

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21

Fish and chips: the origin of human gene families is a predictor of the location of GWAS signals.

Sara V. Good^{1,2^}, Ryan Gotesman,³ Ilya Kisselev¹, Andrew D. Paterson^{3,4^}

¹ Department of Biology, The University of Winnipeg, MB, ² Department of Biological Sciences, The University of Manitoba, Winnipeg, MB, ³ Program in Genetics and Genome Biology, The Hospital for Sick Children, Toronto, ON, Canada; ⁴ Divisions of Biostatistics and Epidemiology, Dalla Lana School of Public Health, University of Toronto, Toronto, ON, Canada

Author emails: s.good@uwinnipeg.ca (SVG), Ryan Gotesman_(RG), Ilya Kisselev, Andrew D.

Paterson Andrew.paterson@sickkids.ca

Author (SVG) ORCID: 0000-0003-0563-7724

Author (ADP) ORCID: 0000-0002-9169-118X

^Joint corresponding authors: (SVG)(ADP)

Running header: Ancestral Linkage Group influence on GWAS

Keywords: genome wide association studies, ancestral genome reconstruction, ancestral linkage groups, macrosynteny, Experimental Factor Ontology, ohnologs, gene-level constraint, GTEEx.

22 Abstract

23 GWAS have identified thousands of loci associated with human complex diseases and traits.
24 How these loci are distributed through the genome has not been systematically evaluated. We
25 hypothesised that the location of GWAS loci differ between ancestral linkage groups (ALGs)
26 related to the paralogy and function of genes. We used data from the NHGRI-EBI GWAS
27 catalog to determine whether the density of GWAS loci relative to HapMap variants in each
28 ALG differed, and whether ALG's were enriched for experimental factor ontological (EFO)
29 terms assigned to the GWAS traits. In a gene-level analyses we explored the characteristics of
30 genes linked to GWAS loci and those mapping to the ALG's. We find that GWAS loci were
31 enriched or deficient in 9 and 7 of the 17 ALG's respectively, while there was no difference in
32 the number of GWAS loci in regions of the human genome unassigned to an ALG. All but 2
33 ALG's were significantly enriched or deficient for one or more EFO terms. Lastly, we find that
34 genes assigned to an ALG are under higher levels of selective constraint, have longer coding
35 sequences and higher median expression in the tissue of highest expression than genes not
36 mapping to an ALG. On the other hand, genes associated with GWAS loci have longer genomic
37 length and exhibit higher levels of selective constraint relative to non-GWAS genes.
38 Collectively, this suggests that understanding the location and ancestral origins of GWAS signals
39 may be informative for the development of tools for variant prioritization and interpretation.

40

41

42

43 Introduction

44 Genome wide association studies (GWAS) have identified thousands of genetic
45 variants that influence human complex traits including both disease and non-clinical phenotypes
46 (Hirschhorn and Daly 2005; Hindorff et al. 2009; Visscher et al. 2017; Buniello et al. 2019).
47 Collectively, these studies indicate that variation in complex traits is often controlled by a few to
48 several dozen loci, i.e. core genes, while the rest, and the majority, of trait heritability is
49 massively polygenic and driven by weak effects of trans-acting expression quantitative trait loci
50 (eQTL) (Shi et al. 2016; Boyle et al. 2017; Liu et al. 2019). In order to prioritize the functional
51 effects of GWAS variants, developing tools to assess the potential impact of variants on
52 Mendelian and complex diseases has become an important area of research. Early tools to
53 prioritize variants utilized orthologous gene conservation across taxa and/or information about
54 the location of missense or nonsense mutations within the 3D structure of proteins (Siepel et al.
55 2005; Kumar et al. 2009; Adzhubei et al. 2013; Mistry et al. 2013; Shihab et al. 2013; Schwarz et
56 al. 2014; Choi and Chan 2015; Ioannidis et al. 2016), while more recent tools incorporate other
57 information such as population allele frequencies, functional genomics data, and other genome
58 annotations (Davydov et al. 2010; Yandell et al. 2011; Hu et al. 2013; Kennedy et al. 2014;
59 Kircher et al. 2014; Schwarz et al. 2014; Gulko et al. 2015; Kichaev et al. 2019).

60 In addition to variant prioritization methods, there are tools that use gene-level
61 information to guide variant analysis or interpretation. For example, based on an empirical
62 approach to calculate the expected number of *de novo* mutations per gene (Samocha et al. 2014),
63 ExAC used information from 60,706 exomes to estimate the probability of loss-of-function
64 (LOF) intolerance (pLI) for genes in the human genome (Lek et al. 2016). The most LOF-
65 intolerant genes (pLI>0.9) predicted by this metric were expected to predict haploinsufficient

66 human dominant disease genes, but pLI conflates estimates of dominance and selection (h and s
67 in population genetics theory), and is thus a measure of the strength of selection of deleterious
68 alleles in heterozygotes (Fuller et al. 2019). Currently, the gnomAD consortium, the
69 continuation of ExAC, recommends using a related metric known as the observed to expected
70 number (oe) of LOF variants in a gene, which is a continuous statistic and provides a better
71 estimate of gene-level purifying selection (Karczewski et al. 2020). Lastly, recently a novel
72 approach to incorporate gene family information (i.e. paralogous genes) to facilitate variant
73 interpretation in neurodevelopmental disorders was developed (Lal et al. 2020).

74 Few studies, however, have examined genome-wide patterns of the location of SNPs
75 associated with diseases or traits from GWAS. Here, we explored the distribution of GWAS
76 variants in the context of ancient conserved blocks of macrosynteny in the human genome.
77 Ancestral genome reconstruction (AGR) has been used to reconstruct the karyotypes of extinct
78 genomes including those of the early chordate (Putnam et al. 2008; Simakov et al. 2020) and
79 vertebrate ancestors (Nakatani et al. 2007; Sacerdot et al. 2018; Nakatani et al. 2021). AGR
80 views genomes as sets of macrosyntenic fragment adjacencies, and employs algorithms
81 enforcing local or global parsimony to reconstruct ancestral karyotypes at internal nodes (Damas
82 et al. 2021). AGR has uncovered that there are regions of strong macrosyntenic conservation
83 across vertebrate genomes that can be traced to the vertebrate ancestor and which are shared with
84 extant invertebrate deuterostome lineages such as scallop and amphioxus (Nakatani et al. 2007;
85 Sacerdot et al. 2018; Simakov et al. 2020; Nakatani et al. 2021).

86 By employing AGR approaches, studies have inferred that there were 17 or 18 ancestral
87 linkage groups (ALG's) in the proto-vertebrate and chordate ancestors (Sacerdot et al. 2018;
88 Simakov et al. 2020; Nakatani et al. 2021). However, a major difference between invertebrate

89 and vertebrate genomes, is that early during vertebrate evolution two rounds (1R and 2R
90 respectively) of whole-genome duplication (WGD) occurred at ~490 mya prior to the divergence
91 of agnathans and at ~438 - 485 mya prior to the divergence of cartilaginous and bony fish
92 (Simakov et al. 2020; Nakatani et al. 2021). This means that the post-2R vertebrate genome has
93 undergone two tetraploidization events relative to the pre-vertebrate ancestor; the non-random
94 distribution of paralogs harbouring subsets of these 2R-derived ohnologs could influence the
95 distribution of GWAS variants. Overall, ~ 35% of all human genes are estimated to have one or
96 more paralogous genes that were duplicated during 2R (so called ohnologs), (Makino et al. 2009;
97 Singh et al. 2015), and ohnologs are enriched for monogenic diseases (Chen et al. 2013).

98 AGR has identified the variation in the length and conservation of homologous syntenic
99 blocks (HSB's) across mammalian genomes, which are interspersed, by definition, with
100 evolutionary breakpoint regions (EBRs) (Pevzner and Tesler 2003; Bourque et al. 2005; Peng et
101 al. 2006; Alekseyev and Pevzner 2007; Becker and Lenhard 2007; Damas et al. 2021). HSB's
102 are blocks of sequence with conserved gene order and orientation in two or more species, while
103 EBR's are regions of nonalignment at a given resolution in between two HSB's (Larkin et al.
104 2009; Damas et al. 2021). To date, there has been no description of the evolutionary
105 characteristics of conserved syntenic blocks in the human genome relative to the pre-vertebrate
106 ALG's or to the amphioxus linkage groups, although such blocks are recognised (Marletaz et al.
107 2018; Simakov et al. 2020). Further, a new extensive genomic, methylation and transcriptomic
108 database for the Mediterranean lancelet (*Branchiostoma lanceolatum*) has provided genome-
109 wide insights into the genomic and regulatory changes associated with early vertebrate WGD
110 events and identified the pre-2R duplicated orthologues of genes in amphioxus with those in
111 diverse human species, including humans (Marletaz et al. 2018). These advances have made it

112 possible to more accurately identify conserved paralogons in the human genome and ohnologous
113 gene family members.

114 Thus, we broadly hypothesized that loci contained within ancient (i.e. pre-2R) conserved
115 linkage blocks in the human genome exhibit different evolutionary or functional characteristics
116 from those not contained in ALGs. Secondly, we hypothesized GWAS variants underlying
117 complex traits may be non-randomly distributed with respect to their location, and exhibit broad
118 functional clustering by ALG. To this end, we examined the following hypotheses: 1) the
119 distribution of GWAS-significant SNPs in the human genome is partially explained by the
120 evolutionary relationship of the genomic regions in which the SNPs are located relative to the
121 pre-2R ancestor, 2) there is enrichment for distinct ontological categories of traits residing in
122 different ALG's, 3) the level of evolutionary constraint on genes will be significantly impacted
123 by whether they are assigned an ALG and 4) the expression of genes associated with enriched
124 ancestry/trait ontological combinations may exhibit clustering in the human genome.

