

1 A human pan-genomic analysis reconfigures the genetic 2 and epigenetic make up of facioscapulohumeral 3 muscular dystrophy

4

5 **Authors:** Valentina Salsi, ^{1†} Matteo Chiara,^{2,3†} Sara Pini,¹ Paweł Kuś,⁴ Lucia Ruggiero,⁵
6 Silvia Bonanno,⁶ Carmelo Rodolico,⁷ Stefano C. Previtali,⁸ Maria Grazia D'Angelo,⁹ Lorenzo
7 Maggi,⁶ Diego Lopercolo,^{10,11} Marek Kimmel,^{4,12} Filippo M. Santorelli,¹³ Graziano Pesole,^{3,14}
8 Rossella G. Tupler^{1,15,16*}

9 **Authors' Affiliations:**

10 1. Department of Biomedical, Metabolic and Neural Sciences, University of Modena and
11 Reggio Emilia, Modena, 41125, Italy. valentina.salsi@unimore.it; pini.saraps@gmail.com;
12 rossella.tupler@unimore.it*

13 2. Dipartimento di Bioscienze, Università degli Studi di Milano, Milan, 20133, Italy.
14 matteo.chiara@unimi.it.

15 3. Institute of Biomembranes, Bioenergetics and Molecular Biotechnologies, Consiglio
16 Nazionale delle Ricerche, Bari, 70126, Italy. matteo.chiara@unimi.it;
17 graziano.pesole@uniba.it.

18 4. Department of Systems Biology and Engineering, Silesian University of Technology,
19 Gliwice, 44-100, Poland. pawel.s.kus@protonmail.com; kimmel@rice.edu.

20 5. Department of Neurosciences, Reproductive and Odontostomatological Sciences,
21 University Federico II of Naples, Naples, 80131, Italy. ruggilucia@gmail.com.

22 6. IRCCS Foundation, C. Besta Neurological Institute, Milan, 20133, Italy.
23 Silvia.Bonanno@istituto-besta.it; Lorenzo.Maggi@istituto-besta.it.

24 7. Department of Clinical and Experimental Medicine, University of Messina, Messina,
25 98125 Italy. carmelo.rodolico@unime.it.

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

- 26 8. INSPE and Division of Neuroscience, IRCCS San Raffaele Scientific Institute, Milan,
27 20133, Italy. previtali.stefano@hsr.it.
- 28 9. Department of Neurorehabilitation, IRCCS Eugenio Medea, Bosisio Parini, 23842, Italy.
29 grazia.dangelo@lanostrafamiglia.it.
- 30 10. Department of Medicine, Surgery and Neurosciences, University of Siena, Siena, 53100
31 Italy. diego.lopergolo@unifi.it.
- 32 11. UOC Neurologia e Malattie Neurometaboliche, Azienda Ospedaliero Universitaria
33 Senese, Policlinico Le Scotte, 53100 Siena, Italy. diego.lopergolo@unifi.it.
- 34 12. Departments of Statistics and Bioengineering, Rice University, Houston, TX, 77005,
35 United States. kimmel@rice.edu.
- 36 13. IRCCS Fondazione Stella Maris, Pisa, 56128, Italy. filippo3364@gmail.com.
- 37 14. Department of Biosciences, Biotechnology and Environment, University of Bari "A.
38 Moro", Bari, 70126, Italy. graziano.pesole@uniba.it.
- 39 15. Department of Molecular, Cell, and Cancer Biology, University of Massachusetts
40 Medical School, Worcester, MA 01605, United States. rossella.tupler@unimore.it *.
- 41 16. Li Weibo Institute for Rare Diseases Research at the University of Massachusetts
42 Medical School, Worcester, MA 01605, United States. rossella.tupler@unimore.it *.

43

44 † Valentina Salsi and Matteo Chiara contributed equally to this study and reserve the right to list themselves
45 as first author.

46

47 ***Corresponding Author:**

48 Rossella Tupler, M.D., Ph.D., Department of Biomedical, Metabolic and Neural Sciences, University of
49 Modena and Reggio Emilia, Email: rossella.tupler@unimore.it).

50

51 **ABSTRACT**

52 Facioscapulohumeral muscular dystrophy (FSHD) is the only human disease associated
53 with epigenetic changes at a macrosatellite array. Almost 95% of FSHD cases carry a
54 reduced number (≤ 10) of tandem 3.3 kilobase repeats, termed D4Z4, on chromosome 4q35;
55 remaining cases bear variants in chromatin remodeling factors, such as SMCHD1,
56 DNMT3B, LRIF1. Reduced CpG methylation is used for the molecular diagnosis of FSHD,
57 but D4Z4-like sequences dispersed in the genome can generate ambiguous results. By
58 analyzing complete haplotype level assemblies from the T2T-CHM13 human genome and
59 86 haploid genomes from the human pangenome project, we uncovered the extensive
60 number of D4Z4-like elements and their widespread inter- and intra-individual variability. An
61 original analytical approach was developed to elucidate this previously unaccounted wealth
62 of D4Z4-like elements and to analyze CpG methylation at D4Z4 in bisulfite-treated DNA
63 from 29 FSHD index cases and 15 relatives. Integrated analysis of clinical phenotype, D4Z4
64 CpG methylation level and gene variants showed that low D4Z4 methylation was associated
65 with variants in the *SMCHD1* gene, but not with the patients' clinical phenotypes. This is the
66 first study showing the relevance of the pangenome and T2T-CHM13 assemblies for
67 investigating the genotype-phenotype correlation in genetic diseases. The extension and
68 the variability of D4Z4-like elements scattered throughout the human genome and the
69 inconsistent association of phenotypes with methylation profiles advocate for a critical
70 revision of FSHD diagnostic tests based on D4Z4 CpG methylation assays and indicate that
71 molecular investigations must be complemented by family studies for the proper
72 interpretation of results.

73

74 Introduction

75

76 Facioscapulohumeral muscular dystrophy (FSHD) is the third most common form of
77 hereditary myopathy with an estimated prevalence from 1/8,000 to 1/20,000^{1,2}. FSHD is the
78 only hereditary human disease associated with reductions in copy number of a 3.3 kb
79 macrosatellite³, termed D4Z4. In the general population, the size of the D4Z4 repeat varies
80 between 11 and 150 units, whereas FSHD patients carry fewer than 11 repeats⁴. Each copy
81 of the 3.3 kb D4Z4 repeat contains an open reading frame (ORF) with two homeobox
82 domains, named *DUX4*-Like (*DUXL*)⁵, and two different classes of GC-rich repetitive
83 DNA^{6,7}: hhspm3, a member of a low copy human repeat family⁸, and LSau, a middle
84 repetitive DNA family^{8,9}. A tandemly arrayed D4Z4 repeat with 98% identity to the 4q35
85 array is located at 10q26¹⁰, but only alleles with reduced D4Z4 copy number on
86 chromosome 4q35 have been associated with FSHD (Figure 1).

87 Two genetically distinct FSHD subtypes, FSHD1 and FSHD2, are currently
88 described. In FSHD1 [OMIM #158900], the reduction of 4q35-D4Z4 repeats below a critical
89 threshold (10 or fewer repeat units, RU)¹¹ is thought to determine epigenetic alterations and
90 the inappropriate expression of nearby genes leading to disease¹²⁻¹⁴. In FSHD2 [OMIM
91 #158901] affected individuals carry two D4Z4 arrays in the healthy range (>10 RU), but bear
92 damaging heterozygous variants in chromatin remodeling factor genes, such as *SMCHD1*
93 (for structural maintenance of chromosomes flexible hinge domain containing 1)¹⁵, *DNMT3B*
94 (for methyltransferase 3B)¹⁶ and *LRIF1* (for ligand-dependent nuclear receptor interacting
95 factor 1)¹⁷. It has thus been proposed that similar to the reduction of D4Z4 RU, the
96 inactivation of these genes alters the epigenetic configuration of the D4Z4 array at 4q35
97 (Figure 1). Recent studies linked the molecular etiology of FSHD with the inappropriate
98 expression of the DUX4 protein from the terminal D4Z4 repeat. Only a specific haplotype
99 labelled 4qA-PAS is considered to be permissive for DUX4 protein expression, due to the

100 presence of a 260-bp sequence, termed pLAM, which provides a 3'-terminal exon (Exon 3)
101 and a polyadenylation signal (PAS) (AUUAAA) in juxtaposition to the terminal copy of *DUX4*
102 ORF (hereafter *ter-DUX4*) on 4q^{18,19} (Figure S1). Although a similar arrangement is
103 observed also on chr 10q26, at this locus the PAS signal is disrupted by a single nucleotide
104 variant (AUCAAA), preventing the polyadenylation of the *ter-DUX4* mRNA. By analogy, 4qB,
105 an alternative (to 4qA-PAS) haplotype at 4q (Figure S1), is considered not permissive for
106 the expression of the *DUX4* protein due to the lack of pLam. Based on these observations
107 (D4Z4 reduced copy number or heterozygous variants in chromatin remodelers, associated
108 with 4qA-PAS haplotype) a reduction of CpG methylation of D4Z4 has been proposed as a
109 faithful indicator of anomalous *DUX4* gene expression and it is used as a proxy of disease
110 status in the clinical setting²⁰.

111 This model is in contrast with clinical and epidemiological data showing reduced
112 penetrance and clinical variability in FSHD families²¹⁻²⁷, as well as with observations
113 obtained by adopting an articulate neurological assessment for the study of genotype-
114 phenotype correlation in FSHD. This latter approach includes the phenotypic classification
115 of probands and relatives by applying the Comprehensive Clinical Evaluation Form
116 (CCEF)²⁸, which beside quantifying the degree of motor disability²⁹, classifies individuals on
117 the basis of detailed phenotypic features which include parameters such as age at onset or
118 site of disease onset beside the detailed description of muscle impairment. Four main
119 descriptive categories have been created. They identify: (1) subjects presenting facial and
120 scapular girdle muscle weakness typical of FSHD (category A, subcategories A1-A3), (2)
121 subjects presenting an incomplete phenotype with muscle weakness limited to scapular
122 girdle or facial muscles (category B, subcategories B1 and B2), (3) asymptomatic or healthy
123 subjects (category C, subcategories C1 and C2), and (4) subjects with myopathic
124 phenotypes presenting clinical features not consistent with FSHD canonical phenotype
125 (category D, subcategories D1 and D2)²⁸. The application of this methodology for the

126 stratification of clinical phenotypes has permitted the quantification of the large phenotypic
127 variability observed in individuals carrying a D4Z4 Reduced Allele (DRA) which is in contrast
128 to the indication that a positive molecular test is the only determining aspect for FSHD
129 diagnosis^{30–37}. Studies also showed that the different categories are associated with
130 different disease trajectories; in particular subjects assessed with a classical FSHD
131 phenotype (Category A) display a steeper functional decline, than subjects with limited
132 muscle impairment (Category B)^{31,35}. Large genotype-phenotype studies also highlighted
133 differences between probands and relatives, including 30-50% of relatives remaining non-
134 penetrant carriers. Consistent with these findings it has been assessed that DRA with 4-8
135 RU have the frequency of a common polymorphism (3% in the general population^{25,38–40}),
136 which is not compatible with the incidence of FSHD; moreover the DRA 4qA-PAS haplotype
137 has a 2% frequency in the human population³⁹. Finally, studies reported a wide variability in
138 D4Z4 CpG methylation levels among FSHD index cases and FSHD families with no clear
139 association between D4Z4 methylation status and disease manifestations or severity^{32,41}.

140 Here, we report the results of comparative analyses of 86 haplotype-level assemblies
141 from the Human Pangenome Project dataset (hereafter PGR)⁴² and the outcome of high-
142 throughput CpG methylation assay of the D4Z4 repeat based on the T2T-CHM13 human
143 genome assembly (hereafter T2T). Our study displays all the potential of complete, high-
144 quality haplotype level genome assemblies for the analysis of D4Z4 arrays in human
145 diseases, confirms the high frequency of the permissive 4qA-PAS haplotype, which is
146 comparable to a common polymorphism and shows that across the 3.3 kb D4Z4 elements
147 only two sequences are 4q/10q-specific and amenable for the study of methylation patterns
148 in FSHD. This refined CpG methylation assay reveals that 4q/10q-specific reduced
149 methylation correlates with the presence of damaging *SMCHD1* heterozygous variants, but
150 not with the clinical status of FSHD2 cases and advocate revisions of previous findings
151 based on the GRCh38 reference (hereafter hg38).

152

153 **Subjects and methods**

154

155 **Participants**

156 The study cohort included 26 subjects (8F, 18M) reported to the Miogen Laboratory of the
157 University of Modena and Reggio Emilia for the diagnosis of FSHD⁴³. These cases carried
158 D4Z4 alleles with 10 or more repeats. In 8 cases the study was extended to one or both
159 parents, for a total of 14 subjects (7F, 7M). Overall, 5 trios and 3 parent-child couples were
160 analyzed. Four subjects carrying a 4U DRA were also included (2F, 2M). In three of these
161 subjects (Family C individuals III-1 and III-3; Family 1252 subject 297/08³²) D4Z4
162 methylation was previously assessed through the BSS assay using primer sets specific for
163 the distal 4qA and 4qA-Long (4qA-L) D4Z4 repeat regions³².

164 The clinical phenotype was determined (Table S1) according to the Comprehensive
165 Clinical Evaluation Form (CCEF), which evaluates the distribution and degree of motor
166 impairment, age at onset of motor impairment and site of muscle weakness at onset, the
167 presence of typical and atypical symptoms²⁴. The CCEF has been developed by the Italian
168 Clinical Network for FSHD to classify subjects on the basis of these clinical features. The
169 CCEF classifies: 1) subjects presenting facial and scapular girdle muscle weakness typical
170 of FSHD (category A, subcategories A1-A3), 2) subjects with muscle weakness limited to
171 scapular girdle or facial muscles (category B subcategories B1, B2), 3)
172 asymptomatic/healthy subjects (category C, subcategories C1, C2), 4) subjects with
173 myopathic phenotype presenting clinical features not consistent with FSHD canonical
174 phenotype (category D, subcategories D1, D2). Degree of motor impairment was also
175 established by applying a standardized evaluation that generates the FSHD score ²⁹.

176 Signed informed consent was obtained from all the subjects prior to the inclusion in
177 the study.

178

179 **Bisulfite sequencing**

180 Bisulfite sequence analysis was performed on high molecular weight genomic DNA
181 obtained from peripheral frozen blood through phenol-chloroform extraction. To assess the
182 methylation level at D4Z4 locus, specific primer sets were designed on hg38 genome
183 assembly using MethPrimer tool⁴⁴. No CpG sites were allowed in the primer sequence. The
184 amplicon selection was based on following criteria: i) the sequences contain SNPs allowing
185 to discriminate between the 4q and 10q D4Z4 repeats; and ii) include the maximum number
186 of CpG sites possible. Amplicons were representative of the whole D4Z4 array and included
187 functional domains such as the D4Z4 binding element (DBE)¹². Table S2 and Table S3,
188 Figure S2 summarize the characteristics of the four selected primer sets. Specific Illumina
189 adapters were added to each set of primers.

190 Bisulfite conversion was performed on 1 µg of genomic DNA by using the EpiTec
191 Bisulfite Kit (Cat N°59104 QIAGEN) following the manufacturer's protocol.

192 PCR amplification of both converted and non-converted DNA was performed using
193 the four selected primer pairs (Table S2). For every subject we generated two pools of
194 amplicons: one from bisulfite converted (BSC), one from non-converted DNA. Illumina
195 paired-end sequencing was then performed by Eurofins genomics on the amplicon pools
196 for a total of 36 Mb (60K read pairs per amplicon with 2 x 300 bp read mode).

197

198 **Bioinformatics analyses**

199 ***Identification and characterization of D4Z4-like sequences and D4Z4 associated*** 200 ***elements in human genome assemblies***

201 The complete sequence of a reference D4Z4 element was obtained from the hg38 assembly
202 in UCSC genome browser (<https://genome.ucsc.edu/>). Sequence similarity searches based
203 on the BLAST program⁴⁵ were performed to annotate D4Z4-like elements and functional

204 elements of the D4Z4 array (see below), in the hg38 reference assembly of the human
205 genome [[GRCh38.p14 - hg38 - Genome - Assembly - NCBI](#)], the T2T assembly [[T2T-
206 CHM13v1.1 - Genome - Assembly - NCBI](#)], and 86 distinct haplotype-level assemblies from
207 the PGR, as available from [https://github.com/human-
208 pangenomics/HPP_Year1_Assemblies](https://github.com/human-pangenomics/HPP_Year1_Assemblies).

209 Results of BLAST sequences similarity searches were stored in simple tabular format
210 and processed by custom Perl scripts. A similarity threshold of 85% and an alignment length
211 threshold of 300 aligned or more aligned residues were used to define D4Z4-like elements.
212 Matches were binned in 5 bins according to their size (<500bp; <1000bp; <1500bp;
213 <2500bp; and <3300bp). Matches of size in-between 3200 and 3300 bp were considered
214 to provide a complete representation of a D4Z4 macrosatellite (full match).

215 D4Z4 arrays were annotated based on the terminal repeat for the presence of qA or
216 qB sequence variants; presence/absence of the pLam sequence and integrity of the
217 AUUAAA polyadenylation signal (PAS). qA/B has been annotated based on the sequence
218 of the respective probes used for the FSHD diagnostic protocol; pLam has been annotated
219 based on the sequence reported in UCSC genome browser (<https://genome.ucsc.edu/>)
220 (Supplemental Material and Methods).

221

222 ***Similarity-based clustering of D4Z4 like elements and terminal DUX4 coding*** 223 ***sequences***

224 D4Z4-like sequences of 3200 or more bp in size (full matches) were extracted from their
225 respective genome assemblies and aligned with Kalign⁴⁶. Conserved alignment blocks were
226 extracted with Gblocks⁴⁷, using the following parameter: Minimum Length Of A Block: 3,
227 Allowed Gap Positions: none. A sequence similarity matrix was computed by means of the
228 EMBOSS distmat program⁴⁸, using the Tajima-Nei correction for multiple substitutions⁴⁹.
229 Phenetic clustering was performed by using the NJ algorithm, as implemented by hclust()

230 function from the cluster package in R⁵⁰. Six distinct groups with high sequence identity
231 were defined for the D4Z4 complete elements. D4Z4 from T2T were used to anchor/assign
232 each group to a chromosome.

233 The same method was applied to cluster *ter-DUX4* coding sequences by sequence
234 similarity. Terminal repeats were identified based on the presence of a qA or qB element *in*
235 *cis* and within 15 Kb from a complete *DUX4* ORF and tentatively assigned to chr 4 or 10
236 based on sequence similarity.

237

238 ***In-silico identification of target regions of primer sets in hg38 and T2T genome*** 239 ***assemblies***

240 Potential target regions of primer sets in the hg38 and T2T genome assemblies were
241 determined by a custom Perl script (Table S4). Sequence similarity searches of primer
242 sequences on the hg38 and T2T reference assemblies were performed by blastn, with
243 default parameters⁴⁵. All matches of 17 bp or longer and with at most 1 mismatch and all
244 perfect matches of 16 bp or longer were recorded. Genomic coordinates were cross-
245 referenced, and potential target genomic regions were identified as those spanned by a 5'
246 and 3' primer pairs, in the correct orientation, and within a distance >50 bp and <600 bp.
247 The same procedure was applied for the identification of potential primer target regions in
248 the BSC space, in this case both the genome and the primer sequences were converted *in-*
249 *silico*.

250

251 **Statistical analyses**

252 Distribution of values were compared by applying the non-parametric, two tailed
253 Kolmogorov Smirnov test, as implemented by the `ks.test()` function from the stats package
254 of the R programming language.

255

256 **Analysis of bisulfite converted reads**

257 ***Quality control and determination of target regions***

258 Reads' quality was assessed by Fastqc^{51,52}. Individual reports were merged by MultiQC⁵³
259 and quality metrics were visually inspected.

260 Non-converted amplicons' reads were aligned to T2T and hg38 human genome
261 assemblies using Bowtie2⁵⁴. Sample-level coverage profiles were computed by bedtools
262 genomecov⁵⁵. Genomic regions covered by 100 or more reads in at least 50% of samples
263 were considered for subsequent analyses. By this approach a total of 172 and 57 candidate
264 target regions were identified in the T2T and hg38 assemblies, respectively. 50/57 and
265 169/172 had at least a high similarity match (17 bp with at most 1 mismatch) to one or more
266 primer sequences.

267 Highly similar candidate regions were collapsed to provide a non-redundant
268 representation of repetitive sequences and allow non-ambiguous mapping of BSC reads.
269 Different sequence identity thresholds were tested through a titration analysis to identify the
270 ideal sequence similarity threshold. The total number of non-ambiguously mapped BSC
271 reads and the total number of CpGs covered by at least 10 BSC reads in at least 50% of
272 samples (testable CpGs) were recorded (Table 1). Considerations based on the number
273 and cumulative size of genomic regions that were merged, and on the number of testable
274 CpGs, prompted us to select 95% identity as the most appropriate threshold. In conclusion,
275 by comparing hg38 and T2T-based analyses of non-converted reads, we defined a total of
276 60 and 9 distinct consensus groups (groups of sequences GS) respectively (Table S5 and
277 Table S6). For every target region the consensus sequence was determined by majority
278 rule consensus, by applying the cons EMBOSS program⁴⁸. These consensus sequences
279 were used for the further analysis of methyl-seq data.

280

281 ***Assessment of methylation levels***

282 BSC reads were aligned to consensus target region sequences and methylation levels were
283 determined by Bismark⁵⁶ +Bowtie2. The bismark_methylation_extractor tool was applied to
284 compute CpG methylation levels. Only CpGs covered by more than 10 reads in at least 25
285 patients were considered for the delineation of methylation profiles.

286

287 ***SMCHD1* variants analysis**

288 FSHD2 probands are routinely tested for the presence of *SMCHD1*, *DNMT3B*, *LRIF1*
289 variants. We collected the identified *SMCHD1* variants in the cohort and performed a
290 reannotation on GRCh37 genome assembly using wAnnovar⁵⁷
291 (<https://wannovar.wglab.org/>; the analysis was performed on the 27.9.22) and filtered. We
292 excluded intron variants and considered only exonic (missense, stop or frameshift variants)
293 and splicing variants with a GnomAD exome all MAF < 10⁻⁴ (GnomAD frequency for PM2
294 pathogenic moderate rule in Varsome)⁵⁸. Among the filtered variants we then excluded the
295 ones which were benign and likely benign according to Varsome ACMG (American College
296 of Medical Genetics and Genomics) prediction^{59,60} (<https://varsome.com/>; the analysis was
297 performed on 27.9.22). As *SMCHD1* gene function has not been thoroughly investigated
298 yet, we decided not to filter out the variants predicted as VUS (Variants of Uncertain
299 Significance). Finally, variant validation through Sanger sequencing was performed on
300 probands and on relatives (Table S7).

