

GigaScience

An improved pig reference genome sequence to enable pig genetics and genomics research --Manuscript Draft--

Manuscript Number:	GIGA-D-19-00374R3	
Full Title:	An improved pig reference genome sequence to enable pig genetics and genomics research	
Article Type:	Research	
Funding Information:	Biotechnology and Biological Sciences Research Council (BBS/E/D/20211550)	Prof Alan L Archibald
	Biotechnology and Biological Sciences Research Council (BBS/E/D/10002070)	Prof Alan L Archibald
	Biotechnology and Biological Sciences Research Council (BB/F021372/1)	Prof Nabeel Affara
	Biotechnology and Biological Sciences Research Council (BB/M011461/1)	Prof Alan L Archibald
	Biotechnology and Biological Sciences Research Council (BB/M011615/1)	Dr Paul Flicek
	Biotechnology and Biological Sciences Research Council (BB/M01844X/1)	Prof Alan L Archibald
	Seventh Framework Programme (KBBE222664)	Not applicable
	Wellcome Trust (WT108749/Z/15/Z)	Dr Paul Flicek
	U.S. Department of Agriculture (8042-31000-001-00-D)	Dr Derek M Bickhart Dr Benjamin D Rosen
	U.S. Department of Agriculture (5090-31000-026-00-D)	Dr Derek M Bickhart
	U.S. Department of Agriculture (3040-31000-100-00-D)	Dr Timothy P L Smith
Abstract:	<p>The domestic pig (<i>Sus scrofa</i>) is important both as a food source and as a biomedical model given its similarity in size, anatomy, physiology, metabolism, pathology and pharmacology to humans. The draft reference genome (Sscrofa10.2) of a purebred Duroc female pig established using older clone-based sequencing methods was incomplete and unresolved redundancies, short range order and orientation errors and associated misassembled genes limited its utility. We present two annotated highly contiguous chromosome-level genome assemblies created with more recent long read technologies and a whole genome shotgun strategy, one for the same Duroc female (Sscrofa11.1) and one for an outbred, composite breed male (USMARCv1.0). Both assemblies are of substantially higher (>90-fold) continuity and accuracy than Sscrofa10.2. These highly contiguous assemblies plus annotation of a further 11 short read assemblies provide an unprecedented view of the genetic make-up of this important agricultural and biomedical model species. We propose that the improved Duroc assembly (Sscrofa11.1) become the reference genome for genomic research in pigs.</p>	
Corresponding Author:	Alan L Archibald UNITED KINGDOM	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		

Corresponding Author's Secondary Institution:	
First Author:	Amanda Warr
First Author Secondary Information:	
Order of Authors:	Amanda Warr
	Nabeel Affara
	Bronwen Aken
	Hamid Beiki
	Derek M Bickhart
	Konstantinos Billis
	William Chow
	Lel Eory
	Heather A Finlayson
	Paul Flicek
	Carlos G Girón
	Darren K Griffin
	Richard Hall
	Greg Hannum
	Thibaut Hourlier
	Kerstin Howe
	David A Hume
	Osagie Izuogu
	Kristi Kim
	Sergey Koren
	Haibou Liu
	Nancy Manchanda
	Fergal J Martin
	Dan J Nonneman
	Rebecca E O'Connor
	Adam M Phillippy
	Gary A Rohrer
	Benjamin D Rosen
	Laurie A Rund
	Carole A Sargent
	Lawrence B Schook
	Steven G Schroeder
	Ariel S Schwartz
	Ben M Skinner
	Richard Talbot
	Elizabeth Tseng

	Christopher K Tuggle
	Mick Watson
	Timothy P L Smith
	Alan L Archibald
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Dear Editors</p> <p>I am pleased to submit a final revised version of the manuscript entitled “An improved pig reference genome sequence to enable pig genetics and genomics research”.</p> <p>The questions and issues raised by the reviewers have been addressed as described previously.</p> <p>The results from the Assemblytics comparisons of the 13 pig genome assemblies available in GigaDB (http://dx.doi.org/10.5524/100732).</p> <p>The figures are integrated into the manuscript and supplementary materials. I have also uploaded copies of each figure image. These images are also available in the associated GigaDB dataset (http://dx.doi.org/10.5524/100732).</p> <p>Yours sincerely</p> <p>Alan L. Archibald</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly</p>	Yes

<p>encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

1 An improved pig reference genome sequence to enable pig genetics and 2 genomics research

3

4 Amanda Warr¹ (amanda.warr@roslin.ed.ac.uk), Nabeel Affara² (na106@cam.ac.uk), Bronwen Aken³
5 (ba1@ebi.ac.uk), Hamid Beiki⁴ (beiki.h.m@gmail.com), Derek M. Bickhart⁵ (derek.bickart@usda.gov),
6 Konstantinos Billis³ (kbillis@ebi.ac.uk), William Chow⁶ (wc2@ebi.ac.uk), Lel Eory¹
7 (lel.eory@roslin.ed.ac.uk), Heather A. Finlayson¹ (heatherfinlayson@gmail.com), Paul Flicek³
8 (flicek@ebi.ac.uk), Carlos G. Girón³ (carlos@ebi.ac.uk), Darren K. Griffin⁷ (d.k.griffin@kent.ac.uk),
9 Richard Hall⁸ (rhall@pacificbiosciences.com), Greg Hannum⁹ (greg@denovium.com), Thibaut
10 Hourlier³ (thibaut@ebi.ac.uk), Kerstin Howe⁶ (kj2@ebi.ac.uk), David A. Hume^{1,†}
11 (david.hume@uq.edu.au), Osagie Izuogu³ (osagie@ebi.ac.uk), Kristi Kim⁸ (kristi.kim07@gmail.com),
12 Sergey Koren¹⁰ (sergey.koren@nih.gov), Haibou Liu⁴ (haiboul2017@gmail.com), Nancy Manchanda¹¹
13 (nancym@iastate.edu), Fergal J. Martin³ (fergal@ebi.ac.uk), Dan J. Nonneman¹²
14 (dan.nonneman@ars.usda.gov), Rebecca E. O'Connor⁷ (r.o'connor@kent.ac.uk), Adam M. Phillippy¹⁰
15 (adam.phillippy@nih.gov), Gary A. Rohrer¹² (gary.rohrer@ars.usda.gov), Benjamin D. Rosen¹³
16 (ben.rosen@usda.gov), Laurie A. Rund¹⁴ (larund@illinois.edu), Carole A. Sargent²
17 (cas1001@cam.ac.uk), Lawrence B. Schook¹⁴ (schook@illinois.edu), Steven G. Schroeder¹³
18 (steven.schroeder@usda.gov), Ariel S. Schwartz⁹ (ariel@denovium.com), Ben M. Skinner²
19 (b.skinner@essex.ac.uk), Richard Talbot¹⁵ (richard.talbot@roslin.ed.ac.uk), Elizabeth Tseng⁸
20 (etseng@pacificbiosciences.com), Christopher K. Tuggle^{4,11} (cktuggle@iastate.edu), Mick Watson¹
21 (mick.watson@roslin.ed.ac.uk), Timothy P. L. Smith^{12*} (tim.smith@ars.usda.gov), Alan L. Archibald^{1*}
22 (alan.archibald@roslin.ed.ac.uk)

23

24 Affiliations

25 ¹The Roslin Institute and Royal (Dick) School of Veterinary Studies, University of Edinburgh, Edinburgh
26 EH25 9RG, U.K.

27 ²Department of Pathology, University of Cambridge, Cambridge CB2 1QP, U.K.

28 ³European Molecular Biology Laboratory, European Bioinformatics Institute, Hinxton, CB10 1SD, U.K.

29 ⁴Department of Animal Science, Iowa State University, Ames, Iowa, U.S.A.

30 ⁵Dairy Forage Research Center, USDA-ARS, Madison, Wisconsin, U.S.A.

31 ⁶Wellcome Sanger Institute, Cambridge, CB10 1SA, U.K.

32 ⁷School of Biosciences, University of Kent, Canterbury CT2 7AF, U.K.

33 ⁸Pacific Biosciences, Menlo Park, California, U.S.A.

34 ⁹Denovium Inc., San Diego, California, U.S.A.

35 ¹⁰Genome Informatics Section, Computational and Statistical Genomics Branch, National Human
36 Genome Research Institute, Bethesda, Maryland, U.S.A.

37 ¹¹Bioinformatics and Computational Biology Program, Iowa State University, Ames, Iowa, U.S.A.

38 ¹²USDA-ARS U.S. Meat Animal Research Center, Clay Center, Nebraska 68933, U.S.A.

39 ¹³Animal Genomics and Improvement Laboratory, USDA-ARS, Beltsville, Maryland, U.S.A

40 ¹⁴Department of Animal Sciences, University of Illinois, Urbana, Illinois, U.S.A.

41 ¹⁵Edinburgh Genomics, University of Edinburgh, Edinburgh EH9 3FL, U.K.

42

43 [†] Current address: Mater Research Institute-University of Queensland, Translational Research
44 Institute, Brisbane, QLD 4102, Australia

45

46 *Corresponding authors: alan.archibald@roslin.ed.ac.uk; tim.smith@ARS.USDA.GOV

47 ORCIDs:

48 Amanda Warr, 0000-0002-7049-1840;

49 Bronwen Aken, 0000-0002-3032-4095;

50 Hamid Beiki, 0000-0002-0516-1431;

51 Derek M Bickhart, 0000-0003-2223-9285;
52 Konstantinos Billis, 0000-0001-8568-4306;
53 William Chow, 0000-0002-9056-201X;
54 LeI Eory, 0000-0001-9358-6239;
55 Paul Flicek, 0000-0002-3897-7955;
56 Carlos G Girón, 0000-0002-0935-7271;
57 Darren K Griffin, 0000-0001-7595-3226;
58 Richard Hall, 0000-0001-6490-8227;
59 Thibaut Hourlier, 0000-0003-4894-7773;
60 Kerstin Howe, 0000-0003-2237-513X;
61 David A. Hume, 0000-0002-2615-1478;
62 Osagie Izuogu, 0000-0003-3116-2558;
63 Sergey Koren, 0000-0002-1472-8962;
64 Fergal J Martin, 0000-0002-1672-050X;
65 Dan J Nonneman, 0000-0001-8073-7843;
66 Rebecca E O'Connor, 0000-0002-4270-970X;
67 Adam M Phillippy, 0000-0003-2983-8934;
68 Gary A Rohrer, 0000-0002-8252-9308;
69 Benjamin D Rosen, 0000-0001-9395-8346;
70 Laurie A Rund, 0000-0003-0761-7196;
71 Lawrence B Schook, 0000-0002-6580-8364;
72 Steven G Schroeder, 0000-0001-9103-5150;
73 Ben M Skinner, 0000-0002-7152-1167;
74 Elizabeth Tseng, 0000-0002-1074-5095;
75 Christopher K Tuggle, 0000-0002-4229-5316;
76 Mick Watson, 0000-0003-4211-0358;

77 Timothy P.L. Smith, 0000-0003-1611-6828;

78 Alan L. Archibald, 0000-0001-9213-1830.

79

80 **Abstract**

81 The domestic pig (*Sus scrofa*) is important both as a food source and as a biomedical model given its
82 similarity in size, anatomy, physiology, metabolism, pathology and pharmacology to humans. The draft
83 reference genome (Sscrofa10.2) of a purebred Duroc female pig established using older clone-based
84 sequencing methods was incomplete, and unresolved redundancies, short range order and
85 orientation errors and associated misassembled genes limited its utility. We present two annotated
86 highly contiguous chromosome-level genome assemblies created with more recent long read
87 technologies and a whole genome shotgun strategy, one for the same Duroc female (Sscrofa11.1) and
88 one for an outbred, composite breed male (USMARCv1.0). Both assemblies are of substantially higher
89 (>90-fold) continuity and accuracy than Sscrofa10.2. These highly contiguous assemblies plus
90 annotation of a further 11 short read assemblies provide an unprecedented view of the genetic make-
91 up of this important agricultural and biomedical model species. We propose that the improved Duroc
92 assembly (Sscrofa11.1) become the reference genome for genomic research in pigs.

93

94 **Keywords**

95 Pig genomes, reference assembly, pig, genome annotation

96

97 **Background**

98 High quality, richly annotated reference genome sequences are key resources and provide important
99 frameworks for the discovery and analysis of genetic variation and for linking genotypes to function.
100 In farmed animal species such as the domestic pig (*Sus scrofa*, NCBI:txid9823) genome sequences have
101 been integral to the discovery of molecular genetic variants and the development of single nucleotide
102 polymorphism (SNP) chips [1] and enabled efforts to dissect the genetic control of complex traits, such
103 as growth, feed conversion, body composition, reproduction, behaviour and responses to infectious
104 diseases [2].

105
106 Genome sequences are not only an essential resource for enabling research but also for applications
107 in the life sciences. Genomic selection, in which associations between thousands of SNPs and trait
108 variation as established in a phenotyped training population are used to choose amongst selection
109 candidates for which there are SNP data but no phenotypes, has delivered genomics-enabled genetic
110 improvement in farmed animals [3] and plants. From its initial successful application in dairy cattle
111 breeding, genomic selection is now being used in many sectors within animal and plant breeding,
112 including by leading pig breeding companies [4, 5].

113
114 The domestic pig (*Sus scrofa*) has importance not only as a source of animal protein but also as a
115 biomedical model. The choice of the optimal animal model species for pharmacological or toxicology
116 studies can be informed by knowledge of the genome and gene content of the candidate species
117 including pigs [6]. A high quality, richly annotated genome sequence is also essential when using gene
118 editing technologies to engineer improved animal models for research or as sources of cells and tissue
119 for xenotransplantation and potentially for improved productivity [7, 8].

120
121 The highly continuous pig genome sequences reported here are built upon a quarter of a century of
122 effort by the global pig genetics and genomics research community including the development of

123 recombination and radiation hybrid maps [9, 10], cytogenetic and Bacterial Artificial Chromosome
124 (BAC) physical maps [11, 12] and a draft reference genome sequence [13].

125

126 The previously published draft pig reference genome sequence (Sscrofa10.2), developed under the
127 auspices of the Swine Genome Sequencing Consortium (SGSC), has a number of significant deficiencies
128 [14-17]. The BAC-by-BAC hierarchical shotgun sequence approach [18] using Sanger sequencing
129 technology can yield a high quality genome sequence as demonstrated by the public Human Genome
130 Project. However, with a fraction of the financial resources of the Human Genome Project, the
131 resulting draft pig genome sequence comprised an assembly, in which long-range order and
132 orientation is good, but the order and orientation of sequence contigs within many BAC clones was
133 poorly supported and the sequence redundancy between overlapping sequenced BAC clones was
134 often not resolved. Moreover, about 10% of the pig genome, including some important genes, were
135 not represented (e.g. *CD163*), or incompletely represented (e.g. *IGF2*) in the assembly [19]. Whilst the
136 BAC clones represent an invaluable resource for targeted sequence improvement and gap closure as
137 demonstrated for chromosome X (SSCX) [20], a clone-by-clone approach to sequence improvement is
138 expensive notwithstanding the reduced cost of sequencing with next-generation technologies.

139

140 The dramatically reduced cost of whole genome shotgun sequencing using Illumina short read
141 technology has facilitated the sequencing of several hundred pig genomes [17, 21, 22]. Whilst a few
142 of these additional pig genomes have been assembled to contig level, most of these genome
143 sequences have simply been aligned to the reference and used as a resource for variant discovery.

144

145 The increased capability and reduced cost of third generation long read sequencing technology as
146 delivered by Pacific Biosciences and Oxford Nanopore platforms, have created the opportunity to
147 generate the data from which to build highly contiguous genome sequences as illustrated recently for
148 cattle [23, 24]. Here we describe the use of Pacific Biosciences (PacBio) long read technology to

149 establish highly continuous pig genome sequences that provide substantially improved resources for
150 pig genetics and genomics research and applications.
151

152 **Results**

153 Two individual pigs were sequenced independently: a) TJ Tabasco (Duroc 2-14) i.e. the sow that was
154 the primary source of DNA for the published draft genome sequence (Sscrofa10.2) [13] and b)
155 MARC1423004 which was a crossbred barrow (i.e. castrated male pig) from a composite population
156 (approximately ½ Landrace, ¼ Duroc and ¼ Yorkshire) at the USDA Meat Animal Research Center. The
157 former allowed us to build upon the earlier draft genome sequence, exploit the associated CHORI-242
158 BAC library resource [25] and evaluate the improvements achieved by comparison with Sscrofa10.2.
159 The latter allowed us to assess the relative efficacy of a simpler whole genome shotgun sequencing
160 and Chicago Hi-Rise scaffolding strategy [26]. This second assembly also provided data for the Y
161 chromosome, and supported comparison of haplotypes between individuals. In addition, full-length
162 transcript sequences were collected for multiple tissues from the MARC1423004 animal, and used in
163 annotating both genomes.

164

165 Sscrofa11.1 assembly

166 Approximately sixty-five fold coverage (176 Gbp) of the genome of TJ Tabasco (Duroc 2-14) was
167 generated using Pacific Biosciences (PacBio) single-molecule real-time (SMRT) sequencing technology.
168 A total of 213 SMRT cells produced 12,328,735 subreads of average length 14,270 bp and with a read
169 N50 of 19,786 bp (Table S1). Reads were corrected and assembled using Falcon (v.0.4.0) [27],
170 achieving a minimum corrected read cutoff of 13 kbp that provided 19-fold genome coverage for input
171 resulting in an initial assembly comprising 3,206 contigs with a contig N50 of 14.5 Mbp.

172

173 The contigs were mapped to the previous draft assembly (Sscrofa10.2) using Nucmer [28]. The long
174 range order of the Sscrofa10.2 assembly was based on fingerprint contig (FPC) [12] and radiation
175 hybrid physical maps with assignments to chromosomes based on fluorescent *in situ* hybridisation
176 data. This alignment of Sscrofa10.2 and the contigs from the initial Falcon assembly of the PacBio data
177 provided draft scaffolds that were tested for consistency with paired BAC and fosmid end sequences

178 and the radiation hybrid map [9]. The draft scaffolds also provided a framework for gap closure using
179 PBJelly [29], or finished quality Sanger sequence data generated from CHORI-242 BAC clones from
180 earlier work [13, 20].

181

182 Remaining gaps between contigs within scaffolds, and between scaffolds predicted to be adjacent on
183 the basis of other available data, were targeted for gap filling with a combination of unplaced contigs
184 and previously sequenced BACs, or by identification and sequencing of BAC clones predicted from
185 their end sequences to span the gaps. The combination of methods filled 2,501 gaps and reduced the
186 number of contigs in the assembly from 3,206 to 705. The assembly, Sscrofa11 (GCA_000003025.5),
187 had a final contig N50 of 48.2 Mbp, only 103 gaps in the sequences assigned to chromosomes, and
188 only 583 remaining unplaced contigs (Table 1). Two acrocentric chromosomes (SSC16, SSC18) were
189 each represented by single, unbroken contigs. The SSC18 assembly also includes centromeric and
190 telomeric repeats (Tables S2, S3; Figs. S1, S2), albeit the former probably represent a collapsed version
191 of the true centromere. The reference genome assembly was completed by adding Y chromosome
192 sequences from other sources (GCA_900119615.2) [20] because TJ Tabasco (Duroc 2-14) was female.
193 The resulting reference genome sequence was termed Sscrofa11.1 and deposited in the public
194 sequence databases (GCA_000003025.6) (Table 1).

195

196 The medium to long range order and orientation of Sscrofa11.1 assembly was assessed by comparison
197 to an existing radiation hybrid (RH) map [9]. The comparison strongly supported the overall accuracy
198 of the assembly (Fig. 1a), despite the fact that the RH map was prepared from a cell line of a different
199 individual. There is one major disagreement between the RH map and the assembly on chromosome
200 3, which will need further investigating. The only other substantial disagreement on chromosome 9,
201 is explained by a gap in the RH map [9]. The assignment and orientation of the Sscrofa11.1 scaffolds
202 to chromosomes was confirmed with fluorescent *in situ* hybridisation (FISH) of BAC clones (Table S4,
203 Fig. S3). The Sscrofa11.1 and USMARCv1.0 assemblies were searched using BLAST [30] with sequences

204 derived from the BAC clones which had been used as probes for the FISH analyses. For most BAC
205 clones these sequences were BAC end sequences [12], but in some cases these sequences were
206 incomplete or complete BAC clone sequences [13, 20]. The links between the genome sequence and
207 the BAC clones used in cytogenetic analyses by fluorescent *in situ* hybridization are summarised in
208 Table S4. The fluorescent *in situ* hybridization results indicate areas where future assemblies might be
209 improved. For example, the Sscrofa11.1 unplaced scaffolds contig 1206 and contig1914 may contain
210 sequences that could be added to end of the long arms of SSC1 and SSC7 respectively.

211

212 The quality of the Sscrofa11 assembly, which corresponds to Sscrofa11.1 after the exclusion of SSCY,
213 was assessed as described previously for the existing Sanger sequence based draft assembly
214 (Sscrofa10.2) [14]. Alignments of Illumina sequence reads from the same female pig were used to
215 identify regions of low quality (LQ; regions with high GC normalised coverage, prevalence of
216 improperly paired reads and prevalence of reads with improper insert sizes) or low coverage (LC;
217 regions with low GC normalised coverage) (Table 2).. The analysis confirms that Sscrofa11 represents
218 a significant improvement over the Sscrofa10.2 draft assembly. For example, the Low Quality Low
219 Coverage (LQLC) proportion of the genome sequence has dropped from 33.07% to 16.3% when
220 repetitive sequence is not masked, and falls to 1.6% when repeats are masked prior to read alignment.

221

222 The remaining LQLC segments of Sscrofa11 have an average GC content of 61.6%. Thus, these regions
223 may represent sequence where short read coverage is low due to the known systematic bias of the
224 short read platform against extreme GC content sequences, rather than deficiencies of the assembly.

225 The Sscrofa11.1 assembly was also assessed visually using gEVAL [31]. The improvement in short range
226 order and orientation as revealed by alignments with isogenic BAC and fosmid end sequences is
227 illustrated for a particularly poor region of Sscrofa10.2 on chromosome 12 (Fig. S4). The problems in
228 this area of Sscrofa10.2 arose from failures to order and orient the sequence contigs and resolve the
229 redundancies between these sequence contigs within BAC clone CH242-147O24 (ENA: FP102566.2).

230 The improved contiguity in Sscrofa11.1 not only resolves these local order and orientation errors, but
231 also facilitates the annotation of a complete gene model for the *ABR* locus. Further examples of
232 comparisons of Sscrofa10.2 and Sscrofa11.1 reveal improvements in contiguity, local order and
233 orientation and gene models (Fig. S5 to S7).

234

235 USMARCv1.0 assembly

236 Approximately sixty-five fold coverage of the genome of the MARC1423004 barrow was generated on
237 a PacBio RSII instrument. The sequence was collected during the transition from P5/C3 to P6/C4
238 chemistry, with approximately equal numbers of subreads from each chemistry. A total of 199 cells of
239 P5/C3 chemistry produced 95.3 Gbp of sequence with mean subread length of 5.1 kbp and subread
240 N50 of 8.2 kbp. A total of 127 cells of P6/C4 chemistry produced 91.6 Gbp of sequence with mean
241 subread length 6.5 kbp and subread N50 of 10.3 kbp, resulting in an overall average subread length,
242 including data from both chemistries, of 6.4 kbp. The reads were assembled using Celera Assembler
243 8.3rc2 [32] and Falcon [27]). The resulting assemblies were compared and the Celera Assembler result
244 was selected based on better agreement with a Dovetail Chicago[®] library [26] (i.e. there was a lower
245 proportion of conflicting links between read pairs from the Chicago[®] library), and was used to create
246 a scaffolded assembly with the HiRise[™] scaffolder consisting of 14,818 contigs with a contig N50 of
247 6.372 Mbp (GenBank accession GCA_002844635.1; Table 1). The USMARCv1.0 scaffolds were
248 therefore completely independent of the existing Sscrofa10.2 or new Sscrofa11.1 assemblies, and they
249 can act as supporting evidence where they agree with those assemblies. However, chromosome
250 assignment of the scaffolds was performed by alignment to Sscrofa10.2, and does not constitute
251 independent confirmation of this ordering. The assignment of these scaffolds to individual
252 chromosomes was confirmed post-hoc by FISH analysis as described for Sscrofa11.1 above. The FISH
253 analysis revealed that several of these chromosome assemblies (SSC1, 5, 6-11, 13-16) are inverted
254 with respect to the cytogenetic convention for pig chromosome (Table S4; Figs. S3, S8 to S10). After

255 correcting the orientation of these inverted scaffolds, there is good agreement between the
256 USMARCv1.0 assembly and the RH map [9] (Fig. 1b, Table S5).

257

258 Sscrofa11.1 and USMARCv1.0 are co-linear

259 The alignment of the two PacBio assemblies reveals a high degree of agreement and co-linearity, after
260 correcting the inversions of several USMARCv1.0 chromosome assemblies (Fig. S11). The agreement
261 between the Sscrofa11.1 and USMARCv1.0 assemblies is also evident in comparisons of specific loci
262 (Figs. S5 to S7) although with some differences (e.g. Fig. S6). The whole genome alignment of
263 Sscrofa11.1 and USMARCv1.0 (Fig. S11) masks some inconsistencies that are evident when the
264 alignments are viewed on a single chromosome-by-chromosome basis (Figs. S8 to S10). It remains to
265 be determined whether the small differences between the assemblies represent errors in the
266 assemblies, or true structural variation between the two individuals (see discussion of the *ERLIN1*
267 locus below).

268

269 Pairwise comparisons amongst the Sscrofa10.2, Sscrofa11.1 and USMARCv1.0 assemblies using the
270 Assemblytics tools [33] revealed a peak of insertions and deletion with sizes of about 300 bp (Figs.
271 S12a to S12c). We assume that these correspond to short interspersed nuclear elements (SINEs). Both
272 Sscrofa11.1 and USMARCv1.0 assemblies have more differences against Sscrofa10.2 (33,347 and
273 44,023 respectively) than against each other (28,733). This is despite the fact that Sscrofa11.1 and
274 Sscrofa10.2 represent the same pig genome. While some differences between Sscrofa10.2 and
275 Sscrofa11.1 may be due to differences in which haplotype has been captured in the assembly, the
276 reduction in low quality and low coverage regions and the dramatic decrease in differences versus
277 USMARCv1.0 leads us to conclude that the majority are improvements in the Sscrofa11.1 assembly.
278 The differences between the Sscrofa11.1 and USMARCv1.0 will represent a mix of true structural
279 differences and assembly errors that will require further research to resolve. The Sscrofa11.1 and
280 USMARCv1.0 assemblies were also compared to 11 Illumina short read assemblies [17] (Table S6).