125

126 **Results**

127 *Identification of human genomic regions that can be conservatively traced to an ancestral*
128 *chordate linkage group.* We compared the ancestral origins of protein coding genes as predicted
129 by Sacercdot et al. (2018) and Simakov et al. (2020) to develop a map of the human genome
130 relative to those predicted to be in the pre-2R linkage groups. After adjusting for differences in
131 linkage group names, we found a high concordance between the assigned linkage groups of
132 genes in the two studies, as well as with those presented by Nakatani et al. (2021) (Table S1A).
133 Starting with the genes in common between the two reconstructions (n=4781), an additional
134 3273 genes were added from the Simakov et al. data that were not included in Sacercdot et al.,

135 14 genes were removed because of inconsistencies between the two models, and another 1795
136 genes were added to the map from the Sacerdot data that were not included in Simakov et al.
137 because they were predicted to be in the same ALG as the flanking genes in the Simakov et al.
138 model. Finally, based on gene family membership information from Marletaz et al. (2018), 560
139 additional genes were assigned to an ALG, giving a total of 10,410 genes being assigned to an
140 ALG (SuppFigS1, Table S1B). The number of genes mapping to each ALG ranged from 330
141 genes on ALG Q to 1162 genes on ALG A (Table S1C). Of the 10,410 genes with assigned ALG
142 ancestry, 7 mapped to the human Y-chromosome and were not included in further analyses due
143 to the lack of GWAS data. As expected, the fewest genes mapped to human chromosome 21
144 (n=73), and the most to chromosome 1 (n=1095), with a significant excess of ALG genes being
145 found of human chromosomes 2, 10, 14, and 15, and a significant deficit on chromosomes 6, 11,
146 19, 21 and X (see Table S1C).

147 In addition to the genes assigned to an ALG, we included 6,674 genes identified by
148 Marletaz et al (2018) as co-orthologs of human genes in amphioxus (“amphioxus orthologs”),
149 but which were not assigned to an ALG based on Sacerdote, Simakov or family information (i.e.
150 genes in Table S1B with a family ID but no assigned ALG). Some of these genes pertain to
151 tandem duplicates or regions of segmental duplication in the human genome (e.g. olfactory
152 receptors) since such genes do not exhibit the requisite conserved synteny over deeper
153 evolutionary time, but, in some cases one of them (the ancestral gene) will be included in the
154 ALGs. Lastly, 2,288 genes were not assigned ancestry based on Sacerdot, Simakov nor did they
155 exhibit orthology with genes in amphioxus and thus their location in the proto-vertebrate genome
156 is unknown (Table S1B, SuppFig. S1).

157 Some regions in the human genome contain > 10Mb of contiguous sequence derived
158 from a single ALG while others show evidence of high rates of fission/fusion over evolutionary
159 time. For example, there are large blocks of ALG A on chromosomes 1, 2, 14 and 20; ALG B on
160 chromosomes 2, 3, 5, 7 and 10; ALG D on chromosomes 16, 18 and 19; and ALG G on
161 chromosome 7 and 12, while chromosomes 1, 5, 9 and 18/19 harbour discontinuous segments of
162 ALGs C, L and/or M (Fig. 1a, Tables S1A, S1D, S1E), chromosomes 2 and 3 blocks of N/P and
163 chromosomes 7 and 12 of E/O. These regions of smaller, intercalated, syntenic blocks of 2 or 3
164 ALG's correspond to the descendants of linkage groups that became fused post - 1R (see
165 Nakatani et al. 2021, and Table S1A), and subsequently underwent local inversions etc. In
166 addition to differences in the level of conservation of contiguous blocks derived from a single
167 ALG's or a post-1R fused ALG, there are differences in the number of ohnolog descendants by
168 ALG which varied from 1.61 for ALG H to 2.07 for ALG B (Table S1F). Given that ohnologs
169 are, by definition, gene duplicates generated during 2R, it is not surprising that the odds of a gene
170 assigned to an ALG being an ohnolog is 4.4 times greater than genes not assigned to an ALG
171 (Odds ratio = 4.4, 95% CI: 4.12-4.71, $z=43.38$, $p<0.0001$).

172 Visualizing the gene families contained within blocks underscores the importance of 2R
173 in generating global syntenic clusters of ancient gene family duplicates. For example, blocks of
174 ALG B are conserved across parts of chromosomes 2, 3, 7, 12 and 17 with large blocks on
175 chromosomes 2 and 7, in which HOX clusters A, B, C and D reside on chromosomes 7, 17, 12
176 and 2 respectively (Fig. 1b, Table S1b). Blocks derived from ALG A are found on chromosomes
177 1, 2, 11, 14, 15, 19 and 20, including almost the entirety of human chromosomes 14: the
178 ohnologs of many gene families are distributed within these blocks (example *AKT*, *RYR*, *GREM*,
179 *SPRED*, *FAM98*, *RASGRP*, *MTA*, *AKT*, *FERMT*, Fig 1c, for a full list see Table S4b). Regions of

180 the human genome derived from the post 1R fused blocks of ALG C and ALG L (Fig 1d, S4b)
181 contain many other well-known ohnologous gene family members (*SKOR*, *SMAD*, *DOCK*, *JAK*,
182 *PDE4*, *PK3R*, *ADAMTSL*, *RAB3*, *TARS*) even when the gene names belie otherwise
183 (*INSR/INSRR/IGF1R*) (see Table S4b for a full list on putative ohnologs).
184 *Mapping and comparing the distribution of GWAS and HapMap variants in the human genome*
185 *relative to the ALG from which they are derived.* The GWAS catalogue contained 128,322
186 entries with 72,041 unique rsID's of which 49,129 could be assigned to an ALG and 22,912 fell
187 in regions of unknown ancestry (31.73%) (Fig. 2a, Table S2A). The density of GWAS
188 significant variants was compared to that available from CEU phase II of the HapMap project,
189 since this has been a common resource for GWAS imputation. To this end, 2,753,189 SNPs from
190 the HapMap project were assigned to ancestral blocks in the human genome: 883,156 SNPs fell
191 in regions of unknown ancestry (32.1%) and 1,870,033 to an ALG (Fig. 2a, SuppFig.2 for higher
192 resolution). A Pearson's Chi-squared test with Yates' continuity correction was performed to
193 assess individual deviation in the assignment of GWAS variants in each ALG relative to the
194 HapMap expected proportion per ALG. All of the ancestral ALG's except ALG E showed either
195 an excess (ALG's A, B, C, G, H, L, M, O and P) or deficit (ALG's E, F, I, J, K, N and Q) of
196 GWAS variants relative to that expected based on HapMap diversity per ALG (Fig. 2b, Table
197 S2b), except for variants falling into unassigned regions of the genome (Fig. 2b, Table S2b).
198 *EFO-analysis of GWAS traits.* Given that there is an excess of GWAS associated variants on
199 some ALG's and a deficit on others relative to HapMap variation, we next asked whether the
200 ontological terms assigned to the GWAS traits were non-randomly distributed with respect to
201 ALG's. Of the ~72,000 GWAS associated SNPs, 50,538 were associated with a primary
202 hierarchical term in the EFO classification (Tables 3Sa). EFO categories with fewer than 300

203 observations were removed or combined with other categories leaving 14 EFO categories:
204 anthropometric, behavioural, biochemical, cardiovascular, digestive-endocrine, hematological,
205 hepatobiliary, lipid, musculoskeletal, nervous, respiratory, immune-infectious, urinary, visual
206 which were collectively associated with 42,099 SNPs. Using a global FDR, chi-square analysis
207 indicated that all EFO traits except urinary system were enriched or underrepresented with at
208 least one but not more than four ALG's (Table S3b). Odds ratio (+/- 95% confidence intervals)
209 were calculated for the 29 associations exhibiting an excess or deficiency of GWAS variants for
210 a given EFO x ALG combination (Table S3c, Figure 3). This revealed the ALG A is enriched
211 for genes in nervous and visual traits, ALG B for musculoskeletal traits, ALG C for nervous
212 system traits, ALG D for lipid and hematological traits, ALG I and Q for cardiovascular traits,
213 and that genes that were unassigned to an ALG were enriched for endocrine and immune system
214 related traits. (Figure 3, Table S3b)). The list of genes associated with each GWAS variants for
215 each of the 14 EFO categories is given in Supplementary Tables 3.

216 *Gene level analyses.* Next, we employed data from gnomAD, GTEx (v8) and Ensembl to explore
217 characteristics of genes: i) assigned to an ALG or ii) associated with GWAS hits (SuppFig. 1).
218 For the analyses of genes assigned to an ALG, genes assigned as unknown and having orthologs
219 in amphioxus were grouped. We identified 11,269 genes associated with the 72,041 unique
220 GWAS variants: unsurprisingly, the number of associations per gene is highly skewed with 4091
221 genes having < 5 associations, and 199 having >30; each gene was counted once for the gene-
222 level analyses (Table S4a). Of the 11,269 GWAS associated genes and the 10,410 genes
223 assigned to an ALG, 11,055 and 10,366 were annotated in the gnomAD V2.1.1 database
224 respectively (n=18,766 total genes at gnomAD, TableS4b). The odds of a GWAS associated
225 gene being assigned to an ALG was 1.64 (CI: 1.54-1.73, z=16.4 p<0.00001, Table S4a), and

226 64.2% of the genes assigned to an ALG were associated with a GWAS, while 52% of the genes
227 unassigned to an ALG were associated with a GWAS. In agreement, with the distribution of
228 GWAS vs HapMap SNPs (Fig. 2b) some ALG's had more GWAS associated genes than
229 expected (ALG's A, B, C, G, M, O and P), while others had fewer (ALG's F, J, K and N)
230 (comparison of Figures 2b and 4A). The genes with the largest number of GWAS signals that
231 were unassigned to an ALG included immune related genes (particularly HLA in the MHC
232 region, and interleukin ligands and/or receptors), a cluster of overlapping genes involved in
233 glucuronidation on chr 2 (*UGT1A4-UGT1A10*), and a variety of endocrine related genes (several
234 related to obesity or lipid metabolism, including *FTO*, *LEPR*, *LPA*)(Table S4b). This is
235 consistent with the EFO analysis that found that immune and endocrine functions are enriched
236 for genes not assigned to an ALG. Notably, orthologs of the ancestral MHC locus have been
237 difficult to identify and, as such, the HLA-loci are not included in the map presented here.
238 However, current data suggests that the MHC locus was present in the ancestral vertebrate
239 genome but acquired its unique function in antigen recognition after the divergence of agnathans
240 (Flajnik and Kasahara 2010; Flajnik 2018; Ohta et al. 2019; Nakatani et al. 2021). The orthologs
241 of these MHC genes would map to ALG M.