301

302 **Results**

303

304 **T2T and PGR reveal the extended variability of D4Z4 repeats throughout the human 305 genome**

306 The extent of D4Z4-like repeat elements in the human genome was investigated by
307 analyzing T2T and 86 haplotype-level genome assemblies from the PGR⁴².

308 Our analyses uncovered hundreds of “complete” (>3200 bp) as well as partial (300
309 – 3200 bp) D4Z4-like sequences (hereafter D4Z4-I), irrespective of the reported geographic
310 origin of the subjects (Table S8). The cumulative size of D4Z4-I spanned in-between 700
311 Kb to 1.5 Mb of sequence in the haplotype-level assemblies of the 43 subjects from the
312 PGR and accounted for approximately 1.2 Mb of sequence in T2T. D4Z4-I were scattered
313 through several chromosomes in T2T and accounted for hundreds of kilobases of
314 sequences on chr 1, chr 4, chr 10, chr 15, chr 21 and chr 22 (Figure 2A). Notably, in hg38
315 only 86.5 Kb of D4Z4-I were observed, representing a more than 10-fold reduction
316 compared with T2T and PGR (Figure 2B). In hg38, D4Z4-I were prevalently associated with
317 the long arms of chr 4 and chr 10 (Figure 2A) where they are arranged as two large arrays
318 of 8 and 10 complete D4Z4 repeats, respectively. PGR and T2T instead harbored a large
319 number (minimum 285, maximum 626, Table S8) of “incomplete” D4Z4-I repeats of < 2.5
320 Kb in size, which are lacking in hg38. Conversely (Figure 2B), the number of complete D4Z4
321 elements included in hg38 was within the range of variability observed in T2T and PGR (26
322 to 121). Interestingly, subjects from AFR (African) ancestry were associated with a higher
323 variability in D4Z4-I copy numbers (Kolmogorov Smirnov test p-value p-value = 0.0006403),
324 while a narrower range was observed in individuals of EAS (Eastern Asia) and AMR (South
325 America) ancestry (Figure S3).

326 Figure 2C shows the variability and the distribution of complete D4Z4-I elements. Six
327 different groups were defined on the basis of sequence identity levels. D4Z4 from T2T were
328 used to anchor/assign each group to a chromosome. The largest clusters (Table 2) and the
329 majority of the sequences (87%) were assigned to either chr 4 (cluster-2) or chr 10 (cluster-
330 1); this notwithstanding, complete D4Z4 elements were also observed on chr 14 (cluster-3),
331 chr 22 (cluster-5 and cluster-6), chr 15 and chr 21 (cluster-4). The number of complete D4Z4
332 elements assigned to distinct clusters in the 86 PGR haplotypes is summarized in Figure
333 2D and Table 2. Chr 4 and chr 10 were associated with a median number of 22 and 19

334 repeats, respectively. The difference was statistically significant according to a Kolmogorov
335 Smirnov test (p-value: 0.02518). Complete D4Z4 repeats on chr 10 displayed lower
336 variability compared with those on chr 4. Substitution rates estimates at haplotype levels
337 suggested a high inter-individual variability, with rates ranging from 0.164 to 0.4312
338 (average of 0.245) substitutions for 100 bp for chr 4 and from 0.015 to 0.21 (average of
339 0.13) for chr 10 (Figure S4 and Figure S5). Interestingly, as shown in Table 2, more than
340 75% of the assemblies carried D4Z4-I assigned to cluster-6 (chr 22); whereas less than 50%
341 of the haplotype assemblies in the PGR had one sequence assigned to cluster-3 (chr 14),
342 cluster-4 (chr 15-21) and cluster-5 (chr 22).

343

344 **The pangenome assemblies show the high frequency of the 4qA-PAS haplotype** 345 **associated with FSHD in the world population**

346 The permissive DRA 4qA-PAS haplotype has been proposed to cause FSHD. However, we
347 previously reported a frequency of 2% for this haplotype³⁸, which is not compatible with the
348 prevalence of the disease⁶¹. We thus analyzed the T2T and PGR to investigate the
349 prevalence of this molecular signature in other populations and to study the conservation
350 and the integrity of *ter-DUX4*.

351 The pLam and telomeric qA and qB sequences (Supplemental Material and
352 Methods) were used to identify and classify terminal D4Z4 repeat in PGR, T2T and hg38. A
353 total of 172 D4Z4 terminal repeats were identified (Table S9) and tentatively assigned to chr
354 4 or chr 10 based on sequence similarity profiles (as described by the clustering of the
355 D4Z4-I sequences in Figure 2C). Haplotypes with inconsistent annotations (2 or more qA/qB
356 probes assigned to the same chromosome and none to the other chromosome) were not
357 included in subsequent analysis^{62,63}. By this approach the complete set of alleles at D4Z4
358 repeat loci at chr 4 and chr 10 was determined for 72 haplotypes from PGR (Table S9).
359 Patterns of conservation of *ter-DUX4* sequences were investigated. Figure 3A shows that

360 *ter-DUX4* assigned to the 10q D4Z4 terminal repeat are found to be highly similar and
361 formed a single cluster; while 4q *ter-DUX4* are more heterogeneous and were partitioned in
362 3 groups (Figure 3A-B). Interestingly, as shown in Figure 3C, in chr 10 and chr4_group 2
363 *ter-DUX4* displayed lower levels of variability (average 0.061 subs per 100 bp) compared to
364 both chr4_group 3 (average 0.12 subs per 100 bp) and chr4_group 1 (average 0.142 subs
365 per 100 bp); Figure 3A also shows that the most variable sequences (chr4_group 1 and
366 chr4_group 3) were preferentially associated with qB alleles. All *ter-DUX4* sequences were
367 found to be intact, and none was associated with frameshifts or premature stop codons.
368 Consistent with previous observations^{18,19}, all the alleles classified as qB did not carry a
369 pLam and were more proximal to the *ter-DUX4* compared to qA (Figure 3D-G). In the same
370 way, (Figure 3E-F) the vast majority of alleles classified as qA (112/119) were associated
371 with a pLam both on chr 4 (41/47) and chr 10 (71/72). Only terminal 4q pLam had valid
372 PAS⁶⁴; while in pLam assigned to 10q the PAS sequences were disrupted by a single
373 nucleotide substitution (AUUAAA→AUCAAA). Perfect conservation of PAS/PAS-like
374 elements was observed both on 4q and 10q terminal D4Z4 repeats. One qB haplotype was
375 assigned to 10q (Table 3).

376 Remarkably, in line with our previous observations³⁸⁻⁴⁰ (Table 3), 4/86 haplotypes
377 (4.6%) carried D4Z4 alleles with 6-8 RU and 4qA-PAS (Figure 3F). None of these
378 haplotypes carried disruptive variants in the *ter-DUX4* or in the pLam sequence element.

379 No statistically significant difference in the total number of chr 4 D4Z4 repeats was
380 observed when qA/qB alleles pLam+/pLam- alleles were compared (Figure 3H-I),
381 suggesting that alleles that are not permissive for the expression of the *ter-DUX4* are not
382 associated with significant differences in D4Z4 copy numbers compared with the permissive
383 qA-PAS haplotype.

384

385 **D4Z4 hypomethylation correlates with *SMCHD1* damaging variants but not with FSHD**
386 **clinical phenotype**

387 Results shown above uncover all the heterogeneity of D4Z4 repeats in the human genome,
388 including the presence of qA D4Z4 alleles with 8 repeats or fewer, and raise potential
389 concerns about the reliability of methylation/diagnostic tests based on the hg38. As bisulfite
390 treated DNA sequencing at targeted regions is commonly used to diagnose FSHD⁶⁵, the
391 uncovered variability of D4Z4-I sequences prompted us to investigate D4Z4 CpG
392 methylation comparing the hg38 and the T2T.

393 A set of primers spanning the entire length of the D4Z4 repeat and a total of 126
394 CpGs (Figure S2 and Table S3) were designed based on hg38 assembly to recapitulate
395 previous findings in the field⁴¹. This experimental design was used to assess a carefully
396 selected cohort of 26 index cases referred for FSHD2 diagnosis, and 14 relatives all carrying
397 10 or more RU. All the participants assayed were searched for damaging variants in the
398 chromatin remodelers *SMCHD1*, *DNMT3B* and *LRIF1*. In addition, 3 FSHD1 index cases
399 and 1 relative carrying 4 RU were studied.

400 Amplicon targeted regions were sequenced both in non-converted and in BSC DNA
401 samples. This approach was undertaken: i. to test the 4q/10q-specificity of the primer sets
402 commonly used and designed based on hg38 assembly⁴¹, ii. to facilitate the alignment of
403 short BSC derived reads. Both non-converted and BSC NGS data were used for sequence
404 alignment based on hg38 and T2T reference genomes (Table S4).

405 A total of 77% and 99.43% non-converted short reads mapped to hg38 and T2T
406 respectively. Moreover, the majority of BSC reads (62% for T2T and 54% for hg38) did not
407 align at 4q primer target regions, potentially suggesting “off target” sequences, and/or high
408 levels of heterogeneity in methylation levels. To identify genomic regions targeted by our
409 assay in an unbiased manner, genomic intervals covered by 100 or more non-converted
410 reads in at least 50% of the samples were recorded. A total of 57 and 172 genomic regions

411 in hg38 and T2T displayed these remarkable levels of coverage. Notably, 50/57 and
412 169/172 of the target regions delineated by coverage analyses had at least a high similarity
413 match (17 bp with at most 1 mismatch) to one or more primer sequences. Based on these
414 observations sequences above 95% sequence similarity were collapsed to prevent
415 inconsistent mapping of BSC short reads across highly similar D4Z4-I elements. By this
416 approach, we defined a total of 60 and 9 distinct consensus groups (groups of sequences,
417 hereafter GS) in T2T and hg38 respectively (Table S5 and Table S6). Consensus
418 sequences were computed for each GS and used for the further analysis of methyl-seq
419 data.

420 This *ad-hoc* bioinformatics workflow revealed remarkable differences in the
421 completeness, number, and breadth of testable CpGs when the GSs derived from hg38 and
422 T2T were compared: hg38 125 CpG; T2T 938 CpG. Global methylation patterns were
423 summarized in the form of a heatmap (Figure S6 and Figure S7). In the T2T-based analysis
424 a total of 4 distinct GS including 81 distinct CpGs had no missing data and complete
425 methylation profiles across all samples, as shown in Figure 4 and Table 4: GS1 (primer D6):
426 28 CpGs; GS2 (primer D1): 32 CpGs; GS5 (chr 13-chr 15): 16 CpGs and GS6 (chr 13-chr
427 15): 5 CpGs). Notably, only GS1 and GS2 included D4Z4 5' 4q/10q-specific regions
428 exclusively, for a total of 60 CpGs, while GS5 and GS6 mapped in D4Z4-I repeats on chr
429 13 and 15. The wide-range distributions of CpG methylation levels, as displayed in Figure
430 S8, highlight a high variability at regions GS1 and GS2, while off-target regions GS5 and
431 GS6 have high methylation levels, with a narrower distribution. Figure S9 shows a general
432 overview of the CpG methylation patterns observed in all the 60 distinct GSs defined in our
433 analysis. Taken together these findings indicate that a significant proportion of the BSC
434 reads recovered by our targeted high throughput assay derive from off-target, hyper-
435 methylated D4Z4-I regions. We also observed that when hg38 is considered as reference,
436 the criteria we established for the T2T-based analysis (namely no missing data and

437 complete methylation profiles across all samples) identified only a single group, GS3 (primer
438 D1) including 26 CpGs (Figure 4, Figure 5 and Table 4). In the light of these observations,
439 methylation profiles inferred from T2T were considered more informative and were selected
440 for subsequent analyses. Only primer D1 had complete data both in T2T and hg38 (Table
441 4); however, CpG methylation levels estimated on hg38 were slightly more elevated
442 compared to those observed on T2T at the equivalent target region (Figure S10, median
443 values: 49.13 hg38, 43.75 T2T). An observation that might be consistent with discrepancies
444 in the mapping of BSC reads on different genome assemblies.

445 Further, our analysis revealed that 4q/10q-specific methylation patterns at GS1 and
446 GS2 stratified samples into 3 groups: low methylation (average CpG methylation 24%),
447 intermediate methylation (average CpG methylation 48%), and high methylation (average
448 CpG methylation 71%). The methylation pattern of GS5 and GS6 was highly uniform across
449 all samples, with an average of 89%.

450 Interestingly, the observed methylation profiles did not correlate with the disease
451 status. As shown in Figure 6A, of 16 samples presenting low methylation 10 (62.5%) derived
452 from cases presenting the classical FSHD features (CCEF category A), 6 (37,5%) derived
453 from cases presenting incomplete (2, CCEF clinical Category B) or complex (4, CCEF
454 clinical category D) phenotypes. Of 16 samples, presenting intermediate methylation 9
455 (56.3%) derived from cases assessed as clinical category A. Of 12 samples presenting high
456 methylation 7 were from Category A cases (58.3%). Consistently, the methylation level
457 distributions assessed in DNA from subjects presenting the CCEF clinical categories A, B
458 and D was not statistically different (Figure S11). Instead, the majority (11/16) of participants
459 in the low methylation group carried Pathogenic, Likely Pathogenic variants or Variants of
460 Uncertain Significance according to ACMG classification⁶⁰ (P/LP/VUS variants) in *SMCHD1*
461 (Fisher exact test p-value 0.0003). This association was confirmed by the comparison of the
462 methylation profile distributions of *SMCHD1* P/LP/VUS variants carriers and non-carriers

463 (Figure 4, Kolmogorov-Smirnov test p -value $\leq 1e-16$). Identical patterns were recovered
464 also when methylation levels on the hg38 assembly were assessed with our approach
465 (Figure 5).

466

467 **Discussion**

468

469 **The pangenome assemblies provide novel hints to evaluate the significance of the** 470 **D4Z4-like sequences in the clinical setting**

471 Thirty years ago, reduction below a critical threshold of the number of repeats belonging to
472 the D4Z4 macrosatellite at 4q35 was causally associated with FSHD. D4Z4 is part a family
473 of repetitive elements characteristic of heterochromatin and it was then known that many
474 D4Z4-I were present in the human genome⁶⁶. At that time, the lack of a complete genome
475 assembly hampered the possibility of fully deciphering the significance of molecular and
476 epigenetic findings in FSHD. Very recently, these limitations have been overcome by the
477 availability of complete, high quality haplotype level genome assemblies. By analyzing the
478 haplotypes from the PGR we collected relevant information on the size, the distribution and
479 the composition of the D4Z4 locus, including distal sequence elements that are considered
480 a hallmark of the FSHD molecular signature. We established that the number of 3.3 kb D4Z4
481 repeats included in each array ranges from 6 RU to 89 RU⁴⁰. We observed that the distance
482 between the qA sequences and the last D4Z4 repeat varies. We also confirmed that qA
483 sequences are more distal compared to qB sequences; the pLam sequence is in association
484 with qA haplotypes¹⁸, and valid PAS is exclusively associated with 4q⁶⁴; chr 10 D4Z4
485 repeats were all marked by qA, with a single exception: a 10qB allele carrying 1 RU. Notably,
486 PGR haplotype analysis confirmed the high frequency of the permissive 4qA-PAS
487 haplotype, which is comparable to a common polymorphism⁶¹. In fact, in 4.6% of cases we
488 identified permissive haplotypes consisting of reduced 4q D4Z4 arrays with ≤ 8 RU and the

489 4qA-PAS. Since all *ter-DUX4* and pLam sequence elements associated with 4qA alleles
490 were found to be conserved and functionally intact, we speculate that disruptive genomic
491 variants at the *ter-DUX4* are unlikely to account for the staggering discrepancy between the
492 frequency of the presumed causative 4qA-PAS allele and the prevalence of FSHD.
493 Furthermore, no skewed association has been detected between 4q D4Z4 arrays copy
494 number and qA/B alleles¹⁹, indicating that non-permissive 4qB alleles and potentially
495 permissive 4qA alleles do not have a different selective pressure for the maintenance of
496 the physiological number of D4Z4 RU.

497 The consistency between data obtained from the analysis of PGR and previous
498 observations^{6,61,66} shows the validity of the PGR to advance knowledge of the molecular
499 signature used to diagnose FSHD on a world-wide scale.

500 The PGR analysis also identified 6 haplotypes with peculiar annotations where 2 or
501 more qA/qB probes were assigned to the same chromosome and none to the other
502 chromosome (Table S9). This condition has been reported before⁶² and is ascribable to the
503 spreading of repetitive and polymorphic D4Z4-I at different loci of the genome^{63,67}.

504 Our comparative analyses of complete D4Z4-I (Figure 2C) highlighted that D4Z4
505 repeats, as well as the *ter-DUX4*, located at chr 4 present a higher variability compared to
506 chr 10. In addition, PGR haplotypes from different ancestries showed different levels of
507 variability in the number of complete D4Z4-I, with a significant difference in D4Z4 copy
508 number between individuals of African ancestry, native Americans and far east Asians.

509

510 **D4Z4 methylation assay must consider the T2T sequence for the identification of** 511 **4q/10q-specific CpGs**

512 Reduced CpG methylation at D4Z4 has been widely studied and proposed as a biomarker
513 for the presence of FSHD. Previous studies of D4Z4 methylation presented some
514 drawbacks: technical limitations and lack of reproducibility led to discordant results

515 regarding which regions are the most representative for D4Z4 methylation status^{20,32,68–74}.
516 In addition, the correlation between D4Z4 methylation level and FSHD clinical status has
517 often been postulated, but different technical settings in the various studies, the lack of clear
518 clinical evaluation of patients and the prominent absence of family studies had hindered the
519 interpretation of results^{20,41,73,74}. In the present work, we comprehensively evaluated the
520 D4Z4 methylation status by considering a high number of reads and excluding technical
521 biases introduced by the previous studies based on hg38 and on standard bioinformatics
522 pipelines. We established (Figure S2) that the primer sets BSS-D1 and BSS-D6 identify
523 CpG methylation signals deriving from 4q/10q-specific regions only, whereas the primer
524 sets BSS-D3 and BSS-D5 amplify D4Z4-I distributed not only at 4q and 10q but also on
525 chromosome 1 and on the short arms of acrocentric chromosomes. Thus, the study of D4Z4
526 CpG methylation can be hazardous if the distribution of D4Z4-I elements and their variability
527 are not correctly taken into account.

528 Despite all this, the application of short read-based assays for the analysis of
529 repetitive elements of relatively large size has limitations. When non-converted reads were
530 aligned to both hg38 and T2T, more than 50% of mapped reads did not align at expected
531 target regions. Nevertheless, the total fraction of mapped reads and the number of testable
532 CpGs were higher when aligned to T2T compared to hg38, confirming that T2T provides a
533 more accurate representation of the human genome. We can anticipate that long read
534 sequencing technologies, such as Oxford Nanopore, will provide relevant technical
535 advances and open new possibilities for the study of repetitive elements' organization and
536 function in the genome. In this respect, Butterfield and colleagues (2023)⁷⁵ applied targeted
537 nanopore sequencing to FSHD patients and healthy subjects and showed that an
538 asymmetric methylation gradient forms in a length-dependent manner at the proximal end
539 of the D4Z4 repeat array reaching saturation approximately at the 10th repeat. They also
540 observed a highly similar methylation pattern at 4q and 10q which was irrespective of the

541 clinical phenotype (FSHD or healthy)⁷⁵. This notwithstanding, at present, long-read
542 sequencing technologies are not mature for systematic large-scale applications in the
543 clinical settings. Instead, the targeted CpG methylation analysis we set up is manageable
544 on a large scale as it follows a simple protocol, is less expensive and highly reproducible.
545 Remarkably, our data are consistent with the results obtained by Butterfield and colleagues.

546

547 **D4Z4 CpG methylation level is not a predictor of FSHD disease**

548 Beside the formal demonstration that primer design is a crucial point in the analysis of D4Z4
549 methylation, our study also revealed that contrary to common knowledge²⁰, not all DNAs
550 from myopathic patients referred as FSHD2 (26 index cases) presented a low methylation
551 level. Our analysis identified three clusters with low, intermediate and high D4Z4 CpG
552 methylation (mean at 24%, 48% and 71% respectively). Figure 6 depicts the results of
553 several comparisons regarding the following parameters: clinical phenotype, D4Z4
554 methylation level and *SMCHD1* mutational status. Figure 6A shows that some genotype-
555 phenotype intersections are not in agreement with the current indications for FSHD
556 diagnosis^{76,77}. Remarkably, we observed the full range of D4Z4 methylation levels in DNA
557 from the 26 participants presenting a classical FSHD phenotype (CCEF Clinical Category
558 A); of those 16 carriers of a damaging *SMCHD1* variant, 11 presented low 4q/10q-specific
559 D4Z4 CpG methylation, 5 displayed intermediate CpG methylation. We also observed 5
560 samples with low CpG methylation with no damaging *SMCHD1* or *DNMT3B* or *LRIF1*
561 variants and presenting classical (2), incomplete (1) or complex (2) phenotypes. underlying
562 the genetic heterogeneity of this molecular phenotype.

563 The variants we identified were distributed all along *SMCHD1* gene body and did not
564 provide information on specific *SMCHD1* region/domain that might account for the reduced
565 methylation at 4q/10q-specific D4Z4 sequences (Figure 6B).

566 To further explore this aspect, we investigated genotype-phenotype correlation in
567 eight trios/parent-child pairs participating to this study (Figure 6C, Table S7). We observed
568 participants presenting the classical FSHD phenotype (CCEF Category A), healthy subjects
569 (CCEF Category C) or complex phenotypes presenting atypical features (CCEF Category
570 D). According to our data, only in families 1 and 26 subjects presenting a classic FSHD
571 phenotype displayed a reduced CpG methylation profile and carried a damaging
572 heterozygous *SMCHD1* variant. In family 22 the classic clinical phenotype, in presence of
573 D4Z4 reduced methylation, was not associated with variants in known chromatin modifiers.
574 In family 8 the probands presented the classic clinical phenotype associated with
575 intermediate D4Z4 methylation level in presence of a damaging *SMCHD1* variant, as the
576 healthy mother. In family 27 the proband presented the Category A phenotype associated
577 with reduced CpG methylation and a damaging *SMCHD1* variant, whereas intermediate
578 D4Z4 methylation level was detected in the healthy father carrying the same *SMCHD1*
579 allele. In family 30 father and daughter presented a complex phenotype with atypical clinical
580 symptoms (Category D1); both reported clinical onset at pelvic girdle and distal lower limbs.
581 They were heterozygous for *SMCHD1* damaging variant and presented reduced D4Z4 CpG
582 methylation. In families 9 both the probands and her father presented Category A phenotype
583 without displaying reduced D4Z4 CpG methylation and carrying WT *SMCHD1* alleles, as
584 the proband in family 11.