281

282 Repetitive sequences, centromeres and telomeres

283 The repetitive sequence content of the Sscrofa11.1 and USMARCv1.0 was identified and
284 characterised. These analyses allowed the identification of centromeres and telomeres for several
285 chromosomes. The previous reference genome (Sscrofa10.2) that was established from Sanger
286 sequence data and a minipig genome (minipig_v1.0, GCA_000325925.2) that was established from
287 Illumina short read sequence data were also included for comparison. The numbers of the different
288 repeat classes and the average mapped lengths of the repetitive elements identified in these four pig
289 genome assemblies are summarised in Figures S13 and S14, respectively.

290

291 Putative telomeres were identified at the proximal ends of Sscrofa11.1 chromosome assemblies of
292 SSC2, SSC3, SSC6, SSC8, SSC9, SSC14, SSC15, SSC18 and SSCX (Fig S1; Table S2). Putative centromeres
293 were identified in the expected locations in the Sscrofa11.1 chromosome assemblies for SSC1-7, SSC9,
294 SSC13 and SSC18 (Fig S2, Table S3). For the chromosome assemblies of each of SSC8, SSC11 and SSC15
295 two regions harbouring centromeric repeats were identified. Pig chromosomes SSC1-12 plus SSCX and
296 SSCY are all metacentric, whilst chromosomes SSC13-18 are acrocentric. The putative centromeric
297 repeats on SSC17 do not map to the expected end of the chromosome assembly.

298

299 Completeness of the assemblies

300 The Sscrofa11.1 and USMARCv1.0 assemblies were assessed for completeness using two tools, BUSCO
301 (Benchmarking Universal Single-Copy Orthologs) [34] and Cogent [35]. BUSCO uses a database of
302 expected gene content based on near-universal single-copy orthologs from species with genomic data,
303 while Cogent uses transcriptome data from the organism being sequenced, and therefore provides an
304 organism-specific view of genome completeness. BUSCO analysis suggests both new assemblies are
305 highly complete, with 93.8% and 93.1% of BUSCOs complete for Sscrofa11.1 and USMARCv1.0

306 respectively, a marked improvement on the 80.9% complete in Sscrofa10.2 and comparable to the
307 human and mouse reference genome assemblies (Table S7).

308

309 Cogent is a tool that identifies gene families and reconstructs the coding genome using full-length,
310 high-quality (HQ) transcriptome data without a reference genome and can be used to check
311 assemblies for the presence of these known coding sequences [35]. PacBio transcriptome (Iso-Seq)
312 data consisting of high-quality isoform sequences from 7 tissues (diaphragm, hypothalamus, liver,
313 skeletal muscle (*longissimus dorsi*), small intestine, spleen and thymus) [36] from the pig whose DNA
314 was used as the source for the USMARCv1.0 assembly were pooled together for Cogent analysis.
315 Cogent partitioned 276,196 HQ isoform sequences into 30,628 gene families, of which 61% had at
316 least 2 distinct transcript isoforms. Cogent then performed reconstruction on the 18,708 partitions.
317 For each partition, Cogent attempts to reconstruct coding 'contigs' that represent the ordered
318 concatenation of transcribed exons as supported by the isoform sequences. The reconstructed contigs
319 were then mapped back to Sscrofa11.1 and contigs that could not be mapped or map to more than
320 one position are individually examined. There were five genes that were present in the Iso-Seq data,
321 but missing in the Sscrofa11.1 assembly. In each of these five cases, a Cogent partition (which consists
322 of 2 or more transcript isoforms of the same gene, often from multiple tissues) exists in which the
323 predicted transcript does not align back to Sscrofa11.1. NCBI-BLASTN of the isoforms from the
324 partitions revealed them to have near perfect hits with existing annotations for *CHAMP1*, *ERLIN1*,
325 *IL1RN*, *MB*, and *PSD4* for other species.

326

327 *ERLIN1* is missing from its predicted location on SSC14 between *CHUK* and *CPN1* gene in Sscrofa11.1.
328 There is good support for the Sscrofa11.1 assembly in the region from the BAC end sequence
329 alignments suggesting this area may represent a true haplotype. Indeed, a copy number variant (CNV)
330 nsv1302227 has been mapped to this location on SSC14 [37] and the *ERLIN1* gene sequences present
331 in BAC clone CH242-513L2 (ENA: CT868715.3) were incorporated into the earlier Sscrofa10.2

332 assembly. However, an alternative haplotype containing *ERLIN1* was not found in any of the
333 assembled contigs from Falcon and this will require further investigation. The *ERLIN1* locus is present
334 on SSC14 in the USMARCv1.0 assembly (30,107,816-30,143,074; note the USMARCv1.0 assembly of
335 SSC14 is inverted relative to Sscrofa11.1). Of eleven short read pig genome assemblies [17] that have
336 been annotated with the Ensembl pipeline (Ensembl release 98, September 2019) [38, 39] *ERLIN1*
337 sequences are present in the expected genomic context in all eleven genome assemblies. As the
338 *ERLIN1* gene is located at the end of a contig in eight of these short read assemblies, it suggests that
339 this region of the pig genome presents difficulties for sequencing and assembly and the absence of
340 *ERLIN1* in the Sscrofa11.1 is more likely to be an assembly error.

341

342 The other 4 genes are annotated in neither Sscrofa10.2 nor Sscrofa11.1. Two of these genes, *IL1RN*
343 and *PSD4*, are present in the original Falcon contigs, however they were trimmed off during the contig
344 QC stage because of apparent abnormal Illumina, BAC and fosmid mapping in the region which was
345 likely caused by the repetitive nature of their expected location on chromosome 3 where a gap is
346 present. The *IL1RN* and *PSD4* genes are present in the USMARCv1.0, albeit their location is anomalous,
347 and are also present in the 11 short read assemblies [17]. *CHAMP1* (ENSSSCG00070014091) is present
348 in the USMARCv1.0 assembly in the sub-telomeric region of the q-arm, after correcting the inversion
349 of the USMARCv1.0 scaffold and is also present in all 11 short read assemblies [17]. After correcting
350 the orientation of the USMARCv1.0 chromosome 11 scaffold there is a small inversion of the distal
351 1.07 Mbp relative to the Sscrofa11.1 assembly; this region harbours the *CHAMP1* gene. The
352 orientation of the Sscrofa11.1 chromosome 11 assembly in this region is consistent with the
353 predictions of the human-pig comparative map [40]. The myoglobin gene (*MB*) is present in the
354 expected location in the USMARCv1.0 assembly flanked by *RASD2* and *RBFOX2*. Partial *MB* sequences
355 are present distal to *RBFOX2* on chromosome 5 in the Sscrofa11.1 assembly. As there is no gap here
356 in the Sscrofa11.1 assembly it is likely that the incomplete *MB* is a result of a misassembly in this
357 region. This interpretation is supported by a break in the pairs of BAC and fosmid end sequences that

358 map to this region of the Sscrofa11.1 assembly. Some of the expected gene content missing from this
359 region of the Sscrofa11.1 chromosome 5 assembly, including *RASD2*, *HMOX1* and *LARGE1* is present
360 on an unplaced scaffold (AEMK02000361.1). Cogent analysis also identified 2 cases of potential
361 fragmentation in the Sscrofa11.1 genome assembly that resulted in the isoforms being mapped to two
362 separate loci, though these will require further investigation. In summary, the BUSCO and Cogent
363 analyses indicate that the Sscrofa11.1 assembly captures a very high proportion of the expressed
364 elements of the genome.

365

366 Improved annotation

367 Annotation of Sscrofa11.1 was carried out with the Ensembl annotation pipeline and released via the
368 Ensembl Genome Browser (Ensembl release 90, August 2017) [38, 41]. Statistics for the annotation
369 as updated in June 2019 (Ensembl release 98, September 2019) are listed in Table 3. This annotation
370 is more complete than that of Sscrofa10.2 and includes fewer fragmented genes and pseudogenes.

371

372 The annotation pipeline utilised extensive short read RNA-Seq data from 27 tissues and long read
373 PacBio Iso-Seq data from 9 adult tissues. This provided an unprecedented window into the pig
374 transcriptome and allowed for not only an improvement to the main gene set, but also the generation
375 of tissue-specific gene tracks from each tissue sample. The use of Iso-Seq data also improved the
376 annotation of UTRs, as they represent transcripts sequenced across their full length from the polyA
377 tract.

378

379 In addition to improved gene models, annotation of the Sscrofa11.1 assembly provides a more
380 complete view of the porcine transcriptome than annotation of the previous assembly (Sscrofa10.2;
381 Ensembl releases 67-89, May 2012 – May 2017) [42] with increases in the numbers of transcripts
382 annotated (Table 3). However, the number of annotated transcripts remains lower than in the human

383 and mouse genomes. The annotation of the human and mouse genomes and in particular the gene
384 content and encoded transcripts has been more thorough as a result of extensive manual annotation.

385

386 Efforts were made to annotate important classes of genes, in particular immunoglobulins and
387 olfactory receptors. For these genes, sequences were downloaded from specialist databases and the
388 literature in order to capture as much detail as possible (see supplementary information for more
389 details).

390

391 These improvements in terms of the resulting annotation were evident in the results of the
392 comparative genomics analyses run on the gene set. The previous annotation had 12,919 one-to-one
393 orthologs with human, while the new annotation of the Sscrofa11.1 assembly has 15,544. Similarly, in
394 terms of conservation of synteny, the previous annotation had 11,661 genes with high confidence
395 gene order conservation scores, while the new annotation has 15,958. There was also a large
396 reduction in terms of genes that were either abnormally short or split when compared to their
397 orthologs in the new annotation.

398

399 The Sscrofa11.1 assembly has also been annotated using the NCBI pipeline [43]. We have compared
400 these two annotations. The Ensembl and NCBI annotations of Sscrofa11.1 are broadly similar (Table
401 S8). There are 17,676 protein coding genes and 1,700 non-coding genes in common. However, 540 of
402 the genes annotated as protein-coding by Ensembl are annotated as non-coding or pseudogenes by
403 NCBI and 227 genes annotated as non-coding by NCBI are annotated as protein-coding (215) or as
404 pseudogenes (12) by Ensembl. The NCBI RefSeq annotation can be visualised in the Ensembl Genome
405 Browser by loading the RefSeq GFF3 track and the annotations compared at the individual locus level.
406 Similarly, the Ensembl annotated genes can be visualised in the NCBI Genome Browser. Despite
407 considerable investment there are also differences in the Ensembl and NCBI annotation of the human
408 reference genome sequence with 20,444 and 19,755 protein-coding genes on the primary assembly,

409 respectively. The MANE (Matched Annotation from NCBI and EMBL-EBI) project was launched to
410 resolve these differences and identify a matched representative transcript for each human protein-
411 coding gene [44]. To date a MANE transcript has been identified for 12,985 genes.

412

413 We have also annotated the USMARCv1.0 assembly using the Ensembl pipeline [38] and this
414 annotation was released via the Ensembl Genome Browser (Ensembl release 97, July 2019) [39](see
415 Table 3 for summary statistics). More recently, we have annotated a further eleven short read pig
416 genome assemblies [17] (Ensembl release 98, September 2019) [39]; see Tables S6 and S11 for
417 summary statistics for the assemblies and annotation, respectively.

418

419 SNP chip probes mapped to assemblies

420 The probes from four commercial SNP chips were mapped to the Sscrofa10.2, Sscrofa11.1 and
421 USMARCv1.0 assemblies. We identified 1,709, 56, and 224 markers on the PorcineSNP60, GGP LD and
422 80K commercial chips that were previously unmapped and now have coordinates on the Sscrofa11.1
423 reference (Table S9). These newly mapped markers can now be imputed into a cross-platform,
424 common set of SNP markers for use in genomic selection. Additionally, we have identified areas of the
425 genome that are poorly tracked by the current set of commercial SNP markers. The previous
426 Sscrofa10.2 reference had an average marker spacing of 3.57 kbp (Stdev: 26.5 kbp) with markers from
427 four commercial genotyping arrays. We found this to be an underestimate of the actual distance
428 between markers, as the Sscrofa11.1 reference coordinates consisted of an average of 3.91 kbp
429 (Stdev: 14.9 kbp) between the same set of markers. We also found a region of 2.56 Mbp that is
430 currently devoid of suitable markers on the new reference.

431

432 A Spearman's rank order (ρ) value was calculated for each assembly (alternative hypothesis: ρ is
433 equal to zero; $p < 2.2 \times 10^{-16}$): Sscrofa10.2: 0.88464; Sscrofa11.1: 0.88890; USMARCv1.0: 0.81260. This
434 rank order comparison was estimated by ordering all of the SNP probes from all chips by their listed

435 manifest coordinates against their relative order in each assembly (with chromosomes ordered by
436 karyotype). Any unmapped markers in an assembly were penalized by giving the marker a "-1" rank in
437 the assembly ranking order.

438

439 In order to examine general linear order of placed markers on each assembly, the marker rank order
440 (y axis; used above in the Spearman's rank order test) was plotted against the rank order of the probe
441 rank order on the manifest file (x axis) (Fig. S15). The analyses revealed some interesting artefacts that
442 suggest that the SNP manifest coordinates for the porcine 60K SNP chip are still derived from an
443 obsolete (Sscrofa9) reference in contrast to all other manifests (Sscrofa10.2). Also, it confirms that
444 several of the USMARCv1.0 chromosome scaffolds are inverted with respect to the canonical
445 orientation of pig chromosomes. The large band of points at the top of the plot corresponds to marker
446 mappings on the unplaced contigs of each assembly. These unplaced contigs often correspond to
447 assemblies of alternative haplotypes in heterozygous regions of the reference animal [24]. Marker
448 placement on these segments suggests that these variants are tracking different haplotypes in the
449 population, which is the desired intent of genetic markers used in Genomic Selection.

450

451 **Discussion**

452 We have assembled a superior, extremely continuous reference assembly (Sscrofa11.1) by leveraging
453 the excellent contig lengths provided by long reads, and a wealth of available data including Illumina
454 paired-end, BAC end sequence, finished BAC sequence, fosmid end sequences, and the earlier curated
455 draft assembly (Sscrofa10.2). The pig genome assemblies USMARCv1.0 and Sscrofa11.1 reported here
456 are 92-fold to 694-fold respectively, more continuous than the published draft reference genome
457 sequence (Sscrofa10.2) [13]. The new pig reference genome assembly (Sscrofa11.1) with its contig
458 N50 of 48,231,277 bp and 506 gaps compares favourably with the current human reference genome
459 sequence (GRCh38.p12) that has a contig N50 of 57,879,411 bp and 875 gaps (Table 1). Indeed,
460 considering only the chromosome assemblies built on PacBio long read data (i.e. Sscrofa11 - the
461 autosomes SSC1-SSC18 plus SSCX), there are fewer gaps in the pig assembly than in human reference
462 autosomes and HSAX assemblies. Most of the gaps in the Sscrofa11.1 reference assembly are
463 attributed to the fragmented assembly of SSCY. The capturing of centromeres and telomeres for
464 several chromosomes (Tables S2, S3; Figs. S1, S2) provides further evidence that the Sscrofa11.1
465 assembly is more complete. The increased contiguity of Sscrofa11.1 is evident in the graphical
466 comparison to Sscrofa10.2 illustrated in Figure 2.

467

468 The improvements in the reference genome sequence (Sscrofa11.1) relative to the draft assembly
469 (Sscrofa10.2) [13] are not restricted to greater continuity and fewer gaps. The major flaws in the BAC
470 clone-based draft assembly were i) failures to resolve the sequence redundancy amongst sequence
471 contigs within BAC clones and between adjacent overlapping BAC clones and ii) failures to accurately
472 order and orient the sequence contigs within BAC clones. Although the Sanger sequencing technology
473 used has a much lower raw error rate than the PacBio technology, the sequence coverage was only 4-
474 6 fold across the genome. The improvements in continuity and quality (Table 2; Figs. S5 to S7) have
475 yielded a better template for annotation resulting in better gene models. The Sscrofa11.1 and
476 USMARCv1.0 assemblies are classed as 4|4|1 and 3|5|1 [10^X : N50 contig (kbp); 10^Y : N50 scaffold

477 (kbp); Z = 1 | 0: assembled to chromosome level] respectively compared to Sscrofa10.2 as 1 | 2 | 1 and
478 the human GRCh38p5 assembly as 4 | 4 | 1 [45].

479

480 The improvement in the complete BUSCO (Benchmarking Universal Single-Copy Orthologs) genes
481 indicates that both Sscrofa11.1 and USMARCv1.0 represent superior templates for annotation of gene
482 models than the draft Sscrofa10.2 assembly and are comparable to the finished human and mouse
483 reference genome sequences (Table S7). Further, a companion bioinformatics analysis of available Iso-
484 seq and companion Illumina RNA-seq data across the nine tissues surveyed has identified a large
485 number (>54,000) of novel transcripts [36]. A majority of these transcripts are predicted to be spliced
486 and validated by RNA-seq data. Beiki and colleagues identified 10,465 genes expressing Iso-seq
487 transcripts that are present on the Sscrofa11.1 assembly, but which are unannotated in current NCBI
488 or Ensembl annotations [36].

489

490 Whilst the alignment of the Sscrofa11.1 and USMARCv1.0 assemblies revealed that several of the
491 USMARCv1.0 chromosome assemblies are inverted relative to Sscrofa11.1 and the cytogenetic map.
492 Such inversions are due to the agnostic nature of genome assembly and post-assembly polishing
493 programs. Unless these are corrected post-hoc by manual curation, they result in artefactual
494 inversions of the entire chromosome. However, such inversions do not generally impact downstream
495 analysis that does not involve the relative order/orientation of whole chromosomes.

496

497 Whether the differences between Sscrofa11.1 and USMARCv1.0 in order and orientation within
498 chromosomes represent assembly errors or real chromosomal differences will require further
499 research. The sequence present at the telomeric end of the long arm of the USMARCv1.0 chromosome
500 7 assembly (after correcting the orientation of the USMARCv1.0 SSC7) is missing from the Sscrofa11.1
501 SSC7 assembly, and currently located on a 3.8 Mbp unplaced scaffold (AEMK02000452.1). This
502 unplaced scaffold harbours several genes including *DIO3*, *CKB* and *NUDT14* whose orthologues map

503 to human chromosome 14 as would be predicted from the pig-human comparative map [40]. This
504 omission will be corrected in an updated assembly in future.

505

506 We demonstrate moderate improvements in the placement and ordering of commercial SNP
507 genotyping markers on the Sscrofa11.1 reference genome which will impact future genomic selection
508 programs. The reference-derived order of SNP markers plays a significant role in imputation accuracy,
509 as demonstrated by a whole-genome survey of misassembled regions in cattle that found a correlation
510 between imputation errors and misassemblies [46]. The gaps in SNP chip marker coverage that we
511 identified will inform future marker selection surveys, which are likely to prioritize regions of the
512 genome that are not currently being tracked by marker variants in close proximity to potential causal
513 variant sites. In addition to the gaps in coverage provided by the commercial SNP chips there are
514 regions of the genome assemblies that are devoid of annotated sequence variation as hitherto
515 sequence variants have been discovered against incomplete genome assemblies. Thus, there is a need
516 to re-analyse good quality re-sequence data against the new assemblies in order to provide a better
517 picture of sequence variation in the pig genome.

518

519 The cost of high coverage whole-genome sequencing (WGS) precludes it from routine use in breeding
520 programs. However, it has been suggested that low coverage WGS followed by imputation of
521 haplotypes may be a cost-effective replacement for SNP arrays in genomic selection [47]. Imputation
522 from low coverage sequence data to whole genome information has been shown to be highly accurate
523 [48, 49]. At the 2018 World Congress on Genetics Applied to Livestock Production Aniek Bouwman
524 reported that in a comparison of Sscrofa10.2 with Sscrofa11.1 (for SSC7 only) for imputation from
525 600K SNP genotypes to whole genome sequence overall imputation accuracy on SSC7 improved
526 considerably from 0.81 (1,019,754 variants) to 0.90 (1,129,045 variants) (Aniek Bouwman, pers.
527 comm). Thus, the improved assembly may not only serve as a better template for discovering genetic
528 variation but also have advantages for genomic selection, including improved imputation accuracy.

529

530 Advances in the performance of long read sequencing and scaffolding technologies, improvements in
531 methods for assembling the sequence reads and reductions in costs are enabling the acquisition of
532 ever more complete genome sequences for multiple species and multiple individuals within a species.
533 For example, in terms of adding species, the Vertebrate Genomes Project [50] aims to generate error-
534 free, near gapless, chromosomal level, haplotyped phase assemblies of all of the approximately 66,000
535 vertebrate species and is currently in its first phase that will see such assemblies created for an
536 exemplar species from all 260 vertebrate orders. At the level of individuals within a species, smarter
537 assembly algorithms and sequencing strategies are enabling the production of high quality truly
538 haploid genome sequences for outbred individuals [24]. The establishment of assembled genome
539 sequences for key individuals in the nucleus populations of the leading pig breeding companies is
540 achievable and potentially affordable. However, 10-30x genome coverage short read data generated
541 on the Illumina platform and aligned to a single reference genome is likely to remain the primary
542 approach to sequencing multiple individuals within farmed animal species such as cattle and pigs [21,
543 51].

544

545 There are significant challenges in making multiple assembled genome resources useful and
546 accessible. The current paradigm of presenting a reference genome as a linear representation of a
547 haploid genome of a single individual is an inadequate reference for a species. As an interim solution
548 the Ensembl team are annotating multiple assemblies for some species such as mouse and dog [52,
549 53, 54]. We have implemented this solution for pig genomes, including eleven Illumina short-read
550 assemblies [17] in addition to the reference Sscrofa11.1 and USMARCv1.0 assemblies reported here
551 (Ensembl release 98, September 2019) [39, 41]. Although these additional pig genomes are highly
552 fragmented (Table S6) with contig N50 values from 32 – 102 kbp, the genome annotation (Table S11)
553 provides a resource to explore pig gene space across thirteen genomes, including six Asian pig

554 genomes. The latter are important given the deep phylogenetic split of about 1 million years between
555 European and Asian pigs [13].

556

557 The current human genome reference already contains several hundred alternative haplotypes and it
558 is expected that the single linear reference genome of a species will be replaced with a new model –
559 the graph genome [55-57]. These paradigm shifts in the representation of genomes present challenges
560 for current sequence alignment tools and the ‘best-in-genome’ annotations generated thus far. The
561 generation of high quality annotation remains a labour-intensive and time-consuming enterprise.
562 Comparisons with the human and mouse reference genome sequences which have benefited from
563 extensive manual annotation indicate that there is further complexity in the porcine genome as yet
564 unannotated (Table 3). It is very likely that there are many more transcripts, pseudogenes and non-
565 coding genes (especially long non-coding genes), to be discovered and annotated on the pig genome
566 sequence [36]. The more highly continuous pig genome sequences reported here provide an improved
567 framework against which to discover functional sequences, both coding and regulatory, and sequence
568 variation. After correction for some contig/scaffold inversions in the USMARCv1.0 assembly, the
569 overall agreement between the assemblies is high and illustrates that the majority of genomic
570 variation is at smaller scales of structural variation. However, both assemblies still represent a
571 composite of the two parental genomes present in the animals, with unknown effects of haplotype
572 switching on the local accuracy across the assembly.

573

574 Future developments in high quality genome sequences for the domestic pig are likely to include: (i)
575 gap closure of Sscrofa11.1 to yield an assembly with one contig per (autosomal) chromosome arm
576 exploiting the isogenic BAC and fosmid clone resource as illustrated here for chromosome 16 and 18;
577 and (ii) haplotype resolved assemblies of a Meishan and White Composite F1 crossbred pig (i.e. the
578 offspring of a Meishan sire and a White Composite dam that is approximately $\frac{1}{2}$ Landrace, $\frac{1}{4}$ Duroc
579 and $\frac{1}{4}$ Yorkshire) currently being sequenced. Beyond this haplotype resolved assemblies for key

580 genotypes in the leading pig breeding company nucleus populations and of miniature pig lines used in
581 biomedical research can be anticipated in the next 5 years. Unfortunately, some of these genomes
582 may not be released into the public domain. The first wave of results from the Functional Annotation
583 of ANimal Genomes (FAANG) initiative [58, 59], are emerging and will add to the richness of pig
584 genome annotation.

585

586 In conclusion, the new pig reference genome (Sscrofa11.1) described here represents a significantly
587 enhanced resource for genetics and genomics research and applications for a species of importance
588 to agriculture and biomedical research.

589

590 **Methods**

591 Additional detailed methods and information on the assemblies and annotation are included in the
592 Supplementary Materials.

593

594 Preparation of genomic DNA

595 DNA was extracted from Duroc 2-14 cultured fibroblast cells passage 16-18 using the Qiagen Blood &
596 Cell Culture DNA Maxi Kit. DNA was isolated from lung tissue from barrow MARC1423004 using a salt
597 extraction method.

598

599 Genome sequencing and assembly

600 Genomic DNAs from the samples described above were used to prepare libraries for sequencing on
601 Pacific Biosciences RS II sequencer (PacBio RS II Sequencing System, RRID:SCR_017988) [60]. For Duroc
602 2-14 DNA P6/C4 chemistry was used, whilst for MARC1423004 DNA a mix of P6/C4 and earlier P5/C3
603 chemistry was used.

604

605 Reads from the Duroc 2-14 DNA were assembled into contigs using the Falcon v0.4.0 assembly pipeline
606 (Falcon, RRID:SCR_016089) following the standard protocol [27]. Quiver v. 2.3.0 [61] was used to
607 correct the primary and alternative contigs. Only the primary pseudo-haplotype contigs were used in
608 the assembly. The reads from the MARC1423004 DNA were assembled into contigs using Celera
609 Assembler v8.3rc2 (Celera assembler, RRID:SCR_010750) [32]. The contigs were scaffolded as
610 described in the results section above.