242 *Gene level constraint.* Genes that trace their ancestry to an ALG and those associated with a
243 GWAS variant may exhibit constrained evolutionary change. Using pLI scores >0.9 as a proxy
244 for haploinsufficiency, we find that genes that trace their ancestry to an ALG (OR=2.58; CI 2.37-
245 2.81, $z=22.00$, $p<0.0001$) or are associated with a GWAS variant (OR = 2.02, CI 1.85-2.2, $Z =$
246 15.9, $p<0.0001$) (Figure 4b) are both enriched for genes that are intolerant to loss of function
247 mutations. To further explore the characteristics of these genes, we performed logistic regression
248 to explore the effect of six variables on ALG or GWAS association: i) gene length ii) coding

249 sequence (cds) length iii) number of coding exons iv) loss-of-function observed/expected upper
250 bound fraction (LOEUF) v) the mean number of nonsynonymous mutations (dN) for a gene
251 between human and mouse and vi) the median expression (in transcripts per million, TPM) in the
252 highest tissue of expression listed in GTEx. After checking for collinearity and removing
253 variables with low explanatory power, the logistic regression showed that genes assigned to any
254 ALG had i) a lower rate of nonsynonymous substitutions ($\beta=-6.08$, $\chi^2 = 644.42$, $p<0.00001$), ii)
255 more coding exons ($\beta=0.088$, $\chi^2 = 558.43$, $p<0.00001$), iii) longer cds length ($\beta=-0.00035$, $\chi^2 =$
256 237.7 , $p<0.0001$), iv) higher median gene expression ($\beta=-0.00010$, $\chi^2 = 14.1$, $p<0.0002$), and iv)
257 lower LOEUF ($\beta=-0.26$, $\chi^2 = 21.26$, $p<0.0002$), (Fig 4b, c, SuppFig3, d, e, respectively). On the
258 other hand, genes associated with GWAS variants had i) longer gene length ($\beta=0.00000084$, $\chi^2 =$
259 644.3 , $p<0.0001$), ii) lower LOEUF scores ($\beta = -0.38$, $\chi^2 = 92.35$, $p<0.0001$), and iii) longer cds
260 length ($\beta=0.00012$, $\chi^2 = 70.60$, $p<0.0001$) (Figures 4 F, H, SuppFig.3 and 4). Collectively, this
261 shows that genes mapping to ALGs evolve more slowly on both shorter (LOEUF) and deeper
262 (dN mouse-human) evolutionary timescales, have more exons and longer coding sequences and
263 higher levels of gene expression than genes not mapping to an ALG. On the other hand, genes
264 associated with GWAS are longer than those not associated with a GWAS, and exhibit greater
265 functional constraint as measured by both pLI and LOEUF but not dN.

266 *Mapping of GWAS associated genes involved in lipid metabolism.* The finding that there are
267 blocks of genes with shared ancestry distributed throughout the genome which are enriched for
268 ohnologs, evolve slowly and show evidence of clustering by functions based on the EFO
269 analysis, suggests that genes that trace their ancestry to an ALG may be involved in core
270 biological processes. As a test of this hypothesis, we took the list of putatively core genes
271 involved in lipid metabolism (Dron and Hegele 2016) and placed them within their human

272 genome context coloured by ancestral ALG and labelled by the tissue of highest expression for
273 each gene in the GTEx catalogue. This reveals that 18 of 19 of the genes involved in lipid
274 metabolism emanate from 5 ALG's (one gene, *APOC2* has unknown origin), and several of them
275 are 2R ohnologs. For example, *APOA1*, *APOA5*, *APOE* and *APOC2* exist in two paralogs on
276 chromosomes 11 and 19 that emanate from ALG A; *LPL*, *LIPC* and *LIPG* are 2R ohnologs on
277 chromosomes 8, 15 and 18 respectively derived from ALG C; *LDLR* and *LDLRAD1* are
278 paralogues derived from ALG L that now reside on chromosomes 1 and 19. Lastly, unrelated
279 genes, *SARIB* and *MTTP* reside on chromosomes 5 and 4 respectively, yet they are all derived
280 from ALG Q in the pre-vertebrate ancestor. Although the remaining four genes – *LCAT* and
281 *CETP* (ALG D), *PCSK9* and *ANGPTL3* (ALG L) are not paralogs, they share current (physical
282 linkage) as well synteny in the pre-vertebrate genome.

283 Similarly, in the EFO analysis, we found that ALG A and C were enriched for nervous
284 and/or visual traits and ALG E for behavioural traits. Given previous suggestions that genes
285 within HSB are enriched for neurodevelopmental functions, we looked for evidence that genes
286 on these three ALG's had brain as the tissue of highest gene expression. For genes without ALG
287 ancestry, 18.5% (1680/9080) have the highest tissue of expression as brain versus 21.5% of those
288 assigned to an ALG (2086/9685) with values ranging from 14.4% on ALG K to 25% on ALG Q.
289 For the three ALG's exhibiting enrichment for brain-related traits, the proportion of genes with
290 highest expression in brain were 24.1% (ALG A), 22.9% (ALGC) and 22.6% (ALG E).
291 Visualization of the density of genes on hg38 (gray histogram), with the density of genes with
292 highest expression in brain (overlaid red histogram), as well as the median TPM expression of
293 brain-related genes on ALG A, C and E (blue dots to the left of the chromosomes) suggests that
294 some regions of the human genome may harbor clusters of brain related traits. For example,

295 large blocks of genes with highest expression in brain derived from ALG A (dark blue bands on
296 chromosomes) are present on chromosomes 1 and 2, from ALG C (burgundy bands) on
297 chromosomes 5 and 15, and from ALG E (pink bands) on chromosome 3 (Figure 5B).

298

299 **Discussion**

300 Although the primary focus of GWAS is to identify causal loci for complex diseases and traits,
301 an understanding of the population and evolutionary processes that maintain heritable variation is
302 essential to their interpretation (Sella and Barton 2019). The analysis of highly studied GWAS
303 traits has uncovered that the genetic basis of many complex phenotypes and diseases is
304 extremely polygenic, highly pleiotropic and, at least in the case of height and body mass index,
305 predominantly influenced by stabilizing (rather than directional) selection (Shi et al. 2016;
306 Simons et al. 2018; Sella and Barton 2019). As more data are collected about the biological
307 function of genes, other factors are coming to the forefront, particularly the role of allelic effects
308 (cis and trans-regulatory elements) on gene expression (Boyle et al. 2017; Liu et al. 2019), which
309 helps to explain the large number of GWAS loci that fall in regulatory regions (Visscher et al.
310 2012; Visscher et al. 2017). At the level of the gene, molecular evolutionary approaches to
311 predicting the functional impact of variants has also revealed that i) biases and differences in the
312 mutational process influence the probability of observing specific variants and ii) genes differ
313 markedly in their average rate of evolutionary change, with slower evolving genes exhibiting
314 greater functional constraint and greater susceptibility to deleterious mutations, a fact that is the
315 basis for many variant interpretation tools (reviewed in (Eilbeck et al. 2017)).

316 At the genome level, approaches in comparative syntenic analyses and ancestral genome
317 reconstructions have identified so-called homologous syntenic blocks (HSBs) and evolutionary
318 breakpoint regions (EBRs) in amniote genomes (Peng et al. 2006; Alekseyev and Pevzner 2007;
319 Kemkemer et al. 2009; Larkin et al. 2009; Alekseyev and Pevzner 2010; Damas et al. 2018;
320 Farre et al. 2019), as well as the imprint of the early vertebrate WGD events (Nakatani et al.
321 2007; Putnam et al. 2008; Sacerdot et al. 2018; Simakov et al. 2020; Nakatani et al. 2021). The
322 location of breakpoints associated with interchromosomal exchange in mammals is non-random,
323 indicating that selection acts at the genome level to maintain combinations of genes and their
324 regulatory elements that are essential for biological organization (Damas et al. 2021). This is in
325 agreement with studies that show that the 3D organization of chromatin in the nucleus and the
326 GC peaks and valleys of topologically associated domains (TADs) strongly influence the number
327 of contacts and inter-chromosomal cross-over events (Berthelot et al. 2015; Jabbari et al. 2019).

328 In this study, we aimed to define the conserved 2-R derived paralogs in the human
329 genome and then examine the distribution of GWAS variants within them. To this end, we
330 generated a map of the human genome with respect to 17 inferred ancestral linkage groups
331 (ALGs) by comparing two reconstructions of the ancestral vertebrate (Putnam et al. 2008;
332 Sacerdot et al. 2018) and one of the ancestral chordate karyotype (Simakov et al. 2020). We
333 relied primarily on the Simakov et al. (2020) model of the ancestral chordate karyotype because
334 their reconstruction is based on a new chromosome level assembly of the amphioxus genome, a
335 basal deuterostome with similar linkage groups to the inferred chordate ancestor. A new
336 reconstruction of the proto-vertebrate karyotype by Nakatani et al., (2021) inferred there to be 18
337 linkage groups in the vertebrate ancestor; their proto-vertebrate chromosomes (pvc's) are highly
338 concordant with those described here, and similar to those present in scallop, an early diverging

339 deuterostome, as well as to the vertebrate ancestor genome inferred by Sacerdot *et al.*, (2018).
340 The additional linkage group they infer is a microchromosome that, based on the Simakov *et al.*
341 (2020) model, corresponds to a small block of genes belonging to ALG A, and which lacked
342 clear synteny with gnathostome or invertebrate genomes (Table S1a, pvc18). Although Nakatani
343 *et al.* (2021) do not provide a list of the human genes that they infer to have putative co-orthologs
344 in the pvc, they provide the linkage groups of 19,969 human genes in the proto-gnathostome
345 genome, which we used to examine the concordance of their gene assignments to those presented
346 here (Table S5a). Of the 10,410 genes assigned to an ALG in this study, 7616 (73.2%) were
347 concordant with and 580 (5.6%) were not included in the Nakatani *et al.* (2021) study, while
348 2214 (21.2%) were assigned to either a different ALG or could have been assigned to the same
349 or different ALG (Table S5b). The ambiguity in some gene assignments is due to fusion of some
350 of the pvc's in the post-1R or post-2R genomes as presented in Nakatani *et al.* (2021) (see Table
351 S1a).

