preferred to analyse the maps in

278 parallel. This allowed us to compare the consistency of the results and, at the same time, helped us

279 to better understand the peculiarities of each CNV set. We saw very consistent trends when we

280 analysed the enrichment of intronic deletions or the differential impact of CNVs depending of the

281 evolutionary age of genes, what show the robustness and generality of our results. At the same time

282 our data suggest that using the different maps in parallel can be a useful way to cross-validate

283 biological findings.

284 Structural variants in the germline DNA constitute an important source of genetic variability that

285 serves as the substrate for evolution. Therefore, dating the evolutionary age of genes allows the

286 study of structural variants that were fixed millions of years ago. We have previously shown that

287 genes of different ages are found in different proportions within current human CNV regions [51].

15

288  Whole young gene loci, contrarily to ancient gene loci, are very variable in copy number and tend

289  to be located in late replicating genomic regions, which are more error-prone and have less precise

290  DNA repair mechanisms than earlier regions [51]. Fixation of duplications or losses of whole genes

291  in these regions can lead to the birth of new genes or to their disappearance (**Fig 5**).

292  Here, we have observed that also gains and losses affecting only part of the coding sequence are

293  also enriched in young genes (**Fig 5**). Such CNVs can disrupt the protein sequence, but they can

294  also eliminate, duplicate or relocate exons or parts of exons, giving to the organism a mechanism to

295  modify young genes.

296  On the other hand, evolutionarily ancient genes, generally depleted of CNVs overlapping with their

297  coding regions, are especially enriched with intronic losses (**Fig 5**). This phenomenon shows that

298  although the protein sequence is usually unaffected, changes in the intronic sequence can modulate

299  the expression of the gene and promote variability in the population. We found in lymphoblastoid

300  cells more differentially expressed genes associated with intronic losses than expected by chance.

301  This association is expected to be even stronger as for many genes the effect of their intronic losses

302  will be only observed in other cell types or tissues.

303  Very interestingly, we observed differences in which genes show changes in gene expression when

304  affected by coding (exonic or whole gene) or purely intronic losses. We see that differences in

305  expression in younger genes are mainly associated to full gene dosage changes or partial disruption

306  of their coding sequence. In contrast, ancient genes that generally are less tolerant to any kind of

307  mutations in their coding sequence, are enriched in intronic deletions which that could be

308  modulating their expression (**Fig 5**). The availability of CNV and population-based gene expression

16

309    data from several tissues will allow to evaluate more accurately what is the impact of coding and

310    non-coding deletions in the whole organism.

311    CNVs can be directly disrupting a regulatory feature or affect the distance, for example, between

312    promoter and enhancer. We found that the presence of enhancers is significantly enriched in

313    introns, agreeing with previous findings in plants [11,44]. In general, genes with complex

314    regulation patterns require more regulatory DNA [52] and introns tend to be longer in tissue-

315    specific and transcription factor genes compared to housekeeping genes [21,53]. Since many

316    enhancers are tissue-specific [54] intronic CNVs might frequently have effects on particular cell

317    types. Therefore, the loss of intronic sequence might be affecting the expression of such genes in a

318    tissue-specific manner.

319    Our results also suggest that non-coding intronic deletions might have a wider impact on population

320    gene expression variability than deletions in non-coding intergenic regions that interact with

321    promoters, given that intronic deletions correlate with gene expression changes more often than

322    expected by chance, while promoter-interacting intergenic regions don't. However, intergenic

323    deletions were associated with genes using promoter-capture HiC data maps derived from a few

324    pooled genomes [46] and we may need to have personal genome interactomes, more tissues and

325    conditions to evaluate more precisely the effect of intergenic deletions on gene expression.

326    Furthermore, with the necessary experimental and analytical future advances, it will be extremely

327    exciting to see how individual copy number variants change the personal landscape of interactions

328    among promoters and other genomic elements.

17

329   We speculate that intronic CNVs might have a previously unsuspected role in shaping gene

330   expression variability in populations with potential important consequences in human evolution and

331   adaptation. After uncovering the relevance of gene length variation in the healthy population by

332   frequent intronic deletions, the next open question will be if any of these common non-coding

333   variants may be associated with disease. In fact, these population-based CNV maps could be useful

334   to identify disease relevant and irrelevant intronic regions. It is now well known that most genome-

335   wide association studies (GWAS) associated SNPs tend to be located in intronic and intergenic

336   regions and the pathogenicity of non-coding CNVs, mostly in upstream promoters, is starting to

337   emerge [55]. Thus, future case-control studies including WGS should also pay attention at

338   potentially important role of purely intronic variation. While exons cover around the 2.8% of the

339   genome, introns cover 35.3%, of the genome (based on the gene set used for this study). WGS

340   studies are starting to focus on distal intergenic enhancers, but intronic regions are commonly

341   ignored in the analyses. A recent analysis of the literature has revealed a substantial amount of

342   pathogenic variants located "deep" within introns (more than 100 bp from exon-intron boundaries)

343   which suggests that the sequence analysis of full introns may help to identify causal mutations for

344   many undiagnosed clinical cases [27]. With the results presented here, we emphasize the

345   importance of sequencing and analysing variants located in introns as they can potentially be as

346   consequential as regulatory elements found in intergenic regions.