611

612 Fluorescence *in situ* hybridisation

613 Metaphase preparations were fixed to slides and dehydrated through an ethanol series (2 mins each
614 in 2×SSC, 70%, 85% and 100% ethanol at RT). Probes were diluted in a formamide buffer (Cytocell)
615 with Porcine Hybloc (Insight Biotech) and applied to the metaphase preparations on a 37°C hotplate

616 before sealing with rubber cement. Probe and target DNA were simultaneously denatured for 2 mins
617 on a 75°C hotplate prior to hybridisation in a humidified chamber at 37°C for 16 h. Slides were washed
618 post hybridisation in 0.4x SSC at 72°C for 2 mins followed by 2x SSC/0.05% Tween 20 at RT for 30 secs,
619 and then counterstained using VECTASHIELD anti-fade medium with DAPI (Vector Labs). Images were
620 captured using an Olympus BX61 epifluorescence microscope with cooled CCD camera and
621 SmartCapture (Digital Scientific UK) system.

622

623 Analysis of repetitive sequences, including telomeres and centromeres

624 Repeats were identified using RepeatMasker (v.4.0.7, RRID:SCR_012954) [62] with a combined repeat
625 database including Dfam (v.20170127) [63] and RepBase (v.20170127) [64]. RepeatMasker was run
626 with “sensitive” (-s) setting using *sus scrofa* as the query species (-- species “*sus scrofa*”). Repeats
627 which showed greater than 40% sequence divergence or were shorter than 70% of the expected
628 sequence length were filtered out from subsequent analyses. The presence of potentially novel
629 repeats was assessed by RepeatMasker using the novel repeat library generated by RepeatModeler
630 (v.1.0.11, RRID:SCR_015027) [62].

631

632 Telomeres were identified by running Tandem Repeat Finder (TRF) [65] with default parameters apart
633 from Mismatch (5) and Minscore (40). The identified repeat sequences were then searched for the
634 occurrence of five identical, consecutive units of the TTAGGG vertebrate motif or its reverse
635 complement and total occurrences of this motif were counted within the tandem repeat. Regions
636 which contained at least 200 identical hexamer units, were >2kbp of length and had a hexamer density
637 of >0.5 were retained as potential telomeres.

638

639 Centromeres were predicted using the following strategy. First, the RepeatMasker output, both
640 default and novel, was searched for centromeric repeat occurrences. Second, the assemblies were
641 searched for known, experimentally verified, centromere specific repeats [66, 67] in the *Sscrofa11.1*

642 genome. Then the three sets of repeat annotations were merged together with BEDTools (BEDTools,
643 RRID:SCR_006646) [68] (median and mean length: 786 bp and 5775 bp, respectively) and putative
644 centromeric regions closer than 500 bp were collapsed into longer super-regions. Regions which were
645 >5kbp were retained as potential centromeric sites.

646

647 Long read RNA sequencing (Iso-Seq)

648 The following tissues were harvested from MARC1423004 at age 48 days: brain (BioSamples:
649 SAMN05952594), diaphragm (SAMN05952614), hypothalamus (SAMN05952595), liver
650 (SAMN05952612), small intestine (SAMN05952615), skeletal muscle – *longissimus dorsi*
651 (SAMN05952593), spleen (SAMN05952596), pituitary (SAMN05952626) and thymus
652 (SAMN05952613). Total RNA from each of these tissues was extracted using Trizol reagent
653 (ThermoFisher Scientific) and the provided protocol. Briefly, approximately 100 mg of tissue was
654 ground in a mortar and pestle cooled with liquid nitrogen, and the powder was transferred to a tube
655 with 1 ml of Trizol reagent added and mixed by vortexing. After 5 minutes at room temperature, 0.2 ml
656 of chloroform was added and the mixture was shaken for 15 seconds and left to stand another
657 3 minutes at room temperature. The tube was centrifuged at 12,000 x g for 15 minutes at 4°C. The
658 RNA was precipitated from the aqueous phase with 0.5 ml of isopropanol. The RNA was further
659 purified with extended DNase I digestion to remove potential DNA contamination. The RNA quality
660 was assessed with a Fragment Analyzer (Advanced Analytical Technologies Inc., IA). Only RNA samples
661 of RQN above 7.0 were used for library construction. PacBio Iso-Seq libraries were constructed per
662 the PacBio Iso-Seq protocol. Briefly, starting with 3 µg of total RNA, cDNA was synthesized by using
663 SMARTer PCR cDNA Synthesis Kit (Clontech, CA) according to the Iso-Seq protocol (Pacific Biosciences,
664 CA). Then the cDNA was amplified using KAPA HiFi DNA Polymerase (KAPA Biotechnologies) for 10 or
665 12 cycles followed by purification and size selection into 4 fractions: 0.8-2 kbp, 2-3 kbp, 3-5 kbp and
666 >5 kbp. The fragment size distribution was validated on a Fragment Analyzer (Advanced Analytical
667 Technologies Inc, IA) and quantified on a DS-11 FX fluorometer (DeNovix, DE). After a second round

668 of large scale PCR amplification and end repair, SMRT bell adapters were separately ligated to the
669 cDNA fragments. Each size fraction was sequenced on 4 or 5 SMRT Cells v3 using P6-C4 chemistry and
670 6 hour movies on a PacBio RS II sequencer (Pacific Bioscience, CA). Short read RNA-Seq libraries were
671 also prepared for all nine tissue using TruSeq stranded mRNA LT kits and supplied protocol (Illumina,
672 CA), and sequenced on a Illumina NextSeq500 platform using v2 sequencing chemistry to generate
673 2 x 75 bp paired-end reads.

674

675 The Read of Insert (ROI) were determined by using *consensustools.sh* in the SMRT-Analysis pipeline
676 v2.0, with reads which were shorter than 300 bp and whose predicted accuracy was lower than 75%
677 removed. Full-length, non-concatemer (FLNC) reads were identified by running the *classify.py*
678 command. The cDNA primer sequences as well as the poly(A) tails were trimmed prior to further
679 analysis. Paired-end Illumina RNA-Seq reads from each tissue sample were trimmed to remove the
680 adaptor sequences and low-quality bases using Trimmomatic (v0.32, RRID:SCR_011848) [69] with
681 explicit option settings: *ILLUMINACLIP:adapters.fa: 2:30:10:1:true LEADING:3 TRAILING:3*
682 *SLIDINGWINDOW: 4:20 LEADING:3 TRAILING:3 MINLEN:25*, and overlapping paired-end reads were
683 merged using the PEAR software (v0.9.6; PEAR, RRID:SCR_003776) [70]. Subsequently, the merged
684 and unmerged RNA-Seq reads from the same tissue samples were *in silico* normalized in a mode for
685 single-end reads by using a Trinity (v2.1.1, RRID:SCR_013048) [71] utility,
686 *insilico_read_normalization.pl*, with the following settings: *--max_cov 50 --max_pct_stdev 100 --*
687 *single*. Errors in the full-length, non-concatemer reads were corrected with the preprocessed RNA-Seq
688 reads from the same tissue samples by using *proofread* (v2.12; Proofread, RRID:SCR_017331) [72].
689 Untrimmed sequences with at least some regions of high accuracy in the *.trimmed.fq* files were
690 extracted based on sequence IDs in *.untrimmed.fq* files to balance off the contiguity and accuracy of
691 the final reads.

692

693 Short read RNA sequencing (RNA-Seq)

694 In addition to the Illumina short read RNA-seq data generated from MARC1423004 and used to correct
695 the Iso-Seq data (see above), Illumina short read RNA-seq data (PRJEB19386) were also generated
696 from a range of tissues from four juvenile Duroc pigs (two male, two female) and used for annotation
697 as described below. Extensive metadata with links to the protocols for sample collection and
698 processing are linked to the BioSample entries under the Study Accession PRJEB19386. The tissues
699 sampled are listed in Table S10. Sequencing libraries were prepared using a ribodepletion TruSeq
700 stranded RNA protocol and 150 bp paired end sequences generated on the Illumina HiSeq 2500
701 platform (Illumina HiSeq 2500 System, RRID:SCR_016383) in rapid mode.

702

703 Annotation

704 The assembled genomes were annotated using the Ensembl pipelines (Ensembl, RRID:SCR_002344)
705 [38] as detailed in the Supplementary materials. The Iso-Seq and RNA-Seq data described above were
706 used to build gene models.

707

708 Mapping SNP chip probes

709 The probes from four commercial SNP chips were mapped to the Sscrofa10.2, Sscrofa11.1 and
710 USMARCv1.0 assemblies using BWA MEM [73] and a wrapper script [74]. Probe sequence was derived
711 from the marker manifest files that are available on the provider websites: Illumina PorcineSNP60 [1,
712 75]; Affymetrix Axiom™ Porcine Genotyping Array [76]; Gene Seek Genomic Profiler Porcine – HD
713 beadChip [77]; Gene Seek Genomic Profiler Porcine v2– LD Chip [77]. In order to retain marker
714 manifest coordinate information, each probe marker name was annotated with the chromosome and
715 position of the marker's variant site from the manifest file. All mapping coordinates were tabulated
716 into a single file, and were sorted by the chromosome and position of the manifest marker site. In
717 order to derive and compare relative marker rank order, a custom Perl script [78] was used to sort and
718 number markers based on their mapping locations in each assembly.

719

720 **Availability of supporting data**

721 The genome assemblies are deposited at NCBI under accession numbers GCA_000003025
722 (Sscrofa11.1) and GCA_002844635.1 (USMARCv1.0). The associated BioSample accession numbers are
723 SAMN02953785 and SAMN07325927, respectively. Iso-seq and RNA-Seq data used for analysis and
724 annotation are available under accession numbers PRJNA351265 and PRJEB19386, respectively.
725 Supporting data and materials are available in the *GigaScience* GigaDB database [100].

726

727 **List of abbreviations**

728 BAC: bacterial artificial chromosome; BLAST: Basic Local Alignment Search Tool; BLASTN: BLAST search
729 of nucleotide database(s); bp: base pairs; BUSCO: Benchmarking Universal Single-Copy Orthologs;
730 BWA: Burrows-Wheeler Aligner; CCD: charged couple device; cDNA: complementary DNA; CNV: copy
731 number variation; DAPI: 4',6-diamidino-2-phenylindole; DNA: Deoxyribonucleic Acid; DNase:
732 deoxyribonuclease; ENA: European Nucleotide Archive; FAANG: Functional Annotation of Animal
733 Genomes; FISH: fluorescent *in situ* hybridisation; FLNC: full-length, non-concatemer; FPC: fingerprint
734 contig; g: relative centrifugal force; Gbp: gigabase pairs; GFF3: genomic feature format, version 3; GC:
735 guanine-cytosine; h: hours; HQ: high quality; ID: identity; Iso-Seq: long read RNA sequencing using
736 PacBio technology; kbp: kilobase pairs; LC: low coverage; LQ: low quality; LQLC: low quality, low
737 coverage; MANE: Matched Annotation from NCBI and EMBL-EBI; Mbp: megabase pairs; mg:
738 milligrams; µg: micrograms; ml: millilitres; mins: minutes; mRNA: messenger RNA; NCBI: National
739 Center for Biotechnology Information; PacBio: Pacific Biosciences; polyA: poly adenine; PCR:
740 polymerase chain reaction; QC: quality control; RefSeq: NCBI Reference Sequence Database; RH:
741 radiation hybrid; RNA: Ribonucleic Acid; RNA-Seq: high-throughput short-read RNA sequencing; ROI:
742 read of interest; RQN: RNA quality number; RT: room temperature; secs: seconds; SGSC: Swine
743 Genome Sequencing Consortium; SINE: Short interspersed nuclear element; SMRT: single-molecule
744 real-time; SNP: single nucleotide **polymorphism**; SSC: saline sodium citrate; SSCn: Sus scrofa
745 chromosome n; Stdev: standard deviation; TRF: Tandem Repeat Finder.

746 **Consent for publication**

747 Not applicable

748 **Competing interests**

749 R.H. K.K. and E.T. are employed by Pacific Biosciences; all other authors declare that they have no
750 competing interests.

751 **Funding**

752 This work was supported by the Biotechnology and Biological Sciences Research Council, Institute
753 Strategic Programme Grant, BBS/E/D/20211550, Alan L Archibald, Mick Watson; the Biotechnology
754 and Biological Sciences Research Council, Institute Strategic Programme Grant, BBS/E/D/10002070,
755 Alan L. Archibald, Mick Watson; the Biotechnology and Biological Sciences Research Council, Response
756 Mode Grant, BB/F021372/1, Nabeel Affara; the Biotechnology and Biological Sciences Research
757 Council, Response Mode Grant, BB/M011461/1, Alan L Archibald; the Biotechnology and Biological
758 Sciences Research Council, Response Mode Grant, BB/M011615/1, Paul Flicek; the Biotechnology and
759 Biological Sciences Research Council, Response Mode Grant, BB/M01844X/1, Alan L Archibald, Mick
760 Watson; EU, FP7 Programme Quantomics, KBBE222664, Alan L Archibald; Wellcome Trust,
761 WT108749/Z/15/Z, Paul Flicek; USDA, CRIS Project, 8042-31000-001-00-D, Derek M Bickhart and
762 Benjamin D Rosen; USDA, CRIS Project, 5090-31000-026-00-D, Derek M Bickhart; USDA, CRIS Project,
763 3040-31000-100-00-D, Timothy P L Smith.

764

765 **Authors' contributions**

766 A.L.A. and T.P.L.S. conceived, coordinated and managed the project; A.L.A., P.F., D.A.H., T.P.L.S. M.W.
767 supervised staff and students performing the analyses; D.J.N., L.R., L.B.S., T.P.L.S. provided biological
768 resources; R.H., K.S.K. and T.P.L.S. generated PacBio sequence data; H.A.F., T.P.L.S. and R.T. generated
769 Illumina WGS and RNA-Seq data; N.A.A., C.A.S., B.M.S. provided SSCY assemblies; D.J.N, and T.P.L.S.
770 generated Iso-Seq data; G.H., R.H., S.K., A.M.P., A.S.S, A.W. generated sequence assemblies; A.W.
771 polished and quality checked Sscrofa11.1; W.C., G.H., K.H., S.K., B.D.R., A.S.S., S.G.S., E.T. performed

772 quality checks on the sequence assemblies; R.E.O'C. and D.K.G. performed cytogenetics analyses; L.E.
773 analysed repeat sequences; H.B., H.L., N.M., C.K.T. analysed Iso-Seq data; D.M.B. and G.A.R. analysed
774 sequence variants; B.A., K.B., C.G.G., T.H., O.I., F.J.M. annotated the assembled genome sequences;
775 A.W. and A.L.A drafted the manuscript; all authors read and approved the final manuscript.

776

777 **Acknowledgments**

778 We are grateful to Chris Tyler-Smith (Wellcome Trust Sanger Institute) for sharing the SSCY sequence
779 data for Sscrofa11.1. In addition to the funding acknowledged above we are grateful for support from
780 the University of Cambridge, Department of Pathology, the European Molecular Biology Laboratory,
781 and the Roslin Foundation. In addition HL and HB were supported by USDA NRSP-8 Swine Genome
782 Coordination funding; SK and AMP were supported by the Intramural Research Program of the
783 National Human Genome Research Institute, US National Institutes of Health; This work used the
784 computational resources of the NIH HPC Biowulf cluster (<https://hpc.nih.gov>); and the Iowa State
785 University Lightning3 and ResearchIT clusters. The Ceres cluster (part of the USDA SCInet Initiative)
786 was used to analyse part of this dataset.

787

788

789 **Supplementary materials**

790 Supplementary materials for this article include:

791 Supplementary Methods and Information

792 Table S1. Pacific Biosciences read statistics.

793 Table S2. Predicted telomeres.

794 Table S3. Predicted centromeres.

795 Table S4. Assigning scaffolds to chromosomes.

796 Table S5 Alignment of Radiation Hybrid maps and genome assemblies.

797 Table S6. Assemblytics comparisons, assembly statistics.

798 Table S7. BUSCO results.

799 Table S8. Annotation statistics. (Ensembl-NCBI comparison)

800 Table S9. Commercial SNP chip probes.

801 Table S10. Tissue samples.

802 Table S11. Ensembl annotation statistics for 13 pig genome assemblies

803 Figure S1. Predicted telomeres.

804 Figure S2. Predicted centromeres.

805 Figure S3. Fluorescent *in situ* hybridisation assignments.

806 Fig. S4. Improvement in local order and orientation and reduction in redundancy.

807 Fig. S5. Assembly comparisons in gEVAL (SSC15).

808 Fig. S6. Assembly comparisons in gEVAL (SSC5).

809 Fig. S7. Assembly comparisons in gEVAL (SSC18).

810 Fig. S8. Order and orientation of SSC18 assemblies.

811 Fig. S9. Order and orientation of SSC7 assemblies.

812 Fig. S10. Order and orientation of SSC8 assemblies.

813 Fig. S11. Assembly alignments.

814 Figure S12. Assemblytics results.

- 815 Figure S13. Counts of repetitive elements in four pig assemblies.
- 816 Figure S14. Average mapped length of repetitive elements in four pig genomes.
- 817 Figure S15. Assembly SNP rank concordance versus reported chromosomal location.
- 818

819 **References**

- 820 1. Ramos AM, Crooijmans RPMA, Affara NA, Amaral AJ, Archibald AL, Beever JE, et al. Design of a
821 high density SNP genotyping assay in the pig using SNPs identified and characterized by next
822 generation sequencing technology. *PLoS One* 2009;**4**:e6524.
- 823 2. Hu ZL, Park CA, ReecyJM, Developmental progress and current status of the Animal QTLdb. *Nucleic*
824 *Acids Res* 2016;**44**:D827–33.
- 825 3. Meuwissen T, Hayes, B, Goddard M, Accelerating Improvement of Livestock with Genomic
826 Selection. *Annu Rev Anim Biosci* 2013;**1**:221–37.
- 827 4. Christensen OF, Madsen P, Nielsen B, Ostersen T, Su G, Single-step methods for genomic
828 evaluation in pigs. *Animal* 2012;**6**:1565–71.
- 829 5. Cleveland M, Hickey JM, Practical implementation of cost-effective genomic selection in
830 commercial pig breeding using imputation. *J Anim Sci* 2013;**91**:3583–92.
- 831 6. Vamathevan JJ, Hall MD, Hasan S, Woollard PM, Xu M, Yang Y, et al. Minipig and beagle animal
832 model genomes aid species selection in pharmaceutical discovery and development. *Toxicol Appl*
833 *Pharmacol* 2013;**270**:149–57.
- 834 7. Klymiuk N, Seeliger F, Bohlooly M, Blutke A, Rudmann DG, Wolf E, Tailored Pig Models for
835 Preclinical Efficacy and Safety Testing of Targeted Therapies. *Toxicol Pathol* 2016;**44**:346–57.
- 836 8. Wells KD, Prather RS, Genome-editing technologies to improve research, reproduction, and
837 production in pigs. *Mol Reprod Dev* 2017;**84**:1012–7.
- 838 9. Servin B, Faraut T, Iannuccelli N, Zelenika D, Milan D, High-resolution autosomal radiation hybrid
839 maps of the pig genome and their contribution to the genome sequence assembly. *BMC Genomics*
840 2012;**13**:585.
- 841 10. Tortereau F, Servin B, Frantz L, Megens HJ, Milan D, Rohrer G, et al. A high density recombination
842 map of the pig reveals a correlation between sex-specific recombination and GC content. *BMC*
843 *Genomics* 2012;**13**:586.

- 844 11. Yerle M, Lahbib-Mansais Y, Mellink C, Goureau A, Pinton P, Echard G, et al. The PiGMaP
845 consortium cytogenetic map of the domestic pig (*Sus scrofa domestica*). *Mamm Genome*
846 1995;**6**:176–86.
- 847 12. Humphray SJ, Scott CE, Clark R, Marron B, Bender C, Camm N, et al. A high utility integrated map
848 of the pig genome. *Genome Biol* 2017;**8**:R139.
- 849 13. Groenen MAM, Archibald AL, Uenishi H, Tuggle CK, Takeuchi Y, Rothschild MF, et al. Analyses of
850 pig genomes provide insight into porcine demography and evolution. *Nature* 2012;**491**:393–8.
- 851 14. Warr A, Robert C, Hume D, Archibald AL, Deeb N, Watson M. Identification of Low-Confidence
852 Regions in the Pig Reference Genome (*Sscrofa* 10.2). *Front Genet* 2015;**6**:338.
- 853 15. O'Connor RE, Fonseka G, Frodsham R, Archibald AL, Lawrie M, Walling GA et al. Isolation of
854 subtelomeric sequences of porcine chromosomes for translocation screening reveals errors in the
855 pig genome assembly. *Anim Genet* 2017;**48**:395–403.
- 856 16. Dawson HD, Chen C, Gaynor B, Shao J, Urban Jr. JF. The porcine translational research database:
857 a manually curated, genomics and proteomics-based research resource. *BMC Genomics*
858 2017;**18**:643.
- 859 17. Li M, Chen L, Tian S, Lin Y, Tang Q, Zhou X, et al. Comprehensive variation discovery and recovery
860 of missing sequence in the pig genome using multiple de novo assemblies. *Genome Res*
861 2017;**27**:865–74.
- 862 18. Schook LB, Beever JE, Rogers J, Humphray S, Archibald A, Chardon P, et al. Swine Genome
863 Sequencing Consortium (SGSC): A strategic roadmap for sequencing the pig genome. *Comp Funct*
864 *Genomics* 2005;**6**:251–5.
- 865 19. Robert C, Fuentes-Utrilla P, Troup K, Loecherbach J, Turner F, Talbot R, et al. Design and
866 development of exome capture sequencing for the domestic pig (*Sus scrofa*). *BMC Genomics*
867 2014;**15**:550.
- 868 20. Skinner BM, Sargent CA, Churcher C, Hunt T, Herrero J, Loveland JE, et al. The pig X and Y
869 Chromosomes: structure, sequence, and evolution. *Genome Res* 2016;**26**:130–9.

- 870 21. Frantz LAF, Schraiber JG, Madsen O, Megens HJ, Cagan A, Bosse M, et al. Evidence of long-term
871 gene flow and selection during domestication from analyses of Eurasian wild and domestic pig
872 genomes. *Nat Genet* 2015;**47**:1141-8.
- 873 22. Groenen MAM. A decade of pig genome sequencing: a window on pig domestication and
874 evolution. *Genet Sel Evol* 2016;**48**:23.
- 875 23. van Dijk EL, Jaszczyszyn Y, Naquin D, Thermes C. The Third Revolution in Sequencing Technology.
876 *Trends Genet* 2018;**34**:666-81.
- 877 24. Koren S, Rhie A, Walenz BP, Dilthey AT, Bickhart DM, Kingan SB, et al. De novo assembly of
878 haplotype-resolved genomes with trio binning. *Nat. Biotechnol* 2018;**36**:1174-82.
- 879 25. CHORI-242: Porcine (*Sus scrofa*) BAC Library <https://bacpacresources.org/library.php?id=124>
880 Accessed 17 Apr 2020.
- 881 26. Putnam NH, O'Connell BL, Stites JC, Rice BJ, Blanchette M, Calef R, et al. Chromosome-scale
882 shotgun assembly using an in vitro method for long-range linkage. *Genome Res* 2016;**26**:342–50.
- 883 27. Chin CS, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, et al. Phased diploid genome
884 assembly with single-molecule real-time sequencing. *Nat Methods* 2016;**13**:1050–4.
- 885 28. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile and open
886 software for comparing large genomes. *Genome Biol* 2004;**5**:R12.
- 887 29. English AC, Richards S, Han Y, Wang M, Lee V, Qu J, et al. Mind the Gap: Upgrading Genomes with
888 Pacific Biosciences RS Long-Read Sequencing Technology. *PLoS One* 2012;**7**:e47768.
- 889 30. Altschul SF, Gish W, Miller W, et al. Basic Local Alignment Search Tool. *J Mol Biol* 1990;**215**:403–
890 10.
- 891 31. Chow W, Brugger K, Caccamo M, Sealy I, Torrance J, Howe K. gEVAL-a web-based browser for
892 evaluating genome assemblies. *Bioinformatics* 2016;**32**:2508–10.
- 893 32. Berlin K, Koren S, Chin CS, Drake JP, Landolin JM, Phillippy AM. Assembling large genomes with
894 single-molecule sequencing and locality-sensitive hashing. *Nat Biotechnol* 2015;**33**:623–30.

- 895 33. Nattestad M, Schatz MC. Assemblytics: A web analytics tool for the detection of variants from an
896 assembly. *Bioinformatics* 2016;**32**:3021-3.
- 897 34. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: Assessing genome
898 assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 2015;**31**:3210–
899 2.
- 900 35. Tseng E. Cogent: COding GENome reconstruction Tool. 2017 <https://github.com/Magdoll/Cogent>
- 901 36. Beiki H, Liu H, Manchanda N, Nonneman D, Smith TPL, Reecy JM et al. Improved annotation of the
902 domestic pig genome through integration of Iso-Seq and RNA-seq data. *BMC Genomics*
903 2019;**20**:344.
- 904 37. Long Y, Su Y, Ai H, Zhang Z, Yang B, Ruan G, et al. A genome-wide association study of copy number
905 variations with umbilical hernia in swine. *Anim Genet* 2016;**47**:298–305.
- 906 38. Cunningham F, Achuthan P, Akanni W, Allen J, Amode MR, Armean IM, et al. Ensembl 2019.
907 *Nucleic Acids Res* 2019;**47**(D1):D745-51.
- 908 39. Ensembl pig strains genome annotation Release 98, September 2019
909 http://www.ensembl.org/Sus_scrofa/Info/Strains
- 910 40. Meyers SN, Rogatcheva MB, Larkin DM, Yerle M, Milan D, Hawken RJ, et al. Piggy-BACing the
911 human genome: II. A high-resolution, physically anchored, comparative map of the porcine
912 autosomes. *Genomics* 2005;**86**:739-52.
- 913 41. Ensembl pig genome annotation Release 98, September 2019
914 http://www.ensembl.org/Sus_scrofa/Info/Index Accessed 15 October 2019.
- 915 42. Ensembl archive http://may2017.archive.ensembl.org/Sus_scrofa/Info/Index Accessed 15 April
916 2020.
- 917 43. NCBI annotation report 106
918 https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Sus_scrofa/106/ Accessed 15 April
919 2020.