352 Based on a thorough study of the orthology of genes between Mediterranean lancelet and
353 vertebrate taxa, we find an additional 6,647 human genes that have a putative ortholog in
354 amphioxus that are not present in conserved ALG's. Many of these genes belong to tandemly
355 duplicated gene clusters which have likely duplicated in the post-2R genome, such that only one
356 gene will map to an ALG. This underscores that the map we generated represents the conserved
357 macrosyntenic clusters in the human genome, rather than all of the genes that originate from the
358 pre-2R genome. On the other hand, some genes that are small and fast evolving may not be
359 represented in the ALG's even if they diversified during 2R. For example, we previously
360 performed detailed phylogenetic and syntenic mapping of the insulin superfamily peptides
361 (Yegorov and Good 2012), and found that the ancestral INS-IGF linked genes trace to ALG O

362 and those of the relaxin like peptides to ALG L. However, these genes are not included in the
363 Sacercdot or Simakov reconstructions. With the availability of improved genome assembly and
364 annotation pipelines, such as that being developed by the Vertebrate Genome Project (Rhie et al.
365 2021), there will be more data with which to assess if genes are “ancient” or have arisen post 2-
366 R by another process such as *de novo* origination, local duplication, and/or a retroposition event
367 (Casola 2018; Van Oss and Carvunis 2019; Cosby et al. 2021). This will allow a deeper
368 understanding of the nature of conserved versus adaptive regions of vertebrate genomes, and
369 facilitate identification of regions in the human genome that would be susceptible to deleterious
370 consequences from structural variants. A motivation for basing our analysis on the ancient pre-
371 2R linkage groups, is that it also facilitates the identification of ohnologs in the human genome
372 which could further help prioritize disease genes (Makino and McLysaght 2010; Chen et al.
373 2013; Lal et al. 2020)

374 We find that approximately half of the protein coding genes in the human genome can be
375 reliably traced to their pre-vertebrate origins and that genes assigned to an ALG have a slower
376 rate of evolution as measured on both deeper (dN) and more recent (LOEUF) timescales.
377 Further, the paralogons harbouring the 2-R derived genes are found in blocks (in which gene
378 memberships but not necessarily gene order is conserved) of varying size with some parts of the
379 human genome showing large regions of conserved ancestry. By defining the chromosomal
380 coordinates of the ALGs, we show that all but one of the ALG’s (ALG E) and unassigned
381 regions, are either enriched or depauperate in GWAS loci relative to that expected based on
382 variation in the human genome identified by HapMap v2. Even though the selection of loci by
383 the HapMap project was subject to ascertainment bias (Clark et al. 2005), the concordance
384 between the GWAS and HapMap variants in unassigned regions, supports the inference that the

385 GWAS signals are non-randomly distributed with respect to ALGs. By further linking the
386 GWAS loci to their closest gene and associated experimental factor ontology (EFO) terms, we
387 determined that all of the EFO terms except one (Urinary) are enriched for ancestry of genes on
388 one to four ALGs, while genes that were unassigned to an ALG are enriched for immune and
389 endocrine system functions. This suggests that the human genome has retained blocks of genes
390 associated with distinct functions that emanate from conserved linkage groups in the vertebrate
391 ancestor. This is concordant with studies in amniotes showing that genes in HSBs are more likely
392 to be involved in conserved processes such as development while those in EBR's are more likely
393 to be involved in adaptive processes (Larkin et al. 2009; Damas et al. 2021).

394 Examination of the distribution of the ALG's in the human genome, indicates high levels of
395 intermixing of some ALG's, such as ALG's C/L/M, ALG's E/O, ALG's N/P and ALG's D/J. Based on
396 the work of Nakatani et al's (2021), it is now clear that these associations are due to the gnathostome
397 specific fusions of duplicates of these ALG's following 1R. They infer that following the 1R
398 tetraploidization event, several of the resulting 36 chromosomes were fused giving rise to a post-1R
399 karyotype of n=25 (see Table S1A), while after 2R, an additional few fusion events occurred leading to a
400 karyotype of n=49, including several microchromosomes. Some of these post-1R fused chromosomes
401 maintained high levels of conservation (e.g. C/M/L and E/O) in the human genome.

402 We found a strong signal that genes that map to an ALG evolve significantly more slowly
403 than genes that do not. Using logistic regression, we find that the largest contributor to ALG vs
404 non-ALG gene assignment was the rate of nonsynonymous substitution of the genes between
405 mouse and human. Further, the level of evolutionary constraint measured LOEUF was higher for
406 genes assigned to ALG's (median LOEUF = 0.38) than those associated with a GWAS loci
407 (median LOEUF= 0.42) and they exhibited greater intolerance to loss of function mutations
408 (pLI). Genes assigned to an ALG also had longer mean cds length and more coding exons

409 (which had a Pearson correlation coefficient of $r^2=0.80$), as well as a higher median transcript
410 expression in the tissue of highest expression based on GTEx v8.0 data. These results are
411 consistent with studies that show that highly expressed genes show greater evolutionary
412 constraint and are more likely to exhibit dosage sensitivity (Duret and Mouchiroud 2000; Pal et
413 al. 2001; Rocha and Danchin 2004; Gout et al. 2010), and further suggests that the genes with
414 deep ancestry in the human genome are likely to be enriched for core genes, as per the
415 conceptual framework of the omnigenetic model (Boyle et al. 2017; Liu et al. 2019)

416 On the other hand, genes associated with GWAS loci, had an approximately twofold
417 higher odds of being haploinsufficient as measured by pLI (OR=2.02) and significantly lower
418 LOEUF scores (median LOEUF=0.42) than non-GWAS genes (median LOEUF=0.57). The pLI
419 statistic has been interpreted as a measure of haploinsufficiency, but this only holds when the
420 selection coefficient is high (i.e. $s=1$), which occurs frequently for Mendelian disease genes, but
421 not other genes (Fuller et al. 2019). When $s < 1$, pLI measures the average effect of loss of
422 function mutations in heterozygotes (i.e., hs). Apart from exhibiting high pLI scores, the best
423 predictor of being a GWAS associated gene was gene length, which is unsurprising since GWAS
424 loci are assigned to the nearest gene.

425 The finding that genes exhibiting deep ancestry in the human genome have strong
426 functional constraint, are highly enriched for ohnologs (OR = 4.4) and have higher median gene
427 expression, suggests that they should be enriched for disease associated loci. Overall, the odds of
428 a gene assigned to an ALG being associated with a GWAS signal was 1.64 (CI:1.54-1.79).
429 Studies have shown that genes associated with neurodevelopmental disorders (NDDs) exhibit
430 greater evolutionary conservation of copy number, i.e. are intolerant to copy number variation,
431 across mammals indicating that they are subject to greater functional constraint (Rice and

432 McLysaght 2017), and that most (perhaps 80%) known monogenic disease genes belong to gene
433 families consisting of ohnologs and/or paralogs (Dickerson and Robertson 2012; Chen et al.
434 2013). Recently, Lal et al., (Lal et al. 2020) examined whether inclusion of gene family
435 information could improve variant interpretation for missense mutations in patients with NDDs,
436 given that NDDs are phenotypically and genetically heterogeneous. They found that by including
437 information about whether variants changed paralog-conserved vs paralog non-conserved
438 residues helped to prioritize missense variants with potential functional consequences. Using the
439 same reasoning, information about membership of genes/ohnologs in ALGs may help prioritize
440 genes influencing the same or similar phenotypic traits. For example, we demonstrated that even
441 though the core genes underlying lipid metabolism are distributed on 11 chromosomes, they
442 emanate from 5 ALG's, harbour 3 ohnologs, and have the highest tissue of expression in a
443 logical tissue (e.g. liver, adrenal gland, spleen, adipose etc). Thus, by adding the information
444 regarding the deep ancestry of genes into our understanding of gene and genome organization,
445 this approach may help to identify regions of susceptibility to structural variants and/or to be
446 another tool to inform variant prioritization.

447 The GWAS catalogue grows at an increasing pace and provides a valuable resource to
448 summarise and interpret the results of GWAS. However, there are, of course, limitations of using
449 the catalogue. Among the more salient caveats is that because the studies deposited into the
450 catalogue were generated using many SNP arrays, different but highly correlated SNPs in LD
451 with each other may appear as independent signals when a single SNP may have been identified
452 if all studies used the same array. Secondly, traits that are easier to measure, such as height or
453 weight, will have much larger study sizes than studies of other traits, and this will provide greater
454 information about these traits than for rare phenotypes. These two caveats will compound each

455 other such that signals at specific genes could be driven by the large number of studies for a trait,
456 measured on a variety of different panels across studies. We tried to minimize this effect by
457 including genes associated with GWAS loci only once regardless of the number of associations.
458 Additionally, we included SNPs in the study based on them having a genome-wide significant p-
459 value ($< 5.0 \times 10^{-8}$). However, the p-value is a crude measure of effect and does not take into
460 account sample size, effect size, nor minor allele frequency into account. Further, the use of the
461 HapMap as a reference for SNP density is most relevant for early GWAS (since early array's
462 typically genotyped variants based on HapMap, and there was limited imputation), but is less
463 appropriate for more recent GWAS which use high-resolution imputation (e.g. TOPMed, (Taliun
464 et al. 2021)) or even whole genome sequencing (WGS). But since each study used different
465 GWAS arrays (often multiple) and in some cases different imputation panels it is not possible to
466 come up with a simple solution for this. Despite these limitations, the data presented here
467 indicate that there are genome-wide influences on the distribution of GWAS variants that can be,
468 partially, explained by the deep and duplicated ancestry of our underwater ancestors.

469 **Methods**

470 We combined information from two studies to create a map of the human (*Homo sapiens*)
471 genome with reference to hypothesized linkage groups in the pre-duplicated ancestral genome.
472 Firstly, we used the reconstruction inferred by Sacerdot et al. (2018); they employed the
473 Algorithm for Gene Order Reconstruction in Ancestors (AGORA) (Berthelot et al. 2015) based
474 on local parsimony to identify contiguous ancestral regions (CARs) which they assembled into
475 17 hypothetical vertebrate ancestor chromosomes (VACs). As part of this, they employed the
476 complete set of phylogenetic trees in Ensembl 69 from 45 extant amniotes to identify ohnologs –
477 genes that were duplicated and retained following 2R in the Amniota ancestor – using criteria

478 such as putative ohnologs must be located at least 2N (90) genes apart on the same CAR or on
479 separate CAR's, and that the estimated time of duplication of duplicates should fall within the
480 1R-2R window (~ 550mya). Secondly, we employed data from the study by Simakov et al.
481 (2020) in which they reconstructed the ancestral chordate linkage groups (CLG's), i.e. an earlier
482 timepoint, based on a new high-depth sequencing of the *Branchiostoma floridae* (Amphioxus)
483 genome. To reconstruct the CLG's, they compared macro-syntenic blocks of genes between
484 Amphioxus with an invertebrate (scallop) against three post-2R vertebrate genomes (chicken,
485 frog and gar). Conservation of macro-synteny was estimated using a synteny discontinuity
486 measure, D, based on the squared Euclidean norm of the difference between left and right
487 windowed averages of the synteny indicator vector between Amphioxus and the four other taxa.
488 In this way, they identified the set of genes and their linkage groups that were present in the
489 chordate ancestor as well as the sets of genes, emanating from a single CLG, that were putatively
490 duplicated in the post-2R vertebrate genome (i.e. ohnologous gene families). They infer that the
491 chordate ancestral genome was also composed of 17 linkage groups. *B. floridae* has 19
492 chromosomes: they find that all of its' chromosomes, except *B. florida* (BFL) linkage groups 2, 3
493 and 4 (which exhibited 4, 1, and 1, breaks respectively) exhibited no discontinuity breaks with
494 later diverging 2R taxa, indicating deeply conserved macro-synteny during chordate evolution,
495 even though local synteny is not well conserved due to copious smaller scale inversions, gene
496 loss and local duplications as well as subsequent chromosomal fission and fusion events during
497 vertebrate evolution.