347   Being intronic deletions so common in the healthy population, it will also be interesting to explore

348   how frequent are purely intronic somatic deletions in cancer and evaluate their potential

349   contribution to the reprogramming of gene regulation of cancer cells. For example, are there

350   somatic deletions of intronic sequences that result in the shortening of oncogenes, favoring their

18

351 higher expression. The role of cancer somatic variants in distal regulatory regions is just starting to

352 be explored [56–60]. As we have shown that intronic regions are significantly enriched in

353 regulatory regions in the human genome, understanding the functional effect somatic intronic

354 deletions in cancer could be an attractive new field of research with high potential for discovery. It

355 has been previously proposed that high-order chromatin architecture is influencing the landscape of

356 chromosomal alterations in cancer [61]. We hypothesize that the high-order genome organisation in

357 healthy cells is applying constraints on where variability can be high or low, allowing high

358 variability anywhere in young genes but only in introns for ancient genes. Therefore, it will be

359 interesting to understand better how these constraints change comparing data from healthy cells

360 with the frequent aneuploidy in tumors, especially in radical re-structuring events originated by

361 chromothripsis [62] .

362 In summary, our data shows that intronic CNVs constitute the most abundant form of CNV in

363 protein-coding genes. These intronic length variation possibly means that the actual size of many

364 genes is not yet fixed in human populations. We show that intronic length variation is particularly

365 frequent in evolutionary old genes, with a significant proportion of them showing associated gene

366 expression changes. This suggests that intronic CNVs might be actively contributing to the

367 evolution of gene regulation in many genes with highly conserved protein sequences. Taken

368 together, our results suggest that copy number variation is shaping gene evolution in different ways

369 depending on the age of genes, duplicating or deleting young genes and fine-tuning the regulation

370 of old genes.

19

# Materials and methods

## Origin and filtering of CNV maps

Whole genome CNV maps were downloaded from 5 different publications [31,33–36]. For our analysis we selected autosomal and not private CNVs. Some extra filters were applied to some maps: In Handsaker et al. we removed CNVs marked as low quality and all the variants from two of the individuals (NA07346 and NA11918) because they were not included in the phased map. From Zarrei's maps we used the stringent map that considered CNVs that appeared in at least 2 individuals and in 2 studies. The complete list of CNVs analysed is available in **S5 Table**.

## Gene structures

Autosomal gene structures and sequences were retrieved from Ensembl [63] (http://www.ensembl.org; version 75) and principal isoforms were determined according to the APPRIS database [64], Ensembl version 74. In order to avoid duplicate identification of introns, intronic regions were defined as regions within introns that aren't coding in any transcript of any gene. When analyzing real introns, in order to avoid duplicate identification of introns, the principal isoform with a higher exonic content was taken. The complete list of genes affected by different types of CNVs is available in **S6 Table**. Genomic sequences were obtained from the primary GRCh37/hg19 assembly, and were used for calculating the GC content of introns and intronic CNVs.

## Dating gene and intron ages

20

390    An age was assigned to all duplicated genes as described in Juan et al. 2013. In the case of

391    singletons gene ages were assigned from the last common ancestor to all the genes in their family

392    according to the gene trees retrieved from ENSEMBL. Singleton's ages can be noisy for genes

393    suffering important alterations as gene fusion/fission events or divergence shifts. As a consequence,

394    these ages should not be interpreted as the age of the oldest region of the gene, but as a restrictive

395    definition of gene age considering a similar gene structure and gene product.

396    The ages (from ancient to recent) and number of genes per age are as follows: FungiMetazoa: 1119,

397    Bilateria: 2892, Chordata: 1152, Euteleostomi: 8230, Sarcopterygii: 182, Tetrapoda: 154, Amniota:

398    408, Mammalia: 375, Theria: 515, Eutheria: 848, Simiiformes: 233, Catarrhini: 170, Hominoidea:

399    106, Hominidae: 64, HomoPanGorilla: 204, HomoSapiens: 500. For some analyses, Primates age

400    groups (Simiiformes to HomoSapiens) were collapsed. For other analyses, we only considered two

401    extreme groups, "ancient" (collapsing groups from FungiMetazoa to Sarcopterygii) and "young"

402    genes (Primates).

403    Intronic regions were assigned the evolutionary age of the gene they belonged to. In the cases when

404    an intron could be assigned to more than one gene, the most recent age was assigned to them.

**Statistical assessment of genome-wide distribution of CNVs**

406    To estimate statistical significance of our results we performed permutation tests. In order to

407    compare the number of overlaps of CNVs with genic functional elements we compared our

408    observed values to a background model. This model was obtained by relocating all the CNVs in the

409    whole genome (except for centromeres and telomeres) 10,000 times.

410   In addition, we generated background models correcting by DNA replication timing. For this, we

411   downloaded DNA replication timing data from 15 cell lines from ENCODE [65,66] and assigned

412   the median value of all cell lines to each 1 Kb window of the genome. Then, we classified the

413   genome in 5 intervals of DNA replication timing and we relocated the CNVs within its interval of

414   replication timing.

415   We compared the location of the CNVs in our datasets and compared with their distribution in the

416   random models in order to calculate enrichments or depletions depending on the intron size and

417   gene age and essentiality.

**Regulatory features**

419   We downloaded a genome-wide set of regions that are likely to be involved in gene regulation from

420   the Ensembl Regulatory Build [67]. We checked if introns are enriched in these regulatory features

421   (promoters, enhancers, promoter flanking regions or insulators) by comparing to a random

422   background model generated by relocating 10,000 times all regulatory features in the genome. P-

423   values are the fraction of random values superior or inferior to the observed values.

424   In order to check for the significance of the overlaps between intronic deletions and regulatory

425   features we relocated 10,000 all intronic deletions within their introns and checked for differences

426   in overlap with regulatory features.