585 Altogether, these results suggest that *SMCHD1* modulates, directly or indirectly,
586 D4Z4 methylation at 4q and 10q, at the same time it seems that deleterious variants at
587 *SMCHD1* are not sufficient for the pathogenesis of FSHD and cannot be considered faithful
588 indicators of FSHD clinical status.

589 In conclusion, our work uncovered the complexity of the genomic setting of D4Z4-I
590 and showed that a molecular diagnostic test for FSHD2 based on *SMCHD1* inactivation and
591 D4Z4 hypomethylation must follow rigorous protocols to avoid biased interpretation.

592 Molecular analysis should be considered together with precise clinical assessment and
593 complemented by family studies for the proper interpretation of results. Moreover, the
594 variability unveiled by our analysis should warn about the risk of relying on a single ideal
595 reference genome sequence in FSHD. In fact, this simplification could strongly narrow our
596 capacity of observing the variability ascribable to each single human genome producing
597 biased interpretations of sequencing data. Finally, our study indicates that a linear approach
598 to diagnose FSHD is obsolete and novel genomic approaches integrated with the precise
599 phenotypic description of patients and their families are needed for the comprehension of
600 the molecular network leading to muscle wasting in FSHD.

601

602 **Declarations**

603

604 **Acknowledgements**

605 We are grateful to Francesca Brocco for contributing to the molecular analysis. In addition,
606 we are indebted to all FSHD patients and their families for participating in this study. We
607 acknowledge the Fondazione di Modena for FAR-FOMO 2021 grant funding this work.

608

609 **Authors' contributions**

610 RGT, VS, MC conceived the study concept and supervised the project; VS, SP conducted
611 the molecular analysis; MC conducted bioinformatics analysis; PK and MK contributed to
612 bioinformatic analysis; LR, SCP, MGD, CR, SB, LM, DL contributed to the patients clinical
613 analysis and sample collection; FMS, GP contributed to results discussion; RGT, VS, MC,
614 SP wrote the paper (original draft); all authors read and approved the final manuscript.

615

616 **Declaration of interests**

617 The authors declare no competing interests.

618

619 **Web sources**

620 The PGR dataset analyzed during the current study is available at

621 https://github.com/human-pangenomics/HPP_Year1_Assemblies.

622

623 **Availability of data and materials**

624 Data generated and/or analyzed during the current study and not included in this
625 published article are available from the corresponding author on reasonable request.

626

627 **Ethics approval and consent to participate**

628 Signed informed consent was obtained from all the subjects prior to the inclusion in the
629 study. The study was approved by the ethics committee of Emilia Romagna Area Vasta
630 Nord 743/2022/OSS/UNIMO SIRER ID 5111.

631

632 **References**

633

- 634 1. Deenen, J.C.W., Arnts, H., Van Der Maarel, S.M., Padberg, G.W., Verschuuren, J.J.G.M.,
635 Bakker, E., Weinreich, S.S., Verbeek, A.L.M., and Van Engelen, B.G.M. (2014). Population-
636 based incidence and prevalence of facioscapulohumeral dystrophy. *Neurology* 83, 1056–
637 1059. 10.1212/WNL.0000000000000797.
- 638 2. Mostacciuolo, M.L., Pastorello, E., Vazza, G., Miorin, M., Angelini, C., Tomelleri, G., Galluzzi,
639 G., and Trevisan, C. Pietro (2009). Facioscapulohumeral muscular dystrophy:
640 Epidemiological and molecular study in a north-east Italian population sample. *Clinical*
641 *Genetics* 75, 550–555. 10.1111/j.1399-0004.2009.01158.x.
- 642 3. Dumbovic, G., Forcales, S. V., and Perucho, M. (2017). Emerging roles of macrosatellite
643 repeats in genome organization and disease development. *Epigenetics* 12, 515–526.
644 10.1080/15592294.2017.1318235.
- 645 4. Lunt, P.W. (1998). 44th ENMC International Workshop: Facioscapulohumeral muscular
646 dystrophy: Molecular studies, 19-21 July 1996, Naarden, The Netherlands. *Neuromuscular*
647 *Disorders* 8, 126–130. 10.1016/S0960-8966(98)00012-1.

- 648 5. Hewitt, J.E., Lyle, R., Clark, L.N., Valleley, E.M., Wright, T.J., Wijmenga, C., Van Deutekom, J.C.
649 t., Francis, F., Sharpe, P.T., Hofker, M., et al. (1994). Analysis of the tandem repeat locus
650 D4Z4 associated with facioscapulohumeral muscular dystrophthyy. *Human Molecular*
651 *Genetics* 3, 1287–1295. 10.1093/hmg/3.8.1287.
- 652 6. Winokur, S.T., Bengtsson, U., Feddersen, J., Mathews, K.D., Weiffenbach, B., Bailey, H.,
653 Markovich, R.P., Murray, J.C., Wasmuth, J.J., Altherr, M.R., et al. (1994). The DNA
654 rearrangement associated with facioscapulohumeral muscular dystrophy involves a
655 heterochromatin-associated repetitive element: Implications for a role of chromatin
656 structure in the pathogenesis of the disease. *Chromosome Research* 2, 225–234.
657 10.1007/BF01553323.
- 658 7. Zhang, X.Y., Loflin, P.T., Gehrke, C.W., Andrews, P.A., and Ehrlich, M. (1987).
659 Hypermethylation of human DNA sequences in embryonal carcinoma cells and somatic
660 tissues but not in sperm. *Nucleic Acids Research* 15, 9429–9449.
661 10.1093/nar/15.22.9429.
- 662 8. Agresti, A., Meneveri, R., Siccardi, A.G., Marozzi, A., Corneo, G., Gaudi, S., and Ginelli, E.
663 (1989). Linkage in human heterochromatin between highly divergent Sau3A repeats and
664 a new family of repeated DNA sequences (HaeIII family). *Journal of Molecular Biology* 205,
665 625–631. 10.1016/0022-2836(89)90308-2.
- 666 9. Meneveri, R., Agresti, A., Marozzi, A., Saccone, S., Rocchi, M., Archidiacono, N., Corneo, G.,
667 Valle, G.D., and Ginelli, E. (1993). Molecular organization and chromosomal location of
668 human GC-rich heterochromatic blocks. *Gene* 123, 227–234. 10.1016/0378-
669 1119(93)90128-P.
- 670 10. Deidda, G., Cacurri, S., Grisanti, P., Vigneti, E., Piazzo, N., and Felicetti, L. (1995). Physical
671 mapping evidence for a duplicated region on chromosome 10qter showing high homology
672 with the facioscapulohumeral muscular dystrophy locus on chromosome 4qter. *European*
673 *Journal of Human Genetics* 3, 155–167. 10.1159/000472291.
- 674 11. Deutekom, J.C.T. van, Wljmenga, C., Tlenhoven, E.A.E. van, Gruter, A.M., Hewitt, J.E.,
675 Padberg, G.W., van Ommen, G.J.B., Hofker, M.H., and Fronts, R.R. (1993). FSHD associated
676 DNA rearrangements are due to deletions of integral copies of a 3.2 kb tandemly repeated
677 unit. *Human Molecular Genetics* 2, 2037–2042. 10.1093/hmg/2.12.2037.
- 678 12. Gabellini, D., Green, M.R., and Tupler, R. (2002). Inappropriate gene activation in FSHD: A
679 repressor complex binds a chromosomal repeat deleted in dystrophic muscle. *Cell* 110,
680 339–348. 10.1016/S0092-8674(02)00826-7.
- 681 13. Cabianca, D.S., Casa, V., Bodega, B., Xynos, A., Ginelli, E., Tanaka, Y., and Gabellini, D.
682 (2012). A long ncRNA links copy number variation to a polycomb/trithorax epigenetic
683 switch in fshd muscular dystrophy. *Cell* 149, 819–831. 10.1016/j.cell.2012.03.035.
- 684 14. Laoudj-Chenivresse, D., Carnac, G., Bisbal, C., Hugon, G., Bouillot, S., Desnuelle, C., Vassetzky,
685 Y., and Fernandez, A. (2005). Increased levels of adenine nucleotide translocator 1 protein
686 and response to oxidative stress are early events in facioscapulohumeral muscular
687 dystrophy muscle. *Journal of molecular medicine (Berlin, Germany)* 83, 216–224.
688 10.1007/s00109-004-0583-7.

- 689 15. Lemmers, R.J.L.F., Tawil, R., Petek, L.M., Balog, J., Block, G.J., Santen, G.W.E., Amell, A.M.,
690 Vliet, P.J. Van Der, Almomani, R., Straasheijm, R., et al. (2012). Digenic inheritance of an
691 SMCHD1 mutation and an FSHD-permissive D4Z4 allele causes facioscapulohumeral
692 muscular dystrophy type 2. *Nat Genet.* 44, 1370–1374. 10.1038/ng.2454.Digenic.
- 693 16. Van Den Boogaard, M.L., Lemmers, R.J.L.F., Balog, J., Wohlgemuth, M., Auranen, M.,
694 Mitsuhashi, S., Van Der Vliet, P.J., Straasheijm, K.R., Van Den Akker, R.F.P., Kriek, M., et al.
695 (2016). Mutations in DNMT3B Modify Epigenetic Repression of the D4Z4 Repeat and the
696 Penetrance of Facioscapulohumeral Dystrophy. *American Journal of Human Genetics* 98,
697 1020–1029. 10.1016/j.ajhg.2016.03.013.
- 698 17. Hamanaka, K., Šikrová, D., Mitsuhashi, S., Masuda, H., Sekiguchi, Y., Sugiyama, A., Shibuya,
699 K., Lemmers, R.J.L.F., Goossens, R., Ogawa, M., et al. (2020). Homozygous nonsense variant
700 in LRIF1 associated with facioscapulohumeral muscular dystrophy. *Neurology* 94, E2441–
701 E2447. 10.1212/WNL.0000000000009617.
- 702 18. Van Geel, M., Dickson, M.C., Beck, A.F., Bolland, D.J., Frants, R.R., Van der Maarel, S.M., De
703 Jong, P.J., and Hewitt, J.E. (2002). Genomic analysis of human chromosome 10q and 4q
704 telomeres suggests a common origin. *Genomics* 79, 210–217. 10.1006/geno.2002.6690.
- 705 19. Lemmers, R.J.L.F., de Kievit, P., Sandkuijl, L., Padberg, G.W., van Ommen, G.J.B., Frants, R.R.,
706 and van der Maarel, S.M. (2002). Facioscapulohumeral muscular dystrophy is uniquely
707 associated with one of the two variants of the 4q subtelomere. *Nature Genetics* 32, 235–
708 236. 10.1038/ng999.
- 709 20. Gould, T., Jones, T.I., and Jones, P.L. (2021). Precise epigenetic analysis using targeted
710 bisulfite genomic sequencing distinguishes fshd1, fshd2, and healthy subjects. *Diagnostics*
711 11, 1–17. 10.3390/diagnostics11081469.
- 712 21. Ricci, G., Scionti, I., Sera, F., Govi, M., D’Amico, R., Frambolli, I., Mele, F., Filosto, M., Vercelli,
713 L., Ruggiero, L., et al. (2013). Large scale genotype-phenotype analyses indicate that novel
714 prognostic tools are required for families with facioscapulohumeral muscular dystrophy.
715 *Brain* 136, 3408–3417. 10.1093/brain/awt226.
- 716 22. K Goto, I Nishino, Y.K.H. (2004). Very low penetrance in 85 Japanese families with
717 facioscapulohumeral muscular dystrophy 1A K. *J Med Genet* 41.
718 10.4324/9780203994559-14.
- 719 23. Wu, Z.Y., Wang, Z.Q., Murong, S.X., and Wang, N. (2004). FSHD in Chinese population:
720 Characteristics of translocation and genotype-phenotype correlation. *Neurology* 63, 581–
721 583. 10.1212/01.WNL.0000133210.93075.81.
- 722 24. Nikolic, A., Ricci, G., Sera, F., Bucci, E., Govi, M., Mele, F., Rossi, M., Ruggiero, L., Vercelli, L.,
723 Ravaglia, S., et al. (2016). Clinical expression of facioscapulohumeral muscular dystrophy
724 in carriers of 1-3 D4Z4 reduced alleles: Experience of the FSHD Italian National Registry.
725 *BMJ Open* 6, 1–10. 10.1136/bmjopen-2015-007798.
- 726 25. Nakagawa, M., Matsuzaki, T., Higuchi, I., Fukunaga, H., Inui, T., Nagamitsu, S., Yamada, H.,
727 Arimura, K., and Osame, M. (1997). Facioscapulohumeral Muscular Dystrophy: Clinical
728 Diversity and Genetic Abnormalities in Japanese Patients. *Internal Medicine* 36, 333–339.
729 10.2169/internalmedicine.36.333.

- 730 26. Sakellariou, P., Kekou, K., Fryssira, H., Sofocleous, C., Manta, P., Panousopoulou, A.,
731 Gounaris, K., and Kanavakis, E. (2012). Mutation spectrum and phenotypic manifestation
732 in FSHD Greek patients. *Neuromuscular Disorders* 22, 339–349.
733 10.1016/j.nmd.2011.11.001.
- 734 27. Salort-Campana, E., Nguyen, K., Bernard, R., Jouve, E., Solé, G., Nadaj-Pakleza, A.,
735 Niederhauser, J., Charles, E., Ollagnon, E., Bouhour, F., et al. (2015). Low penetrance in
736 facioscapulohumeral muscular dystrophy type 1 with large pathological D4Z4 alleles: A
737 cross-sectional multicenter study. *Orphanet Journal of Rare Diseases* 10, 4–11.
738 10.1186/s13023-014-0218-1.
- 739 28. Ricci, G., Ruggiero, L., Vercelli, L., Sera, F., Nikolic, A., Govi, M., Mele, F., Daolio, J., Angelini,
740 C., Antonini, G., et al. (2016). A novel clinical tool to classify facioscapulohumeral muscular
741 dystrophy phenotypes. *Journal of Neurology* 263, 1204–1214. 10.1007/s00415-016-
742 8123-2.
- 743 29. Lamperti, C., Fabbri, G., Vercelli, L., D’Amico, R., Frusciante, R., Bonifazi, E., Fiorillo, C.,
744 Borsato, C., Cao, M., Servida, M., et al. (2010). A standardized clinical evaluation of patients
745 affected by facioscapulohumeral muscular dystrophy: The FSHD clinical score. *Muscle and*
746 *Nerve* 42, 213–217. 10.1002/mus.21671.
- 747 30. Ruggiero, L., Mele, F., Manganelli, F., Bruzzese, D., Ricci, G., Vercelli, L., Govi, M., Vallarola,
748 A., Tripodi, S., Villa, L., et al. (2020). Phenotypic Variability Among Patients With D4Z4
749 Reduced Allele Facioscapulohumeral Muscular Dystrophy. *JAMA network open* 3,
750 e204040. 10.1001/jamanetworkopen.2020.4040.
- 751 31. Vercelli, L., Mele, F., Ruggiero, L., Sera, F., Tripodi, S., Ricci, G., Vallarola, A., Villa, L., Govi,
752 M., Maranda, L., et al. (2021). A 5-year clinical follow-up study from the Italian National
753 Registry for FSHD. *Journal of Neurology* 268, 356–366. 10.1007/s00415-020-10144-7.
- 754 32. Nikolic, A., Jones, T.I., Govi, M., Mele, F., Maranda, L., Sera, F., Ricci, G., Ruggiero, L., Vercelli,
755 L., Portaro, S., et al. (2020). Interpretation of the epigenetic signature of
756 facioscapulohumeral muscular dystrophy in light of genotype-phenotype studies.
757 *International Journal of Molecular Sciences* 21. 10.3390/ijms21072635.
- 758 33. Ricci, G., Cammish, P., Siciliano, G., Tupler, R., Lochmuller, H., and Evangelista, T. (2019).
759 Phenotype may predict the clinical course of facioscapulohumeral muscular dystrophy.
760 *Muscle and Nerve* 59, 711–713. 10.1002/mus.26474.
- 761 34. Ricci, G., Mele, F., Govi, M., Ruggiero, L., Sera, F., Vercelli, L., Bettio, C., Santoro, L., Mongini,
762 T., Villa, L., et al. (2020). Large genotype–phenotype study in carriers of D4Z4 borderline
763 alleles provides guidance for facioscapulohumeral muscular dystrophy diagnosis.
764 *Scientific Reports* 10, 1–12. 10.1038/s41598-020-78578-7.
- 765 35. He, J.J., Lin, X.D., Lin, F., Xu, G.R., Xu, L.Q., Hu, W., Wang, D.N., Lin, H.X., Lin, M.T., Wang, N., et
766 al. (2018). Clinical and genetic features of patients with facial-sparing
767 facioscapulohumeral muscular dystrophy. *European Journal of Neurology* 25, 356–364.
768 10.1111/ene.13509.
- 769 36. Wang, Z., Qiu, L., Lin, M., Chen, L., Zheng, F., Lin, L., Lin, F., Ye, Z., Lin, X., He, J., et al. (2022).
770 Prevalence and disease progression of genetically-confirmed facioscapulohumeral
771 muscular dystrophy type 1 (FSHD1) in China between 2001 and 2020: a nationwide

- 772 population-based study. *The Lancet regional health. Western Pacific* 18, 100323.
773 [10.1016/j.lanwpc.2021.100323](https://doi.org/10.1016/j.lanwpc.2021.100323).
- 774 37. Salsia, V., Vatteemi, G.N.A., and Tupler, R.G. (2023). The FSHD jigsaw: are we placing the
775 tiles in the right position? *Curr Opin Neurol*. [10.1097/WCO.0000000000001176](https://doi.org/10.1097/WCO.0000000000001176).
- 776 38. Wohlgenuth, M., Lemmers, R.J., Van der Kooi, E.L., Van der Wielen, M.J., Van Overveld, P.G.,
777 Dauwerse, H., Bakker, E., Frants, R.R., Padberg, G.W., and Van der Maarel, S.M. (2003).
778 Possible phenotypic dosage effect in patients compound heterozygous for FSHD-sized
779 4q35 alleles. *Neurology* 61, 909–913. [10.1212/WNL.61.7.909](https://doi.org/10.1212/WNL.61.7.909).
- 780 39. Scionti, I., Greco, F., Ricci, G., Govi, M., Arashiro, P., Vercelli, L., Berardinelli, A., Angelini, C.,
781 Antonini, G., Cao, M., et al. (2012). Large-scale population analysis challenges the current
782 criteria for the molecular diagnosis of fascioscapulohumeral muscular dystrophy.
783 *American Journal of Human Genetics* 90, 628–635. [10.1016/j.ajhg.2012.02.019](https://doi.org/10.1016/j.ajhg.2012.02.019).
- 784 40. Van Overveld, P.G.M., Lemmers, R.J.F.L., Deidda, G., Sandkuijl, L., Padberg, G.W., Frants,
785 R.R., and Van Der Maarel, S.M. (2000). Interchromosomal repeat array interactions
786 between chromosomes 4 and 10: A model for subtelomeric plasticity. *Human Molecular*
787 *Genetics* 9, 2879–2884. [10.1093/hmg/9.19.2879](https://doi.org/10.1093/hmg/9.19.2879).
- 788 41. Salsi, V., Magdinier, F., and Tupler, R. (2020). Does DNA methylation matter in FSHD?
789 *Genes* 11. [10.3390/genes11030258](https://doi.org/10.3390/genes11030258).
- 790 42. Wang, T., Antonacci-Fulton, L., Howe, K., Lawson, H.A., Lucas, J.K., Phillippy, A.M., Popejoy,
791 A.B., Asri, M., Carson, C., Chaisson, M.J.P., et al. (2022). The Human Pangenome Project: a
792 global resource to map genomic diversity. *Nature* 604, 437–446. [10.1038/s41586-022-](https://doi.org/10.1038/s41586-022-04601-8)
793 [04601-8](https://doi.org/10.1038/s41586-022-04601-8).
- 794 43. Bettio, C., Salsi, V., Orsini, M., Calanchi, E., Magnotta, L., Gagliardelli, L., Kinoshita, J.,
795 Bergamaschi, S., and Tupler, R. (2021). The Italian National Registry for FSHD: an
796 enhanced data integration and an analytics framework towards Smart Health Care and
797 Precision Medicine for a rare disease. *Orphanet Journal of Rare Diseases* 16, 1–21.
798 [10.1186/s13023-021-02100-z](https://doi.org/10.1186/s13023-021-02100-z).
- 799 44. Li, L.C., and Dahiya, R. (2002). MethPrimer: Designing primers for methylation PCRs.
800 *Bioinformatics* 18, 1427–1431. [10.1093/bioinformatics/18.11.1427](https://doi.org/10.1093/bioinformatics/18.11.1427).
- 801 45. Altschup, S.F., Gish, W., Pennsylvania, T., and Park, U. (1990). Basic Local Alignment
802 Search Tool 2Department of Computer Science. 403–410.
- 803 46. Lassmann, T. (2020). Kalign 3: Multiple sequence alignment of large datasets.
804 *Bioinformatics* 36, 1928–1929. [10.1093/bioinformatics/btz795](https://doi.org/10.1093/bioinformatics/btz795).
- 805 47. Castresana, J. (2000). Selection of conserved blocks from multiple alignments for their use
806 in phylogenetic analysis. *Molecular Biology and Evolution* 17, 540–552.
807 [10.1093/oxfordjournals.molbev.a026334](https://doi.org/10.1093/oxfordjournals.molbev.a026334).
- 808 48. Rice, P., Longden, L., and Bleasby, A. (2000). EMBOSS: The European Molecular Biology
809 Open Software Suite. *Trends in Genetics* 16, 276–277. [10.1016/S0168-9525\(00\)02024-2](https://doi.org/10.1016/S0168-9525(00)02024-2).
- 810 49. Tajima, F., and Nei, M. (1984). Estimation of evolutionary distance between nucleotide
811 sequences. *Molecular Biology and Evolution* 1, 269–285.