- 920 44. MANE (Matched Annotation from NCBI and EMBL-EBI)
921 <https://www.ensembl.org/info/genome/genebuild/mane.html> Accessed 15 April 2020.
- 922 45. gEVAL: Genome Evaluation Browser. <https://geval.sanger.ac.uk/> Accessed 15 October 2019.
- 923 46. Utsunomiya ATH, Santos DJ, Boison SA, Utsunomiya YT, Milanese M, Bickhart DM, et al. Revealing
924 misassembled segments in the bovine reference genome by high resolution linkage disequilibrium
925 scan. *BMC Genomics* 2016;**17**:705.
- 926 47. Hickey JM. Sequencing millions of animals for genomic selection 2.0. *J Anim Breed Genet*
927 2013;**130**:331-2.
- 928 48. Le SQ, Durbin R. SNP detection and genotyping from low-coverage sequencing data on multiple
929 diploid samples. *Genome Res* 2011;**21**:952-60.
- 930 49. Li Y, Sidore C, Kang HM, Boehnke M, Abecasis GR. Low-coverage sequencing: Implications for
931 design of complex trait association studies. *Genome Res* 2011;**21**:940-51.
- 932 50. Vertebrate Genomes Project <https://vertebrategenomesproject.org/>
- 933 51. Daetwyler HD, Capitan A, Pausch H, Stotthard P, van Binsbergen R, Brøndum RF, et al. Whole-
934 genome sequencing of 234 bulls facilitates mapping of monogenic and complex traits in cattle.
935 *Nat Genet* 2014;**46**:858-65.
- 936 52. Lilue J, Doran AG, Fiddes IT, Abrudan M, Armstrong J, Bennett R, et al. Sixteen diverse laboratory
937 mouse reference genomes define strain-specific haplotypes and novel functional loci. *Nat Genet*
938 2018;**50**:1574-83.
- 939 53. Ensembl mouse strain annotation http://www.ensembl.org/Mus_musculus/Info/Strains
940 Accessed 15 April 2020.
- 941 54. Ensembl dog breed annotation http://www.ensembl.org/Canis_familiaris/Info/Strains Accessed
942 15 April 2020.
- 943 55. Baier U, Beller T, Ohlebusch E. Graphical pan-genome analysis with compressed suffix trees and
944 the Burrows-Wheeler transform. *Bioinformatics* 2015;**32**:497–504.

- 945 56. Chaisson MJP, Wilson RK, Eichler EE. Genetic variation and the de novo assembly of human
946 genomes. *Nat Rev Genet* 2015;**16**:627–40.
- 947 57. Garrison E, Sirén J, Novak AM, Hickey G, Eizenga JM, Dawson ET, et al. Variation graph toolkit
948 improves read mapping by representing genetic variation in the reference. *Nat Biotechnol*
949 2018;**36**:875-9.
- 950 58. Andersson L, Archibald AL, Bottema CD, Brauning R, Burgess SC, Burt DW, et al. Coordinated
951 international action to accelerate genome-to-phenome with FAANG, the Functional Annotation
952 of Animal Genomes project. *Genome Biol* 2015;**16**:57.
- 953 59. Foissac S, Djebali S, Munyard K, Vialaneix N, Rau A, Muret K, et al. Multispecies annotation of
954 transcriptome and chromatin structure in domesticated animals. *BMC Biol* 2019;**17**:108
- 955 60. Pendleton M, Sebra R, Pang AA, Ummat A, Franzen O, Rausch T, et al. Assembly and diploid
956 architecture of an individual human genome via single-molecule technologies. *Nat Methods*
957 2015;**12**:780–6.
- 958 61. Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, et al. Nonhybrid, finished
959 microbial genome assemblies from long-read SMRT sequencing data. *Nat Methods* 2013;**10**:563-
960 9.
- 961 62. RepeatMasker <http://www.repeatmasker.org> Accessed 17 April 2020
- 962 63. Hubley R, Finn RD, Clements J, Eddy SR, Jones TA, Bao W, et al. The Dfam database of repetitive
963 DNA families. *Nucleic Acids Res* 2016;**44**:D81–9.
- 964 64. Bao W, Kojima KK, Kohany O. Repbase Update, a database of repetitive elements in eukaryotic
965 genomes. *Mob DNA* **6**(1). doi: 10.1186/s13100-015-0041-9. (2015).
- 966 65. Benson G. Tandem repeats finder: A program to analyze DNA sequences. *Nucleic Acids Res*
967 1999;**27**:573–80.
- 968 66. Miller JR, Hindkjær J, Thomsen PD. A chromosomal basis for the differential organization of a
969 porcine centromere-specific repeat. *Cytogenet Cell Genet* 1993;**62**:37–41.

- 970 67. Riquet J, Mulsant P, Yerle M, Cristobal-Gaudy MS, Le Tissier P, Milan D, et al. Sequence analysis
971 and genetic mapping of porcine chromosome 11 centromeric S0048 marker. *Cytogenet Cell Genet*
972 1996;**74**:127-32.
- 973 68. Quinlan AR, Hall IM. BEDTools: A flexible suite of utilities for comparing genomic features,
974 *Bioinformatics* 2010;**26**:841–2.
- 975 69. Bolger AM, Lohse M, Usadel B. Trimmomatic: A flexible trimmer for Illumina sequence data.
976 *Bioinformatics* 2014;**30**:2114–20.
- 977 70. Zhang J, Kobert K, Flouri T, Stamatakis A. PEAR: A fast and accurate Illumina Paired-End reAd
978 mergeR. *Bioinformatics* 2014;**30**:614–20.
- 979 71. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Full-length transcriptome
980 assembly from RNA-Seq data without a reference genome. *Nat Biotechnol* 2011;**29**:644–52.
- 981 72. Hackl T, Hedrich R, Schultz J, Förster F. Proovread: Large-scale high-accuracy PacBio correction
982 through iterative short read consensus. *Bioinformatics* 2014;**30**:3004–11.
- 983 73. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.
984 *Bioinformatics* 2009;**25**:1754-60.
- 985 74. Bickhart D. 2019. Align and order SNP probes.
986 https://github.com/njdbickhart/perl_toolchain/blob/master/assembly_scripts/alignAndOrderSn
987 [pProbes.pl](https://github.com/njdbickhart/perl_toolchain/blob/master/assembly_scripts/alignAndOrderSn)
- 988 75. Illumina PorcineSNP60 Beadchip [https://emea.illumina.com/products/by-type/microarray-](https://emea.illumina.com/products/by-type/microarray-kits/porcine-snp60.html)
989 [kits/porcine-snp60.html](https://emea.illumina.com/products/by-type/microarray-kits/porcine-snp60.html)
- 990 76. Axiom™ Porcine Genotyping Array
991 <https://www.thermofisher.com/order/catalog/product/550588#/550588>
- 992 77. Gene Seek Genomic Profiler Porcine <https://genomics.neogen.com/uk/ggp-porcine>
- 993 78. Bickhart D. 2019. Pig genome SNP sort rank order.
994 https://github.com/njdbickhart/perl_toolchain/blob/master/assembly_scripts/pigGenomeSNPS
995 [ortRankOrder.pl](https://github.com/njdbickhart/perl_toolchain/blob/master/assembly_scripts/pigGenomeSNPS)

- 996 79. Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, et al. Integrative
997 genomics viewer. *Nat Biotechnol* 2011;**29**:24–6.
- 998 80. Koren S, Walenz BP, Berlin K, JMILLER JR, Bergman NH, Phillippy AM. Canu: Scalable and accurate
999 long-read assembly via adaptive κ -mer weighting and repeat separation. *Genome Res*
1000 2017;**27**:722–36.
- 1001 81. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
1002 *ArXiv:1303.3997v1 [q-bio.GN]* (2013).
- 1003 82. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: An integrated tool
1004 for comprehensive microbial variant detection and genome assembly improvement. *PLoS ONE*
1005 2014;**9**:e112963.
- 1006 83. Anderson SI, Lopez-Corrales NL, Gorick B, Archibald AL. A large-fragment porcine genomic library
1007 resource in a BAC vector. *Mamm Genome* 2000;**11**:811–4.
- 1008 84. Morgulis A, Gertz EM, Schäffer AA, Agarwala R. A Fast and Symmetric DUST Implementation to
1009 Mask Low-Complexity DNA Sequences. *J Comp Biol* 2006;**13**:1028-40.
- 1010 85. Down TA, Hubbard TJP. Computational detection and location of transcription start sites in
1011 mammalian genomic DNA. *Genome Res* 2002;**12**:458–61.
- 1012 86. Lowe TM, Eddy SR. TRNAscan-SE: A program for improved detection of transfer RNA genes in
1013 genomic sequence. *Nucleic Acids Res* 1996;**25**:955–64.
- 1014 87. Burge C, Karlin S. Prediction of complete gene structures in human genomic DNA. *J Mol Biol*
1015 1997;**268**:78–94.
- 1016 88. O’Leary NA, Wright MW, Brister JR, Ciufu S, Haddad D, McVeigh R, et al. Reference sequence
1017 (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation.
1018 *Nucleic Acids Res* 2016;**44**:D733-45.
- 1019 89. Slater GSC, Birney E. Automated generation of heuristics for biological sequence comparison. *BMC*
1020 *Bioinformatics* 2005;**6**:31.

- 1021 90. Frankish A, Diekhans M, Ferreira AM, Johnson R, Jungreis I, Loveland J, et al. GENCODE reference
1022 annotation for the human and mouse genomes. *Nucleic Acids Res* 2019;**47**:D766-73.
- 1023 91. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 2018;**34**:3094-3100.
- 1024 92. She R, Chu JS, Uyar B, Wang J, Wang K, Chen N. genBlastG: Using BLAST searches to build
1025 homologous gene models. *Bioinformatics* 2011;**27**:2141-3.
- 1026 93. Lefranc MP, Giudicelli V, Duroux P, Jabado-Michaloud J, Folch G, Aouinti S, et al. IMGT®, the
1027 international ImMunoGeneTics information system® 25 years on. *Nucleic Acids Res*
1028 2015;**43**:D413-22.
- 1029 94. Lewis SE, Searle SMJ, Harris N, Gibson M, Iyer V, Richter J, et al. Apollo: a sequence annotation
1030 editor. *Genome Biol* 2002;**3**:research0082.1.
- 1031 95. Griffiths-Jones S, Bateman A, Marshall M, Khanna A, Eddy SR. Rfam: An RNA family database',
1032 *Nucleic Acids Res* 2003;**31**:439-41.
- 1033 96. Griffiths-Jones S, Grocock RJ, van Dongen S, Bateman A, Enright AJ. miRBase: microRNA
1034 sequences, targets and gene nomenclature. *Nucleic Acids Res* 2006;**34**:D140-4.
- 1035 97. Lorenz R, Bernhart SH, Höner Zu Siederdisen C, Tafer H, Flamm C, Stadler PF, et al. ViennaRNA
1036 Package 2.0. *Algorithms Mol Biol* 2011;**6**:26.
- 1037 98. Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics*
1038 2013;**29**:2933-5.
- 1039 99. Papatheodorou I, Fonesca NA, Keays M, Tang YA, Barrera E, Bazant W, et al. Expression Atlas:
1040 Gene and protein expression across multiple studies and organisms. *Nucleic Acids Res*
1041 2018;**46**:D246-51.
- 1042 100. Warr A; Affara N; Aken B; Beiki H; Bickhart DM; Billis K; Chow W; Eory L; Finlayson HA; Flicek
1043 P; Girón CG; Griffin DK; Hall R; Hannum G; Hourlier T; Howe K; Hume DA; Izuogu O; Kim K; Koren
1044 S; Liu H; Manchanda N; Martin FJ; Nonneman DJ; O'Connor RE; Phillippy AM; Rohrer GA; Rosen
1045 BD; Rund LA; Sargent CA; Schook LB; Schroeder SG; Schwartz AS; Skinner BM; Talbot R; Tseng E;
1046 Tuggle CK; Watson M; Smith TPL; Archibald AL (2020): Supporting data for "An improved pig

1047 reference genome sequence to enable pig genetics and genomics research" *GigaScience*

1048 Database. <http://dx.doi.org/10.5524/100732>

1049

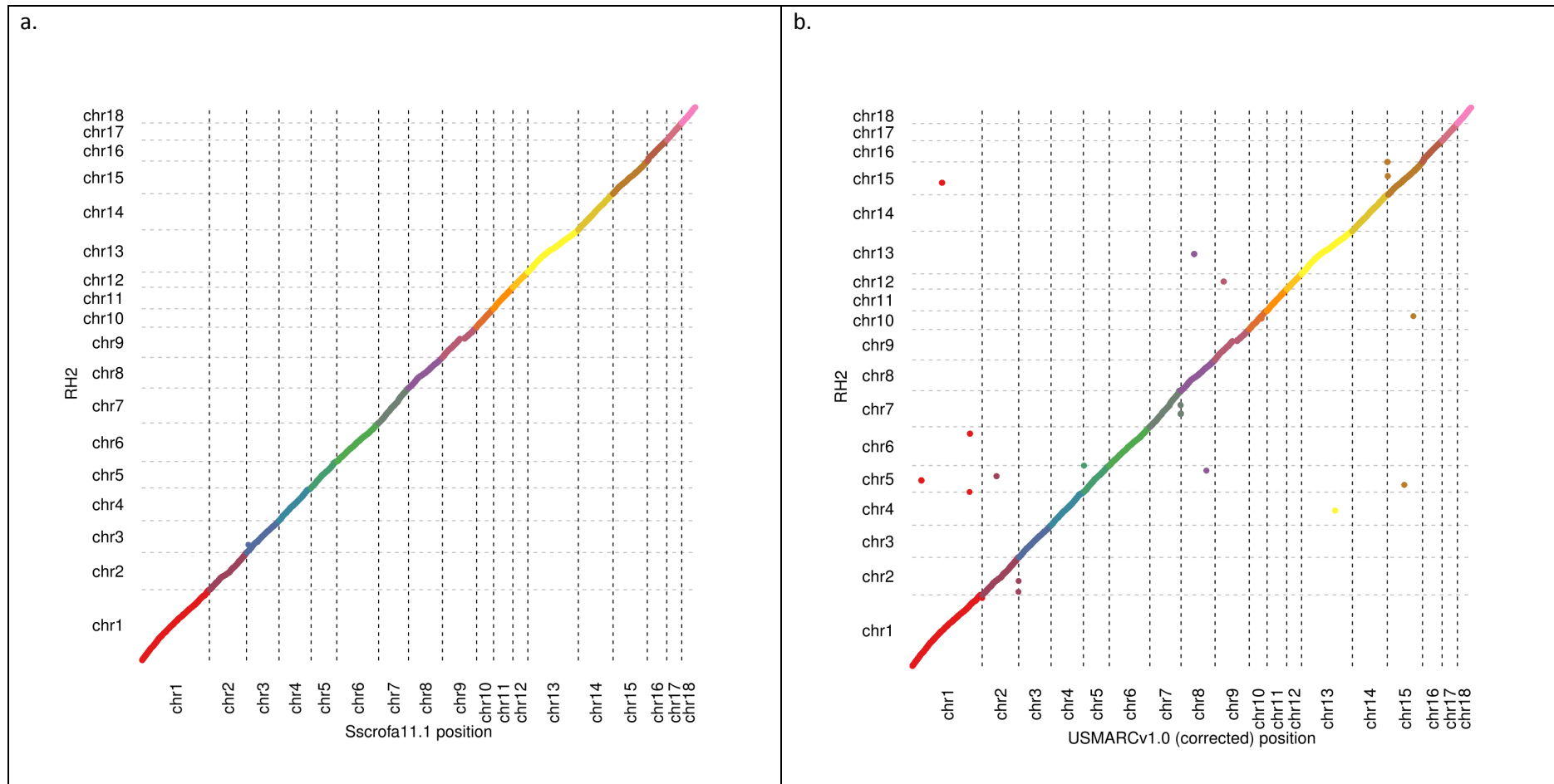
1050

1051

1052

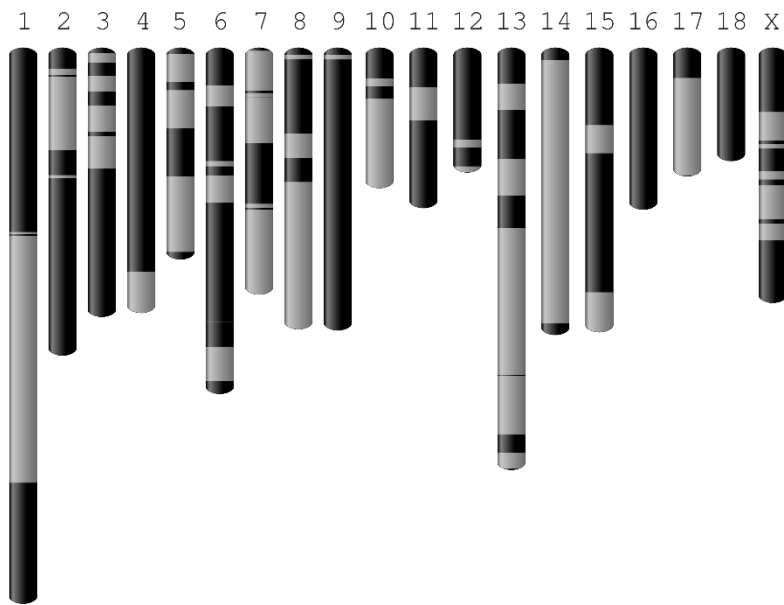
1053

1054 **Figure 1. Assemblies and radiation hybrid map alignments.** Plots illustrating co-linearity between radiation hybrid map and a) Sscrofa11.1 and b) USMARCv1.0
1055 assemblies (autosomes only).

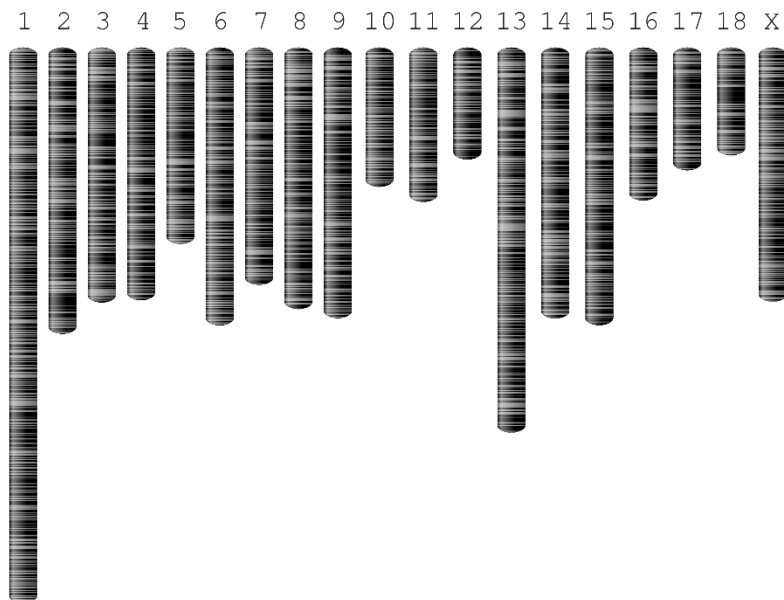


1056

1057 **Figure 2. Visualisation of improvements in assembly contiguity.** Graphical visualisation of contigs
1058 for Sscrofa11 (top) and Sscrofa10.2 (bottom) as alternating dark and light grey bars.



1059



1060

1061

1062

1063 **Table 1. Assembly statistics.** Summary statistics for assembled pig genome sequences and comparison with current human reference genome. (source: NCBI,

1064 <https://www.ncbi.nlm.nih.gov/assembly/>; * includes mitochondrial genome.

Assembly	Sscrofa10.2	Sscrofa11	Sscrofa11.1	USMARCv1.0	GRCh38.p12
Total sequence length	2,808,525,991	2,456,768,445	2,501,912,388	2,755,438,182	3,099,706,404
Total ungapped length	2,519,152,092	2,454,899,091	2,472,047,747	2,623,130,238	2,948,583,725
Number of scaffolds	9,906	626	706	14,157	472
Gaps between scaffolds	5,323	24	93	0	349
Number of unplaced scaffolds	4,562	583	583	14,136	126
Scaffold N50	576,008	88,231,837	88,231,837	131,458,098	67,794,873
Scaffold L50	1,303	9	9	9	16
Number of unspanned gaps	5,323	24	93	0	349
Number of spanned gaps	233,116	79	413	661	526
Number of contigs	243,021	705	1,118	14,818	998
Contig N50	69,503	48,231,277	48,231,277	6,372,407	57,879,411
Contig L50	8,632	15	15	104	18
Number of chromosomes*	*21	19	*21	*21	24

1065

1066

1067 **Table 2. Summary of quality statistics for SSC1-18, SSCX.** Quality measures and terms as defined [14].

1068

	Mean (Sscrofa11)	Std (Sscrofa11)	Bases (Sscrofa11)	% genome (Sscrofa11)	% genome (Sscrofa10.2)
High Coverage	50	7	119,341,205	4.9	2.6
Low Coverage (LC)	50	7	185,385,536	7.5	26.6
% Properly paired	86	6.8	95,508,007	3.9	4.95
% High inserts	0.3	1.6	40,835,320	1.72	1.52
% Low inserts	8.2	4.3	114,793,298	4.7	3.99
Low quality (LQ)	-	-	284,838,040	11.6	13.85
Total LQLC	-	-	399,927,747	16.3	33.07
LQLC windows that do not intersect RepeatMasker regions			39,918,551	1.6	

1069

1070

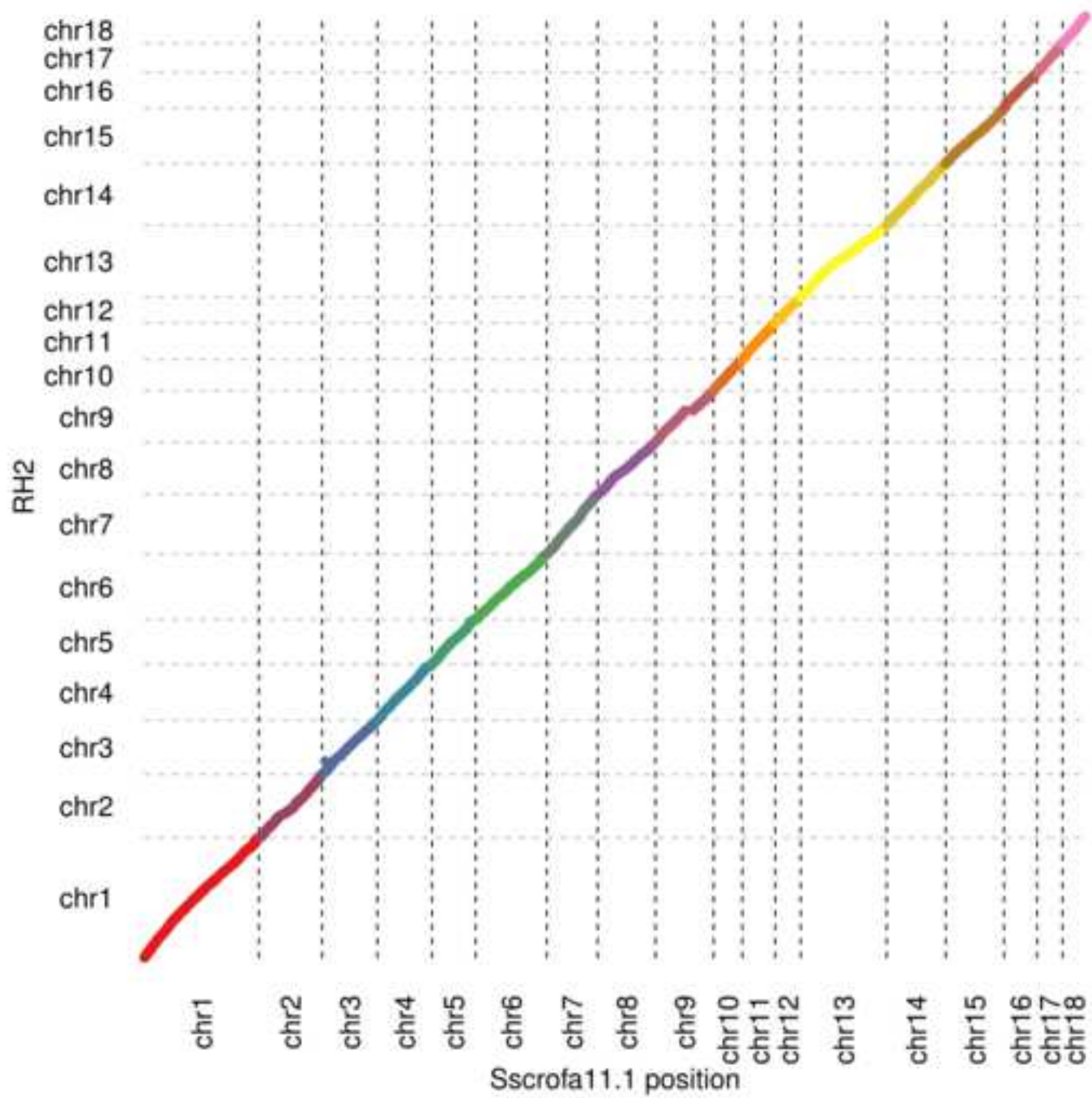
1071 **Table 3. Annotation statistics.** Ensembl annotation of pig (Sscrofa10.2, Sscrofa11.1, USMARCv1.0), human (GRCh38.p12) and mouse (GRCm38.p6)
 1072 assemblies.