**Gene expression analysis**

428   We used available RNA-seq data at Geuvadis [43] that was derived from lymphoblastoid cell lines

429   for 445 individuals who were sequenced by the 1000 Genomes Project and for whom we have the

22

430 intronic deletions in the largest CNV map [36]. We focused our analyses on the 763 genes that have

431 only one intronic deletion in the population with at least two individuals affected in the Geuvadis

432 dataset. For each of these genes we classified the PEER normalized gene expression levels [68] in

433 two groups depending if the individual carried or not the intronic deletion and performed Student's

434 t-tests. We corrected for multiple testing with p.adjust R function (Benjamini-Hochberg method). In

435 addition, we randomized the individuals with the intronic deletions 10,000 times and calculated the

436 expected percentages of significantly differentially expressed genes.

437

438

## References

445    1.    Lynch M. The Origins of Genome Architecture. Sinauer Associates Incorporated; 2007.

446    2.    Chorev M, Carmel L. The function of introns. Front Genet. 2012;3: 55.

447    3.    Jo B-S, Choi SS. Introns: The Functional Benefits of Introns in Genomes. Genomics Inform.

448          2015;13: 112.

450    4.    Gilbert W. The Exon Theory of Genes. Cold Spring Harb Symp Quant Biol. 1987;52: 901–

451          905.

452    5.    Schwartz SH, Silva J, Burstein D, Pupko T, Eyras E, Ast G. Large-scale comparative analysis

453          of splicing signals and their corresponding splicing factors in eukaryotes. Genome Res.

454          2008;18: 88–103.

455    6.    Gelfman S, Burstein D, Penn O, Savchenko A, Amit M, Schwartz S, et al. Changes in exon-

456          intron structure during vertebrate evolution affect the splicing pattern of exons. Genome Res.

457          2012;22: 35–50.

458    7.    Kirkconnell KS, Magnuson B, Paulsen MT, Lu B, Bedi K, Ljungman M. Gene length as a

459          biological timer to establish temporal transcriptional regulation. Cell Cycle. 2017;16: 259–270.

460    8.    Le Hir H, Nott A, Moore MJ. How introns influence and enhance eukaryotic gene expression.

461          Trends Biochem Sci. 2003;28: 215–220.

462    9.    Nott A, Meislin SH, Moore MJ. A quantitative analysis of intron effects on mammalian gene

24

463    expression. RNA. 2003;9: 607–617.

464    10.  Gaffney DJ, Keightley PD. Genomic selective constraints in murid noncoding DNA. PLoS

465    Genet. 2006;2: e204.

466    11.  Rose AB. Intron-mediated regulation of gene expression. Curr Top Microbiol Immunol.

467    2008;326: 277–290.

468    12.  Le Hir H, Nott A, Moore MJ. How introns influence and enhance eukaryotic gene expression.

469    Trends Biochem Sci. 2003;28: 215–220.

470    13.  Seoighe C, Korir PK. Evidence for intron length conservation in a set of mammalian genes

471    associated with embryonic development. BMC Bioinformatics. 2011;12 Suppl 9: S16.

472    14.  Eisenberg E, Levanon EY. Human housekeeping genes, revisited. Trends Genet. 2013;29:

473    569–574.

474    15.  Castillo-Davis CI, Mekhedov SL, Hartl DL, Koonin EV, Kondrashov FA. Selection for short

475    introns in highly expressed genes. Nat Genet. 2002;31: 415–418.

476    16.  Heyn P, Kircher M, Dahl A, Kelso J, Tomancak P, Kalinka AT, et al. The earliest transcribed

477    zygotic genes are short, newly evolved, and different across species. Cell Rep. 2014;6: 285–

478    292.

479    17.  Artieri CG, Fraser HB. Transcript length mediates developmental timing of gene expression

480    across Drosophila. Mol Biol Evol. 2014;31: 2879–2889.

25

481   18.   Swinburne IA, Silver PA. Intron delays and transcriptional timing during development. Dev

482         Cell. 2008;14: 324–330.

483   19.   Jeffares DC, Penkett CJ, Bähler J. Rapidly regulated genes are intron poor. Trends Genet.

484         2008;24: 375–378.

485   20.   Keane PA, Seoighe C. Intron Length Coevolution across Mammalian Genomes. Mol Biol

486         Evol. 2016;33: 2682–2691.

487   21.   Eisenberg E, Levanon EY. Human housekeeping genes are compact. Trends Genet. 2003;19:

488         362–365.

489   22.   Vinogradov AE. Compactness of human housekeeping genes: selection for economy or

490         genomic design? Trends Genet. 2004;20: 248–253.

491   23.   Carmel L, Koonin EV. A universal nonmonotonic relationship between gene compactness and

492         expression levels in multicellular eukaryotes. Genome Biol Evol. 2009;1: 382–390.

493   24.   Agrawal A, Hamvas A, Cole FS, Wambach JA, Wegner D, Coghill C, et al. An intronic

494         ABCA3 mutation that is responsible for respiratory disease. Pediatr Res. 2012;71: 633–637.

495   25.   Lo Y-F, Nozu K, Iijima K, Morishita T, Huang C-C, Yang S-S, et al. Recurrent deep intronic

496         mutations in the SLC12A3 gene responsible for Gitelman's syndrome. Clin J Am Soc

497         Nephrol. 2011;6: 630–639.

498   26.   Nurnberg ST, Zhang H, Hand NJ, Bauer RC, Saleheen D, Reilly MP, et al. From Loci to

499         Biology: Functional Genomics of Genome-Wide Association for Coronary Disease. Circ Res.

500     NIH Public Access; 2016;118: 586.

501     27. Vaz-Drago R, Custódio N, Carmo-Fonseca M. Deep intronic mutations and human disease.