- 812 50. R Core Team (2020). R: A language and environment for statistical computing. R
813 Foundation for Statistical Computing, Vienna, Austria. <https://www.r-project.org/>.
- 814 51. Andrews, S. (2010). FastQC: A Quality Control Tool for High Throughput Sequence Data.
815 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
- 816 52. FastQC (2015). <https://qubeshub.org/resources/fastqc>.
- 817 53. Ewels, P., Lundin, S., and Max, K. (2016). Data and text mining MultiQC : summarize
818 analysis results for multiple tools and samples in a single report. *32*, 3047–3048.
819 10.1093/bioinformatics/btw354.
- 820 54. Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2.
821 *Nature Methods* *9*, 357–359. 10.1038/nmeth.1923.
- 822 55. Quinlan, A.R., and Hall, I.M. (2010). BEDTools : a flexible suite of utilities for comparing
823 genomic features. *26*, 841–842. 10.1093/bioinformatics/btq033.
- 824 56. Krueger, F., and Andrews, S.R. (2011). Bismark: A flexible aligner and methylation caller
825 for Bisulfite-Seq applications. *Bioinformatics* *27*, 1571–1572.
826 10.1093/bioinformatics/btr167.
- 827 57. Chang, X.; Wang, K. (2008). wANNOVAR: annotating genetic variants for personal
828 genomes via the web. *Journal of Human Genetics* *6*, 2166–2171. 10.1136/jmedgenet-
829 2012-100918.wANNOVAR.
- 830 58. Chen, K., Hu, J., Moore, D.L., Liu, R., Kessans, S.A., Breslin, K., Lucet, I.S., Keniry, A., Leong,
831 H.S., Parish, C.L., et al. (2015). Genome-wide binding and mechanistic analyses of Smchd1-
832 mediated epigenetic regulation. *Proceedings of the National Academy of Sciences of the*
833 *United States of America* *112*, E3535–E3544. 10.1073/pnas.1504232112.
- 834 59. Kopanos, C., Tsiolkas, V., Kouris, A., Chapple, C.E., Albarca Aguilera, M., Meyer, R., and
835 Massouras, A. (2019). VarSome: the human genomic variant search engine. *Bioinformatics*
836 *35*, 1978–1980. 10.1093/bioinformatics/bty897.
- 837 60. Richards S, Aziz N, Bale S, Bick D, Das S, Gastier-Foster J, Grody WW, Hegde M, Lyon E,
838 Spector E, Voelkerding K, R.H., and Committee, A.L.Q.A. (2015). Standards and guidelines
839 for the interpretation of sequence variants: a joint consensus recommendation of the
840 American College of Medical Genetics and Genomics and the Association for Molecular
841 Pathology. *Genet Med* *17*, 405–424. 10.1038/gim.2015.30.
- 842 61. Scionti, I., Greco, F., Ricci, G., Govi, M., Arashiro, P., Vercelli, L., Berardinelli, A., Angelini, C.,
843 Antonini, G., Cao, M., et al. (2012). Large-scale population analysis challenges the current
844 criteria for the molecular diagnosis of fascioscapulohumeral muscular dystrophy.
845 *American Journal of Human Genetics* *90*, 628–635. 10.1016/j.ajhg.2012.02.019.
- 846 62. Delourme, M., Charlene, C., Gerard, L., Ganne, B., Perrin, P., Vovan, C., Bertaux, K., Nguyen,
847 K., and Magdinier, F. (2023). Complex 4q35 and 10q26 Rearrangements.
848 10.1212/NXG.000000000200076.
- 849 63. Cacurri, S., Piazzo, N., Deidda, G., Vigneti, E., Galluzzi, G., Colantoni, L., Merico, B., Ricci, E.,
850 and Felicetti, L. (1998). Sequence homology between 4qter and 10qter loci facilitates the

- 851 instability of subtelomeric KpnI repeat units implicated in facioscapulohumeral muscular
852 dystrophy. *American Journal of Human Genetics* 63, 181–190. 10.1086/301906.
- 853 64. Lemmers, R.J.L.F., van der Vliet, P.J., van der Gaag, K.J., Zuniga, S., Frants, R.R., de Knijff, P.,
854 and van der Maarel, S.M. (2010). Worldwide Population Analysis of the 4q and 10q
855 Subtelomeres Identifies Only Four Discrete Interchromosomal Sequence Transfers in
856 Human Evolution. *American Journal of Human Genetics* 86, 364–377.
857 10.1016/j.ajhg.2010.01.035.
- 858 65. de Greef, J.C., Lemmers, R.J.L.F., van Engelen, B.G.M., Sacconi, S., Venance, S.L., Frants, R.R.,
859 Tawil, R., and van der Maarel, S.M. (2009). Common epigenetic changes of D4Z4 in
860 contraction-dependent and contraction-independent FSHD. *Human Mutation* 30, 1449–
861 1459. 10.1002/humu.21091.
- 862 66. Hewitt, J.E., Lyle, R., Clark, L.N., Valleley, E.M., Wright, T.J., Wijmenga, C., Van Deutekom, J.C.
863 t., Francis, F., Sharpe, P.T., Hofker, M., et al. (1994). Analysis of the tandem repeat locus
864 D4Z4 associated with facioscapulohumeral muscular dystrophthhy. *Human Molecular*
865 *Genetics* 3, 1287–1295. 10.1093/hmg/3.8.1287.
- 866 67. Pini, S., Napoli, F.M., Tagliafico, E., La Marca, A., Bertucci, E., Salsi, V., and Tupler, R. (2023).
867 De novo variants and recombination at 4q35: Hints for preimplantation genetic testing in
868 facioscapulohumeral muscular dystrophy. *Clinical Genetics* 103, 242–246.
869 10.1111/cge.14250.
- 870 68. Gaillard, M.C., Roche, S., Dion, C., Tasmadjian, A., Bouget, G., Salort-Campana, E., Vovan, C.,
871 Chaix, C., Broucqsault, N., Morere, J., et al. (2014). Differential DNA methylation of the
872 D4Z4 repeat in patients with FSHD and asymptomatic carriers. *Neurology* 83, 733–742.
873 10.1212/WNL.0000000000000708.
- 874 69. Lemmers, R.J.L.F., Goeman, J.J., Van der Vliet, P.J., Van Nieuwenhuizen, M.P., Balog, J., Vos-
875 Versteeg, M., Camano, P., Ramos Arroyo, M.A., Jerico, I., Rogers, M.T., et al. (2015). Inter-
876 individual differences in CpG methylation at D4Z4 correlate with clinical variability in
877 FSHD1 and FSHD2. *Human Molecular Genetics* 24, 659–669. 10.1093/hmg/ddu486.
- 878 70. Jones, T.I., King, O.D., Himeda, C.L., Homma, S., Chen, J.C.J., Beermann, M. Lou, Yan, C.,
879 Emerson, C.P., Miller, J.B., Wagner, K.R., et al. (2015). Individual epigenetic status of the
880 pathogenic D4Z4 macrosatellite correlates with disease in facioscapulohumeral muscular
881 dystrophy. *Clinical Epigenetics* 7, 1–22. 10.1186/s13148-015-0072-6.
- 882 71. Calandra, P., Cascino, I., Lemmers, R.J.L.F., Galluzzi, G., Teveroni, E., Monforte, M., Tasca, G.,
883 Ricci, E., Moretti, F., Van Der Maarel, S.M., et al. (2016). Allele-specific DNA
884 hypomethylation characterises FSHD1 and FSHD2. *Journal of Medical Genetics* 53, 348–
885 355. 10.1136/jmedgenet-2015-103436.
- 886 72. Dion, C., Roche, S., Laberthonnière, C., Broucqsault, N., Mariot, V., Xue, S., Gurzau, A.D.,
887 Nowak, A., Gordon, C.T., Gaillard, M.C., et al. (2019). SMCHD1 is involved in de novo
888 methylation of the DUX4-encoding D4Z4 macrosatellite. *Nucleic Acids Research* 47, 2822–
889 2839. 10.1093/nar/gkz005.
- 890 73. Hiramuki, Y., Kure, Y., Saito, Y., Ogawa, M., Ishikawa, K., Yoshimura, M.M., Oya, Y.,
891 Takahashi, Y., Kim, D.S., Arai, N., et al. (2022). Simultaneous measurement of the size and
892 methylation of chromosome 4qA - D4Z4 repeats in facioscapulohumeral muscular

- 893 dystrophy by long - read sequencing. *Journal of Translational Medicine*, 1–12.
894 [10.1186/s12967-022-03743-7](https://doi.org/10.1186/s12967-022-03743-7).
- 895 74. Erdmann, H., Scharf, F., Gehling, S., Benet-Pagès, A., Jakubiczka, S., Becker, K., Seipelt, M.,
896 Kleefeld, F., Knop, K.C., Prott, E.-C., et al. (2023). Methylation of the 4q35 D4Z4 repeat
897 defines disease status in facioscapulohumeral muscular dystrophy. *Brain* *146*, 1388–1402.
898 [10.1093/brain/awac336](https://doi.org/10.1093/brain/awac336).
- 899 75. Butterfield, R.J., Dunn, D.M., Duvall, B., Moldt, S., and Weiss, R.B. (2023). Deciphering D4Z4
900 CpG methylation gradients in facioscapulohumeral muscular dystrophy using nanopore
901 sequencing. *bioRxiv : the preprint server for biology*. [10.1101/2023.02.17.528868](https://doi.org/10.1101/2023.02.17.528868).
- 902 76. Lemmers, R.J.L.F., Wohlgemuth, M., Van Der Gaag, K.J., Van Der Vliet, P.J., Van Teijlingen,
903 C.M.M., De Knijff, P., Padberg, G.W., Frants, R.R., and Van Der Maarel, S.M. (2007). Specific
904 sequence variations within the 4q35 region are associated with facioscapulohumeral
905 muscular dystrophy. *American Journal of Human Genetics* *81*, 884–894. [10.1086/521986](https://doi.org/10.1086/521986).
- 906 77. Tawil, R., Kissel, J.T., Heatwole, C., Pandya, S., Gronseth, G., and Benatar, M. (2015).
907 Evidence-based guideline summary: Evaluation, diagnosis, and management of
908 facioscapulohumeral muscular dystrophy. *Neurology* *85*, 357–364.
909 [10.1212/WNL.0000000000001783](https://doi.org/10.1212/WNL.0000000000001783).

910

911 **Figure titles and legends**

912

913 **Figure 1. Scheme of the proposed models for FSHD1 and FSHD2 pathogenesis**

914 Tandemly arrayed D4Z4 repeats (triangles) and the 4q35 haplotype elements are
915 represented along with the molecular signatures proposed for FSHD1 and FSHD2
916 pathogenesis.

917

918 **Figure 2. Stratification and variability of D4Z4-I elements**

919 (A) D4Z4-I elements in hg38 and T2T. Cumulative size per chromosome of D4Z4-I elements.
920 Only chromosomes for which at least 1 Kb of D4Z4-I sequence was identified, in any of the
921 2 assemblies, are reported. Size (in Kb) is shown on the Y axis. Different colors are used
922 for T2T and hg38 and to mark complete D4Z4-I (see legend). (B) Cumulative size (in Kb) of
923 D4Z4-I elements in the hg38, T2T and PGR haplotype level assemblies. Left: cumulative
924 size of D4Z4-I elements. Right: cumulative size of full D4Z4-I elements. Data is displayed in

925 the form of a barplot. Values are on the Y axis, labels on the X axis. For the human
926 pangenome haplotype level assemblies, 3 distinct bars are used to indicate the minimum,
927 maximum and average values. Color codes according to the legend. (C) Dendrogram of
928 complete D4Z4-l elements. Red rectangles and color codes are used to delineate the 5 main
929 groups of D4Z4-l elements, as identified by sequence similarity-based clustering. (D) Total
930 estimated number of complete D4Z4 elements per sequence similarity group. Distributions
931 of values are displayed in the form of a boxplot, with groups indicated on the X axis and
932 copy numbers on the Y axis. The top-right panel provides zoom on groups chr 14, chr 15-
933 21, chr 22-1 and chr 22-2.

934

935 **Figure 3. Characterization of D4Z4 alleles**

936 (A) Heatmap of *ter-DUX4*. Scaled identity levels are displayed by the green color gradient
937 on the right. The dendrogram on the left shows clustering of *ter-DUX4* based on sequence
938 identity. Colored vertical bars are used to indicate: presence/absence of pLam; classification
939 as qA or qB; sequence similarity cluster (see top). Color codes are explained directly under
940 each bar. (B) Numerosity of sequence similarity clusters. The barplot shows (y-axis) the
941 total number of *ter-DUX4* in each cluster (x-axis). (C) Boxplot of substitution rates.
942 Substitution rates (substitutions by 100 bp) by cluster. Distribution of values are shown as
943 a boxplot. Clusters are on the y-axis, values on the x-axis. (D) Distance of qA and qB
944 elements from *ter-DUX4*. Data are represented in the form of a histogram. Distances (in Kb)
945 are indicated on the x-axis. qA and qB elements are marked with different colors (see
946 legend). (E-G) Annotation of terminal D4Z4 repeats. Two barplots are displayed for the 3
947 (out of 4) clusters to which 5 or more sequences have been assigned. Each barplot reports,
948 for qA and qB alleles respectively: the total count (tot); the number pLam-like elements
949 (pLam); the number pLam-like elements with a valid PAS (plamPolyA); the total number of
950 haplotypes with 8 or less D4Z4 repeats. Values are indicated on the y-axis, labels on the x-

951 axis. Distinct sequence similarity clusters are indicated by a corresponding label directly
952 under each plot. (G-H) Total number of complete D4Z4-I elements in H) qA and qB
953 haplotypes; I) haplotypes with (pLam+) and without (pLam-). Distribution of values are
954 displayed in the form of a boxplot, with values on the X-axis and groups on the Y axis. Color
955 codes are consistent throughout Figure 2.

956

957 **Figure 4. CpG Methylation profiles of cohort subjects inferred on T2T**

958 A total of 82 CpGs for which complete data are available for all the study subjects are shown.
959 (A) Heatmap of % CpG methylation: each row represents the CpG methylation pattern of a
960 single individual. Columns represent the methylation pattern of each analyzed CpG. The
961 dendrogram shows the clustering of the subjects based on the observed methylation
962 profiles. Colored vertical bars are used to display CCEF clinical category status, and the
963 presence/absence of deleterious genomic variants in *SMCHD1*. Color codes are illustrated
964 directly under each bar. The colored bar at the top demarcates each of the 4 distinct 95%
965 sequence identity genomic region clusters with complete data for all the subjects. (B)
966 Boxplots of methylation levels distributions by subject. The color-scale used for the heatmap
967 and the boxplot is shown at the bottom. (C) Violin-plot of average methylation levels,
968 restricted to CpGs in regions GS1 and GS6, in subjects with/without deleterious genetic
969 variants in *SMCHD1*.

970

971 **Figure 5. Methylation profiles of cohort subjects inferred on hg38**

972 A total of 32 distinct CpGs, for which complete data are available for all the study subjects,
973 are displayed. (A) Heatmap of % CpG methylation: individuals are reported on the rows and
974 CpGs on the columns. The dendrogram shows a clustering of the subjects based on the
975 observed methylation profiles. Colored vertical bars are used to display CCEF grades, and
976 the presence/absence of deleterious genomic variants in *SMCHD1*. Color codes are

977 illustrated directly under each bar. The colored bar at the top demarcates each of the 4
978 distinct 95% sequence identity genomic regions clusters with complete data for all the
979 subjects. (B) Boxplots of methylation levels distributions by subject. The color-scale used
980 for the heatmap and the boxplot is shown at the bottom. (C) Violin-plot of average
981 methylation levels in subjects with/without deleterious genetic variants in *SMCHD1*. Only
982 CpGs in primer-targeted regions are considered.

983

984 **Figure 6. Genotype-phenotype correlation**

985 (A) Venn diagrams reporting the number of subjects presenting different combinations of
986 4q/10q-specific D4Z4 methylation level, *SMCHD1* mutational status and FSHD clinical
987 category. (B) ACMG classification, median 4q/10q-specific D4Z4 methylation of variant
988 carriers, variant type and location are reported for each *SMCHD1* variant. When the variant
989 is carried by more than one subject, median D4Z4 methylation status of all the subjects is
990 reported. (C) Distribution of 4q/10q-specific D4Z4 methylation, *SMCHD1* mutational status
991 and FSHD clinical category in the 8 trios/parent-child couples included in the study.
992 Circles/squares are colored according to the median 4q/10q-specific D4Z4 methylation
993 level; the color is fading as it represents the variability of CpG methylation pattern in each
994 individual. WT= subject carrying no *SMCHD1* variants, VAR= subject carrying a damaging
995 *SMCHD1* variants.

996

997

998

999

1000

1001

1002 **Table titles and legends**

1003

1004 **Table 1. Titration analysis for the identification of the optimal sequence similarity**
1005 **threshold**

Identity-level (%)	# regions	size (Kb)	uniquely mapped BSC reads (%)	# testable CpGs
99	145	67,41	5,43%	237
98	103	59,13	12,13%	411
97	87	54,01	26,74%	873
96	73	46,14	32,47%	997
95	60	40,03	49,17%	1176
94	56	38,71	51,01%	1083
93	53	37,55	53,43%	891
92	46	36,49	58,47%	766
91	45	33,02	72,71%	637
90	42	32,75	75,41%	451

1006 Identity-level (%): sequence identity threshold. #regions: number of distinct regions defined
1007 by the threshold. size (Kb): total size of the regions in Kb. uniquely mapped BSC reads (%):
1008 percentage of uniquely mapping BSC reads. # testable CpGs: total number of distinct CpGs
1009 for which more than 10 reads were available for more than 25 subjects. Italic characters
1010 indicate the selected sequence similarity threshold.

1011

1012

1013

1014

1015

1016

1017 **Table 2. Characterization of complete D4Z4 elements clusters**

Cluster	Chromosome	% sequences	Median number of complete repeats (per haplotype)	Average substitution rate per 100 bp	Substitution rates estimates at haplotype level
1	Chr 10	87%	19	0.13	0.015-0.21
2	Chr 4		22	0.245	0.164-0.4312
3	Chr 14	2.3%	0.5	1.73	NA
4	Chr 15, Chr 21	3.1%	1	1.81	NA
5	Chr 22-1	1.5%+6.1%	0	0.798	NA
6	Chr 22-2		2	1.24	NA

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

Table 3- Comparison between PGR analysis results and previous knowledge	
Haplotype	Previous knowledge
<p>4qA</p>	<p>- 6 to 89 RU alleles;</p> <p>- The majority of alleles classified as qA are associated with a pLam both on 4q and 10q;</p> <p>- AUUAAA PAS at 4q;</p> <p>- 4.6% individuals carry a 4qA-PAS alleles with <8RU;</p> <p>- qA sequence is not always at the same distance from the last D4Z4 repeat (from 4.7 Kb to 8.3 Kb).</p>
<p>4qB</p>	<p>- qB alleles do not carry a pLam sequence;</p> <p>- qB sequence is more proximal than qA (2.5 Kb and 7.3 Kb respectively);</p> <p>- No differences in D4Z4 repeat numbers between 4qA and 4qB alleles.</p>
<p>10qA</p> <p>10qB</p>	<p>- AUCAAA PAS at 10q;</p> <p>- One 10qB allele was identified.</p>

Table 3- Comparison between PGR analysis results and previous knowledge

1032 **Table 4: GS containing CpGs covered by at least 10 BSC reads in at least 50% of**
1033 **samples and no missing data across all samples**

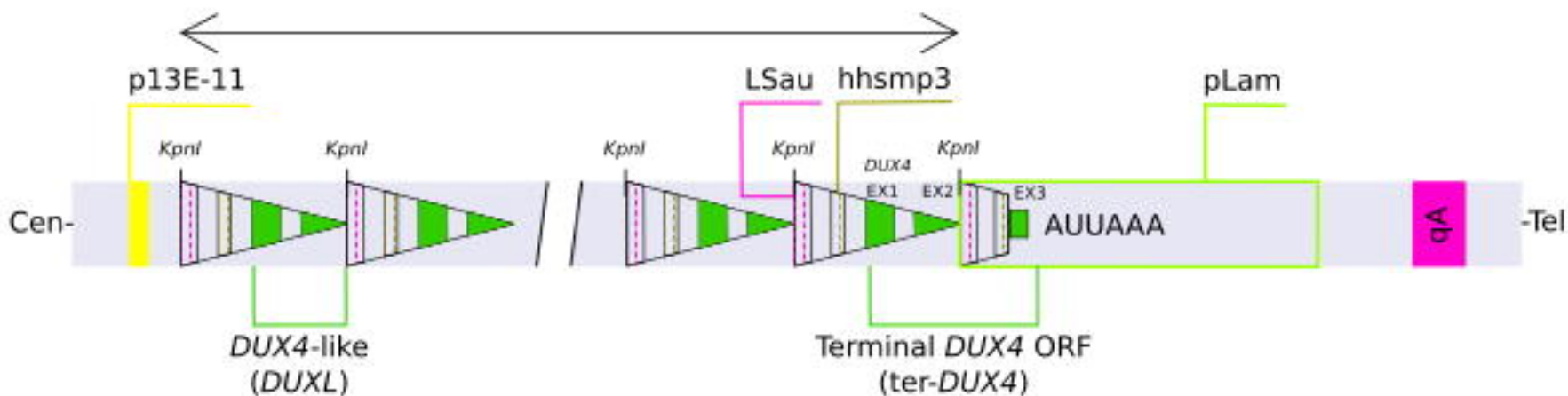
Assembly	Group	Chromosome	N° CpGs
T2T	GS1	Chr 4, Chr 10	28
	GS2		32
	GS5	Chr 13, Chr 15	16
	GS6		5
hg38	GS3	Chr 4, Chr 10	26

1034

FSHD

Prevalence 1:8.000 - 1:20.000

D4Z4 repeated units (RU)