498 *Construction of a map of the human genome based on the ancestral chordate genome.* We
499 downloaded the complete set of genes traced to the pre-vertebrate ancestral genome by Sacercdot
500 et al. (2018) from (<ftp://ftp.biologie.ens.fr/pub/dyogen/genomicus/69.10/> last accessed, October

501 [2020](#)). In total, 5,051 gene families were assigned to one of the ancestral linkage groups 1-17 and
502 an additional 4773 were assigned to the ancestral genome but the linkage relationships were not
503 defined. 4,968 of these families had at least one descendent in the human genome, with an
504 average number of 1.69 paralogues per family, such that 8,389 protein coding genes from the
505 human genome were assigned to one of the 17 VACs (number 1-17). Of these 8,389 genes,
506 definitive gene names and locations were identified for 8,242 genes that were annotated by
507 ENSEMBL based on build 38 of the human genome including all 22 autosomes and the X-
508 chromosome (Table S1b).

509 Next, we downloaded the information from Simakov et al. (2020) who assigned ancestry
510 of human genes relative to one of the 17 hypothetical CLG's (number A-Q)
511 (<https://bitbucket.org/viemet/public/src/master/public>, date accessed=November 2020). In total,
512 they include 6082 gene families with representatives in the human genome, with an average
513 number of 1.5 paralogues per family resulting in 9106 human protein coding genes. After
514 removal of duplicates and selecting genes that were annotated in ENSEMBL, we obtained 5853
515 gene families representing 8065 genes. Using these two gene lists, we created a consensus map
516 of the human genome relative to the ancestral linkage groups (ALGs) by using the Simakov-
517 CLG ancestry assignment as a reference point and then adding or removing genes based on
518 consistency with their assigned ancestral linkage group in Sacercdot et al. (2018). Specifically,
519 we added genes included in the Sacercdot-VAC data if the genes fell within a contiguous stretch
520 of the human genome mapping to equivalent regions in both studies. On the other hand, in cases
521 where the two reconstructions differed in the assigned ancestry of a gene, gene ancestry was
522 treated as unknown. For breakpoints, i.e. where there was a change in the ancestry of an ALG,
523 the beginning and end of the block ID were assigned to be 10 kb from the nearest (start or end)

524 of a gene and intervening genomic regions were treated as having unknown ancestry. (Map:
525 CLG_Map_locations_23-new.csv). To provide a third level of information, we used information
526 about family membership based on orthologous gene comparisons between the human and
527 amphioxus genomes as shown in Marletaz et al. (Marletaz et al. 2018). In cases in which the
528 Marletaz study demonstrated that genes families with up to four members situated on different
529 human chromosomes or more than 100Kb apart on the same chromosomes were putative
530 ohnologs, the inferred mapped position based on other members of the gene family was inferred.
531 This added additional ALG membership for 560 genes. Using the list of orthologous genes
532 between human and amphioxus based in the Marletaz study (2018), genes with putative
533 orthologs in amphioxus but which were not mapped to an ALG based on the above criteria were
534 categorized as “amphioxus orthologs”. Lastly, genes for which no ALG nor amphioxus ortholog
535 could be identified were classified as having “unknown” ancestry.

536 *Mapping GWAS and HapMap variants to the ancestral linkage groups.* Next, we
537 downloaded the complete list of single nucleotide polymorphisms (SNPs) that have been
538 identified in Genome Wide Association Studies (GWAS) and deposited into the GWAS
539 catalogue from <https://www.ebi.ac.uk/gwas/> (file: GWAS_catalog_v1.0.2-
540 associations_e100_r2020-07-14.tsv, last accessed on July 14, 2020). Using their genomic
541 coordinates on build 38, we mapped all GWAS SNPs with a unique rsID and an association p-
542 value $< 5 * 10^{-8}$ to the ALG map generated above and then counted the number of SNPs falling
543 into each ALG (Buniello et al. 2019). To compare the density of human genetic variants
544 identified in GWAS to overall diversity in the genome, we used Phase II r24 of autosomal and
545 X-chromosome variants from the HapMap Utah residents (CEPH- yes; Centre d'Etude du
546 polymorphisme humain) with northern and western European ancestry (CEU) obtained from

547 ftp://ftp.ncbi.nlm.nih.gov/hapmap/genotypes/2008-10_phaseII/ and employed liftover
548 (<genome.ucsc.org>) to update the variant coordinates from hg36 to hg38. A total of 2,753,012
549 HapMap SNPs were available (for both autosomes and X, 275 did not lift over) from which we
550 calculated i) the number of SNPs in 200 kb windows and ii) the number of HapMap SNPs
551 mapping to each CLG.

552 *Ontological classification of GWAS catalogue phenotypes.* The GWAS catalogue reports
553 associations for a wide range of phenotypes including both disease (e.g. Type 2 diabetes) and
554 non-clinical phenotypes (e.g. hair color). All curated traits included in the catalogue are mapped
555 to terms from the Experimental Factor Ontology (<https://www.ebi.ac.uk/efo/>). In most cases,
556 traits are mapped to a singular phenotypic class, although some traits map to 2-3 EFO terms,
557 with most mapping to one of the 27 higher level terms but sometimes to a sub-classification
558 term. We selected only variants associated with the highest order EFO terms per category, and
559 assessed association between EFO traits and ALG ancestry using a global chi-square analyses set
560 with an FDR < 0.05. For those traits that rejected the null hypothesis based on the global FDR, a
561 2x2 Odd's ratio with 95% confidence intervals were calculated for those traits/ancestry
562 combinations whose standardized residuals were $> \pm 3.84$.

563 *Gene-level analyses: Characteristics of genes with assigned ancestry or associated with a GWAS*
564 *variant.* We tabulated the possible protein-coding genes and the frequency of association for
565 each GWAS variant from the GWAS catalogue. For SNPs that are intra-genic, the GWAS
566 catalogue reports a single putative associated gene, the one flanking the focal SNP, while for
567 inter-genic SNPs, the catalogue typically reports a minimum of two protein coding genes (those
568 immediately up and down-stream of the variant), plus additional non-coding or splice variants.
569 We retained a singular gene for intra-genic SNPs and no more than two protein-coding variants

570 per SNP, and counted the number of times associations were observed. To assess characteristics
571 of genes associated with GWAS or that have an ALG ancestry, we merged this list and the list of
572 genes that were assigned to an ALG (tabulated by ENSEMBL ID's), with a database of gene-
573 level metrics of all human protein coding genes curated at gnomAD v 2.1.1
574 (gnomad.v2.1.1.lof_metrics.by_gene) (last accessed, July 15, 2020).

575 ***Gene-level Constraint.*** The degree of evolutionary constraint of GWAS-associated genes and
576 those with ancestry assigned to an ALG was compared using two metrics of the relative
577 intolerance to loss of function of mutations for genes i) associated with GWAS studies ii)
578 mapping to an ALG and iii) those with unassigned ancestry. Firstly, we compared the
579 distribution of the statistic, “loss-of-function observed/expected upper bound fraction (LOEUF)”.
580 This statistic is akin to a measure of haploinsufficiency on a continuous scale, in which genes
581 with a lower value evolve under strong selection against predicted loss-of-function variation, and
582 those with values = 1 are not constrained. Secondly, we employed a dichotomous categorization
583 of genes as being haploinsufficient if the probability of being loss of function intolerant (pLi)
584 was >0.9. Values of LOEUF and pLI were obtained from gnomAD V 2.1.1 (Karczewski et al.
585 2021) (last accessed, July 15, 2020).

586 ***Expression of genes on ALG's enriched for EFO terms.*** Given the finding that genes mapping to
587 different ALG's were associated with different traits based on the EFO classification, we
588 examined the relationship between ALG enrichment and gene expression in tissues relevant to
589 EFO enrichment. To this end, we obtained the median transcripts per million (TPM) for all
590 genes and tissues included in GTEx V8 (gtexportal.org, last date accessed, November 15, 2020),
591 and recorded the tissue with the highest median TPM per gene. Tissues were grouped into
592 biologically meaningful groups a) Adipose (sub-cutaneous and visceral (omentum)) b) Blood-

593 Vascular (Whole Blood, arteries (Coronary, Tibial), heart (left ventricle; atrial appendage) c)
594 Brain (all sections) d) Cells (EBV-transformed lymphocytes and cultured Fibroblasts) e)
595 Digestive System (Esophagous, stomach, colon, intestine, pancreas) f) Female Reproductive
596 system (breast, uterus, cervix, vagina, fallopian tube) g) Lung h) Muscle i) Nerve Tibial, j) Skin
597 (exposed and not-exposed) k) spleen l) Male Reproduction (Testis, prostate) m) Glands (Thyroid,
598 Pituitary, Adrenal and Salivary) n) urinary (Bladder, Kidney). Overall patterns of gene
599 expression by ALG and tissue/EFO terms were explored; in particular, the relationship between
600 genes mapping to ALG's that were enriched for the EFO terms "brain" and "visual traits" were
601 assessed for expression of genes in nervous tissues. Additionally, differences in gene expression
602 for genes with and without ALG-ancestry were examined by GTEx grouped tissue. Lastly, the
603 ALG ancestry and tissue of expression for genes identified as being core genes in lipid
604 metabolism (Table 1) by Dron et al., (2016) were extracted (Dron and Hegele 2016).