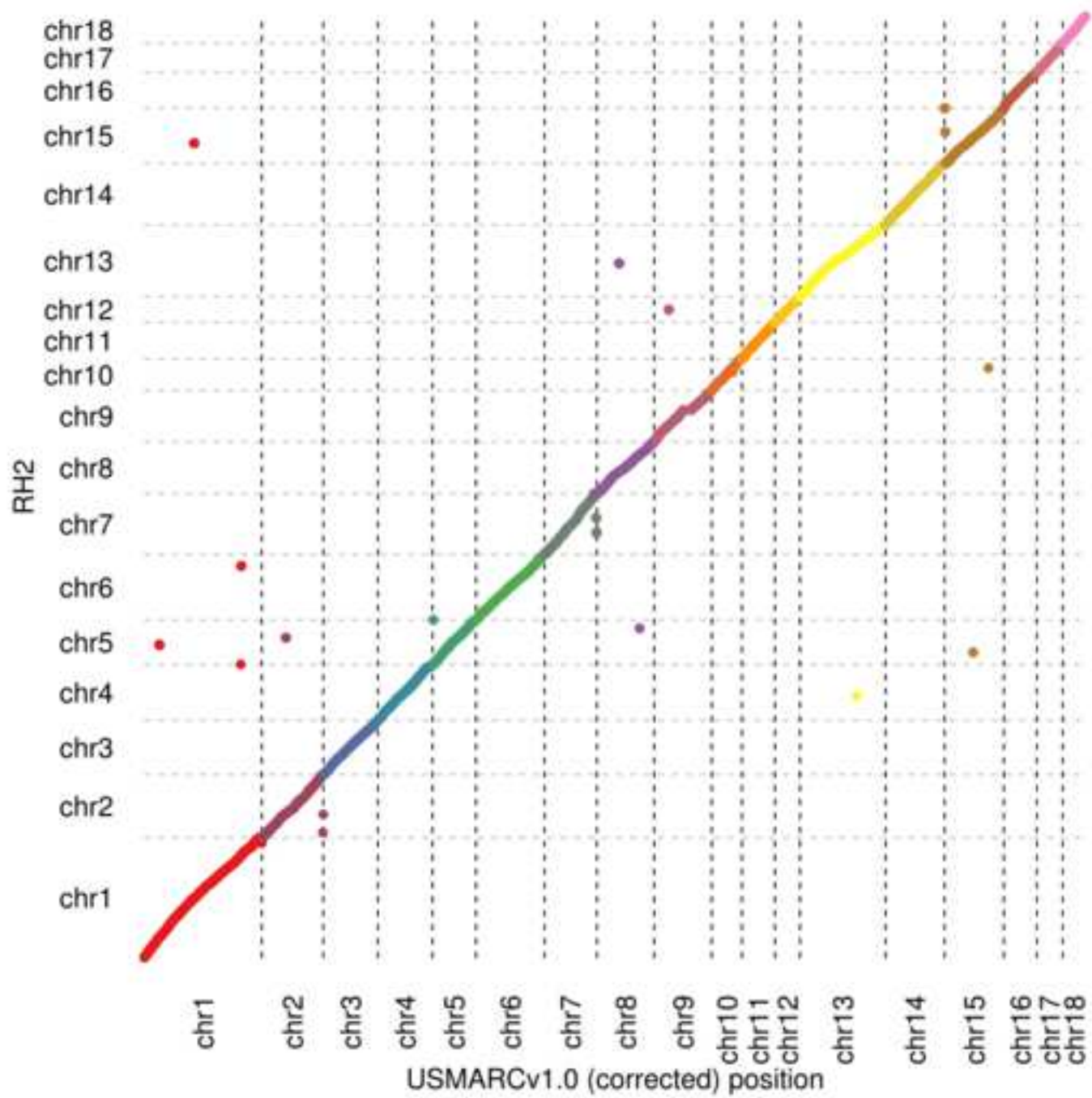
1073

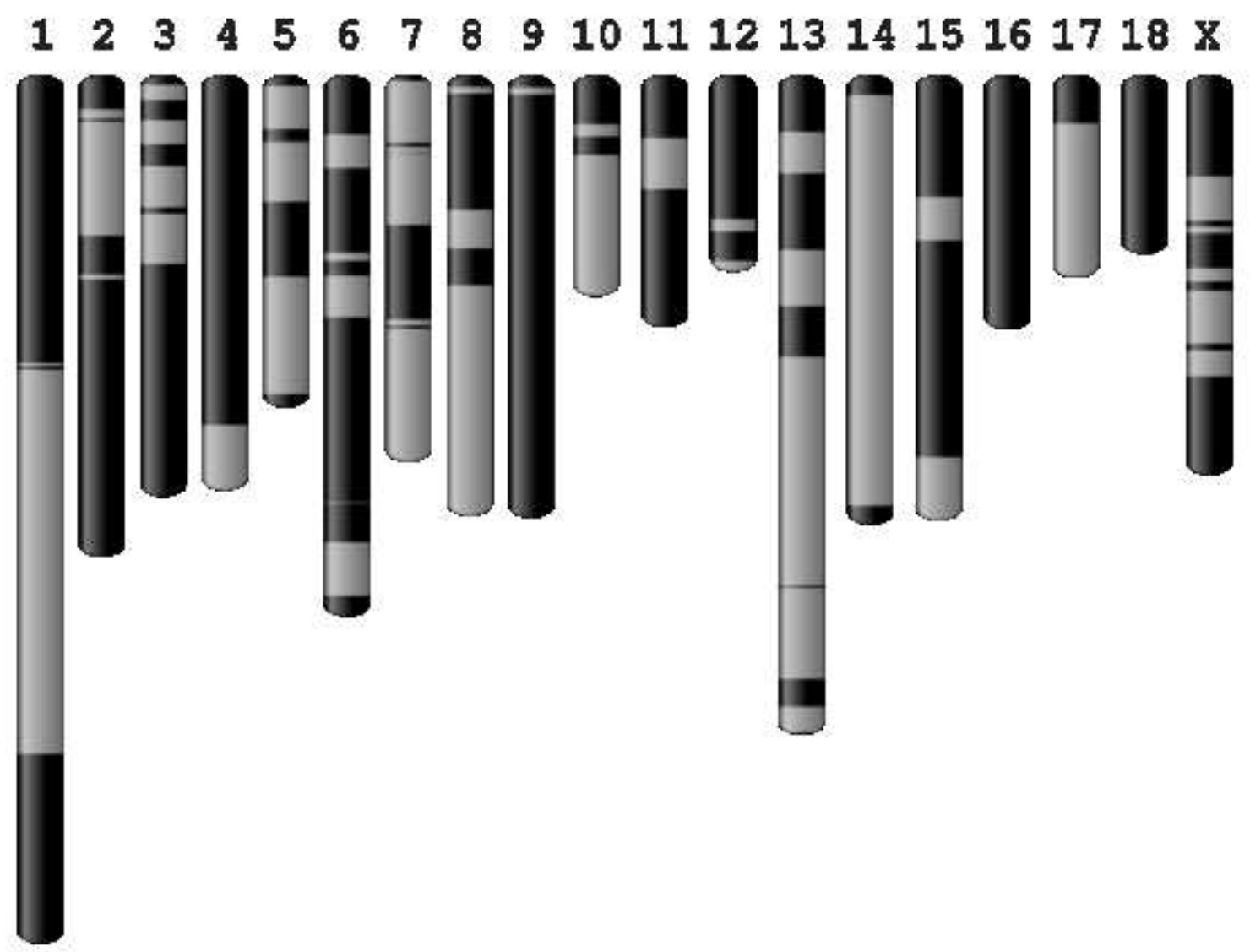
	Sscrofa10.2	Sscrofa11.1	USMARCv1.0	GRCh38.p13	GRCm38.p6
	Ensembl (Release 89)	Ensembl (Release 98)	Ensembl (Release 97)	Ensembl (Release 98)	Ensembl (Release 98)
Coding genes	21,630 (Incl. 10 read through)	21,301	21,535	20,444 incl 667 read through	22,508 incl 270 read through
Non-coding genes	3,124	8,971	6,113	23,949	16,078
small non-coding genes	2,804	2,156	2,427	4,871	5,531
long non-coding genes	135 (incl 1 read through)	6,798	3,307	16,857 incl 304 read through	9,985 incl 75 read through
misc. non-coding genes	185	17	379	2,221	562
Pseudogenes	568	1,626	674	15,214 incl 8 read through	13,597 incl 4 read through
Gene transcripts	30,585	63,041	58,692	227,530	142,446
Genscan gene predictions	52,372	46,573	152,168	51,756	57,381
Short variants	60,389,665	64,310,125		665,834,144	83,761,978
Structural variants	224,038	224,038		6,013,113	791,878