FSHD 1

95% cases

FSHD 2

5% cases

D4Z4 reduced allele (DRA)
 ≤ 10 RU

+
4qA-PAS haplotype

↓
Terminal D4Z4
reduced methylation

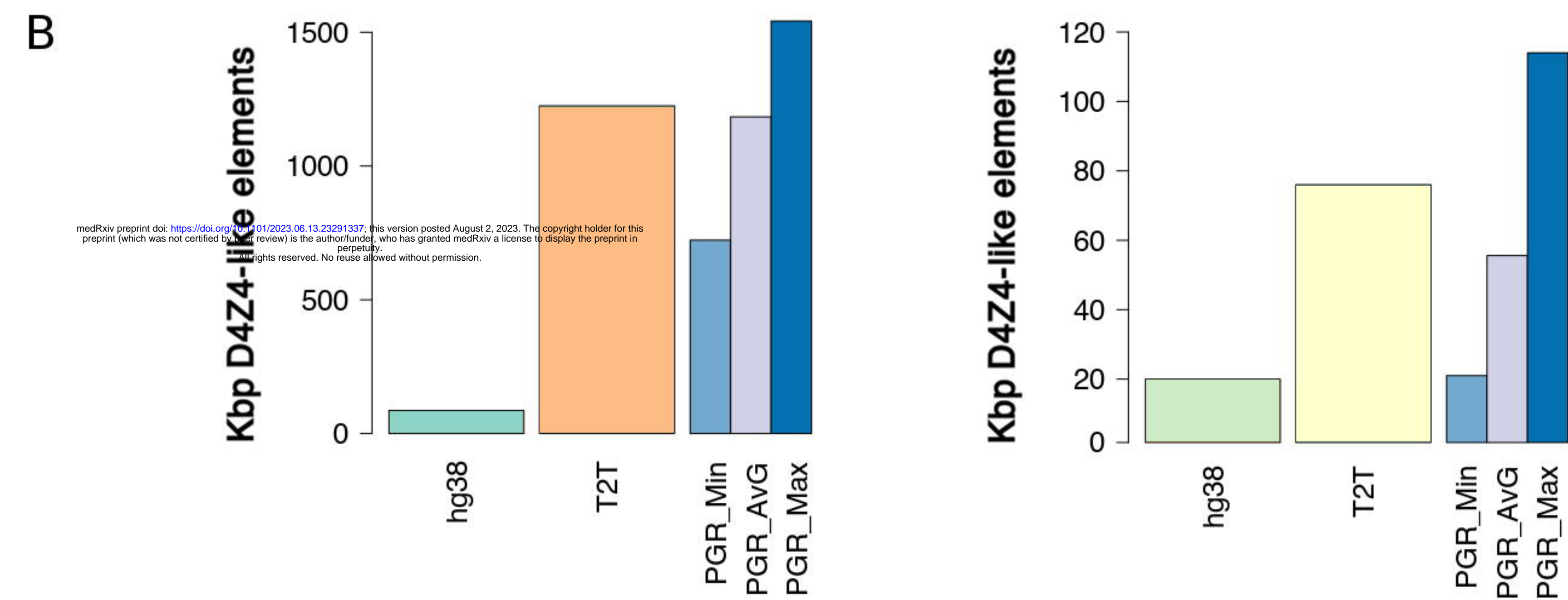
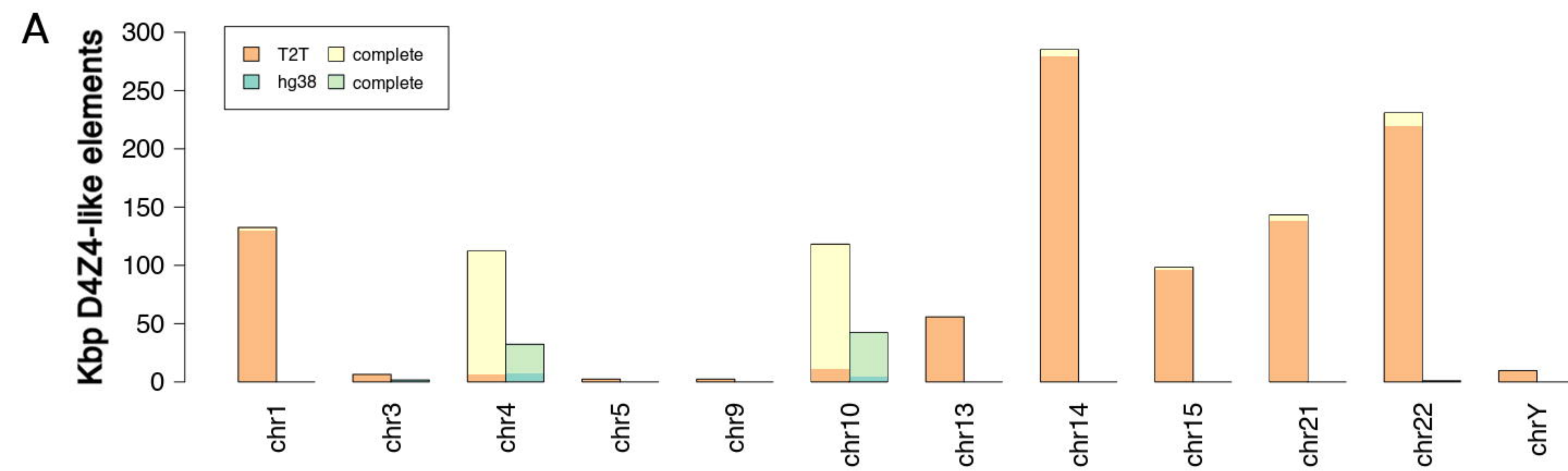
↓
Terminal DUX4 ORF (ter-DUX4) toxic expression

> 10 RU

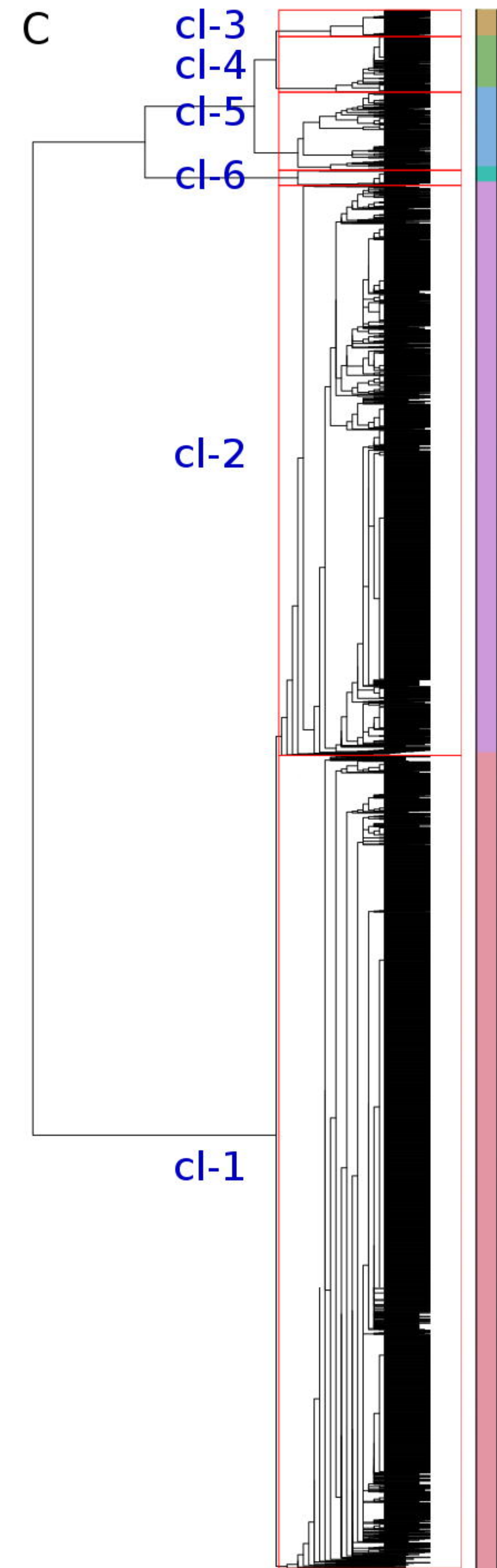
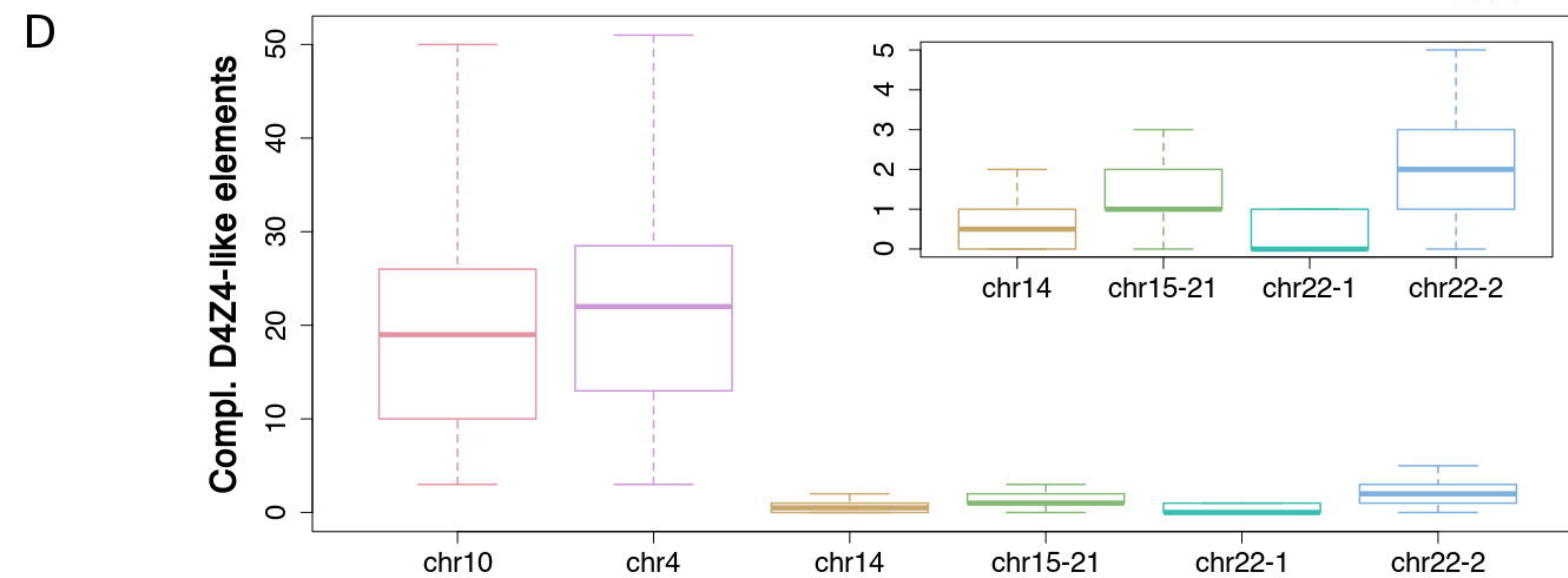
+
4qA-PAS haplotype

+
Variants in chromatin remodeling genes
(SMCHD1/DNMT3B/LRIF-1)

↓
D4Z4 repeats
reduced methylation



medRxiv preprint doi: <https://doi.org/10.1101/2023.06.13.23291337>; this version posted August 2, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. All rights reserved. No reuse allowed without permission.