605 *Statistical and graphical Analysis.* Fisher's Exact tests were used to compare associations of
606 genes with different classifications (e.g. associated with a GWAS or not associated, mapping to
607 an ALG or not-mapping). Odd's ratio (OR) and 95% confidence intervals calculated from 2x2
608 tables according to Altman (1991) for rows with complete data (no NA's in any cell counted for
609 the 2x2 table). A global G-test and individual Pearson Chi-Square tests with Yates continuity
610 correction were performed to assess the difference in observed (GWAS) vs expected (HapMap)
611 proportion of SNPs mapping to the ALGs. Kolmogorv-Smirnov tests were performed to compare
612 the distribution of continuous variables (e.g. LOEUF) between discrete classes. Logistic
613 regression was used to assess the influence of six gene traits or measures of evolutionary
614 constraint (all uncorrelated with pearson's $r^2 = 0.81$ or less) on being associated with an ALG or
615 GWAS variant, and removing terms that had low explanatory value and using an FDR = 0.05. T-

616 tests or Mann-Whitney U tests were used to compare the means of normally and the ranks of
617 non-normally distributed variables respectively. Chromosome plots were generated with R
618 packages BioCircos (Cui et al. 2016) and ChromPlots (Orostica and Verdugo 2016).

619 **Availability of data and materials**

620 The source data and R and python codes used for all analyses are available at
621 <https://github.com/saravictoriagood/GWAS-origins>

622 **Competing interests.**

623 The authors declare that they have no competing interests.

624 **Acknowledgements**

625 The authors would like to thank Jaanus Sarvoli for comments on an earlier version of this
626 manuscript. The research was supported by a CIHR Strategic Training in Advanced Genetic
627 Epidemiology (STAGE) fellowship to SVG, by an NSERC – DG to SVG, and CIHR project
628 grant to ADP.

629

630

LITERATURE CITED

- 631 Adzhubei I, Jordan DM, Sunyaev SR. 2013. Predicting functional effect of human missense mutations
632 using PolyPhen-2. *Curr Protoc Hum Genet* **Chapter 7**: Unit7 20.
- 633 Alekseyev MA, Pevzner PA. 2007. Are there rearrangement hotspots in the human genome? *PLoS*
634 *Comput Biol* **3**: e209.
- 635 Alekseyev MA, Pevzner PA. 2010. Comparative genomics reveals birth and death of fragile regions in
636 mammalian evolution. *Genome Biol* **11**: R117.
- 637 Becker TS, Lenhard B. 2007. The random versus fragile breakage models of chromosome evolution: a
638 matter of resolution. *Mol Genet Genomics* **278**: 487-491.
- 639 Berthelot C, Muffato M, Abecassis J, Roest Crolius H. 2015. The 3D organization of chromatin explains
640 evolutionary fragile genomic regions. *Cell Rep* **10**: 1913-1924.
- 641 Bourque G, Zdobnov EM, Bork P, Pevzner PA, Tesler G. 2005. Comparative architectures of mammalian
642 and chicken genomes reveal highly variable rates of genomic rearrangements across different
643 lineages. *Genome Res* **15**: 98-110.
- 644 Boyle EA, Li YI, Pritchard JK. 2017. An Expanded View of Complex Traits: From Polygenic to Omnigenic.
645 *Cell* **169**: 1177-1186.
- 646 Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, McMahon A, Morales J,
647 Mountjoy E, Sollis E et al. 2019. The NHGRI-EBI GWAS Catalog of published genome-wide
648 association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res* **47**: D1005-
649 D1012.
- 650 Casola C. 2018. From De Novo to "De Nono": The Majority of Novel Protein-Coding Genes Identified with
651 Phylostratigraphy Are Old Genes or Recent Duplicates. *Genome Biol Evol* **10**: 2906-2918.
- 652 Chen WH, Zhao XM, van Noort V, Bork P. 2013. Human monogenic disease genes have frequently
653 functionally redundant paralogs. *PLoS Comput Biol* **9**: e1003073.
- 654 Choi Y, Chan AP. 2015. PROVEAN web server: a tool to predict the functional effect of amino acid
655 substitutions and indels. *Bioinformatics* **31**: 2745-2747.
- 656 Clark AG, Hubisz MJ, Bustamante CD, Williamson SH, Nielsen R. 2005. Ascertainment bias in studies of
657 human genome-wide polymorphism. *Genome Res* **15**: 1496-1502.
- 658 Cosby RL, Judd J, Zhang R, Zhong A, Garry N, Pritham EJ, Feschotte C. 2021. Recurrent evolution of
659 vertebrate transcription factors by transposase capture. *Science* **371**.
- 660 Cui Y, Chen X, Luo H, Fan Z, Luo J, He S, Yue H, Zhang P, Chen R. 2016. BioCircos.js: an interactive Circos
661 JavaScript library for biological data visualization on web applications. *Bioinformatics* **32**: 1740-
662 1742.
- 663 Damas J, Corbo M, Lewin HA. 2021. Vertebrate Chromosome Evolution. *Annu Rev Anim Biosci* **9**: 1-27.
- 664 Damas J, Kim J, Farre M, Griffin DK, Larkin DM. 2018. Reconstruction of avian ancestral karyotypes
665 reveals differences in the evolutionary history of macro- and microchromosomes. *Genome Biol*
666 **19**: 155.
- 667 Davydov EV, Goode DL, Sirota M, Cooper GM, Sidow A, Batzoglou S. 2010. Identifying a high fraction of
668 the human genome to be under selective constraint using GERP++. *PLoS Comput Biol* **6**:
669 e1001025.
- 670 Dickerson JE, Robertson DL. 2012. On the origins of Mendelian disease genes in man: the impact of gene
671 duplication. *Mol Biol Evol* **29**: 61-69.
- 672 Dron JS, Hegele RA. 2016. Genetics of Lipid and Lipoprotein Disorders and Traits. *Curr Genet Med Rep* **4**:
673 130-141.
- 674 Duret L, Mouchiroud D. 2000. Determinants of substitution rates in mammalian genes: expression
675 pattern affects selection intensity but not mutation rate. *Mol Biol Evol* **17**: 68-74.

- 676 Eilbeck K, Quinlan A, Yandell M. 2017. Settling the score: variant prioritization and Mendelian disease.
677 *Nat Rev Genet* **18**: 599-612.
- 678 Farre M, Kim J, Proskuryakova AA, Zhang Y, Kulemzina AI, Li Q, Zhou Y, Xiong Y, Johnson JL, Perelman PL
679 et al. 2019. Evolution of gene regulation in ruminants differs between evolutionary breakpoint
680 regions and homologous synteny blocks. *Genome Res* **29**: 576-589.
- 681 Flajnik MF. 2018. A cold-blooded view of adaptive immunity. *Nat Rev Immunol* **18**: 438-453.
- 682 Flajnik MF, Kasahara M. 2010. Origin and evolution of the adaptive immune system: genetic events and
683 selective pressures. *Nature Reviews Genetics* **11**: 47-59.
- 684 Fuller ZL, Berg JJ, Mostafavi H, Sella G, Przeworski M. 2019. Measuring intolerance to mutation in human
685 genetics. *Nat Genet* **51**: 772-776.
- 686 Gout JF, Kahn D, Duret L, Paramecium Post-Genomics C. 2010. The relationship among gene expression,
687 the evolution of gene dosage, and the rate of protein evolution. *PLoS Genet* **6**: e1000944.
- 688 Gulko B, Hubisz MJ, Gronau I, Siepel A. 2015. A method for calculating probabilities of fitness
689 consequences for point mutations across the human genome. *Nat Genet* **47**: 276-283.
- 690 Hindorff LA, Sethupathy P, Junkins HA, Ramos EM, Mehta JP, Collins FS, Manolio TA. 2009. Potential
691 etiologic and functional implications of genome-wide association loci for human diseases and
692 traits. *Proc Natl Acad Sci U S A* **106**: 9362-9367.
- 693 Hirschhorn JN, Daly MJ. 2005. Genome-wide association studies for common diseases and complex
694 traits. *Nat Rev Genet* **6**: 95-108.
- 695 Hu H, Huff CD, Moore B, Flygare S, Reese MG, Yandell M. 2013. VAAST 2.0: improved variant
696 classification and disease-gene identification using a conservation-controlled amino acid
697 substitution matrix. *Genet Epidemiol* **37**: 622-634.
- 698 Ioannidis NM, Rothstein JH, Pejaver V, Middha S, McDonnell SK, Baheti S, Musolf A, Li Q, Holzinger E,
699 Karyadi D et al. 2016. REVEL: An Ensemble Method for Predicting the Pathogenicity of Rare
700 Missense Variants. *Am J Hum Genet* **99**: 877-885.
- 701 Jabbari K, Chakraborty M, Wiehe T. 2019. DNA sequence-dependent chromatin architecture and nuclear
702 hubs formation. *Scientific Reports* **9**: 14646.
- 703 Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alfoldi J, Wang Q, Collins RL, Laricchia KM, Ganna A,
704 Birnbaum DP et al. 2020. The mutational constraint spectrum quantified from variation in
705 141,456 humans. *Nature* **581**: 434-443.
- 706 Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alfoldi J, Wang Q, Collins RL, Laricchia KM, Ganna A,
707 Birnbaum DP et al. 2021. Author Correction: The mutational constraint spectrum quantified
708 from variation in 141,456 humans. *Nature* **590**: E53.
- 709 Kemkemer C, Kohn M, Cooper DN, Froenicke L, Hogel J, Hameister H, Kehrer-Sawatzki H. 2009. Gene
710 synteny comparisons between different vertebrates provide new insights into breakage and
711 fusion events during mammalian karyotype evolution. *BMC Evol Biol* **9**: 84.
- 712 Kennedy B, Kronenberg Z, Hu H, Moore B, Flygare S, Reese MG, Jorde LB, Yandell M, Huff C. 2014. Using
713 VAAST to Identify Disease-Associated Variants in Next-Generation Sequencing Data. *Curr Protoc*
714 *Hum Genet* **81**: 6 14 11-25.
- 715 Kichaev G, Bhatia G, Loh PR, Gazal S, Burch K, Freund MK, Schoech A, Pasaniuc B, Price AL. 2019.
716 Leveraging Polygenic Functional Enrichment to Improve GWAS Power. *Am J Hum Genet* **104**: 65-
717 75.
- 718 Kircher M, Witten DM, Jain P, O'Roak BJ, Cooper GM, Shendure J. 2014. A general framework for
719 estimating the relative pathogenicity of human genetic variants. *Nat Genet* **46**: 310-315.
- 720 Kumar P, Henikoff S, Ng PC. 2009. Predicting the effects of coding non-synonymous variants on protein
721 function using the SIFT algorithm. *Nat Protoc* **4**: 1073-1081.
- 722 Lal D, May P, Perez-Palma E, Samocha KE, Kosmicki JA, Robinson EB, Moller RS, Krause R, Nurnberg P,
723 Weckhuysen S et al. 2020. Gene family information facilitates variant interpretation and