502         Hum Genet. 2017; doi:10.1007/s00439-017-1809-4

503     28. Arlt MF, Wilson TE, Glover TW. Replication stress and mechanisms of CNV formation. Curr

504         Opin Genet Dev. 2012;22: 204–210.

505     29. Conrad DF, Bird C, Blackburne B, Lindsay S, Mamanova L, Lee C, et al. Mutation spectrum

506         revealed by breakpoint sequencing of human germline CNVs. Nat Genet. 2010;42: 385–391.

507     30. Pang AW, MacDonald JR, Pinto D, Wei J, Rafiq MA, Conrad DF, et al. Towards a

508         comprehensive structural variation map of an individual human genome. Genome Biol.

509         2010;11: R52.

510     31. Zarrei M, MacDonald JR, Merico D, Scherer SW. A copy number variation map of the human

511         genome. Nat Rev Genet. 2015;16: 172–183.

512     32. Pirooznia M, Goes FS, Zandi PP. Whole-genome CNV analysis: advances in computational

513         approaches. Front Genet. 2015;6: 138.

514     33. Abyzov A, Li S, Kim DR, Mohiyuddin M, Stütz AM, Parrish NF, et al. Analysis of deletion

515         breakpoints from 1,092 humans reveals details of mutation mechanisms. Nat Commun.

516         2015;6: 7256.

517     34. Handsaker RE, Van Doren V, Berman JR, Genovese G, Kashin S, Boettger LM, et al. Large

518         multiallelic copy number variations in humans. Nat Genet. 2015;47: 296–303.

27

519    35.  Sudmant PH, Mallick S, Nelson BJ, Hormozdiari F, Krumm N, Huddleston J, et al. Global

520        diversity, population stratification, and selection of human copy-number variation. Science.

521        2015;349: aab3761.

522    36.  Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, et al. An

523        integrated map of structural variation in 2,504 human genomes. Nature. 2015;526: 75–81.

524    37.  Hart T, Chandrashekhar M, Aregger M, Steinhart Z, Brown KR, MacLeod G, et al. High-

525        Resolution CRISPR Screens Reveal Fitness Genes and Genotype-Specific Cancer Liabilities.

526        Cell. 2015;163: 1515–1526.

527    38.  Wang T, Birsoy K, Hughes NW, Krupczak KM, Post Y, Wei JJ, et al. Identification and

528        characterization of essential genes in the human genome. Science. 2015;350: 1096–1101.

529    39.  Blomen VA, Májek P, Jae LT, Bigenzahn JW, Nieuwenhuis J, Staring J, et al. Gene

530        essentiality and synthetic lethality in haploid human cells. Science. 2015;350: 1092–1096.

531    40.  Juan D, Rico D, Marques-Bonet T, Fernández-Capetillo O, Valencia A. Late-replicating CNVs

532        as a source of new genes. Biol Open. 2014;3. doi:10.1242/bio.20147815

533    41.  Chen W-H, Trachana K, Lercher MJ, Bork P. Younger genes are less likely to be essential

534        than older genes, and duplicates are less likely to be essential than singletons of the same age.

535        Mol Biol Evol. 2012;29: 1703–1706.

536    42.  Petrovski S, Wang Q, Heinzen EL, Allen AS, Goldstein DB. Genic intolerance to functional

537        variation and the interpretation of personal genomes. PLoS Genet. 2013;9: e1003709.

538    43. Lappalainen T, Sammeth M, Friedländer MR, 't Hoen PAC, Monlong J, Rivas MA, et al.

539         Transcriptome and genome sequencing uncovers functional variation in humans. Nature.

540         2013;501: 506–511.

541    44. Majewski J, Ott J. Distribution and characterization of regulatory elements in the human

542         genome. Genome Res. 2002;12: 1827–1836.

543    45. Schoenfelder S, Furlan-Magaril M, Mifsud B, Tavares-Cadete F, Sugar R, Javierre B-M, et al.

544         The pluripotent regulatory circuitry connecting promoters to their long-range interacting

545         elements. Genome Res. 2015;25: 582–597.

546    46. Javierre BM, Burren OS, Wilder SP, Kreuzhuber R, Hill SM, Sewitz S, et al. Lineage-Specific

547         Genome Architecture Links Enhancers and Non-coding Disease Variants to Target Gene

548         Promoters. Cell. 2016;167: 1369–1384.e19.

549    47. Korbel JO, Urban AE, Affourtit JP, Godwin B, Grubert F, Simons JF, et al. Paired-end

550         mapping reveals extensive structural variation in the human genome. Science. 2007;318: 420–

551         426.

552    48. McCarroll SA, Kuruvilla FG, Korn JM, Cawley S, Nemesh J, Wysoker A, et al. Integrated

553         detection and population-genetic analysis of SNPs and copy number variation. Nat Genet.

554         Nature Publishing Group; 2008;40: 1166–1174.

555    49. Conrad DF, Pinto D, Redon R, Feuk L, Gokcumen O, Zhang Y, et al. Origins and functional

556         impact of copy number variation in the human genome. Nature. 2010;464: 704–712.

29

557   50. Sudmant PH, Kitzman JO, Antonacci F, Alkan C, Malig M, Tsalenko A, et al. Diversity of

558        human copy number variation and multicopy genes. Science. 2010;330: 641–646.

559   51. Juan D, Rico D, Marques-Bonet T, Fernandez-Capetillo O, Valencia A. Late-replicating CNVs

560        as a source of new genes. Biol Open. 2013;3: 231–231.

561   52. Gaffney DJ, Keightley PD. Genomic selective constraints in murid noncoding DNA. PLoS

562        Genet. 2006;2: e204.

563   53. Kirkconnell KS, Magnuson B, Paulsen MT, Lu B, Bedi K, Ljungman M. Gene length as a

564        biological timer to establish temporal transcriptional regulation. Cell Cycle. 2017;16: 259–270.