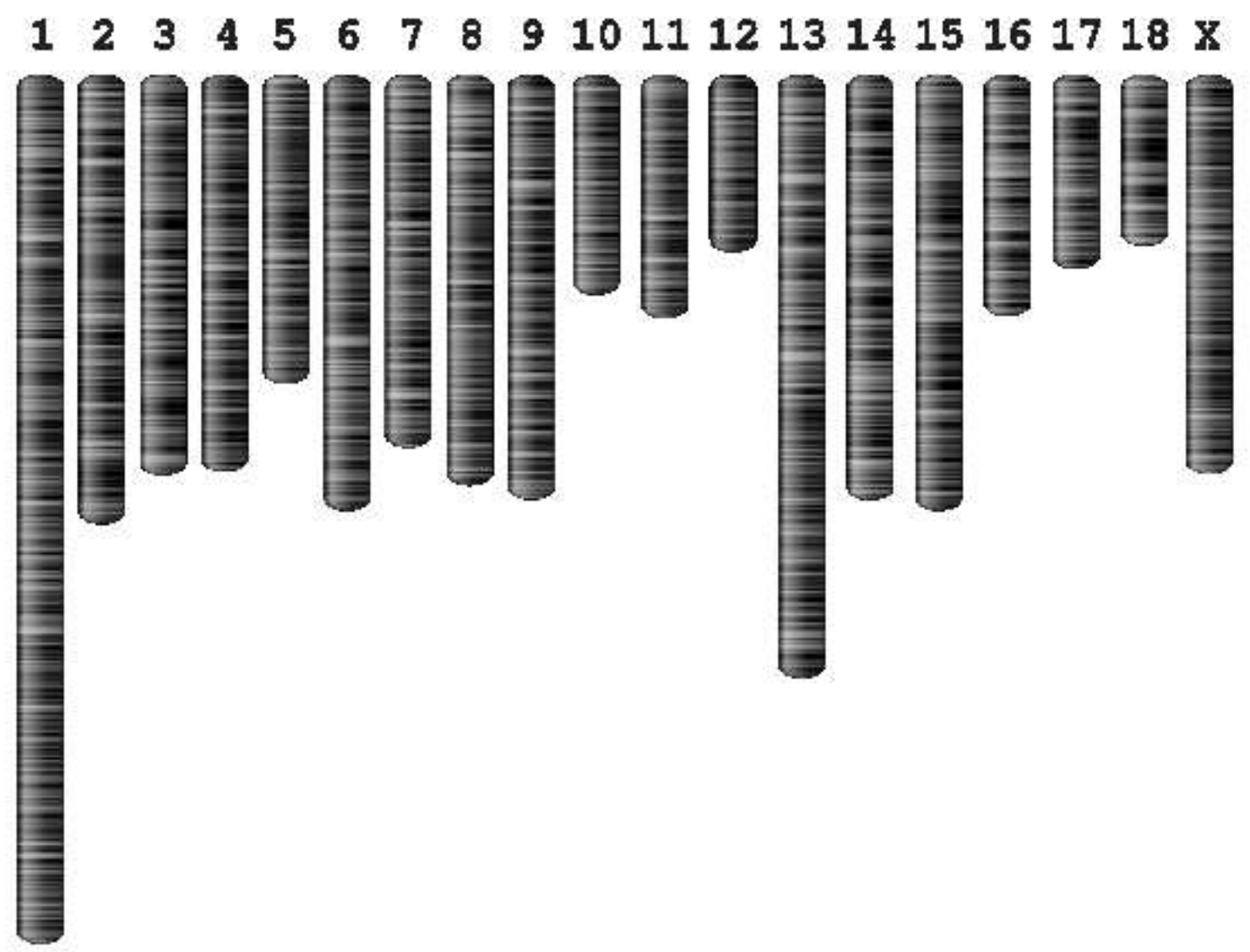
1074

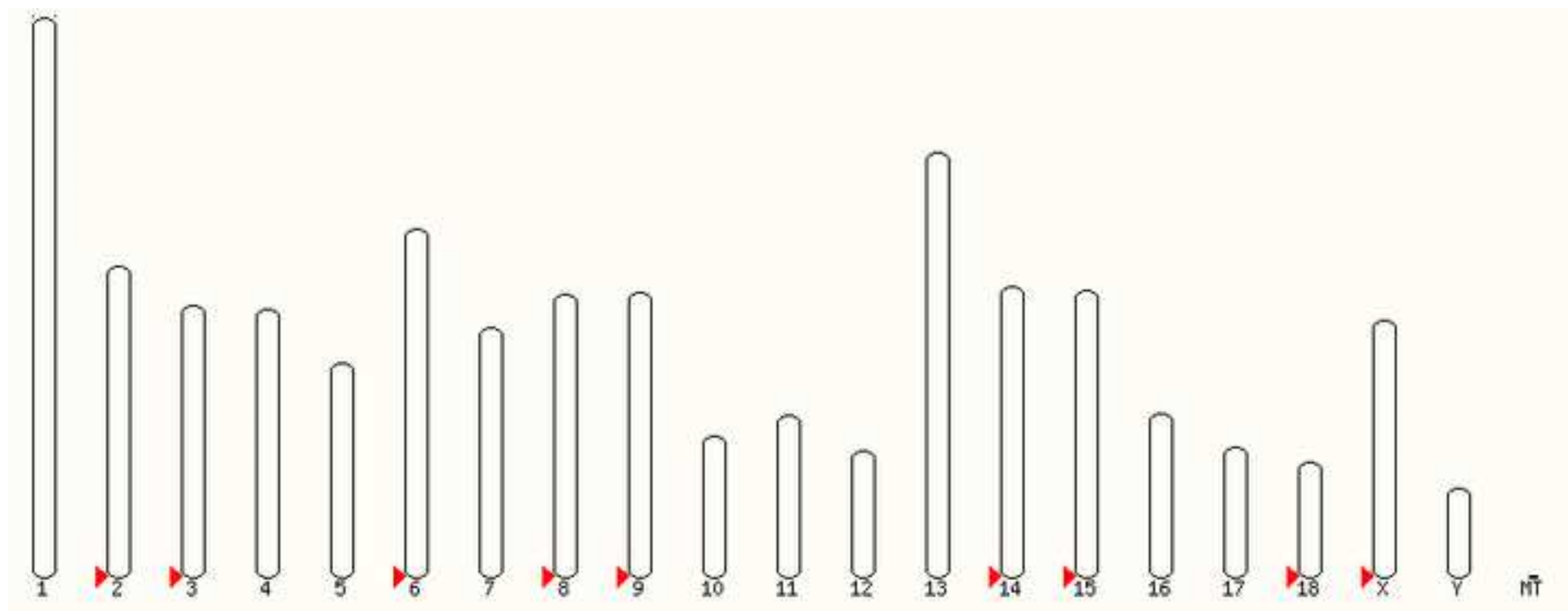
1075

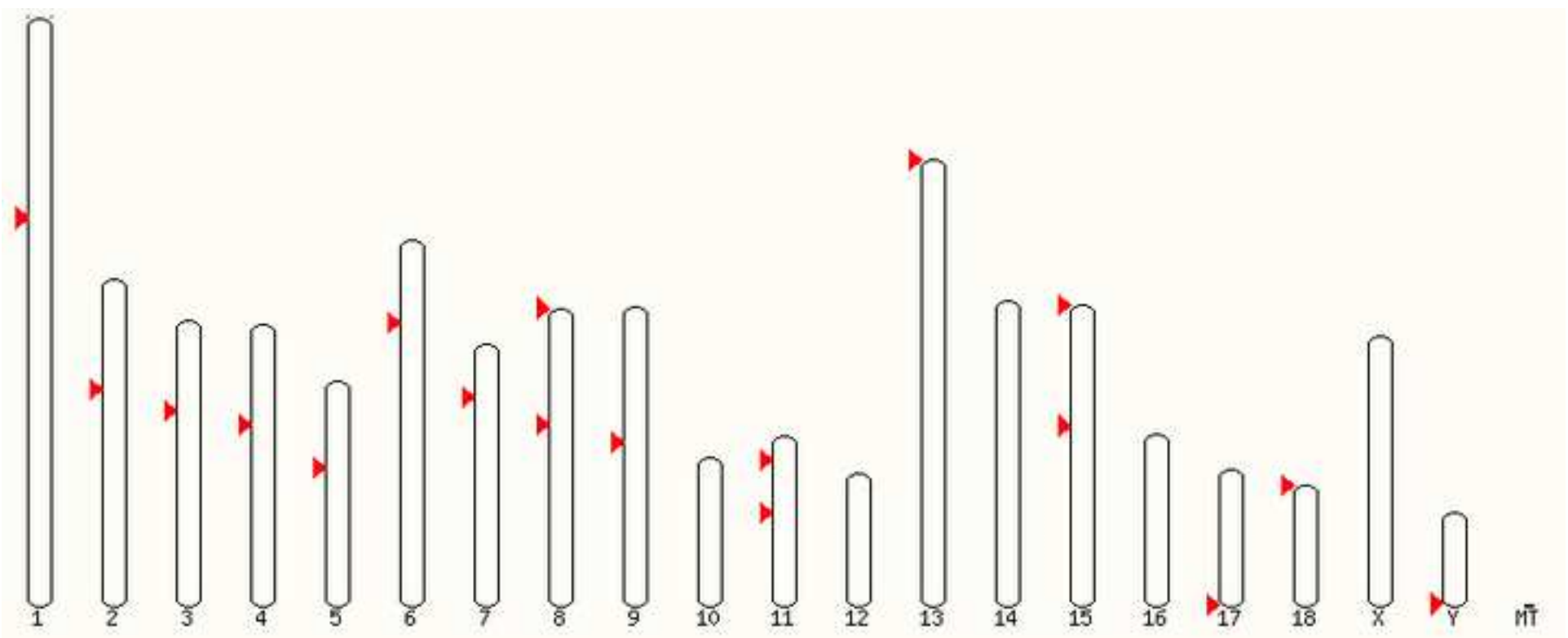


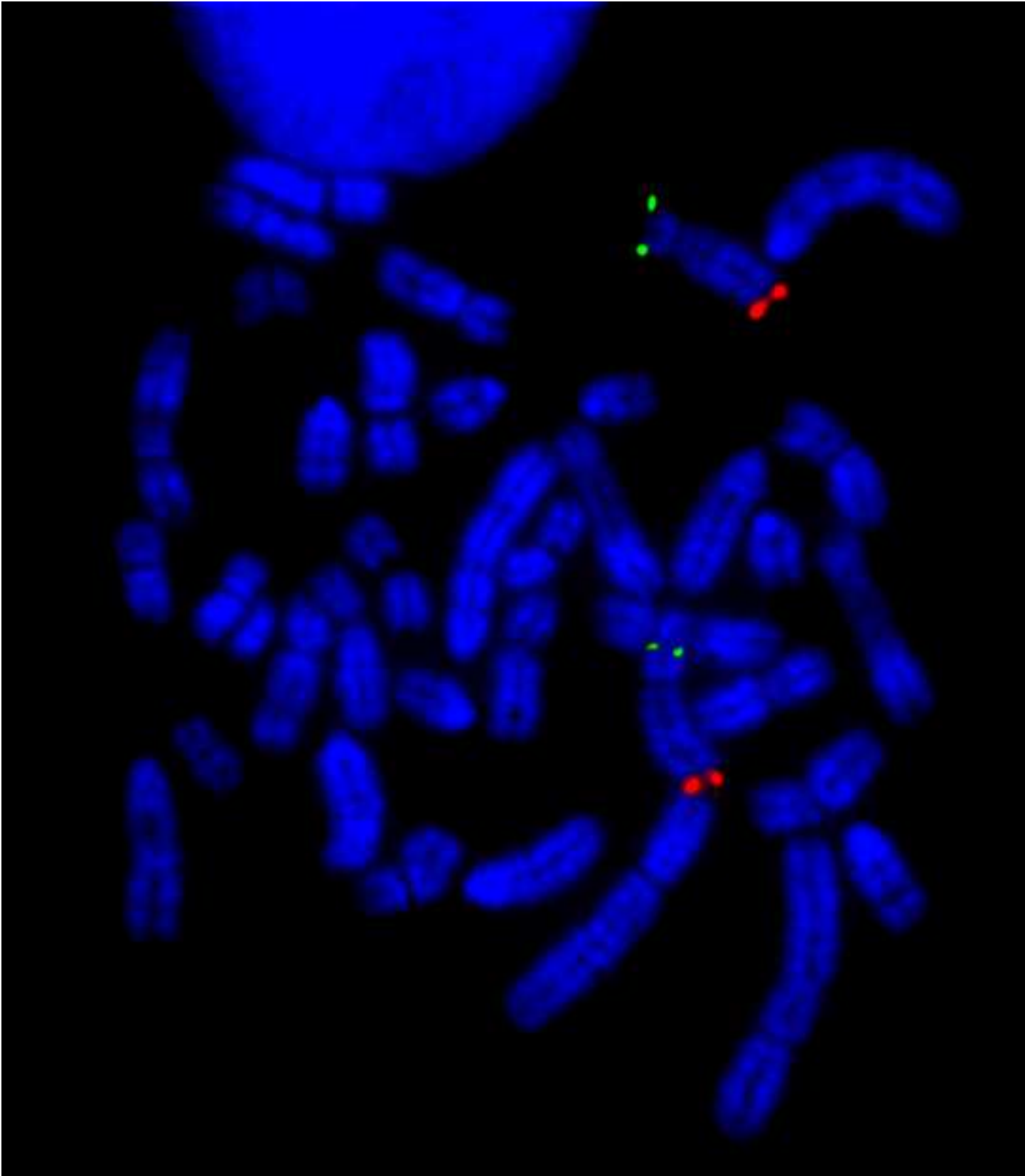




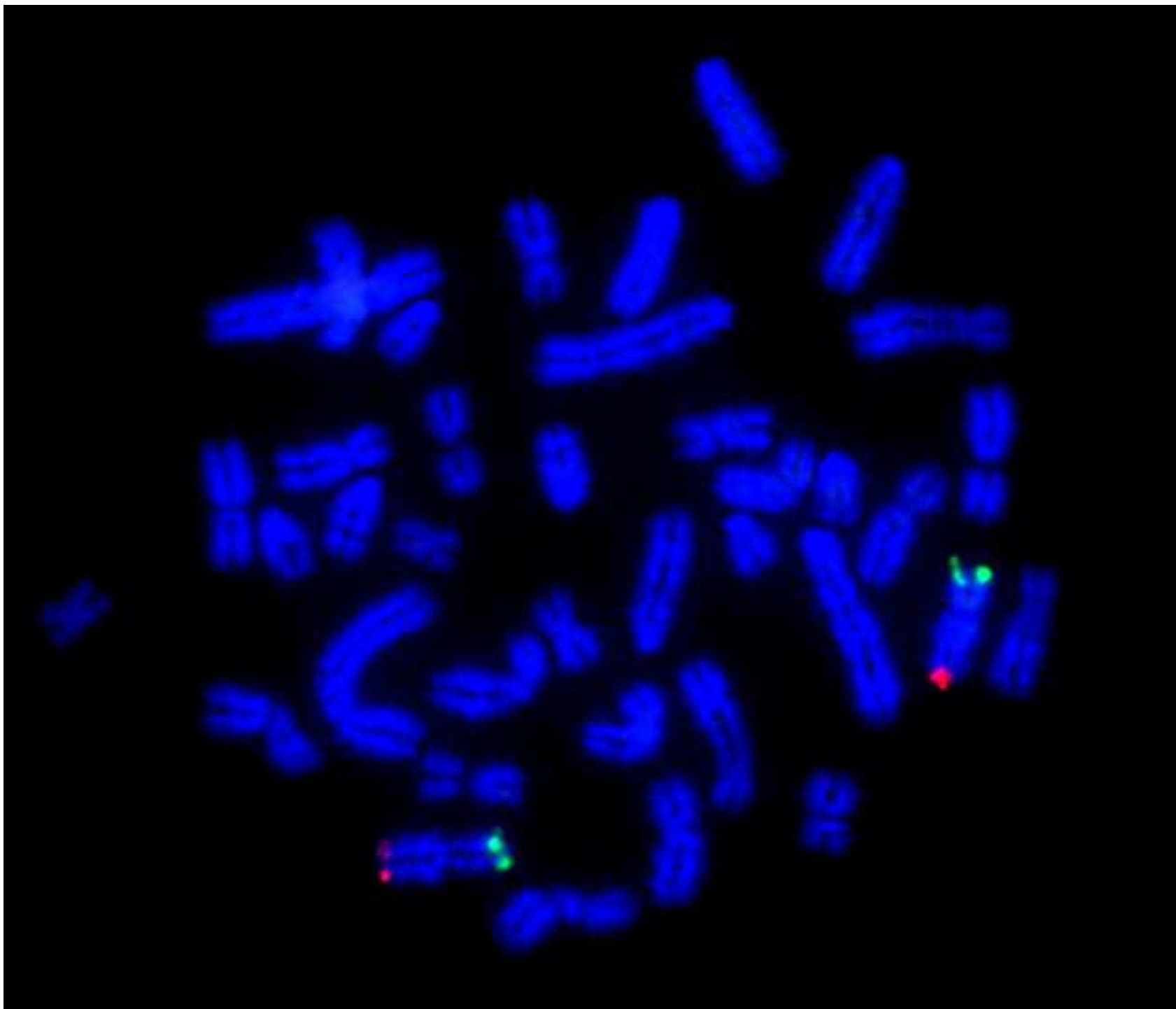


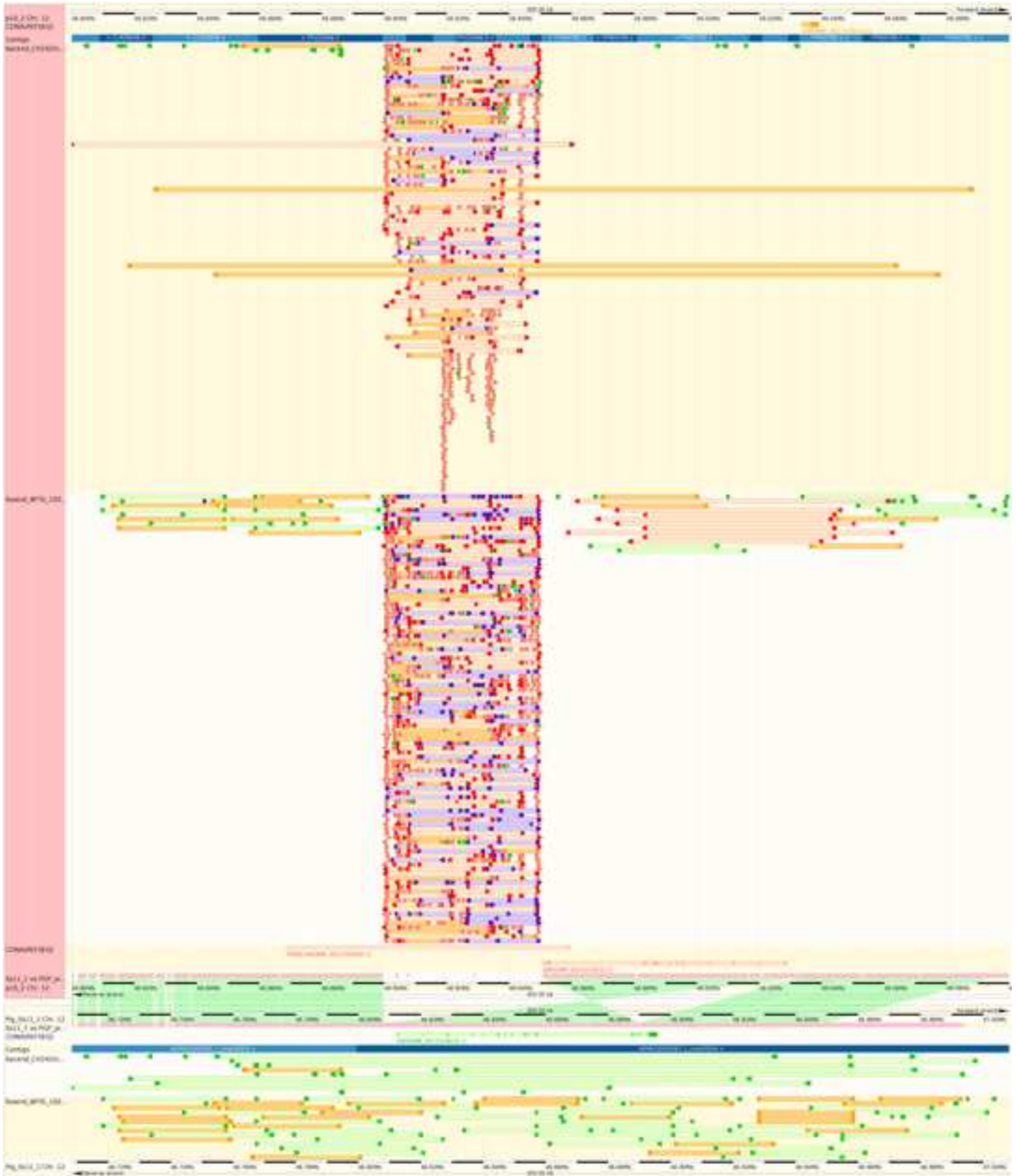


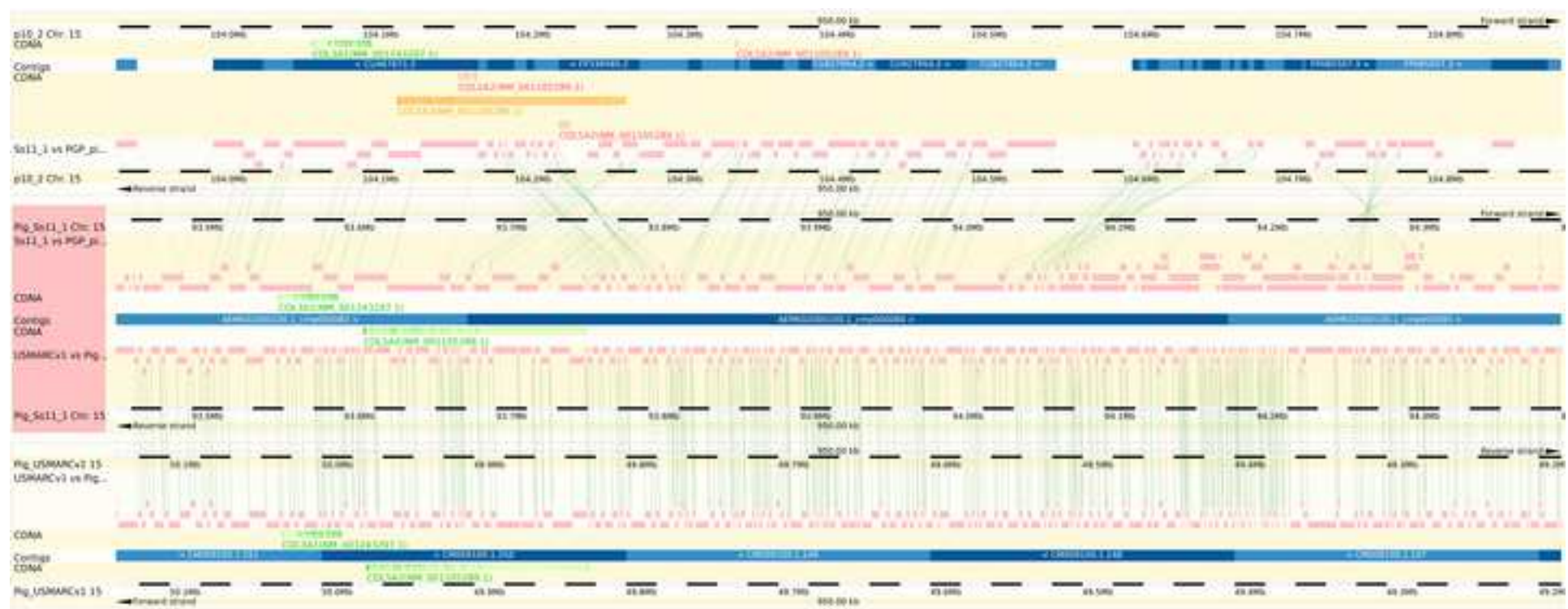




s19.Mge

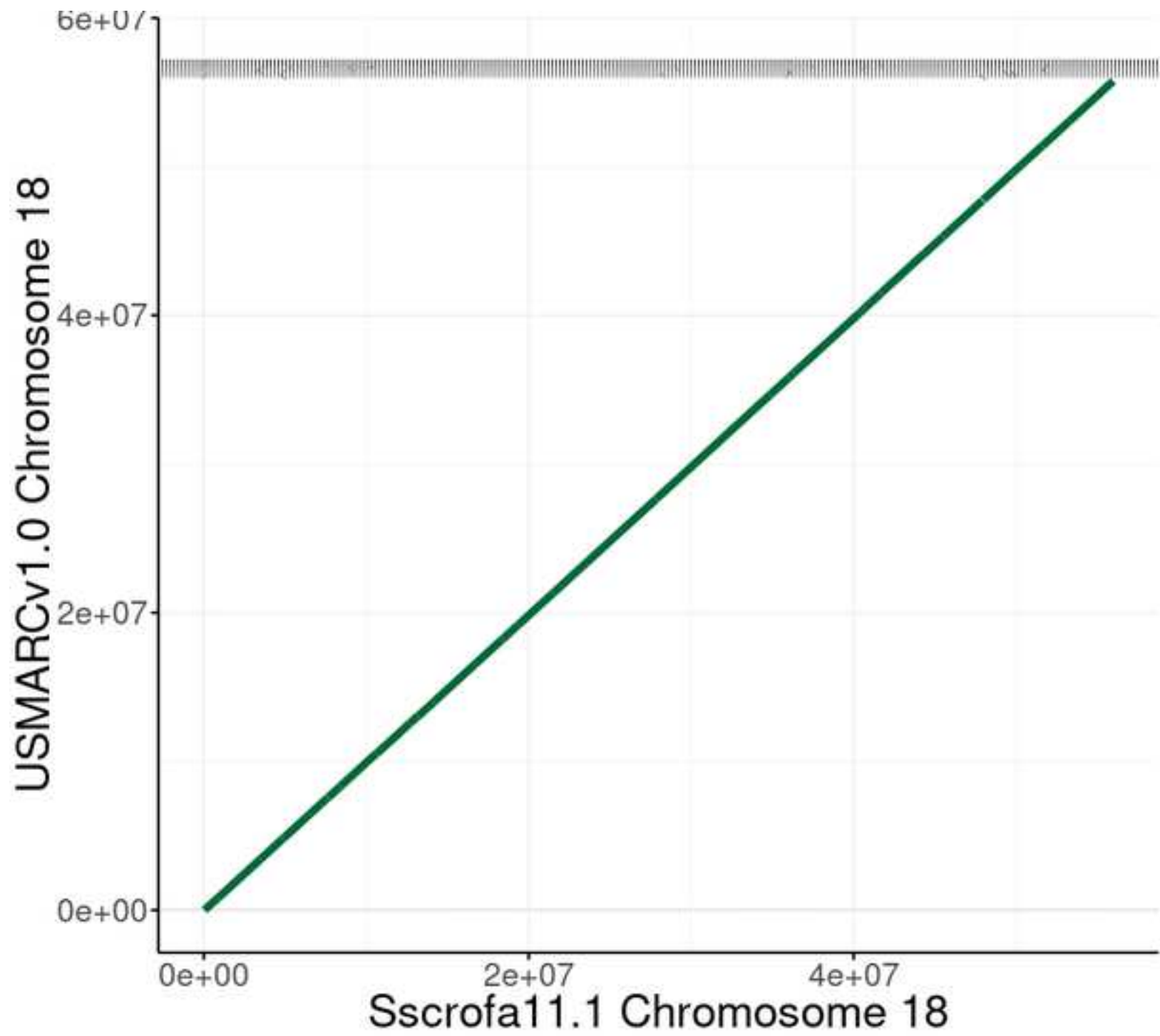


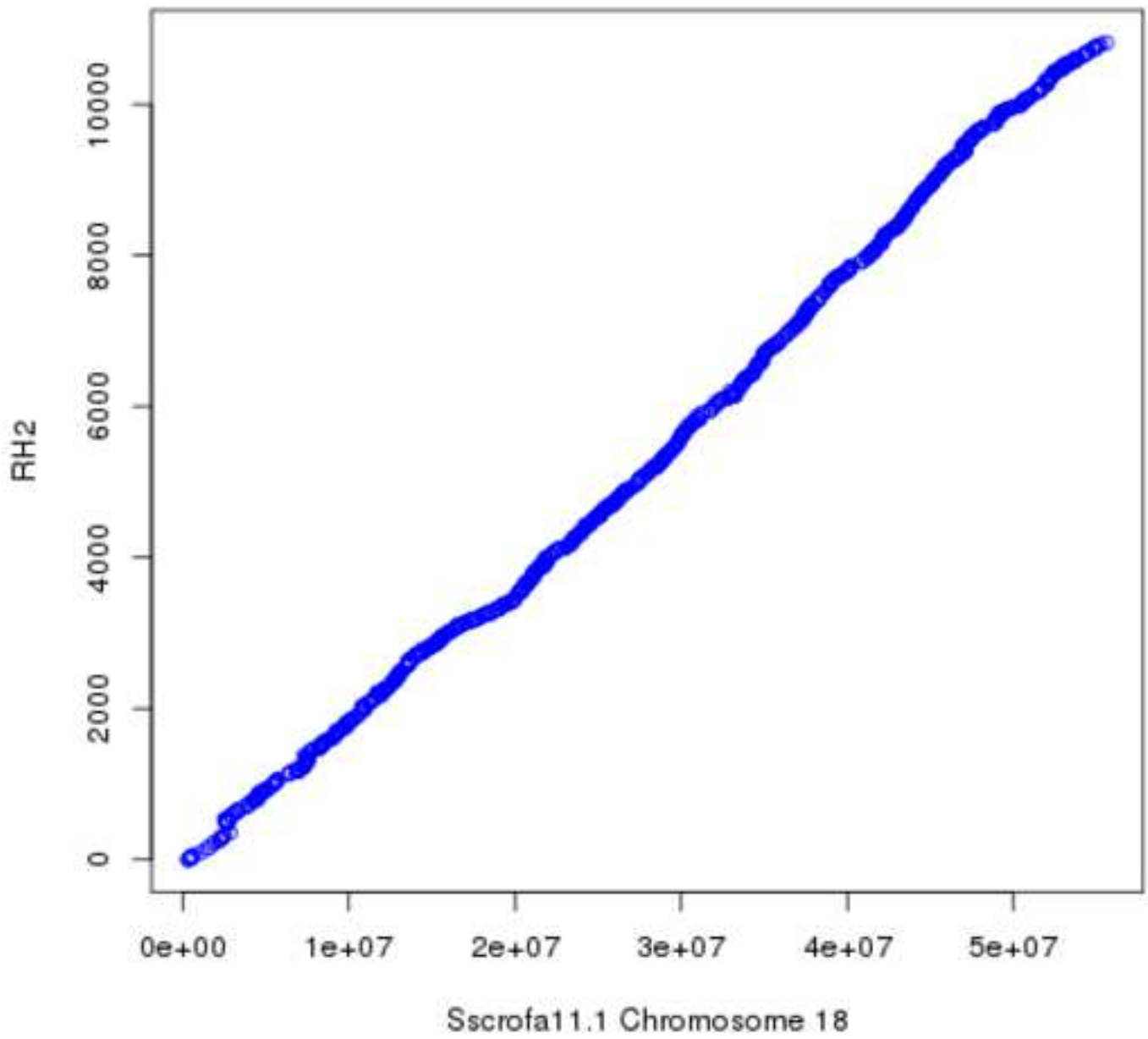


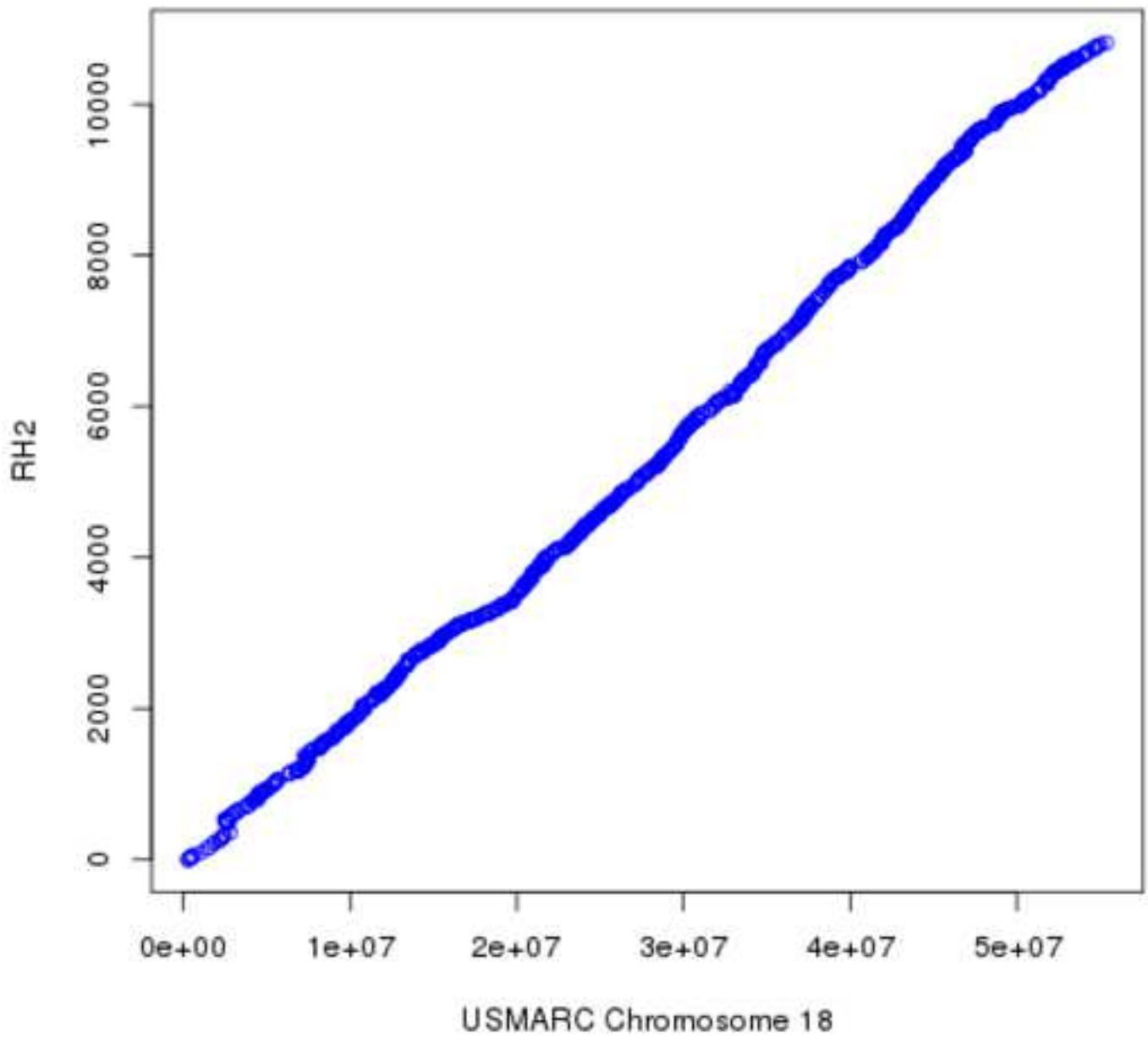


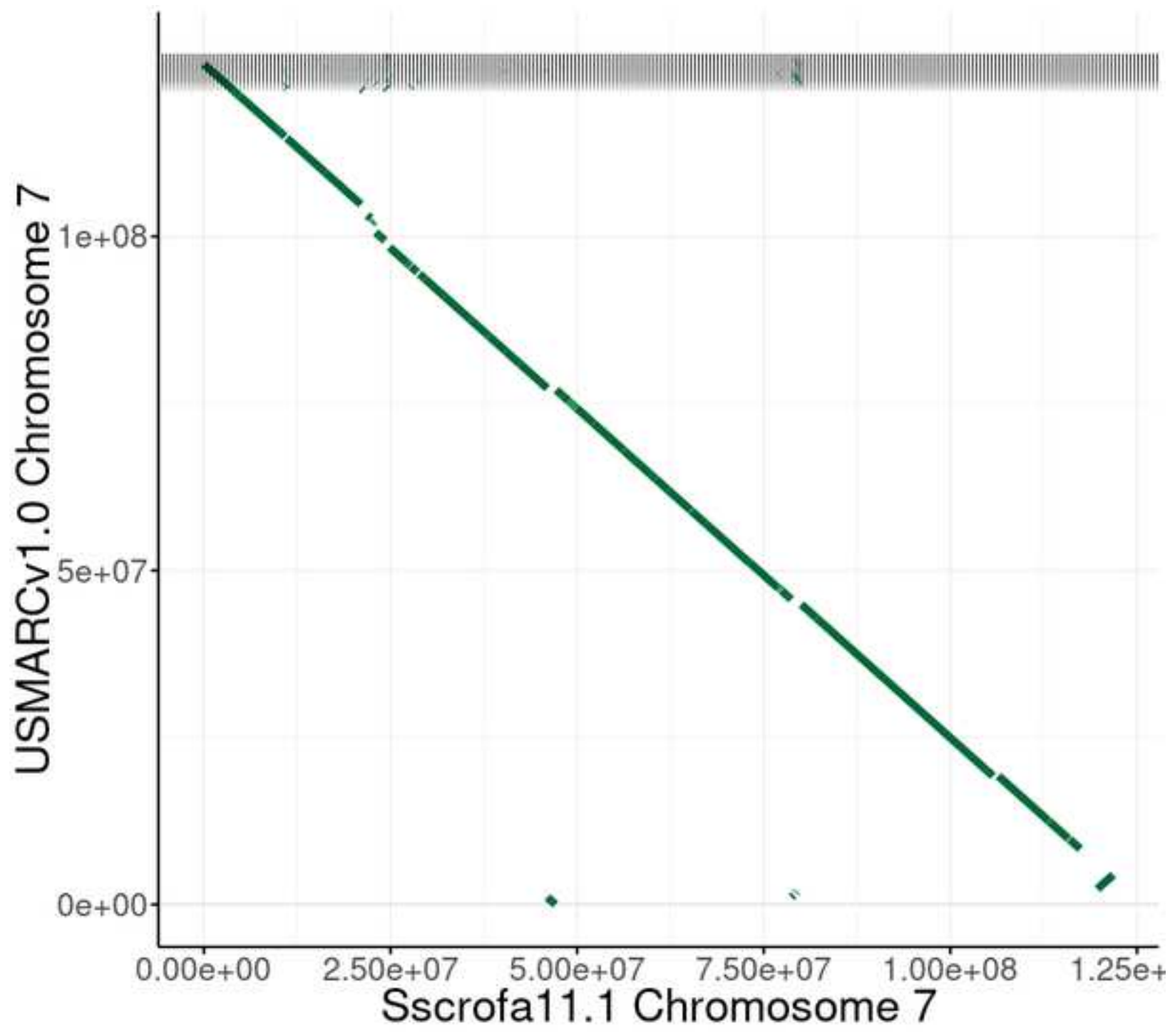


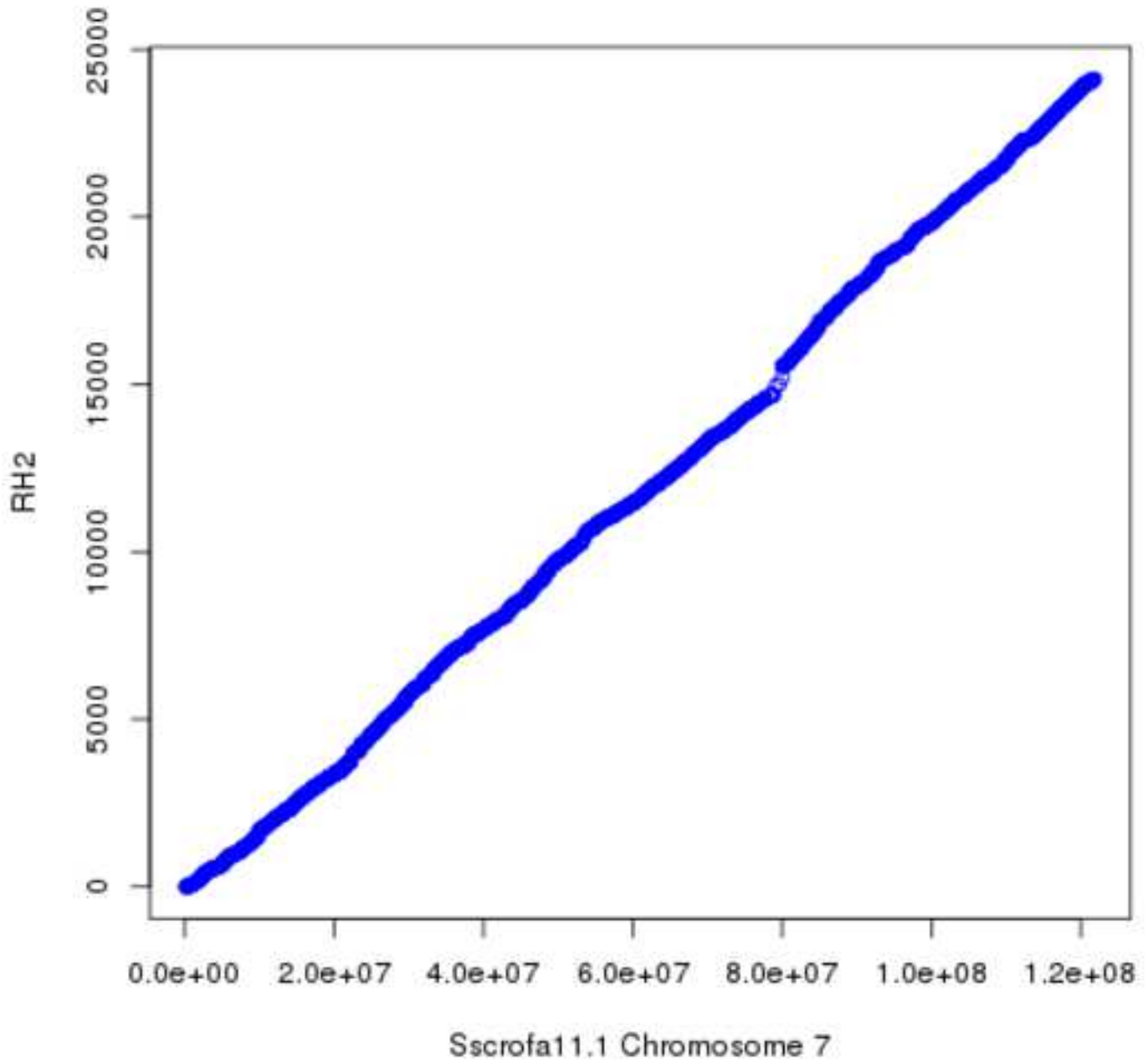


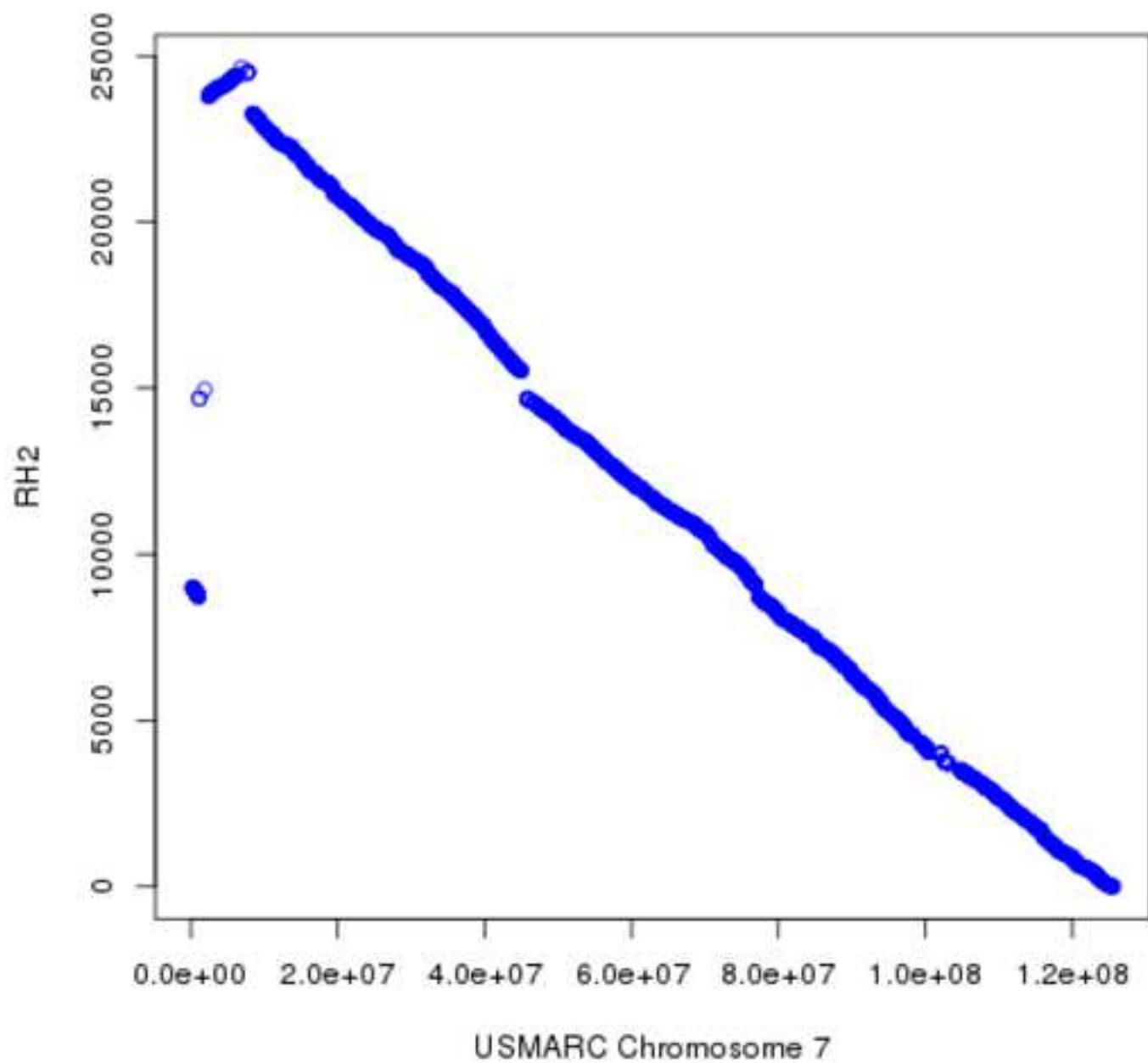


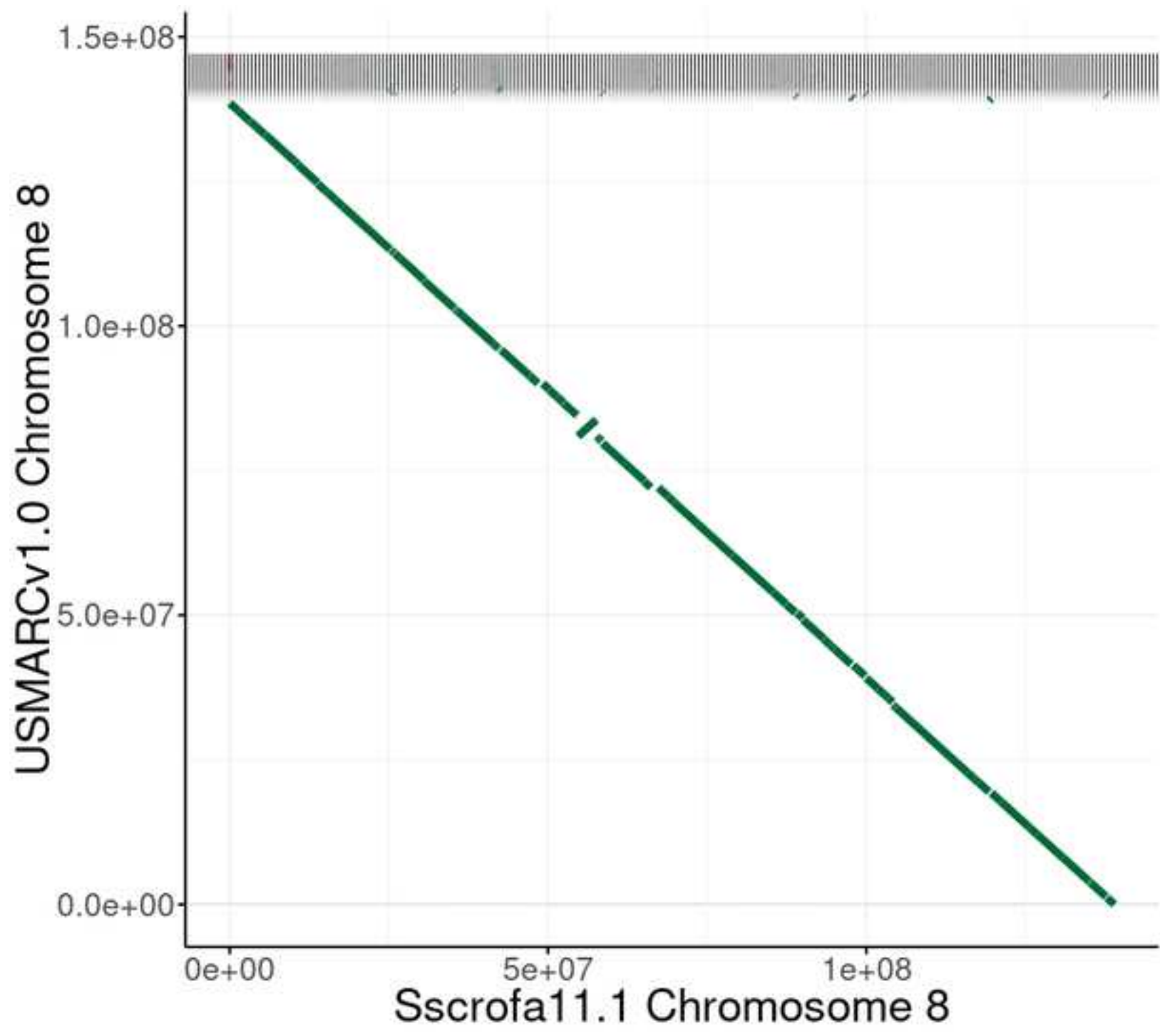


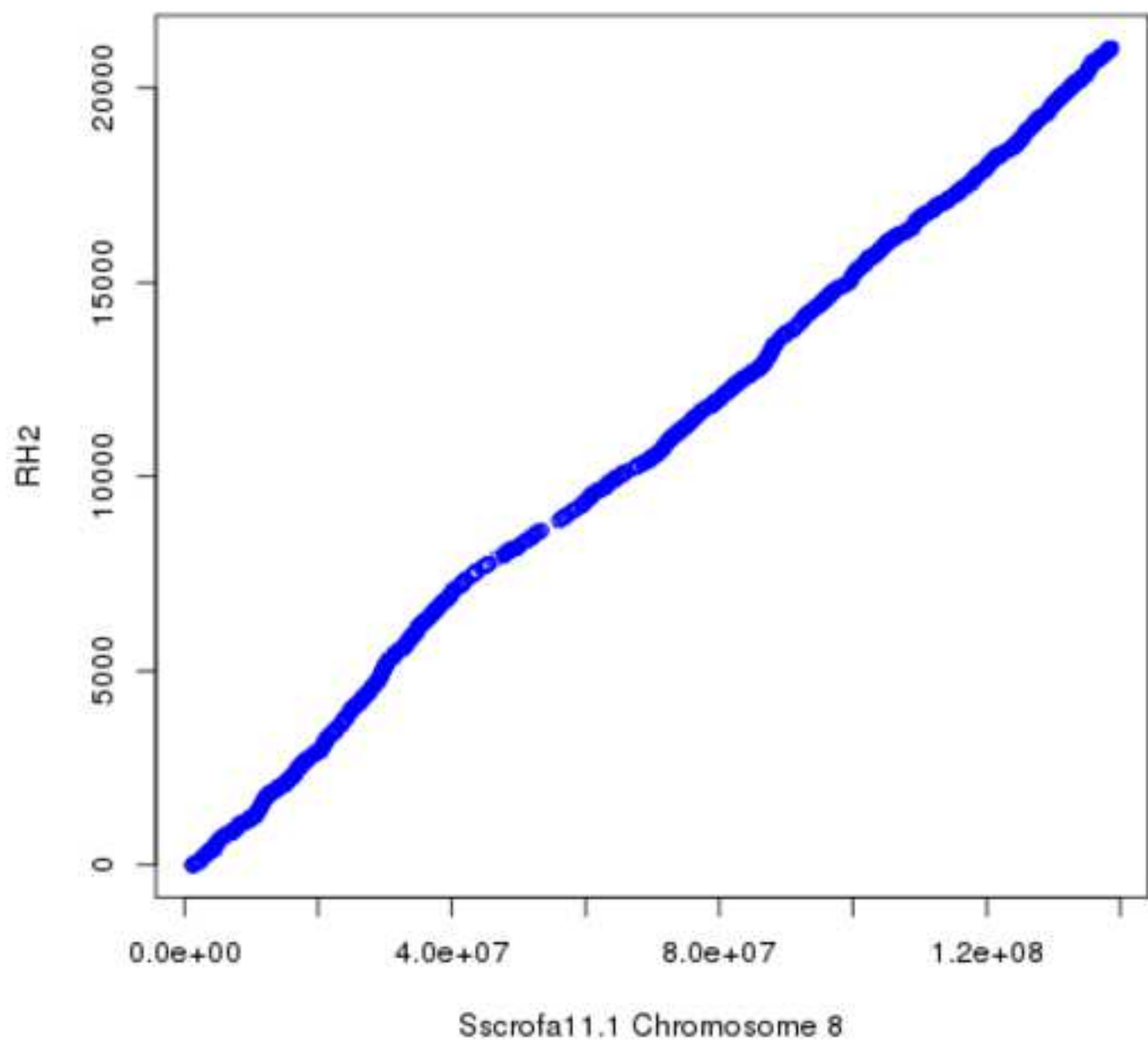


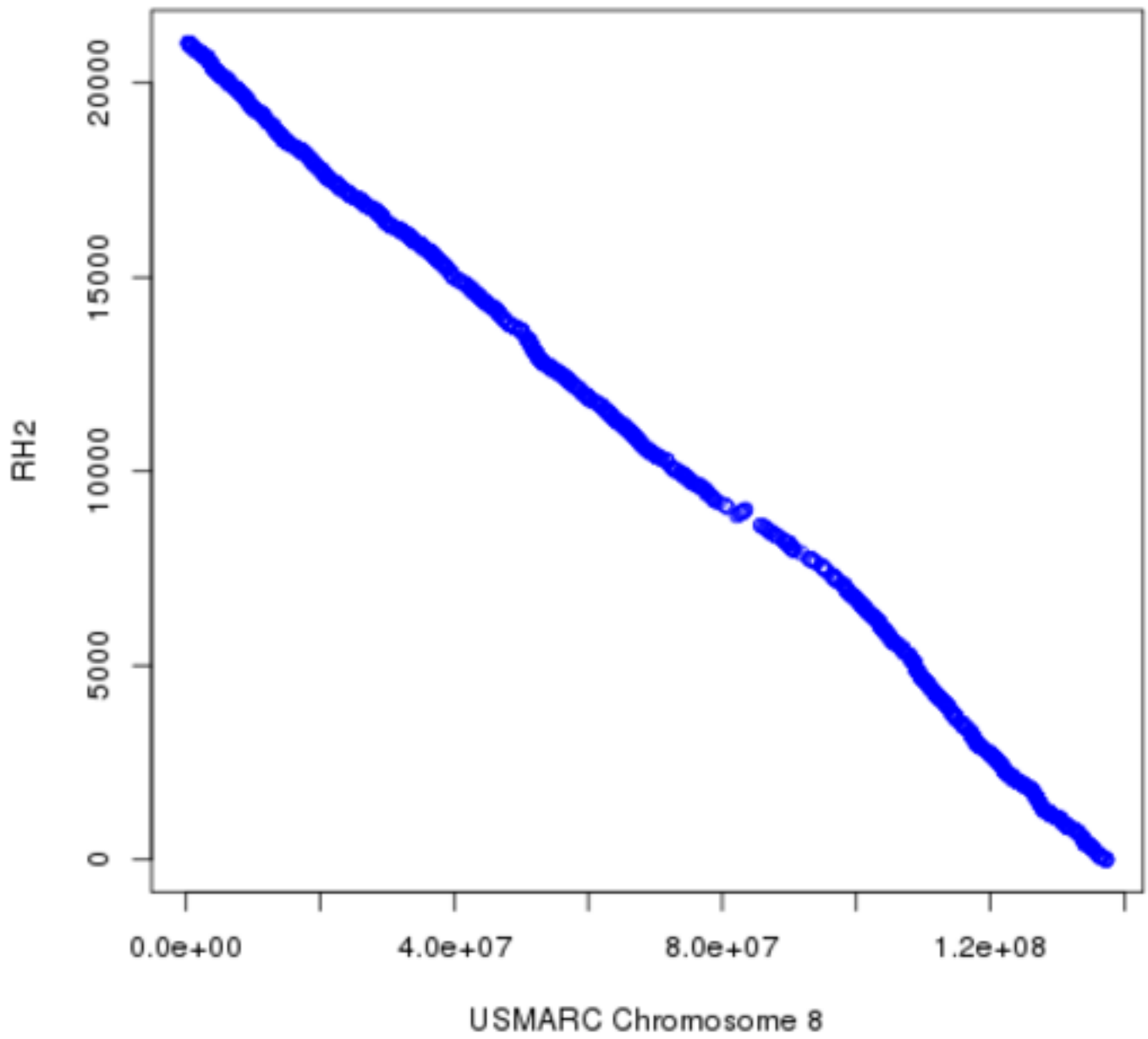


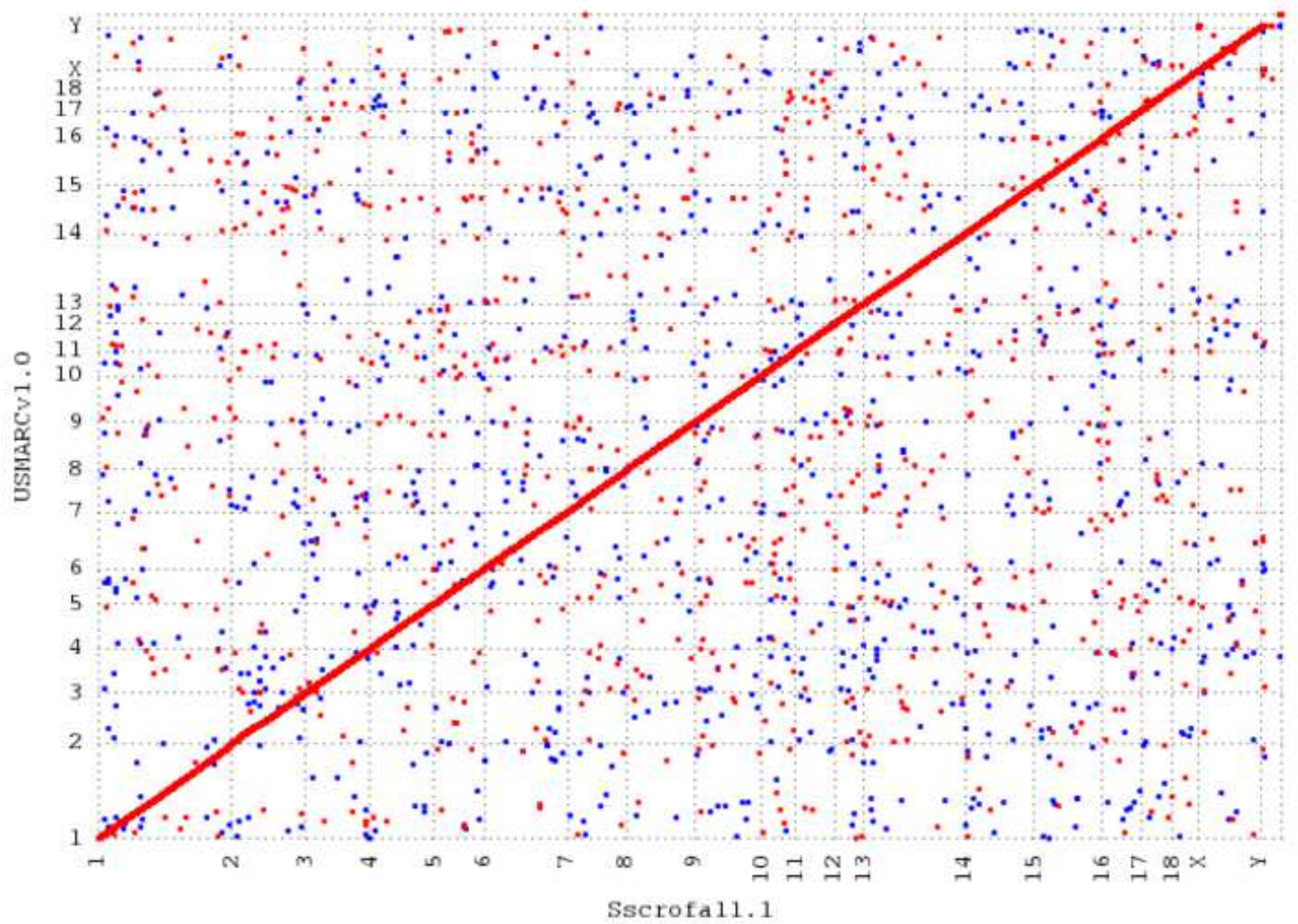