- 724 identification of disease-associated genes in neurodevelopmental disorders. *Genome Med* **12**:
725 28.
- 726 Larkin DM, Pape G, Donthu R, Auvil L, Welge M, Lewin HA. 2009. Breakpoint regions and homologous
727 synteny blocks in chromosomes have different evolutionary histories. *Genome Res* **19**: 770-777.
- 728 Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, O'Donnell-Luria AH, Ware JS, Hill AJ,
729 Cummings BB et al. 2016. Analysis of protein-coding genetic variation in 60,706 humans. *Nature*
730 **536**: 285-291.
- 731 Liu X, Li YI, Pritchard JK. 2019. Trans Effects on Gene Expression Can Drive Omnigenic Inheritance. *Cell*
732 **177**: 1022-1034 e1026.
- 733 Makino T, Hokamp K, McLysaght A. 2009. The complex relationship of gene duplication and essentiality.
734 *Trends in genetics : TIG* **25**: 152-155.
- 735 Makino T, McLysaght A. 2010. Ohnologs in the human genome are dosage balanced and frequently
736 associated with disease. *Proc Natl Acad Sci U S A* **107**: 9270-9274.
- 737 Marletaz F, Firbas PN, Maeso I, Tena JJ, Bogdanovic O, Perry M, Wyatt CDR, de la Calle-Mustienes E,
738 Bertrand S, Burguera D et al. 2018. Amphioxus functional genomics and the origins of vertebrate
739 gene regulation. *Nature* **564**: 64-70.
- 740 Mistry J, Finn RD, Eddy SR, Bateman A, Punta M. 2013. Challenges in homology search: HMMER3 and
741 convergent evolution of coiled-coil regions. *Nucleic Acids Res* **41**: e121.
- 742 Nakatani Y, Shingate P, Ravi V, Pillai NE, Prasad A, McLysaght A, Venkatesh B. 2021. Reconstruction of
743 proto-vertebrate, proto-cyclostome and proto-gnathostome genomes provides new insights
744 into early vertebrate evolution. *Nature communications* **12**: 4489.
- 745 Nakatani Y, Takeda H, Kohara Y, Morishita S. 2007. Reconstruction of the vertebrate ancestral genome
746 reveals dynamic genome reorganization in early vertebrates. *Genome Res* **17**: 1254-1265.
- 747 Ohta Y, Kasahara M, O'Connor TD, Flajnik MF. 2019. Inferring the "Primordial Immune Complex": Origins
748 of MHC Class I and Antigen Receptors Revealed by Comparative Genomics. *J Immunol* **203**:
749 1882-1896.
- 750 Orostica KY, Verdugo RA. 2016. chromPlot: visualization of genomic data in chromosomal context.
751 *Bioinformatics* **32**: 2366-2368.
- 752 Pal C, Papp B, Hurst LD. 2001. Highly expressed genes in yeast evolve slowly. *Genetics* **158**: 927-931.
- 753 Peng Q, Pevzner PA, Tesler G. 2006. The fragile breakage versus random breakage models of
754 chromosome evolution. *PLoS Comput Biol* **2**: e14.
- 755 Pevzner P, Tesler G. 2003. Human and mouse genomic sequences reveal extensive breakpoint reuse in
756 mammalian evolution. *Proc Natl Acad Sci U S A* **100**: 7672-7677.
- 757 Putnam NH, Butts T, Ferrier DE, Furlong RF, Hellsten U, Kawashima T, Robinson-Rechavi M, Shoguchi E,
758 Terry A, Yu JK et al. 2008. The amphioxus genome and the evolution of the chordate karyotype.
759 *Nature* **453**: 1064-1071.
- 760 Rhie A McCarthy SA Fedrigo O Damas J Formenti G Koren S Uliano-Silva M Chow W Functammasan A
761 Kim J et al. 2021. Towards complete and error-free genome assemblies of all vertebrate species.
762 *Nature* **592**: 737-746.
- 763 Rice AM, McLysaght A. 2017. Dosage sensitivity is a major determinant of human copy number variant
764 pathogenicity. *Nat Commun* **8**: 14366.
- 765 Rocha EP, Danchin A. 2004. An analysis of determinants of amino acids substitution rates in bacterial
766 proteins. *Mol Biol Evol* **21**: 108-116.
- 767 Sacerdot C, Louis A, Bon C, Berthelot C, Roest Crollius H. 2018. Chromosome evolution at the origin of
768 the ancestral vertebrate genome. *Genome biology* **19**: 166-166.
- 769 Samocha KE, Robinson EB, Sanders SJ, Stevens C, Sabo A, McGrath LM, Kosmicki JA, Rehnstrom K,
770 Mallick S, Kirby A et al. 2014. A framework for the interpretation of de novo mutation in human
771 disease. *Nat Genet* **46**: 944-950.

- 772 Schwarz JM, Cooper DN, Schuelke M, Seelow D. 2014. MutationTaster2: mutation prediction for the
773 deep-sequencing age. *Nat Methods* **11**: 361-362.
- 774 Sella G, Barton NH. 2019. Thinking About the Evolution of Complex Traits in the Era of Genome-Wide
775 Association Studies. *Annual review of genomics and human genetics* **20**: 461-493.
- 776 Shi H, Kichaev G, Pasaniuc B. 2016. Contrasting the Genetic Architecture of 30 Complex Traits from
777 Summary Association Data. *Am J Hum Genet* **99**: 139-153.
- 778 Shihab HA, Gough J, Cooper DN, Stenson PD, Barker GL, Edwards KJ, Day IN, Gaunt TR. 2013. Predicting
779 the functional, molecular, and phenotypic consequences of amino acid substitutions using
780 hidden Markov models. *Hum Mutat* **34**: 57-65.
- 781 Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, Clawson H, Spieth J, Hillier LW,
782 Richards S et al. 2005. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast
783 genomes. *Genome Res* **15**: 1034-1050.
- 784 Simakov O, Marletaz F, Yue JX, O'Connell B, Jenkins J, Brandt A, Calef R, Tung CH, Huang TK, Schmutz J et
785 al. 2020. Deeply conserved synteny resolves early events in vertebrate evolution. *Nat Ecol Evol*
786 **4**: 820-830.
- 787 Simons YB, Bullaughey K, Hudson RR, Sella G. 2018. A population genetic interpretation of GWAS
788 findings for human quantitative traits. *PLoS Biol* **16**: e2002985.
- 789 Singh PP, Arora J, Isambert H. 2015. Identification of Ohnolog Genes Originating from Whole Genome
790 Duplication in Early Vertebrates, Based on Synteny Comparison across Multiple Genomes. *PLoS*
791 *Comput Biol* **11**: e1004394.
- 792 Taliun D Harris DN Kessler MD Carlson J Szpiech ZA Torres R Taliun SAG Corvelo A Gogarten SM Kang HM
793 et al. 2021. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature*
794 **590**: 290-299.
- 795 Van Oss SB, Carvunis AR. 2019. De novo gene birth. *PLoS Genet* **15**: e1008160.
- 796 Visscher PM, Brown MA, McCarthy MI, Yang J. 2012. Five years of GWAS discovery. *Am J Hum Genet* **90**:
797 7-24.
- 798 Visscher PM, Wray NR, Zhang Q, Sklar P, McCarthy MI, Brown MA, Yang J. 2017. 10 Years of GWAS
799 Discovery: Biology, Function, and Translation. *Am J Hum Genet* **101**: 5-22.
- 800 Yandell M, Huff C, Hu H, Singleton M, Moore B, Xing J, Jorde LB, Reese MG. 2011. A probabilistic disease-
801 gene finder for personal genomes. *Genome Res* **21**: 1529-1542.
- 802 Yegorov S, Good S. 2012. Using Paleogenomics to Study the Evolution of Gene Families: Origin and
803 Duplication History of the Relaxin Family Hormones and Their Receptors. *Plos One* **7**: e32923.
- 804
- 805

Figure Legends

806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826

827
828
829
830
831
832
833
834
835
836
837
838
839
840
841

842
843
844
845
846
847
848

849
850

Fig. 1 Circular maps of the human genome (minus Y-chromosome) partitioned by ALG. Outermost circle, ideogram of C-banding patterns; second ring, blocks of the human genome coloured by the 17 ancestral ALGs; third ring, location of GWAS SNPs in the GWAS catalogue coloured by ancestral ALG ($p < 10^{-08}$). Innermost, ribbons connecting regions of the human genome with $> 2\text{Mbp}$ of contiguous sequence belonging to each ALG (top left). Circular Chromosome maps showing broad syntenic regions of ALG B (top right), ALG A (bottom left) and ALG's C & L (bottom right) with examples of ohnologous gene families membership indicated (outermost ring), ideogram with gene individual gene names, GWAS SNP density (third ring) and ribbons connecting regions with $> 1\text{Mb}$ of contiguous sequence.

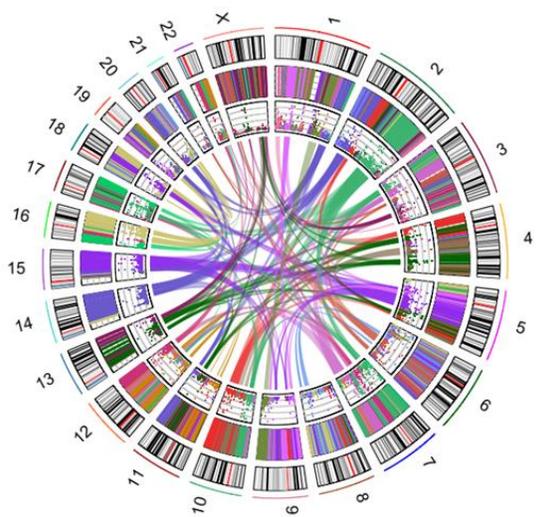
Fig. 2 a) Density of 49,129 GWAS significant variants coloured by ALG and number of HapMap SNPs in 200kb windows (gray line). The x-axis denotes the position in million of base pairs (Mbp) and the Y-axis the density of SNPs on a log-scale. b) Proportion of HapMap and GWAS significant SNPs by ALG, excluding those unassigned to a ALG.

Fig. 3 Odd's ratio and 95% confidence intervals for the association of the EFO traits assigned to 42,999 GWAS variants in the GWAS catalogue. Traits showing an excess ($\text{OR} > 1$) or deficit ($\text{OR} < 1$) of variants by ALG were those in which the standardized residuals from a chi-square analysis were $= \pm 3.841$.