565   54. Pennacchio LA, Loots GG, Nobrega MA, Ovcharenko I. Predicting tissue-specific enhancers

566        in the human genome. Genome Res. 2007;17: 201–211.

567   55. Zhang F, Lupski JR. Non-coding genetic variants in human disease. Hum Mol Genet. 2015;24:

568        R102–10.

569   56. Northcott PA, Lee C, Zichner T, Stütz AM, Erkek S, Kawauchi D, et al. Enhancer hijacking

570        activates GFI1 family oncogenes in medulloblastoma. Nature. 2014;511: 428–434.

571   57. Puente XS, Beà S, Valdés-Mas R, Villamor N, Gutiérrez-Abril J, Martín-Subero JI, et al. Non-

572        coding recurrent mutations in chronic lymphocytic leukaemia. Nature. 2015;526: 519–524.

573   58. Zhang X, Choi PS, Francis JM, Imielinski M, Watanabe H, Cherniack AD, et al. Identification

574        of focally amplified lineage-specific super-enhancers in human epithelial cancers. Nat Genet.

575        2016;48: 176–182.

576   59.  Khurana E, Fu Y, Chakravarty D, Demichelis F, Rubin MA, Gerstein M. Role of non-coding

577        sequence variants in cancer. Nat Rev Genet. 2016;17: 93–108.

578   60.  Weischenfeldt J, Dubash T, Drainas AP, Mardin BR, Chen Y, Stütz AM, et al. Pan-cancer

579        analysis of somatic copy-number alterations implicates IRS4 and IGF2 in enhancer hijacking.

580        Nat Genet. 2017;49: 65–74.

581   61.  Fudenberg G, Getz G, Meyerson M, Mirny LA. High order chromatin architecture shapes the

582        landscape of chromosomal alterations in cancer. Nat Biotechnol. 2011;29: 1109–1113.

583   62.  Zhang C-Z, Leibowitz ML, Pellman D. Chromothripsis and beyond: rapid genome evolution

584        from complex chromosomal rearrangements. Genes Dev. 2013;27: 2513–2530.

585   63.  Yates A, Akanni W, Amode MR, Barrell D, Billis K, Carvalho-Silva D, et al. Ensembl 2016.

586        Nucleic Acids Res. 2016;44: D710–6.

587   64.  Rodriguez JM, Maietta P, Ezkurdia I, Pietrelli A, Wesselink J-J, Lopez G, et al. APPRIS:

588        annotation of principal and alternative splice isoforms. Nucleic Acids Res. 2013;41: D110–7.

589   65.  Thurman RE, Day N, Noble WS, Stamatoyannopoulos JA. Identification of higher-order

590        functional domains in the human ENCODE regions. Genome Res. 2007;17: 917–927.

591   66.  Hansen RS, Thomas S, Sandstrom R, Canfield TK, Thurman RE, Weaver M, et al. Sequencing

592        newly replicated DNA reveals widespread plasticity in human replication timing. Proc Natl

593        Acad Sci U S A. 2010;107: 139–144.

594   67.  Zerbino DR, Wilder SP, Johnson N, Juettemann T, Flicek PR. The ensembl regulatory build.

595       Genome Biol. 2015;16: 56.

596   68.  Stegle O, Parts L, Durbin R, Winn J. A Bayesian framework to account for complex non-

597       genetic factors in gene expression levels greatly increases power in eQTL studies. PLoS

598       Comput Biol. 2010;6: e1000770.

32

599    **Figure captions**

600    **Fig 1. Types of CNVs in the different datasets.** (A) CNVs can overlap entire genes or fractions of

601    genes. CNVs overlapping with exons of a gene (exonic CNVs) and CNVs found within introns

602    (intronic CNVs). (B-D) Number of whole gene, exonic and intronic CNV events, showing the

603    different proportions of CNV gains, losses and gain and loss CNVs. (E-F) Venn diagrams showing

604    the number of genes with whole gene, exonic and intronic losses (E) or gains (F) in the three maps

605    with the higher number of deletions reported (Zarrei, Sudmant-Nature and Abyzov). Circle sizes are

606    proportional to the number of genes affected.

607    **Fig 2. Evolutionary age of affected genes.** Percentage of genes from each gene evolutionary age

608    that contain deletions overlapping with exons, including partial and whole gene CNVs (A) or

609    intronic deletions (B). The light blue line represents the expected value, calculated as the mean of

610    the genes in the 10,000 random permutations. Red asterisks mark the significantly enriched groups

611    of genes, while black asterisks mark gene age groups with fewer deletions than expected ($P < 0.05$).

612    Plot (C) shows, from all the genes overlapping with deletions after aggregating the three maps,

613    what is the proportion of genes that have all or part of their exons affected by deletions and what is

614    the percentage of genes with intronic deletions only. Bar width is proportional to the percentage of

615    genes from each evolutionary age that is affected by deletions of any kind, which spans from 18.5%

616    (Mammalia) to 49.8% (HomoPanGorilla). The equivalent figure for each separate map is shown in

617    S7 Fig.

618    **Fig 3. Changes in intron and gene size.** (A) Proportion of the reference intron that has been

619    observed as deleted in any of the studies. (B) Proportion of the whole intronic content of a gene that

33

620 has been observed as deleted. (C) Change in gene length by intronic deletions. (D) Example of gene

621 with a substantial change in gene size with a single intronic deletion.

622 **Fig 4. Differential expression.** (A) Number of genes with whole gene, exonic or intronic deletions

623 or with intergenic deletions in a region in long-range contact with it. (B) Differentially expressed

624 genes (DEGs) for each type of deletion. The colored bars show the observed number of DEGs per

625 group. The white bars represent the median random number of DEGs when randomizing 10,000

626 times the individuals with and without the deletion. Significant differences ($P < 0.05$) are marked

627 with * ($P < 0.0001$ in all cases) and error bars show median absolute deviation.