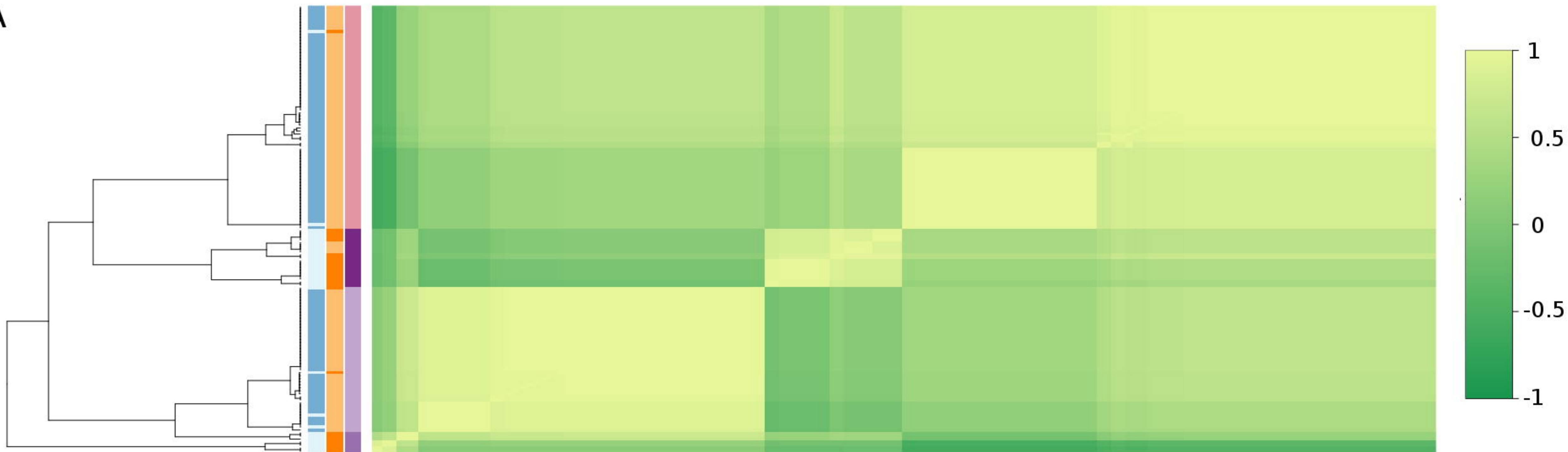


● chr10
 ● chr4_gr1
 ● chr4_gr2
 ● chr4_gr3

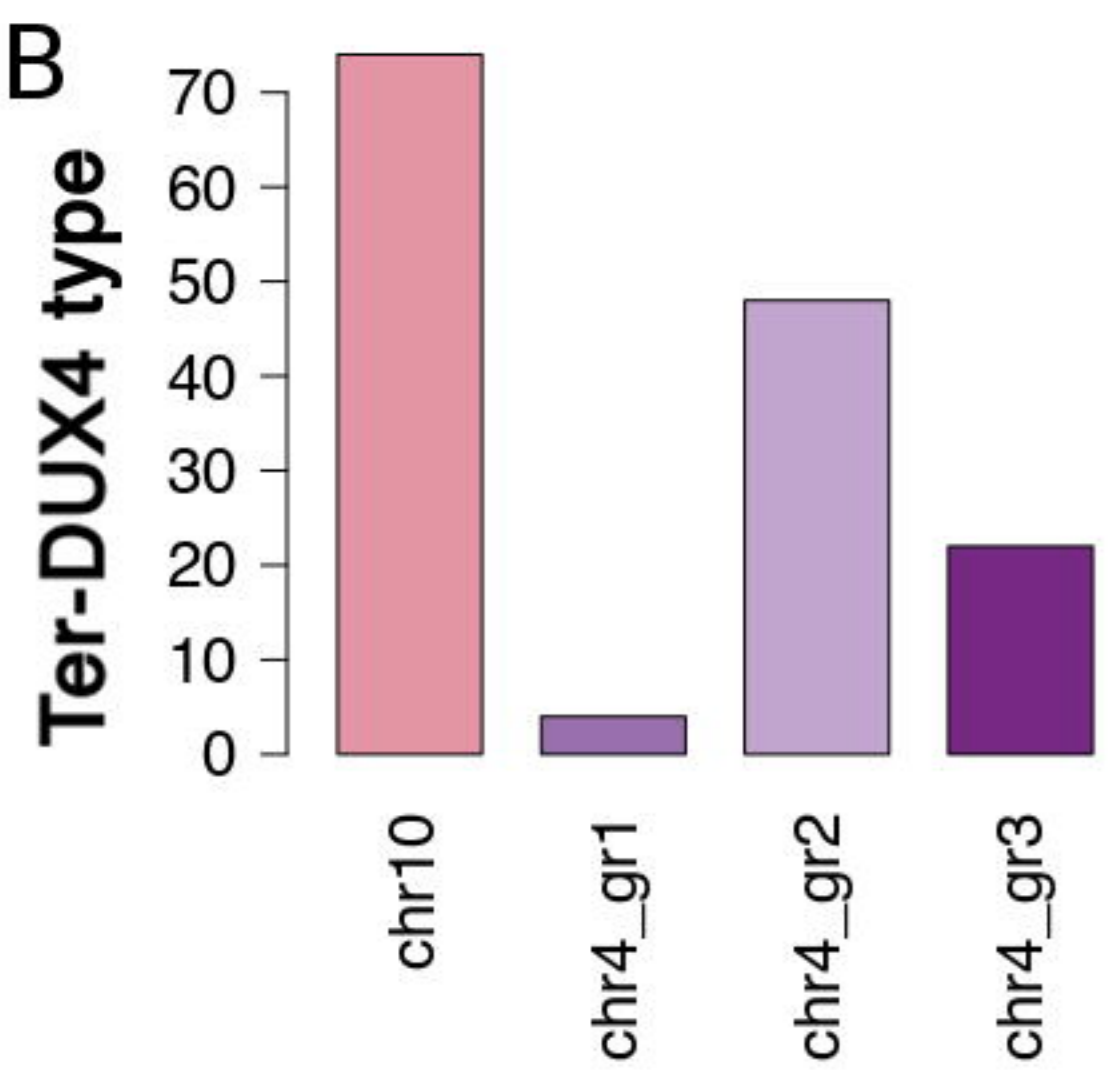
 ● plam+
 ● plam-

 ● qA
 ● qB

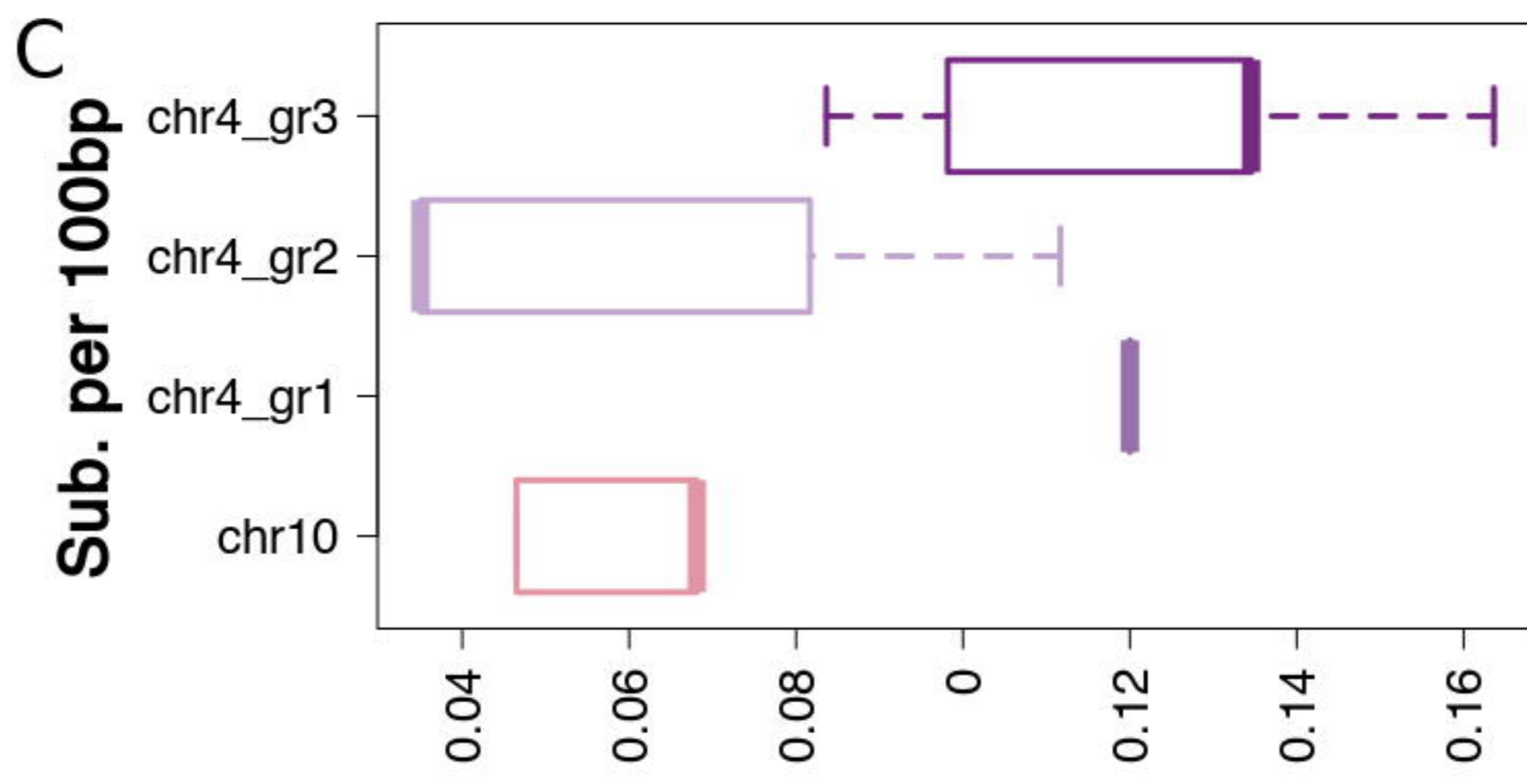
A



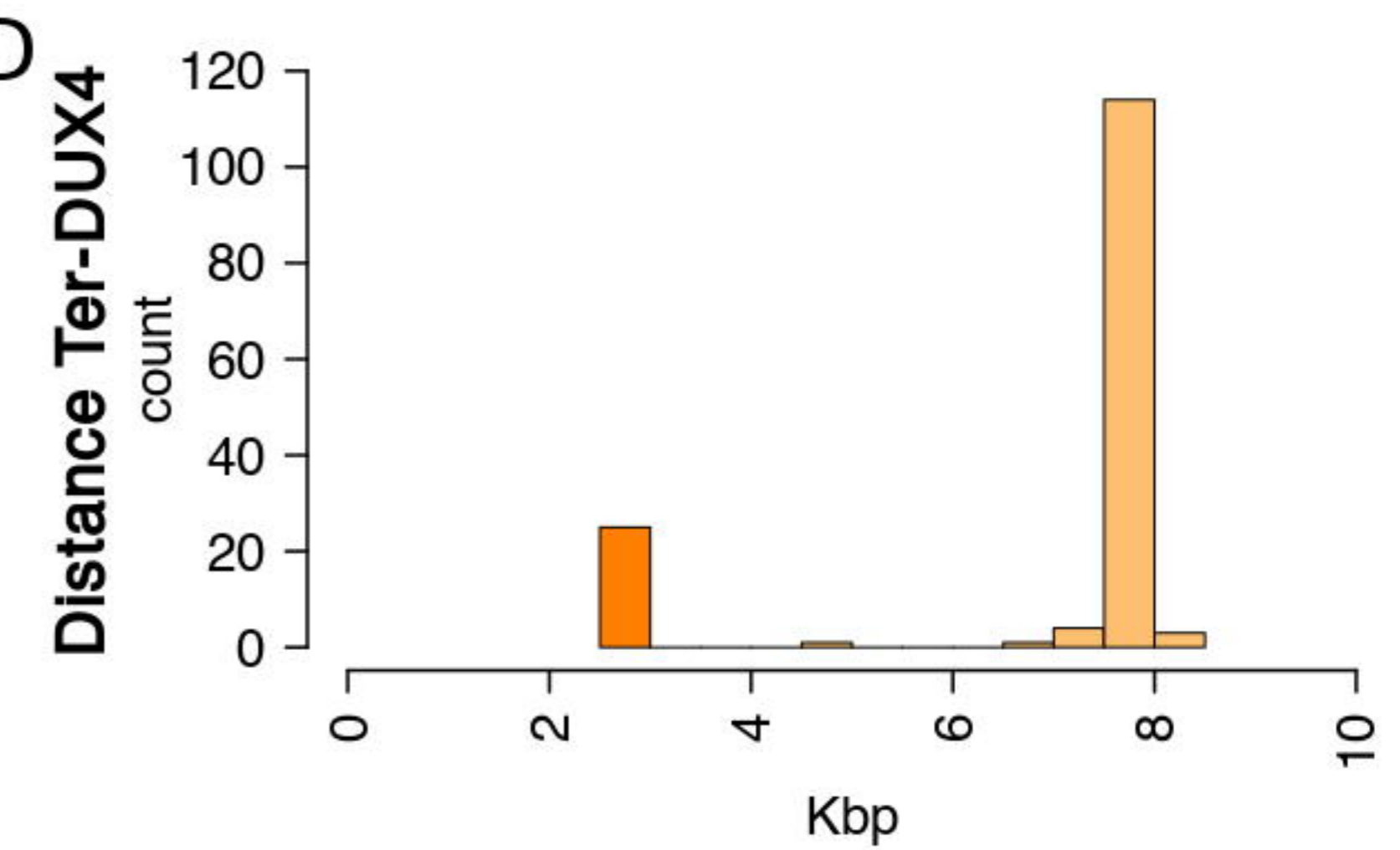
B



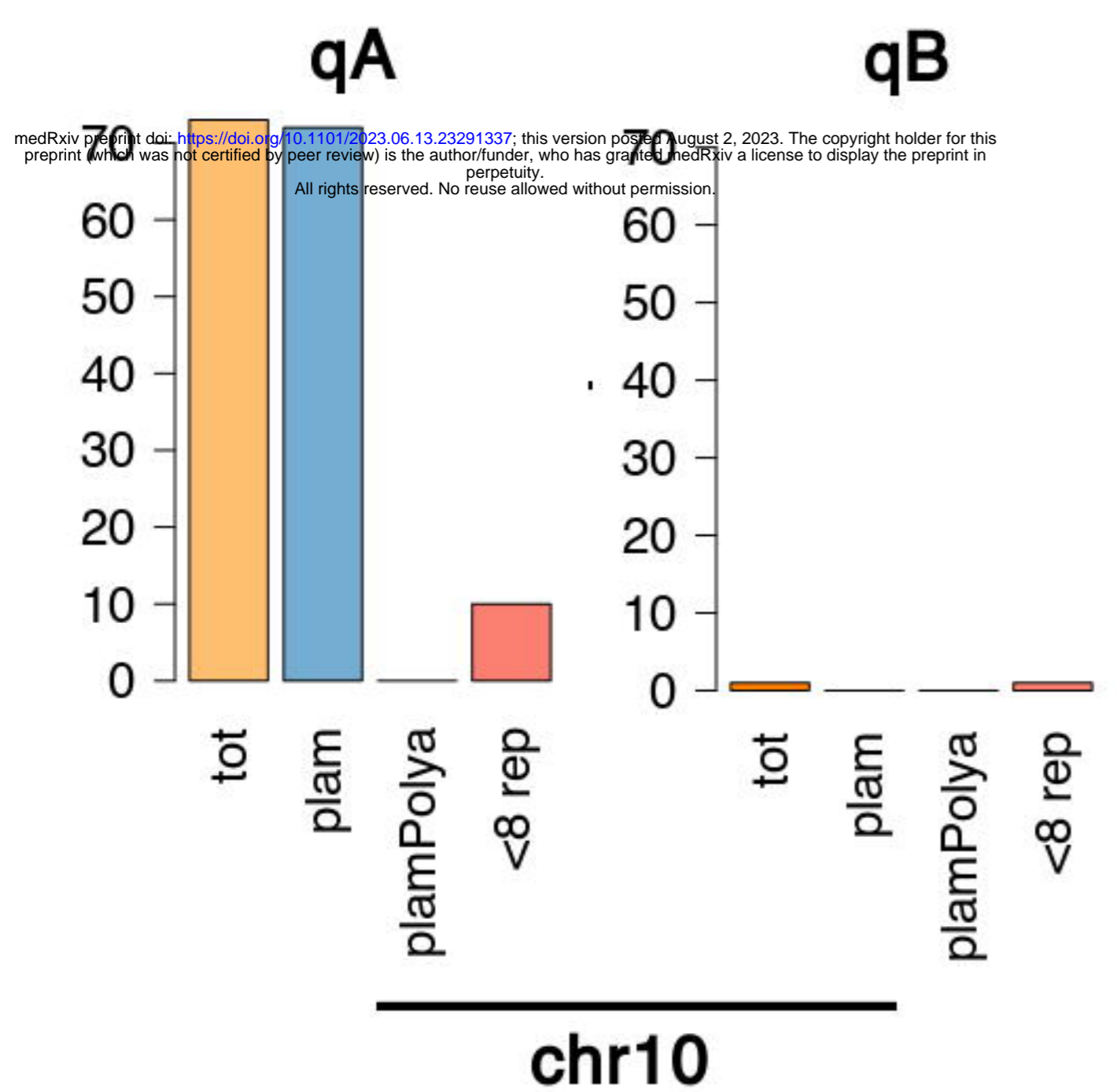
C



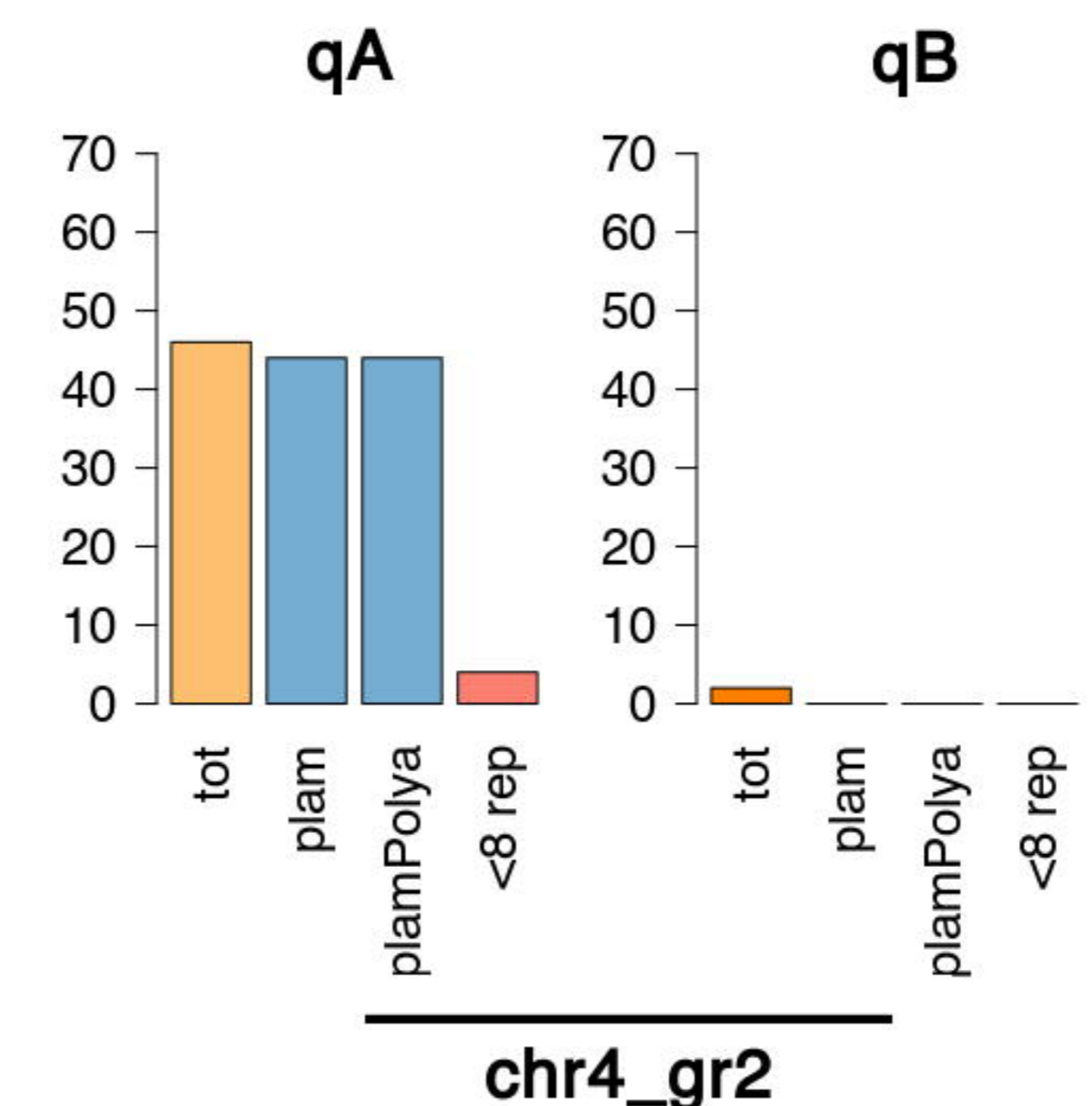
D



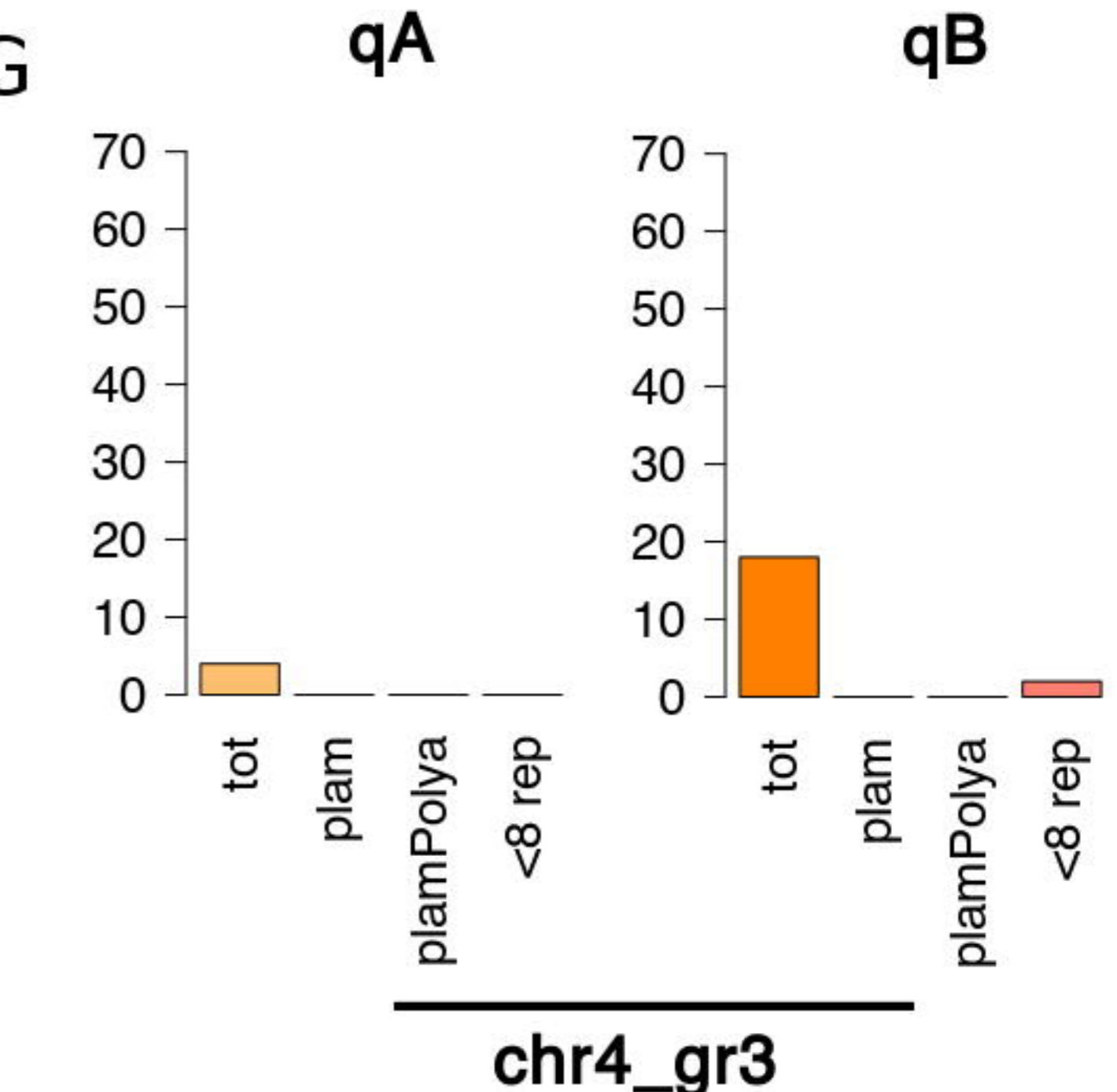
E