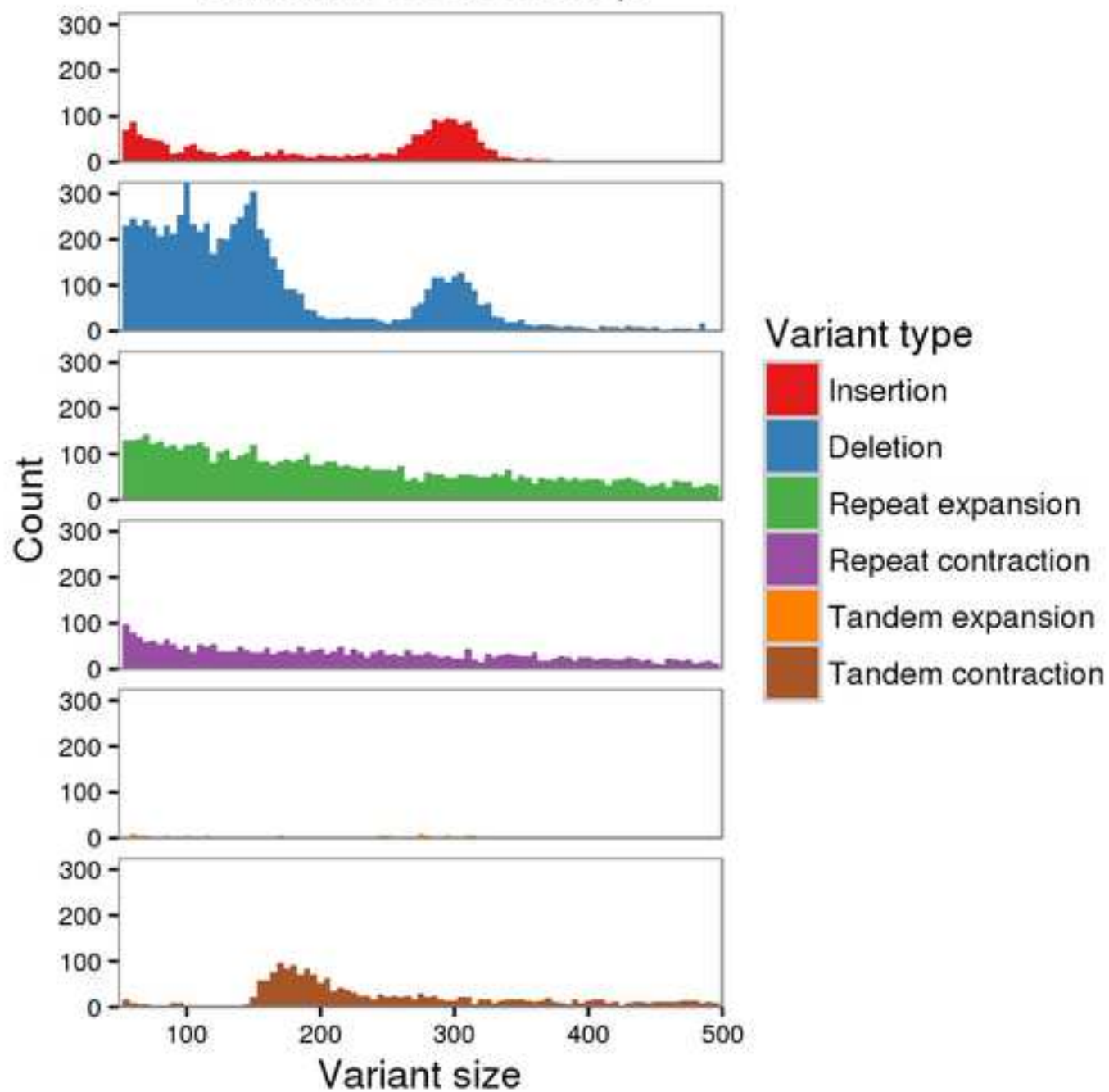


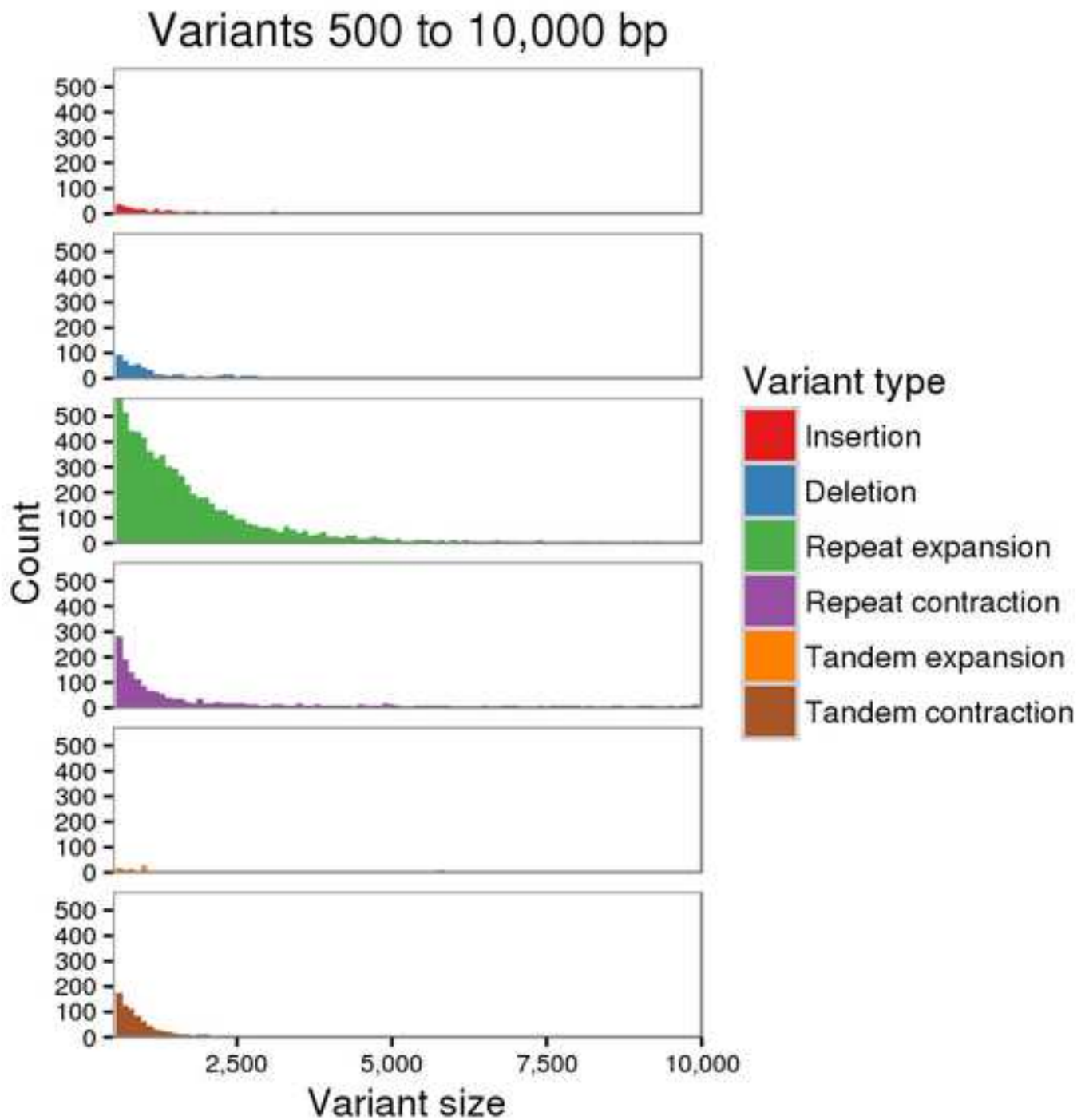


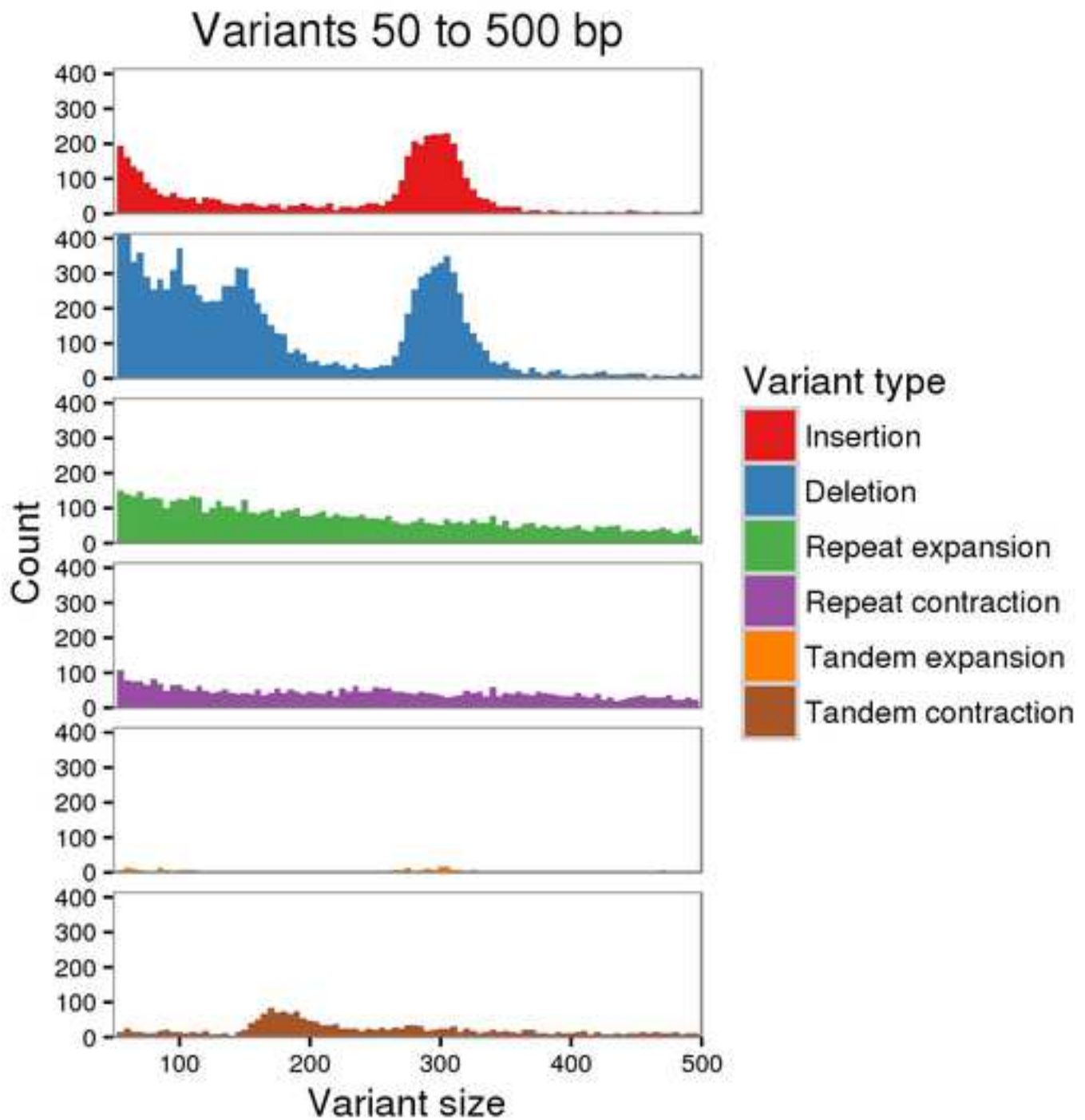


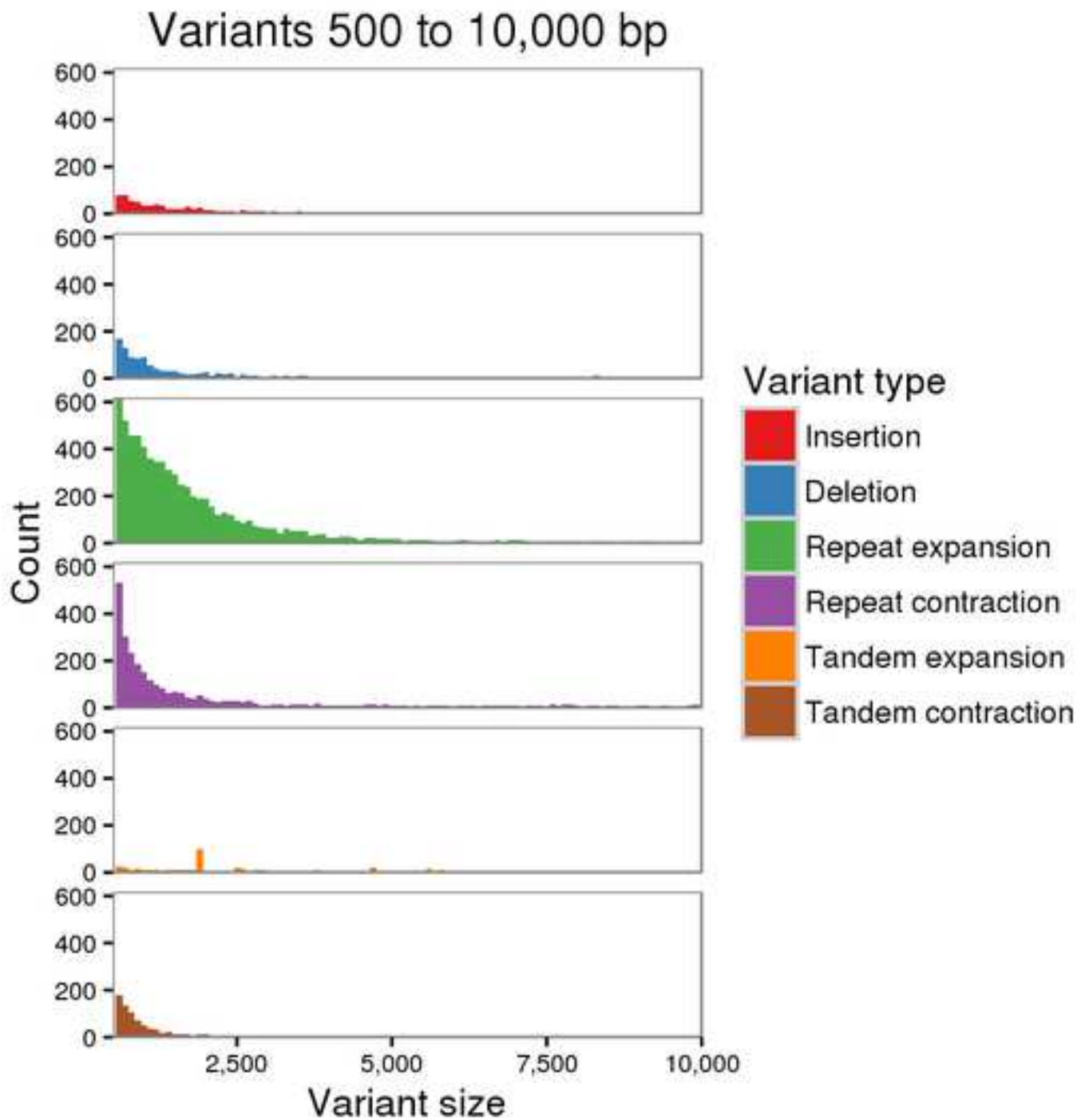


Variants 50 to 500 bp

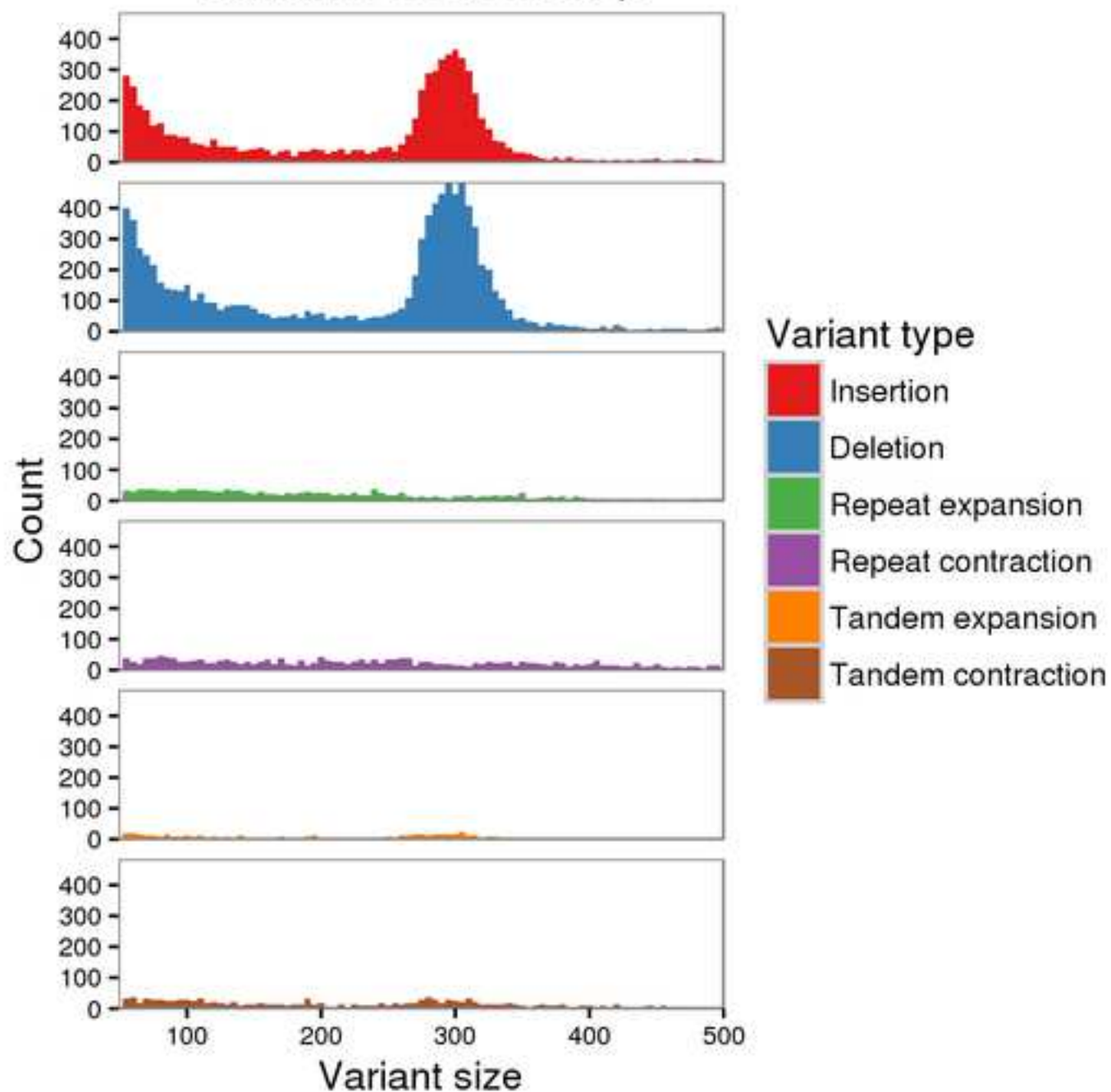


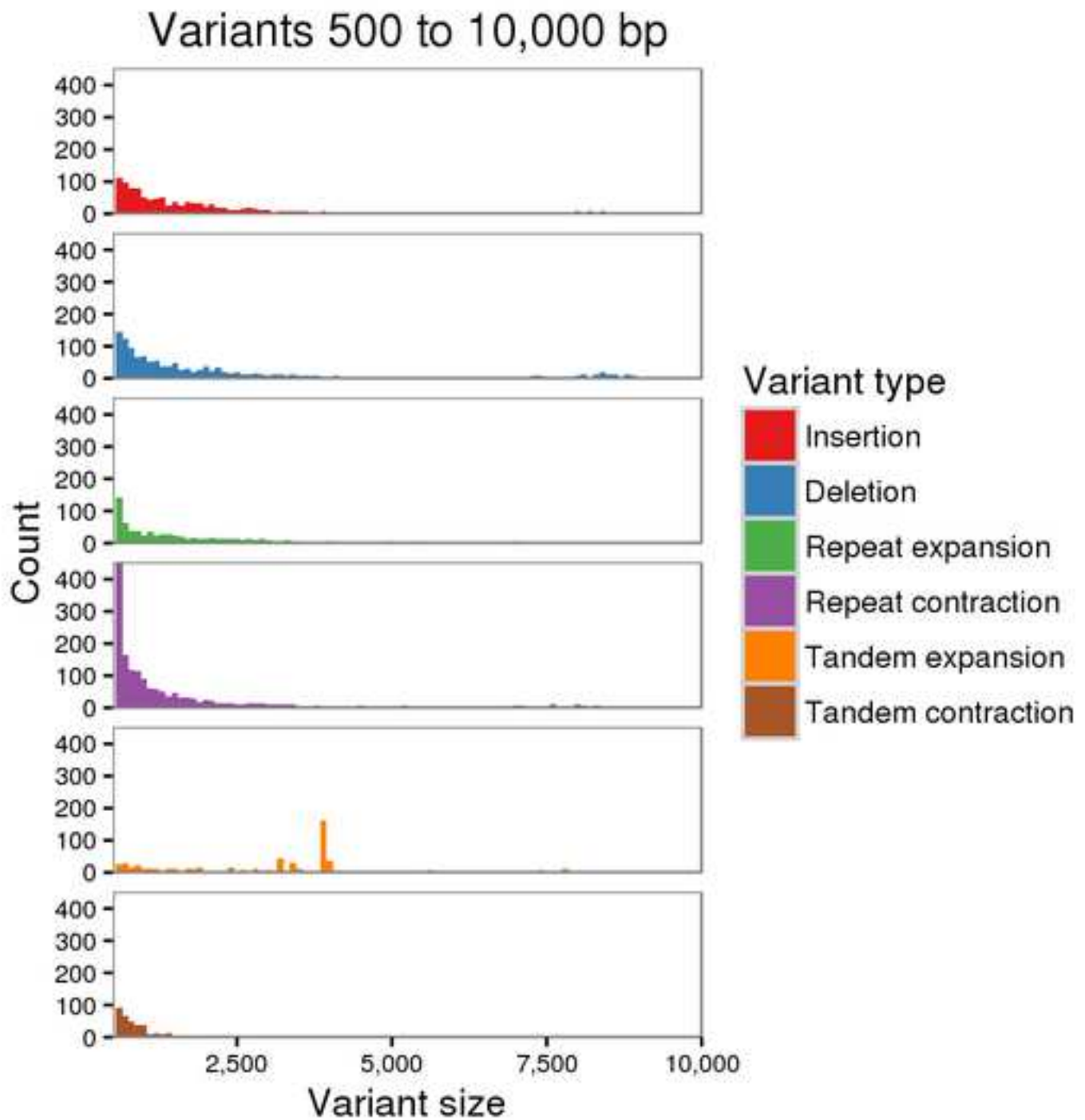


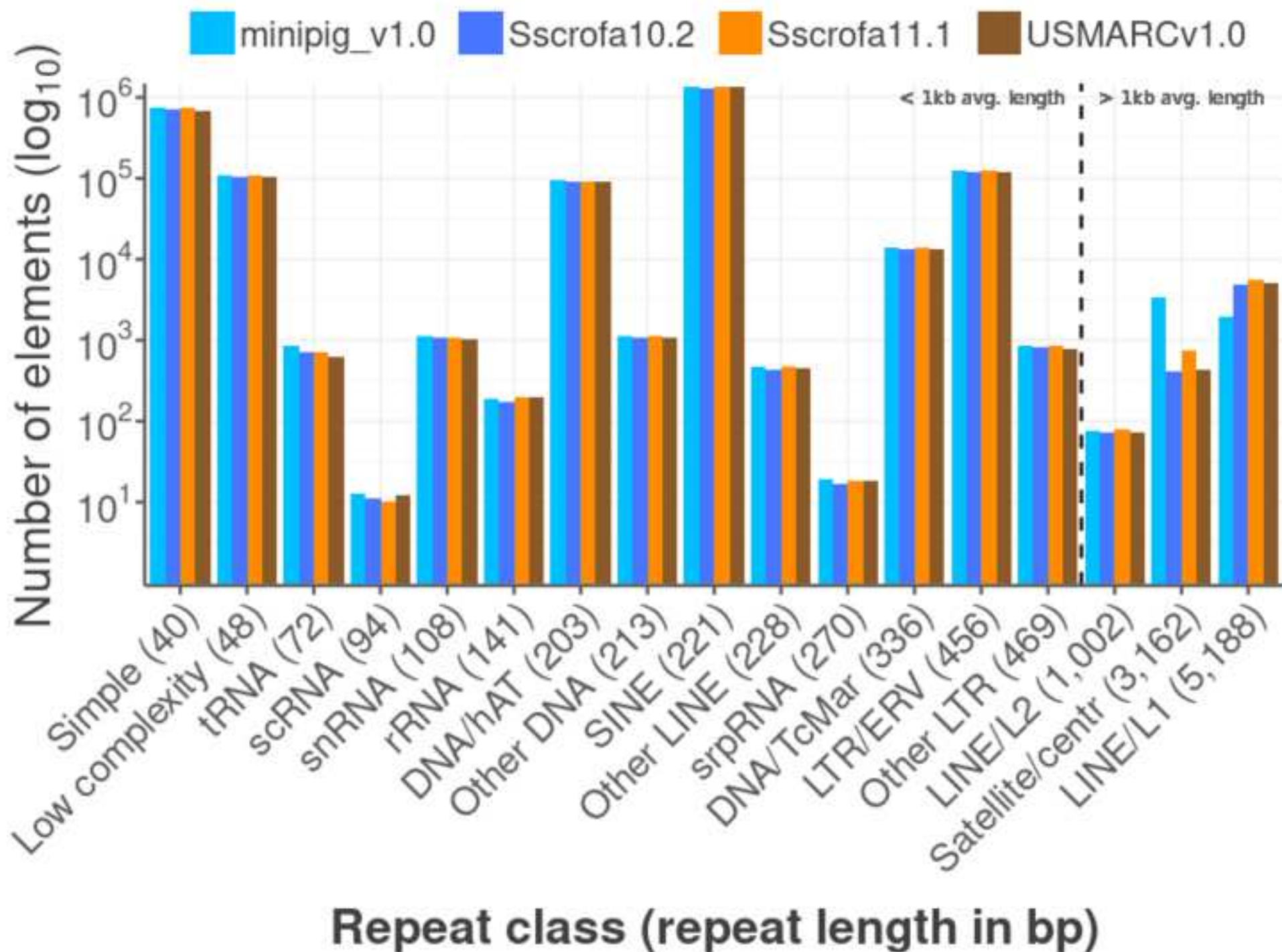


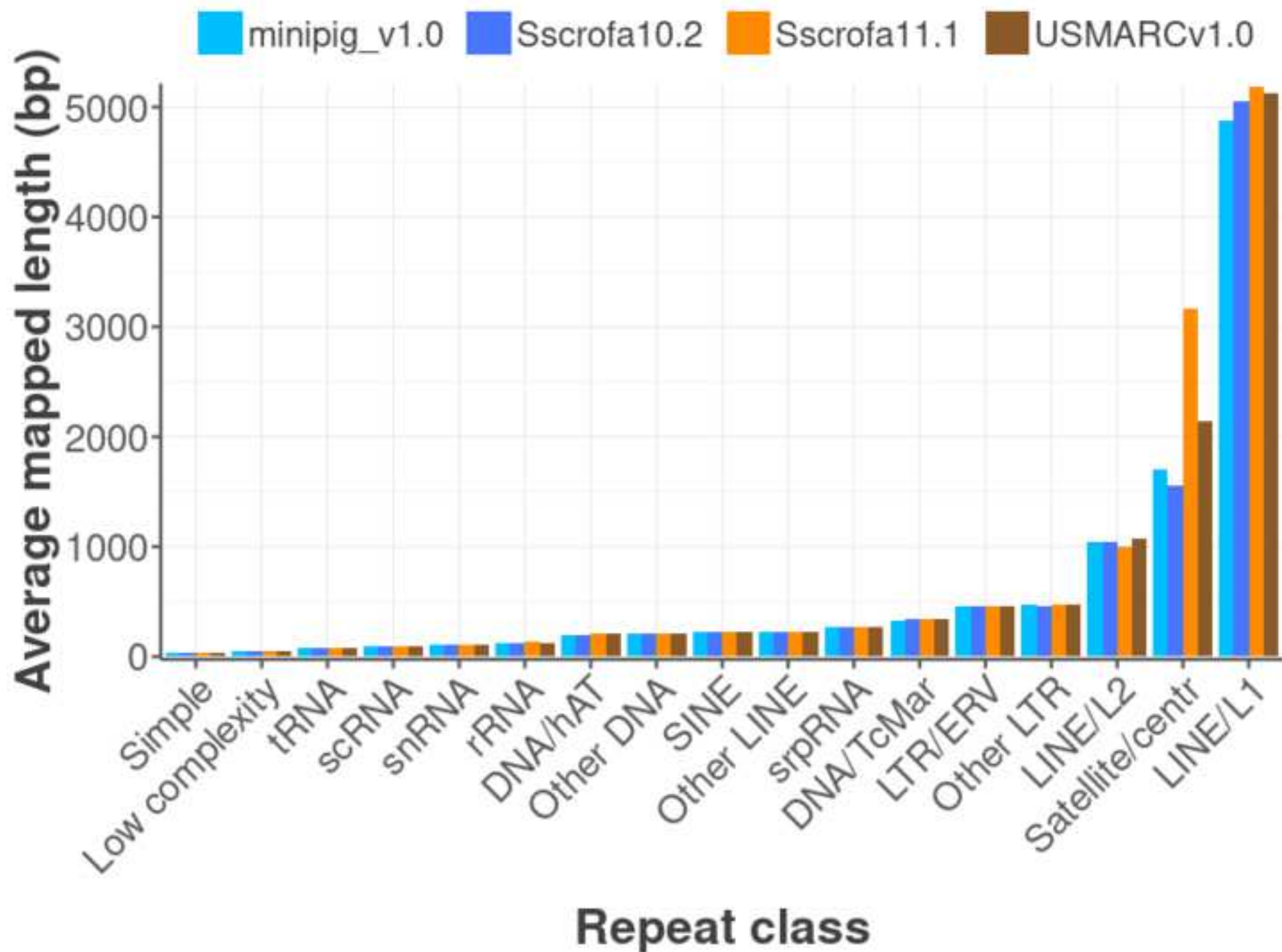


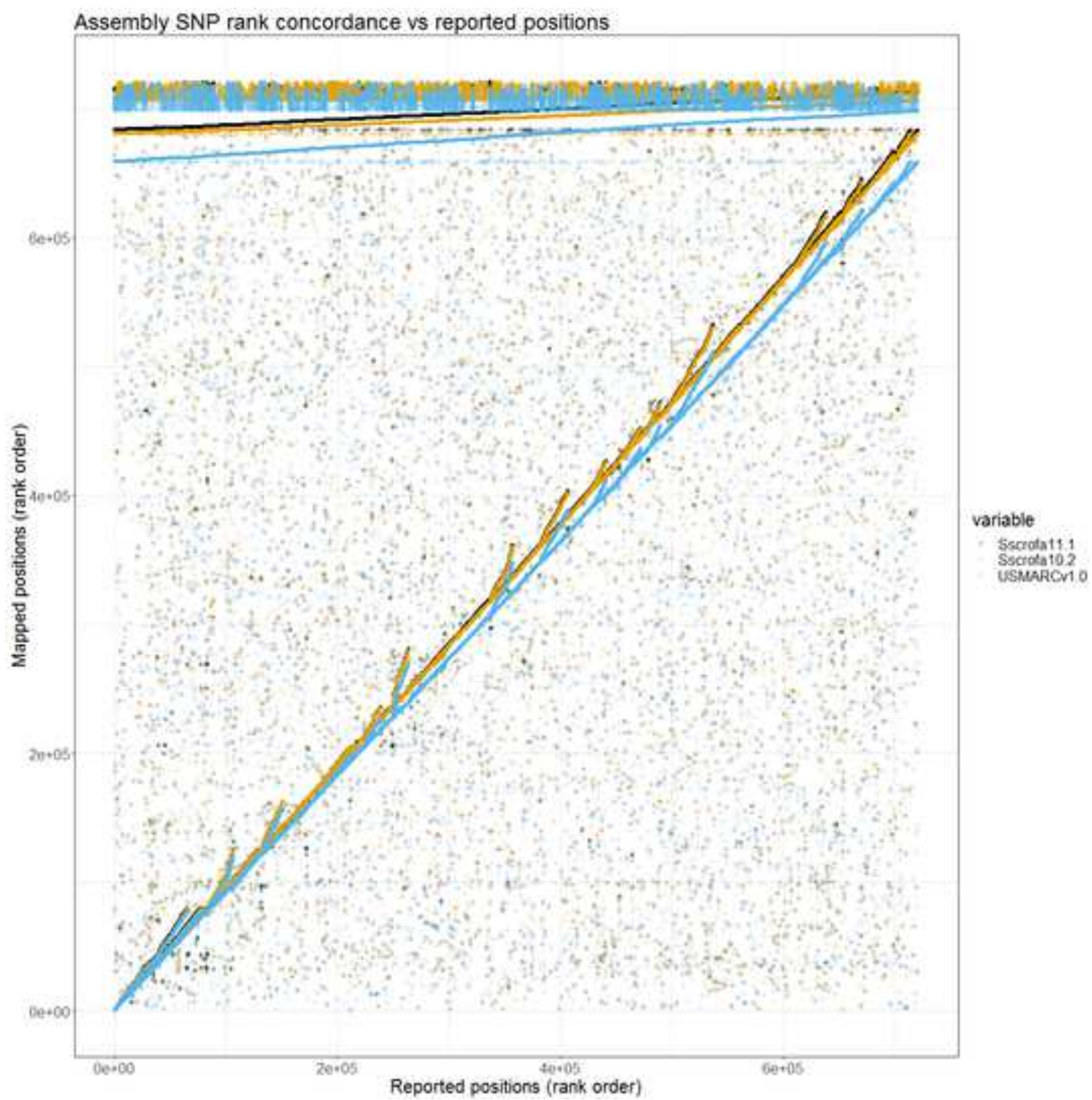
Variants 50 to 500 bp













[Click here to access/download](#)

Supplementary Material

Pig_genomes_suppl_21042020.docx





[Click here to access/download](#)

Supplementary Material

[TableS11_Pig_strains_annotate.xlsx](#)





THE UNIVERSITY of EDINBURGH
The Royal (Dick) School
of Veterinary Studies

THE ROSLIN INSTITUTE
The University of Edinburgh
Easter Bush
Midlothian
EH25 9RG
Telephone: +44 (0)131 651 9100
www.roslin.ed.ac.uk

Dear Editors

I am pleased to submit a final revised version of the manuscript entitled "An improved pig reference genome sequence to enable pig genetics and genomics research".

The questions and issues raised by the reviewers have been addressed as described previously.

The results from the Assemblytics comparisons of the 13 pig genome assemblies available in GigaDB (<http://dx.doi.org/10.5524/100732>).

The figures are integrated into the manuscript and supplementary materials. I have also uploaded copies of each figure image. These images are also available in the associated GigaDB dataset (<http://dx.doi.org/10.5524/100732>).

Yours sincerely

Alan L. Archibald