628 **Fig 5. Impact of CNVs on genes and their evolution.** Evolutionarily ancient and young genes

629 accumulate different kinds of structural variants. While young genes are enriched in coding

630 deletions (which alter gene dosage or disrupt the protein, sometimes affecting gene expression),

631 ancient genes have highly conserved coding sequence but an enrichment of deletions within their

632 introns. As we have shown, these changes in introns are sometimes associated with changes in gene

633 expression, showing that although the protein is highly conserved, the expression of it can change

634 from an individual to another due to changes in regulation.

635

636 **Supporting information**

637 **S1 Fig. Comparison of datasets.** Only variants in autosomes are considered and private events are

638 excluded. (A) Number and type of CNVs per dataset. (B) Autosomal Mb that are CNV. Gray part

639 of the bars corresponds to the CNV Mbs that are shared among maps. Colored parts of the bars are

34

640    map-specific CNV regions. (C) Width distribution of gains and losses in each map bean lines and

641    overall line are means). (D) Number of subjects and number of populations of origin used for

642    building of each filtered map.

643    **S2 Fig. Overlap of genes affected among datasets.** Number of genes overlapping with CNVs (A)

644    or found completely within a CNV (B) by number of maps in which the overlap is observed.

645    **S3 Fig. Enrichment of deletions in introns of different sizes.** Number of deletions in each size

646    bin. The light blue line represents the expected value, calculated as the mean of the affected introns

647    in the 10,000 random permutations. Red asterisks mark the bins significantly enriched with intronic

648    deletions ($P < 0.05$).

649    **S4 Fig: Differential effect of CNVs on protein coding genes of different evolutionary ages.** The

650    proportion of genes with CNVs (gains, losses and gain and loss CNVs) affecting their coding region

651    in each gene age class is shwon. Red stars show a significantly higher percentage of genes with

652    coding compared to 10,000 randomizations of the CNVs over the genome ($P < 0.05$). Black stars

653    show in which age groups there is a depletion of coding CNVs ($P < 0.05$).

654    **S5 Fig. Differential effect of intronic deletions in introns and protein coding genes of different**

655    **evolutive ages.** Figure (A) shows the percentage of introns over 1.5kb that contain intronic

656    deletions. Introns shared among genes of different gene ages were assigned the youngest age.

657    Figure (B) represents the percentage of genes with introns over 1.5kb that have intronic deletions.

658    In all plots, red stars show a significantly higher percentage of introns or genes with intronic

659    deletions compared to 10,000 randomizations of the deletions over the genome ($P < 0.05$). Black

660    stars show in which age groups there is a depletion of intronic deletions ($P < 0.05$).

35

661 **S6 Fig. Essential genes.** (A) Percentage of essential genes per evolutionary age. (B) Intron sizes of

662 non-essential and essential genes. (C) Percentage of non-essential (solid bars) and essential (empty

663 bars) genes of different evolutionary ages that have intronic deletions. Red stars show a

664 significantly higher percentage of genes with intronic deletions compared to 10,000 randomizations

665 of the deletions over the genome ($P < 0.05$).

666 **S7 Fig. Effect of the different types of deletions on all evolutionary ages.** Proportion of genes

667 with deletions that have the whole locus deleted, only part of their exons (exonic) affected by

668 deletions or intronic deletions only. (Equivalent to Fig. 2C, but here separated by CNV map.)

669 **S8 Fig. GC content in introns and intronic deletions.** (A) Bean-plots showing the different GC

670 distribution between the flanking exons of introns with or without deletions, separated by intron

671 size bins (with equal number of introns per bin). (B) GC content distributions in introns with or

672 without deletions, separated by intron size bins. Significance is considered for p-values $< 0.05$.

673 Beans show the estimated density of each distribution; horizontal lines show the mean values of

674 each side of the bean and the dashed horizontal line line represents the overall average of all values.

675 **S9 Fig. Examples of introns with a drop of GC content.** X-axis represents the coordinates of the

676 intron with its flanking exons (black boxes). Y-axis shows the GC content, calculated with sliding

677 200 bp windows. The deleted region is highlighted in grey.

678 **S10 Fig. Overlap of genes with different types of deletions.** Venn diagram showing the overlaps

679 for genes carrying different types of mutations (A) and for differentially expressed genes (DEGs)

680 (B).

681    **S1 Table. Number of individuals in each map, project the variants belong to and methods**

682    **used for CNV detection.**

683    **S2 Table. Enrichment of intronic deletions in non-essential and essential genes.**

684    **S3 Table. Overlap of intronic deletions with regulatory features.**

685    **S4 Table. Differentially expressed genes.**

686    **S5 Table. Filtered CNV maps used in this study.**
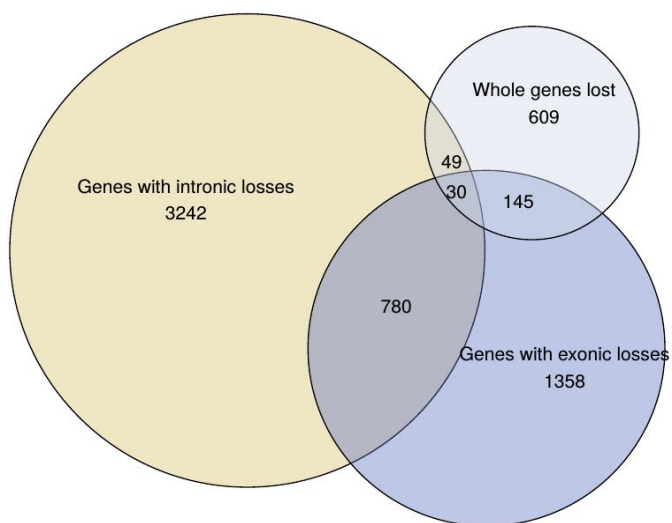
687    **S6 Table. Genes affected by CNVs, including their ages.**