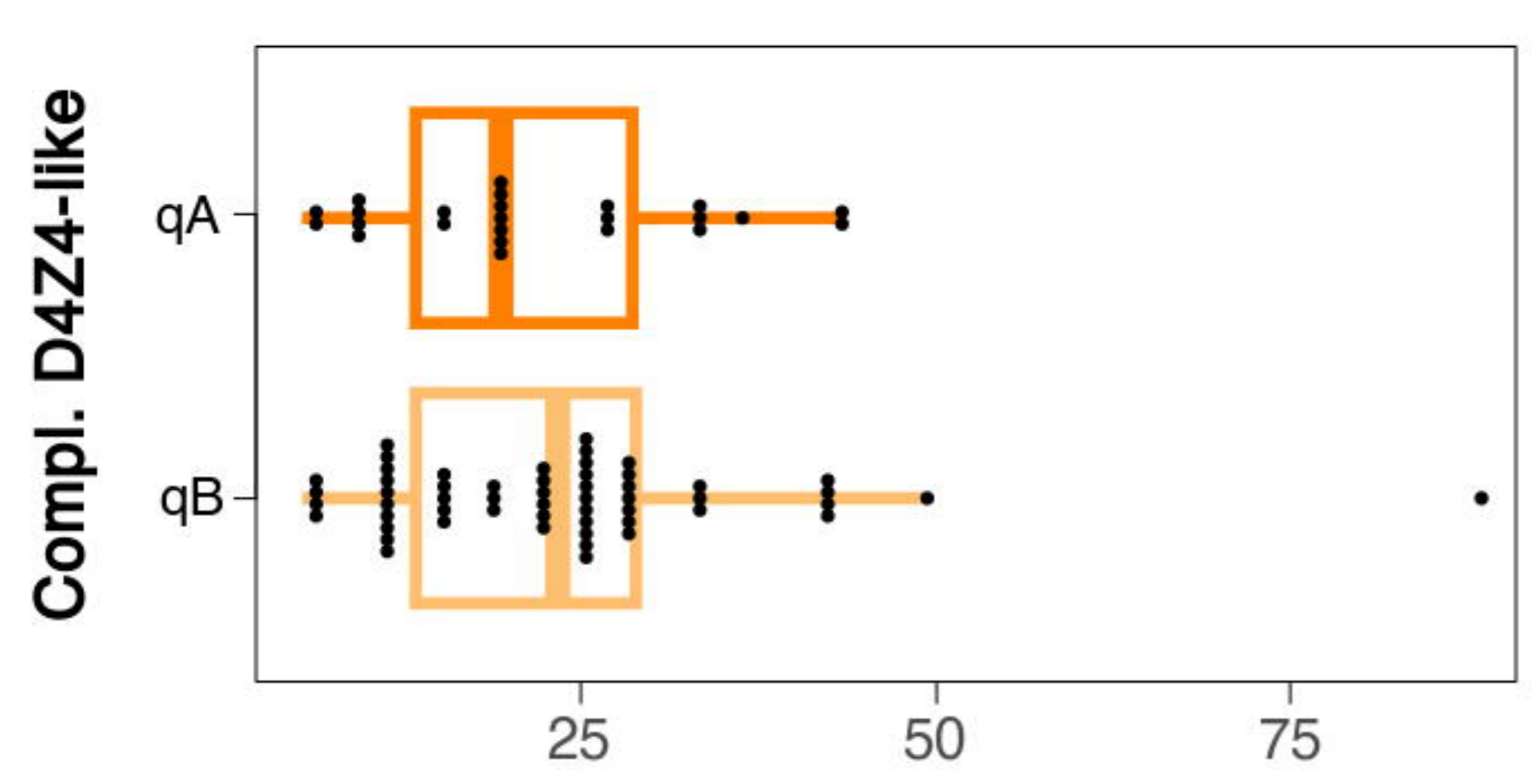
F



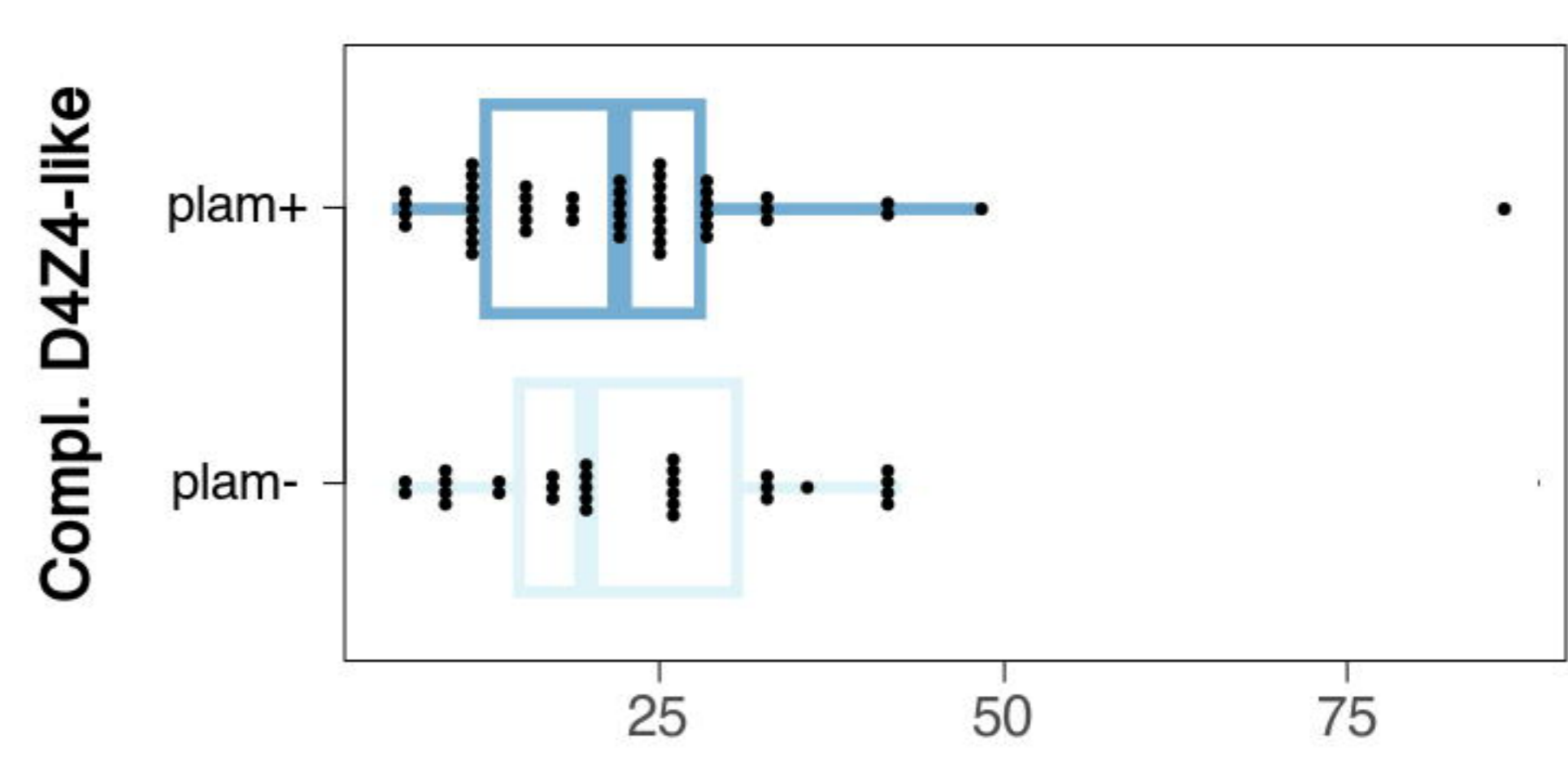
G

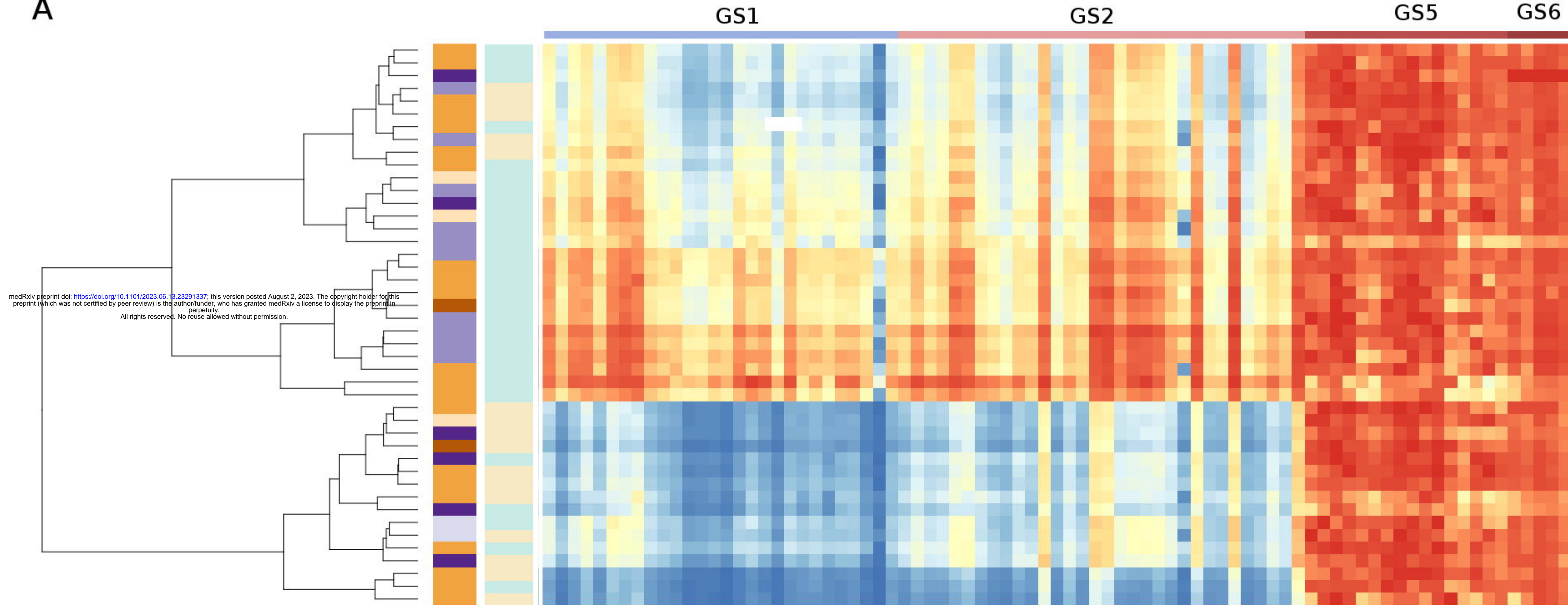
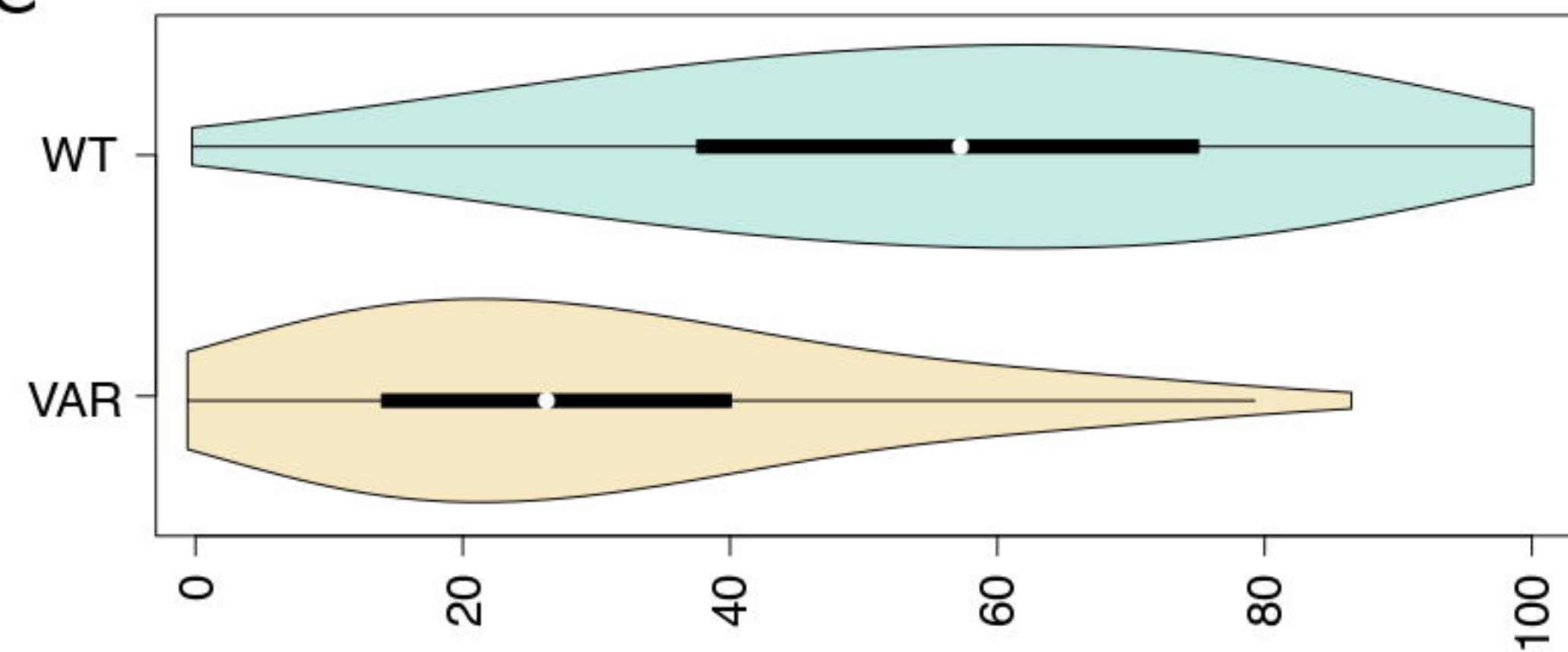
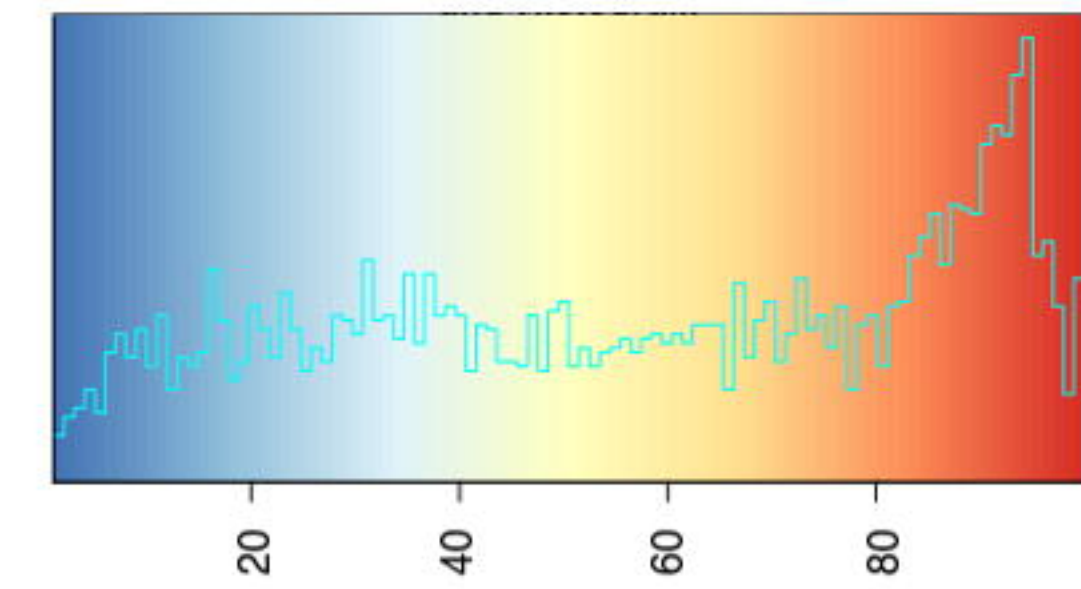
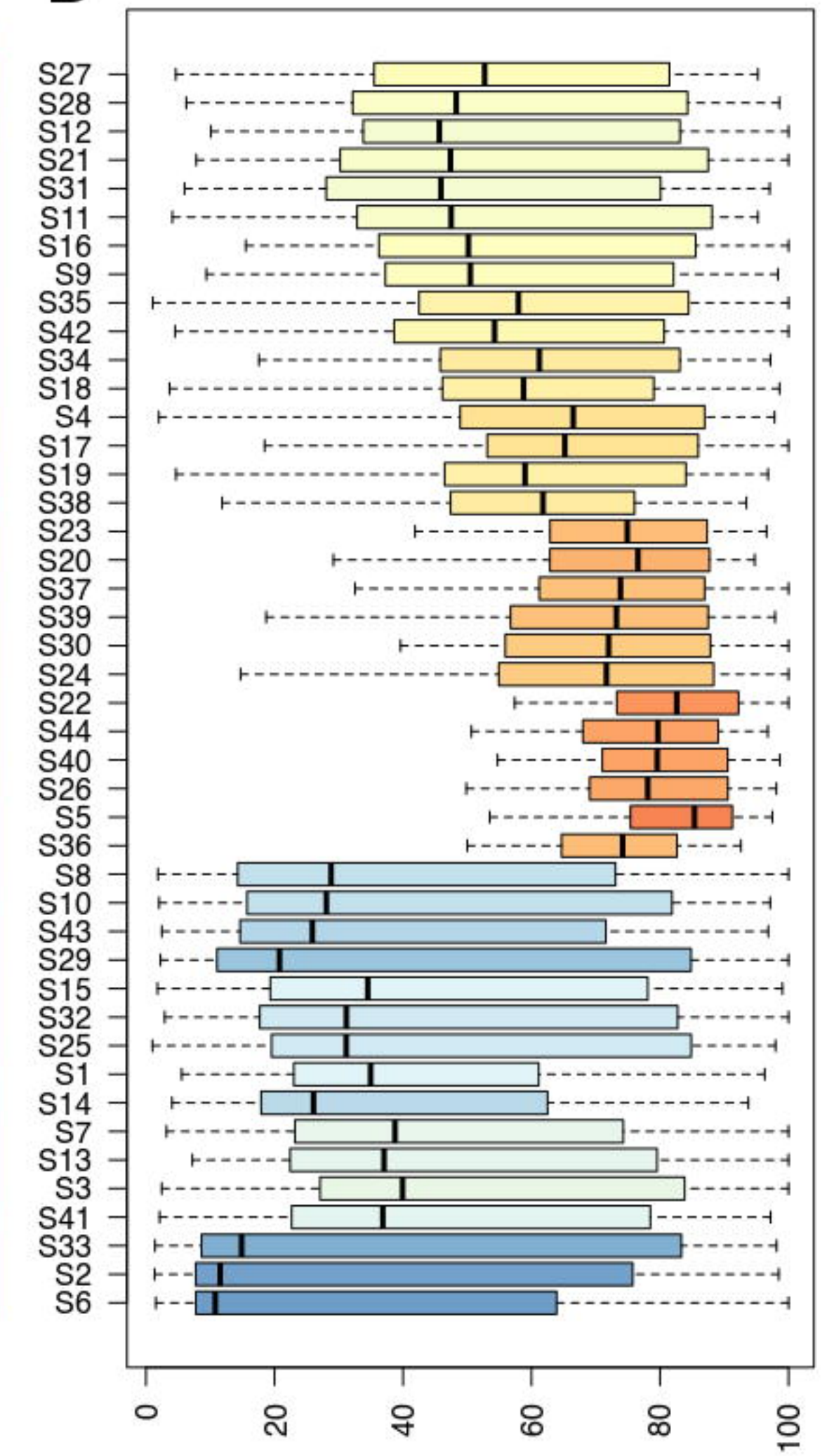


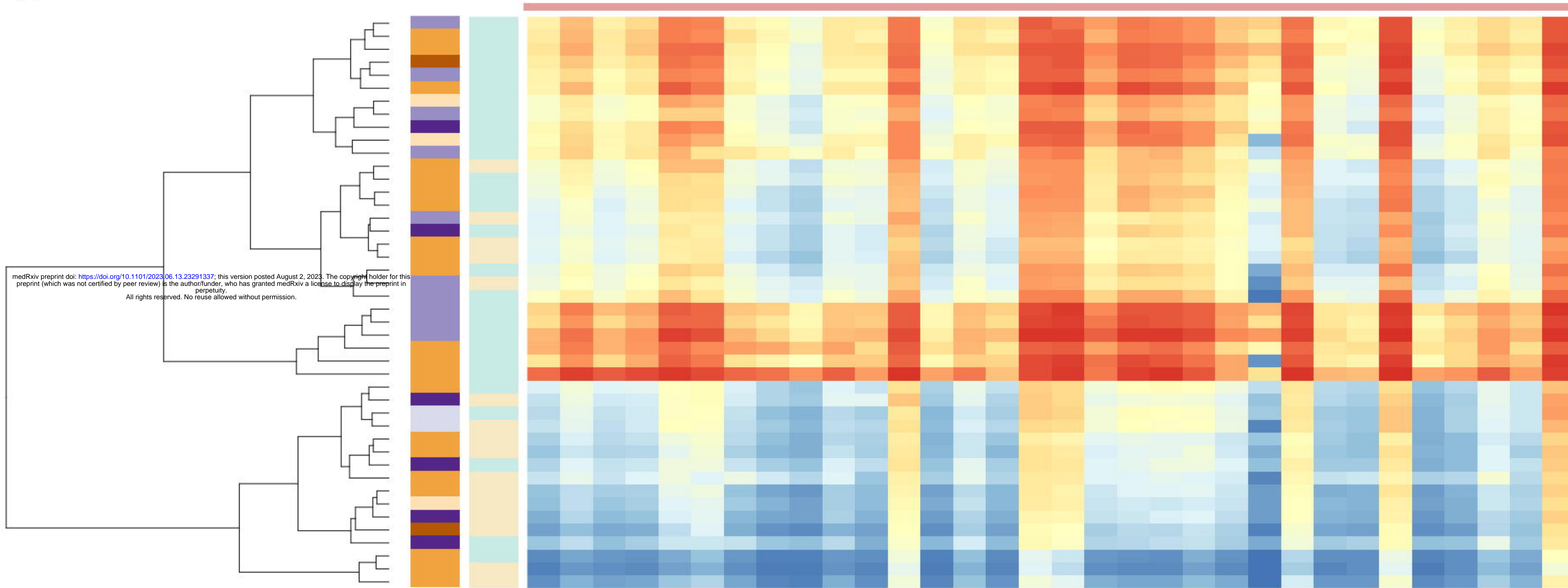
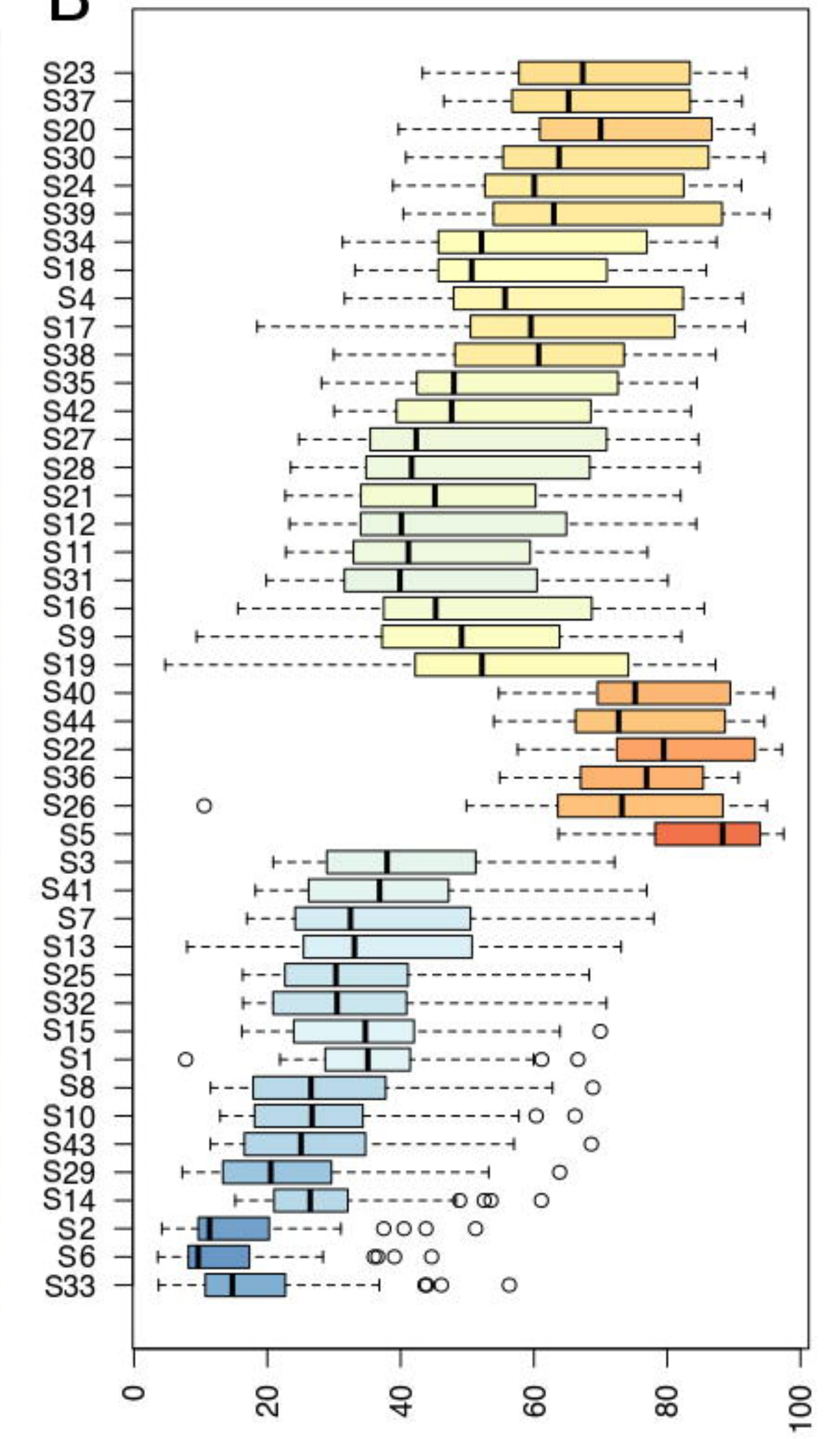
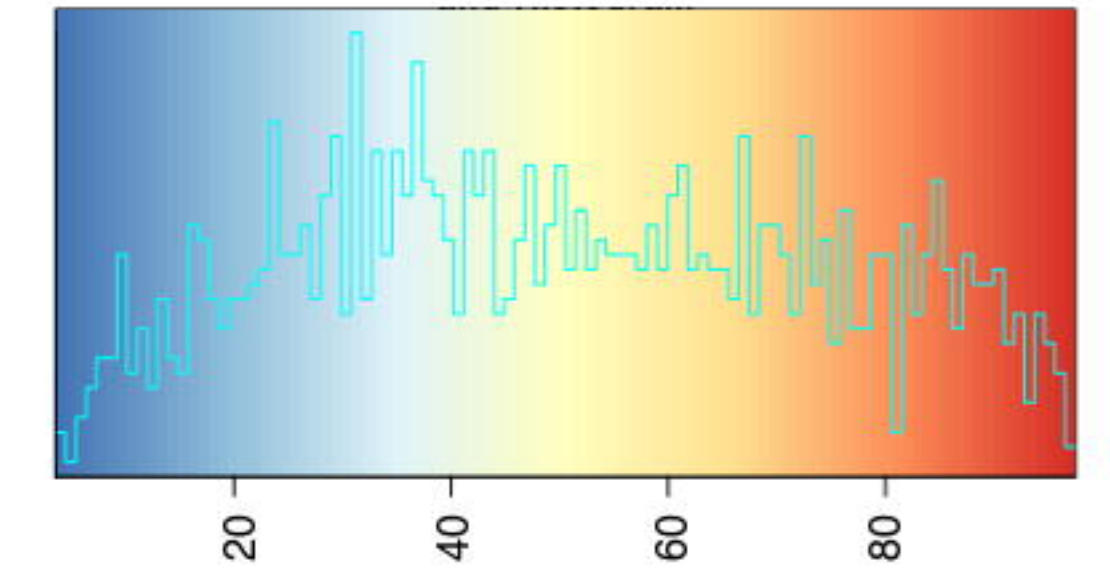
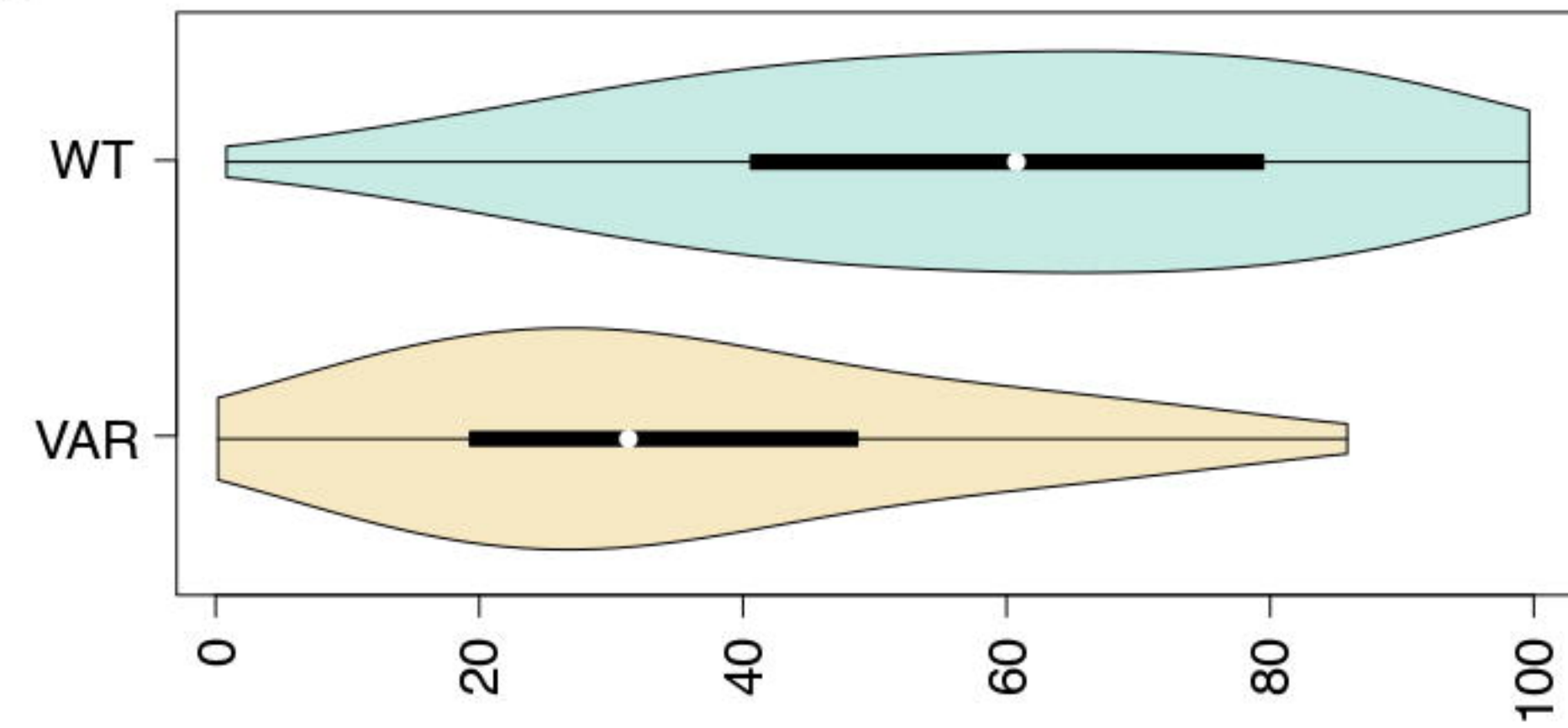
H