37

**A)**

Whole gene CNVs | Exonic CNVs | Intronic CNVs

CNV

**B)** Whole gene CNVs

**C)** Exonic CNVs

**D)** Intronic CNVs

**E)** Genes with losses

Genes with intronic losses 3242
Whole genes lost 609
Genes with exonic losses 1358
49
30
145
780

**F)** Genes with gains

Genes with exonic gains 1215
Whole genes gained 758
Genes with intronic gains 347
147
122
1

**A) Coding deletions**

Sudmant (Nature) — Zarrei — Abyzov

**B) Intronic deletions**

Sudmant (Nature) — Zarrei — Abyzov

**C)**

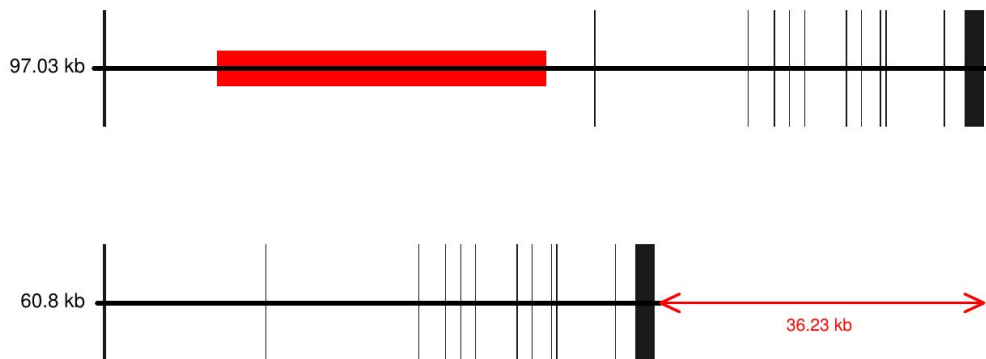Whole gene / Exonic / Intronic

**A) Proportion of intron lost**
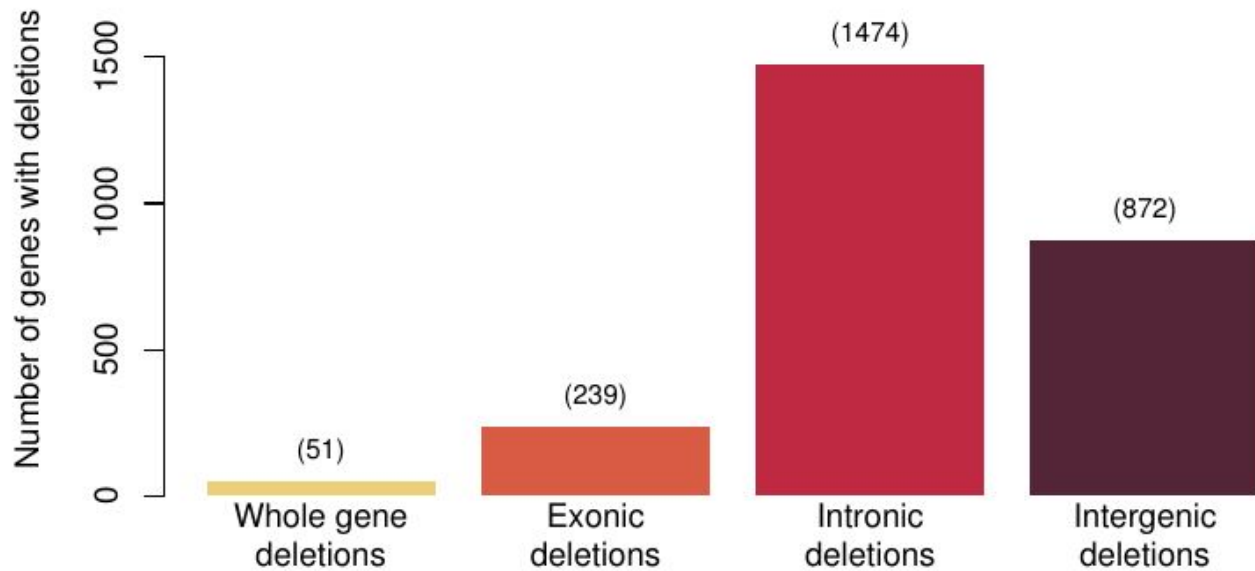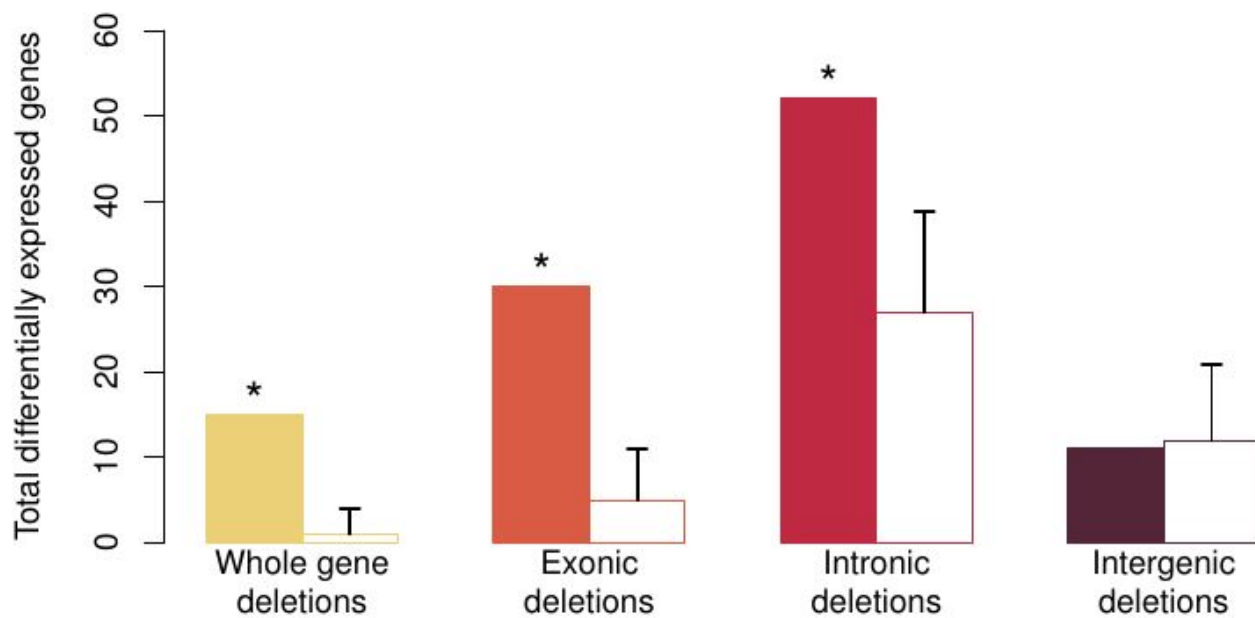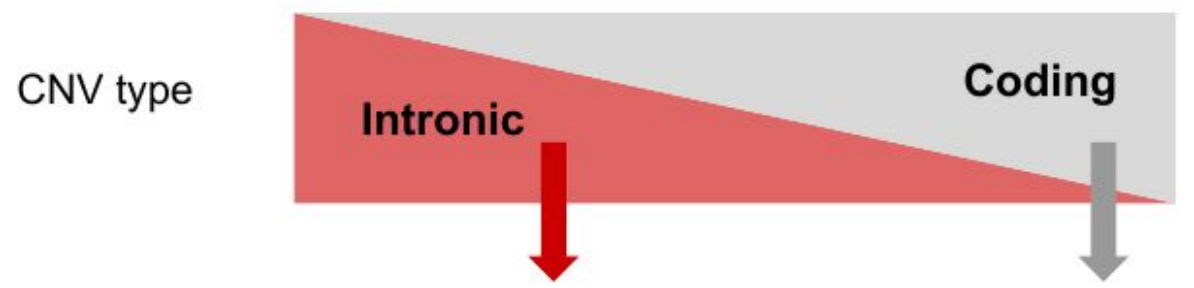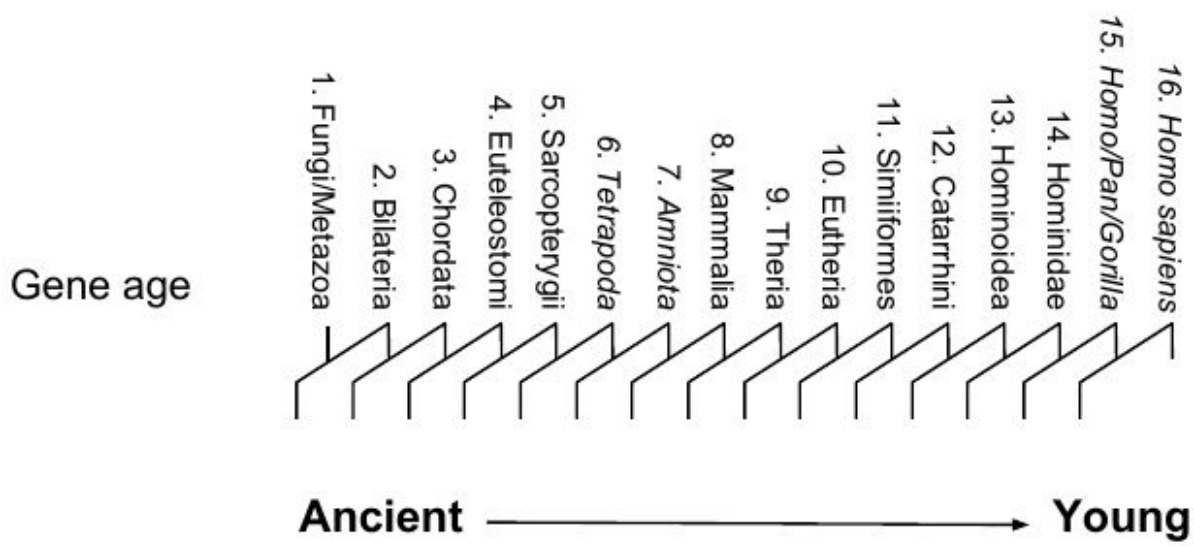
**B) Proportion of intronic content lost per gene**

**C) Change in gene length by intronic losses**

**D) *SLC1A* gene**

**A)**

**B)**

**Gene age**

1. Fungi/Metazoa
2. Bilateria
3. Chordata
4. Euteleostomi
5. Sarcopterygii
6. Tetrapoda
7. Amniota
8. Mammalia
9. Theria
10. Eutheria
11. Simiiformes
12. Catarrhini
13. Hominoidea
14. Hominidae
15. *Homo/Pan/Gorilla*
16. *Homo sapiens*

**Ancient** → **Young**

**Essentiality**

**CNV type**

**Intronic**

**Coding**

**CNV effect**

Transcription variability:
(Timing, expression
levels, splicing)

Birth and death of genes
Reshaping of proteins