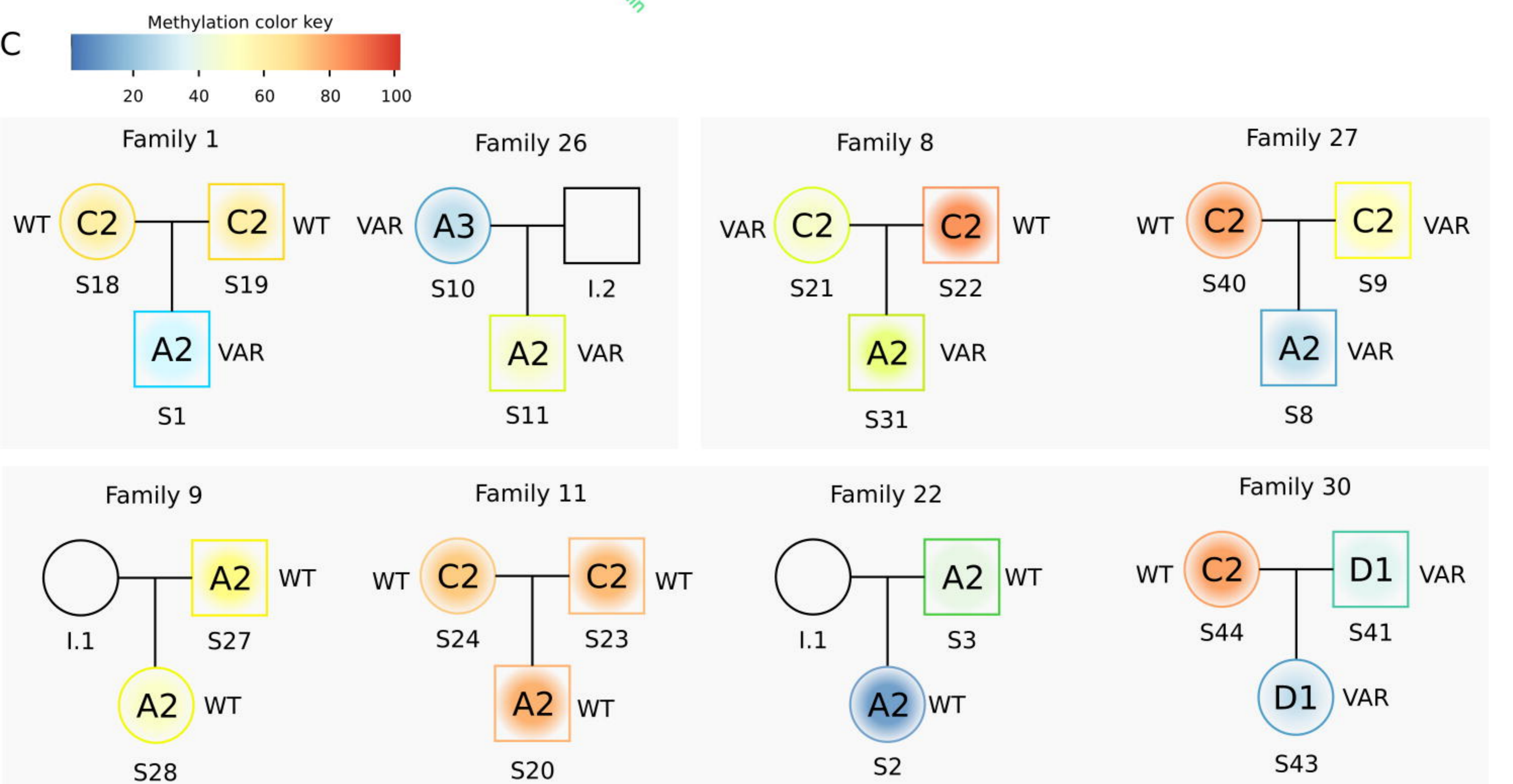
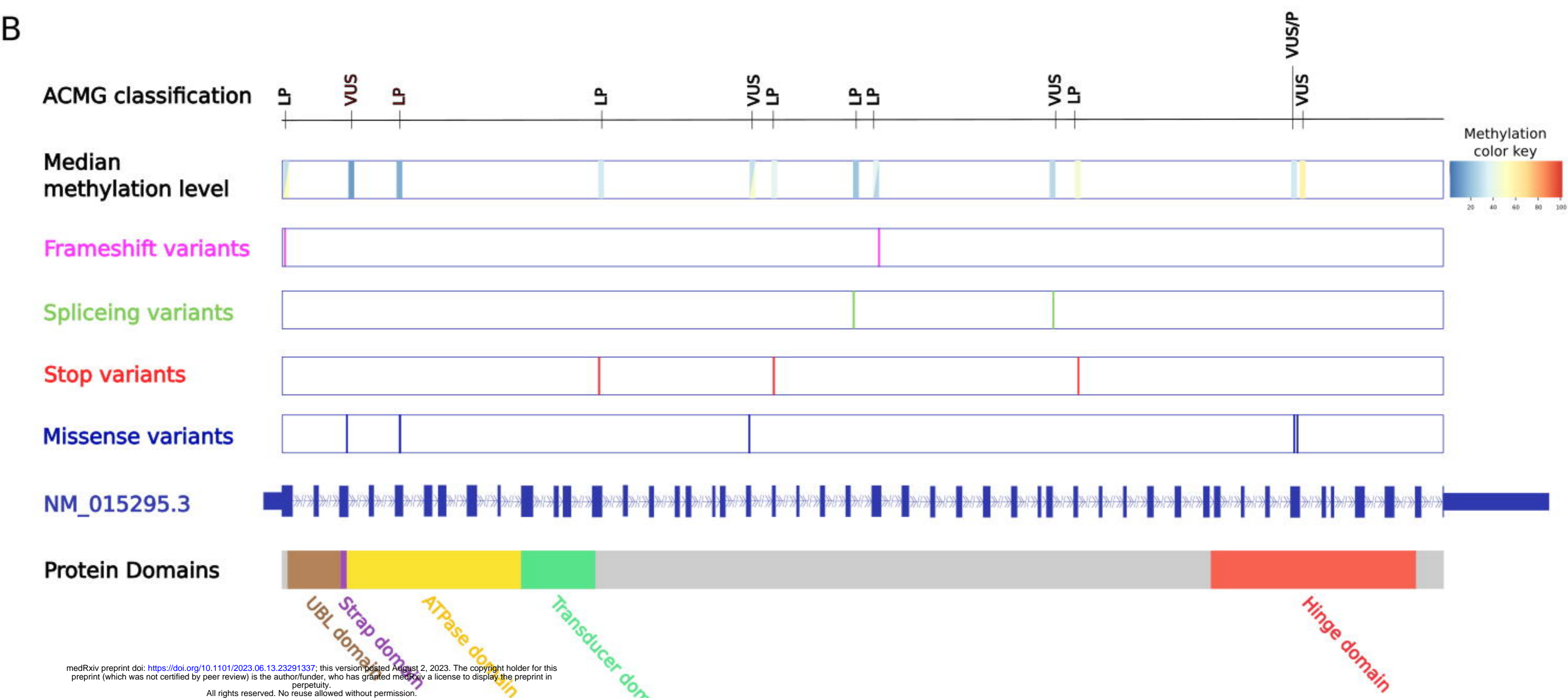
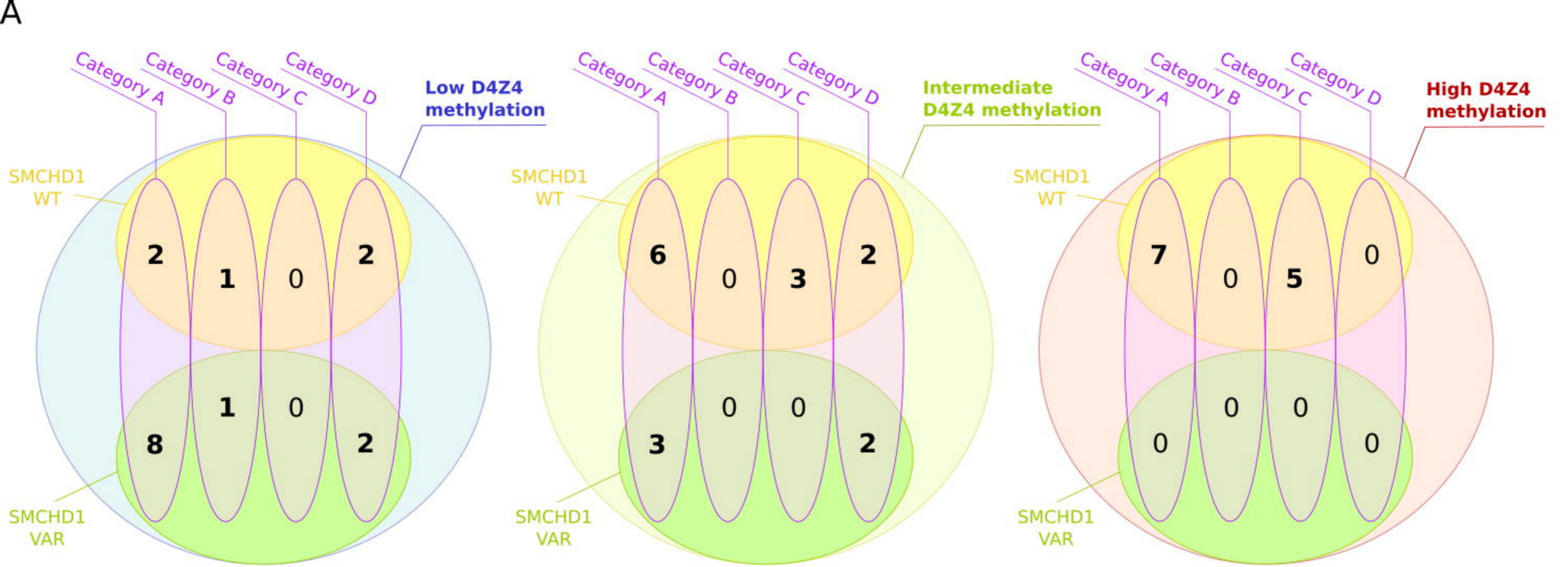


I

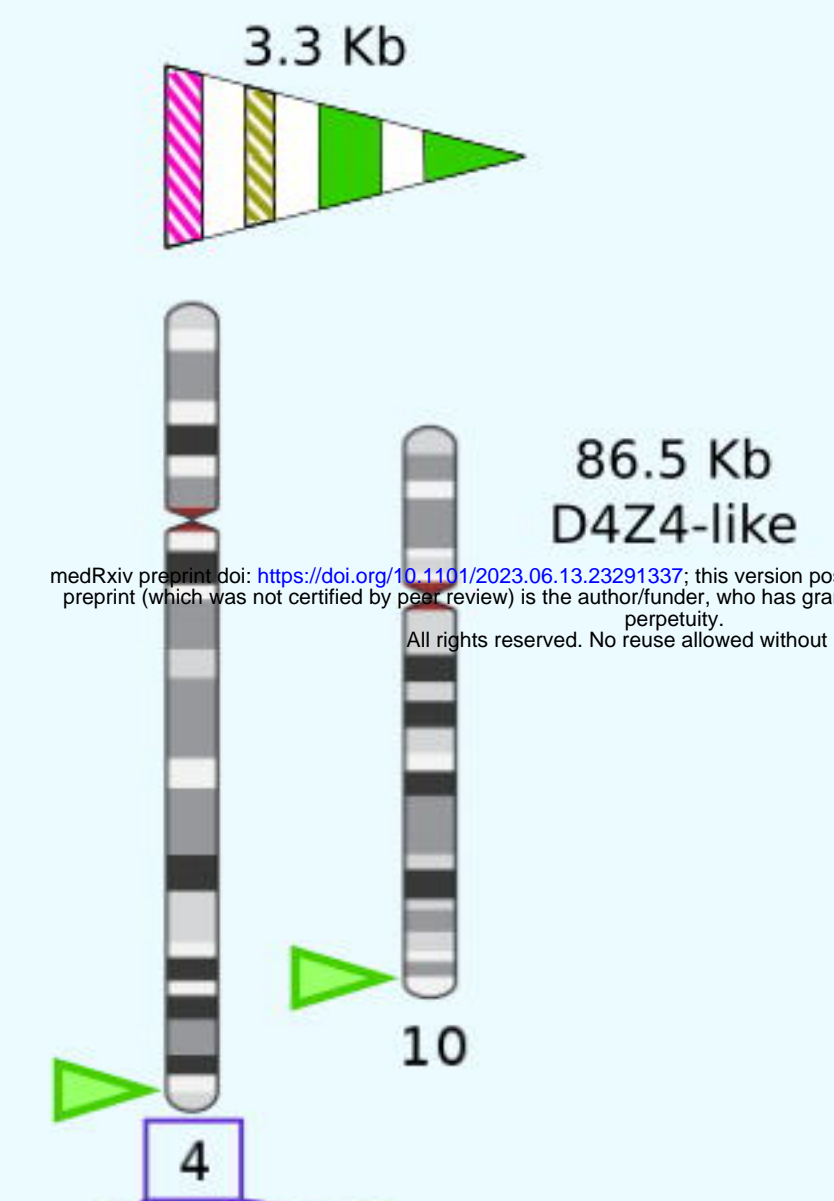
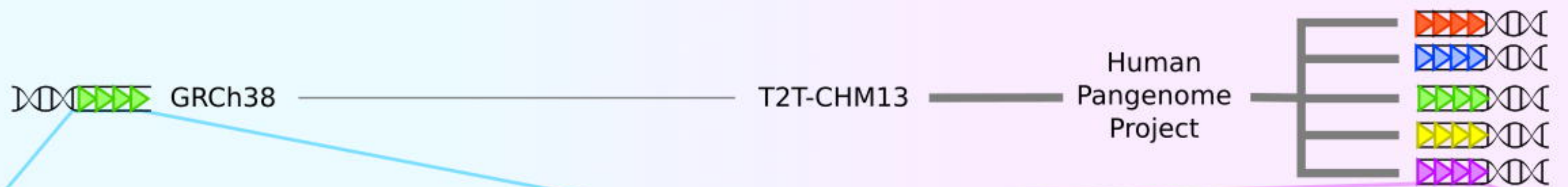


A**C****B**

A**B****C**

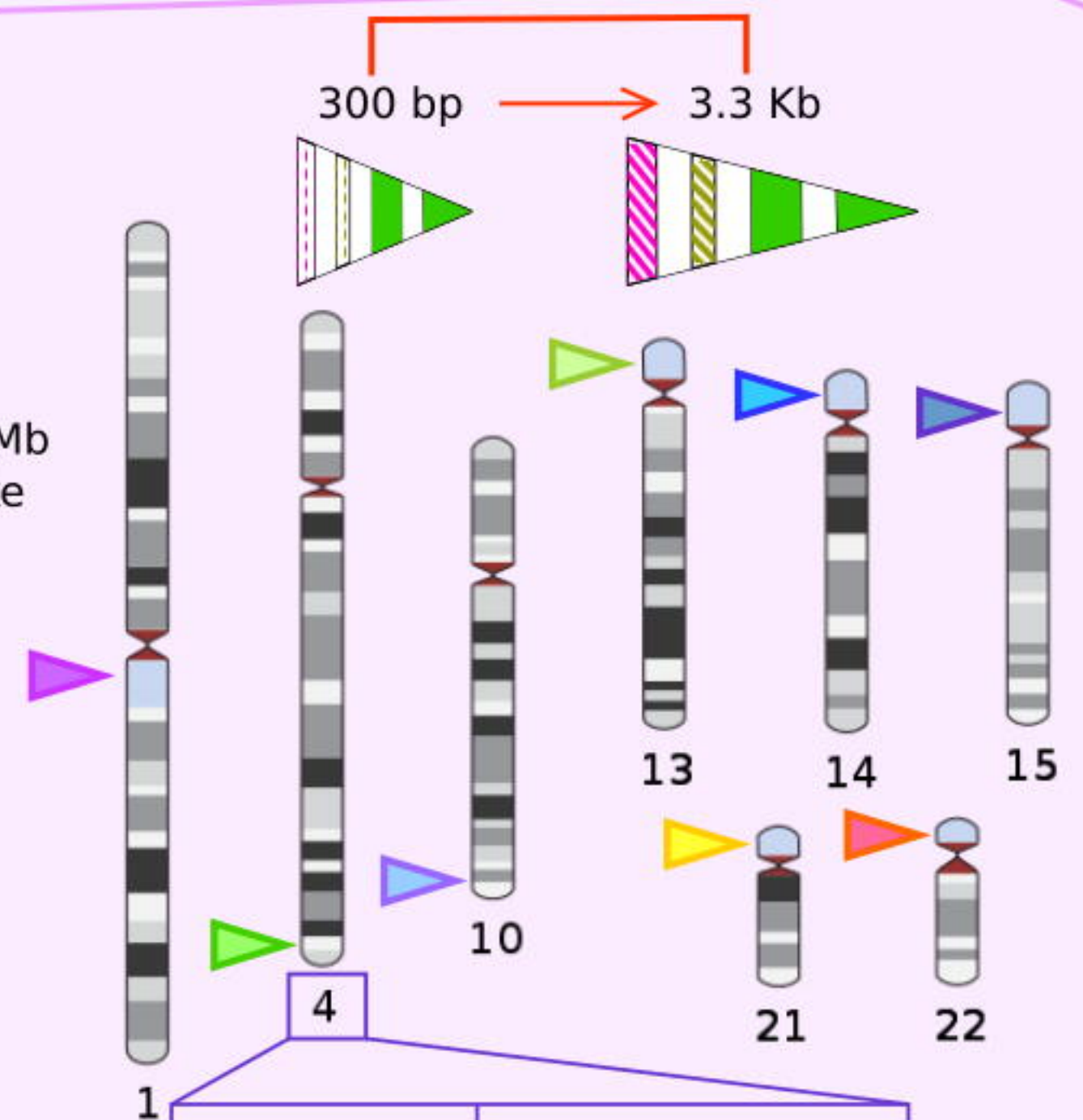


medRxiv preprint doi: <https://doi.org/10.1101/2023.06.13.23291337>; this version posted August 2, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. All rights reserved. No reuse allowed without permission.

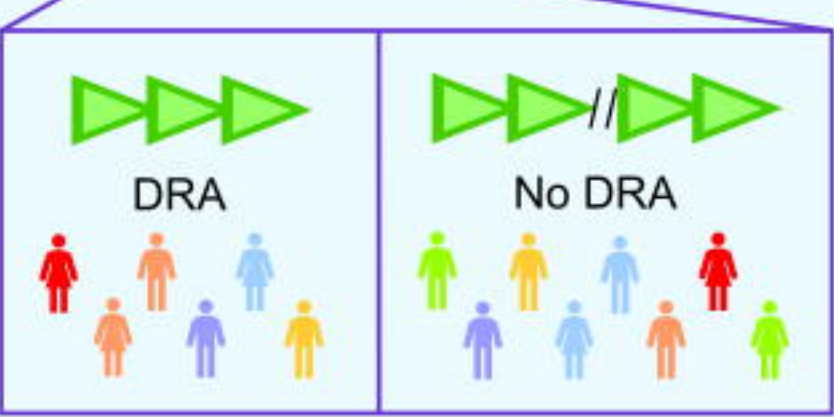


Genetic variability

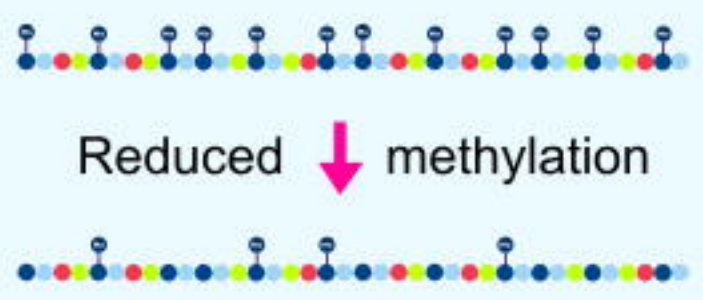
0.7 - 1.5 Mb D4Z4-like



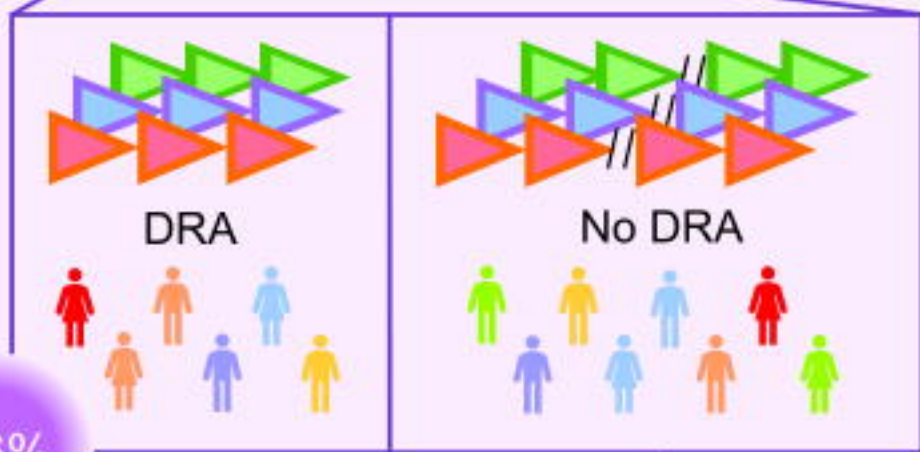
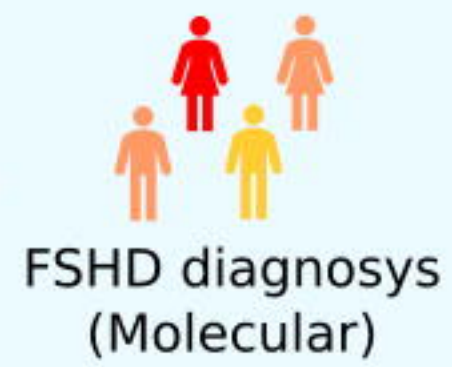
Phenotypic variability



+4qA +SMCHD1 variants



Transcriptional derepression

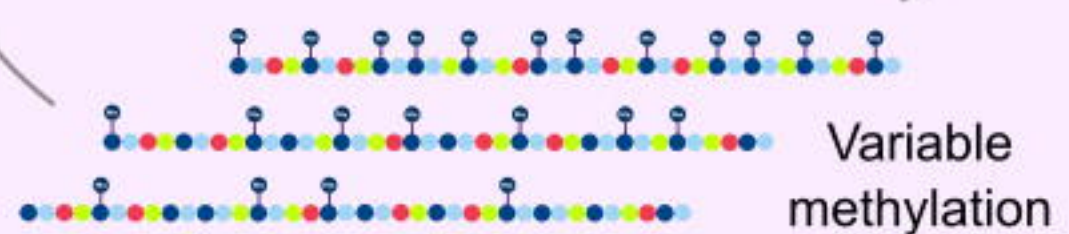


FSHD diagnosis (Clinical)

Other factors

+4qA +/-SMCHD1 variants

Transcriptional derepression



medRxiv preprint doi: <https://doi.org/10.1101/2023.06.13.23291337>; this version posted August 2, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. All rights reserved. No reuse allowed without permission.