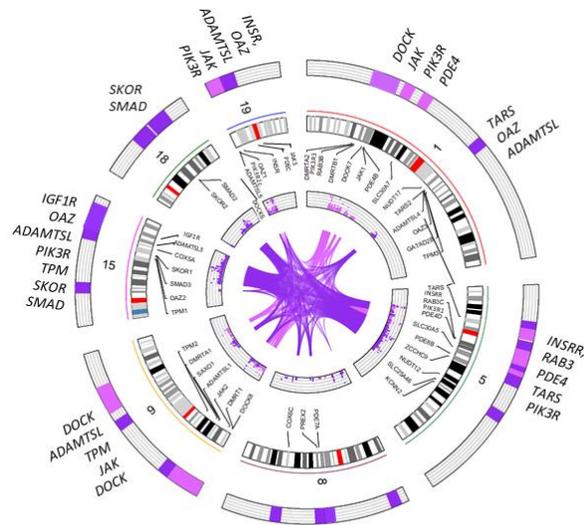
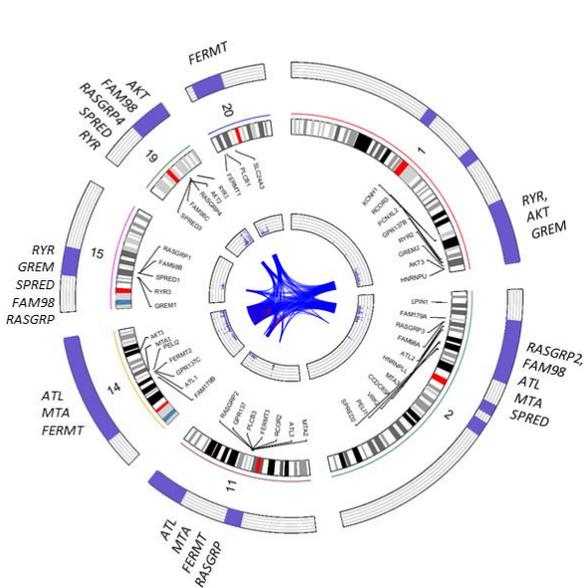
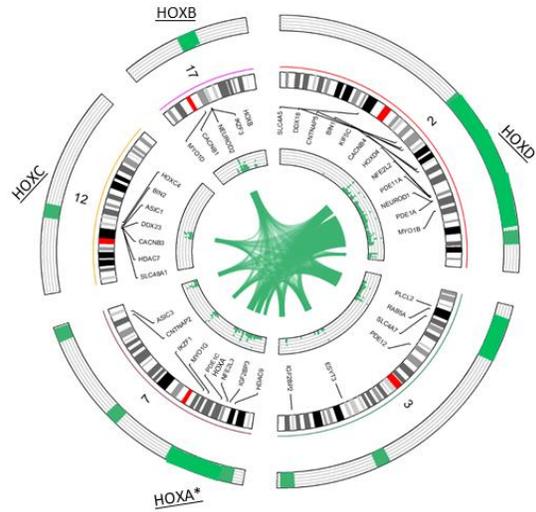
Fig. 4 a) Weighted proportion of genes on a ALG that are GWAS. The dotted line marks 63.7, the overall percent of GWAS associated genes that map to any CLG. Total number of genes is 19,705. b) Odds ratio for genes that are associated with a GWAS variant or assigned to an ALG have a $\text{pLI} > 0.9$, a measure of haploinsufficiency c) the mean number of nonsynonymous mutations (dN) of genes between human and mouse ($\text{dN}_{\text{ALG}} = 0.071$, $\text{dN}_{\text{unassigned}} = 0.132$, $t = -31.95$, $p < 2.2 \times 10^{-16}$) d) distribution of the number of coding exons for genes that map to an ALG or are unassigned (mean/median ALG: 12.5/10, mean/median unassigned: 7.5/5, $W = 60266263$, $p\text{-value} < 2.2 \times 10^{-16}$) e) Median expression (in transcripts per million, TPM) of genes assigned to an ALG (median = 54.0) or unassigned (median = 33.3) in the tissue of highest expression obtained from GTEx (Mann-Whitney U -test $W = 52949771$, $p < 2.2 \times 10^{-16}$) f) Distribution of the length of genes associated with a GWAS vs those not associated (mean length of GWAS associated genes 94069.35 vs 33134 bp not associated, $t = 36.819$, $\text{df} = 14156$, $p\text{-value} < 2.2 \times 10^{-16}$). g) Distribution of the loss-of-function observed/expected upper bound fraction (LOEUF) was significantly lower for genes that map ALG (Kolmogorov-Sminov $D = 0.186$, $p < 2.2 \times 10^{-16}$) and h) for genes associated with a GWAS variant ($D = 0.133$, $D =$, $p < 2.2 \times 10^{-16}$).

Fig. 5 a) Circular map depicting the chromosomes with ALG assignment (outer ring), ideogram (second ring), and tissue of highest expression (innermost ring) for 19 genes associated with core processes in lipid metabolism (see text for details). b) The human chromosomes housing blocks derived from ALG A, C and E (coloured by ALG), enriched for brain-related functions. Median expression (blue dots, left of chromosomes) of genes on ALG A, C or E whose highest tissue of expression in GTEx is brain. Overall gene density (gray histogram) and gene density in brain (red histogram) – right side of chromosomes.

1
2



3



4



5

6

7

8 **Fig. 1**

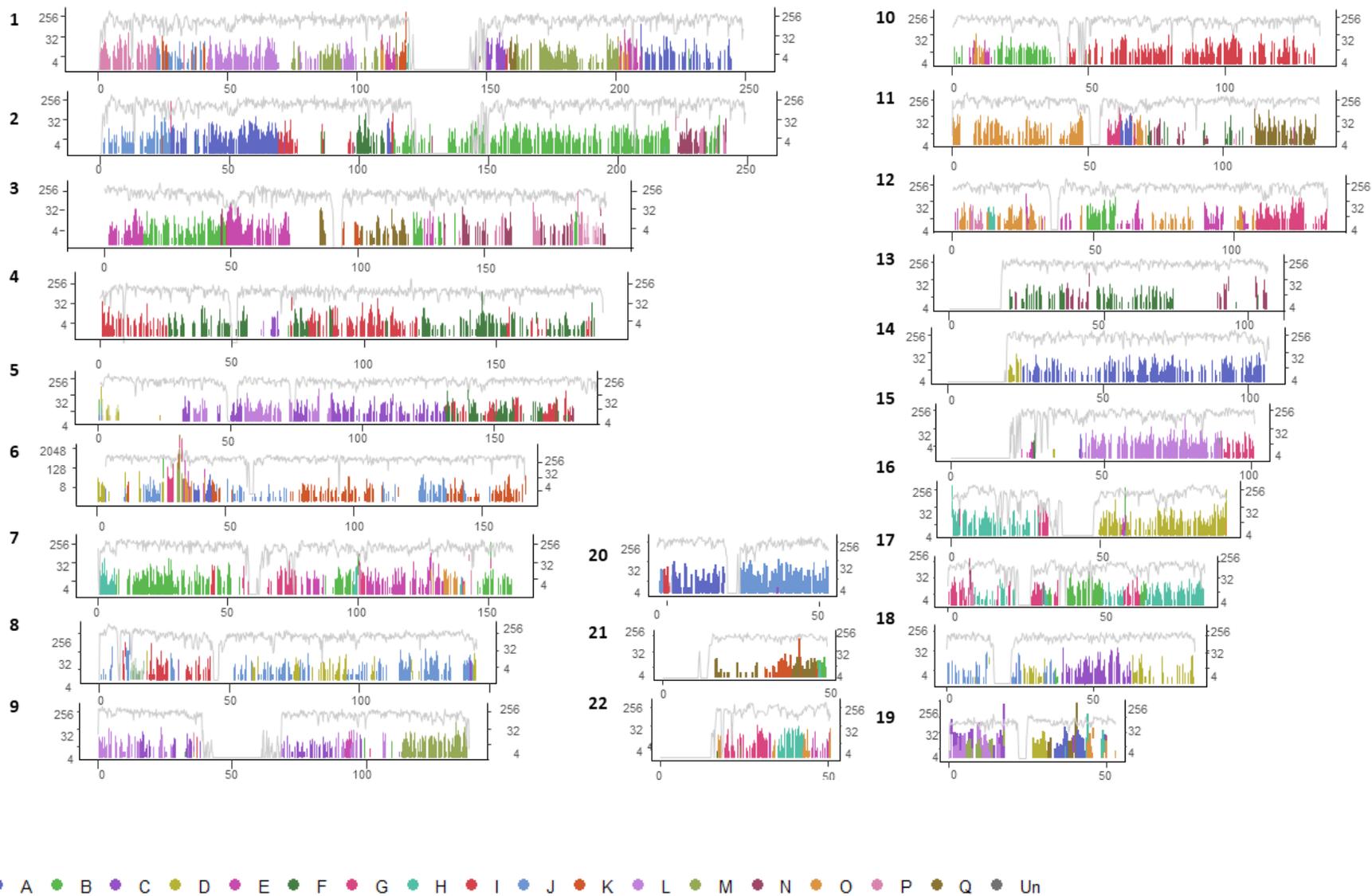
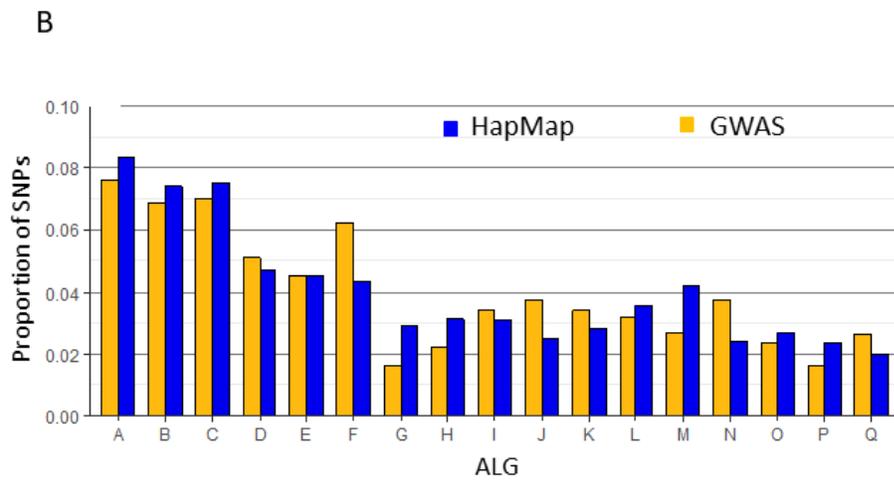
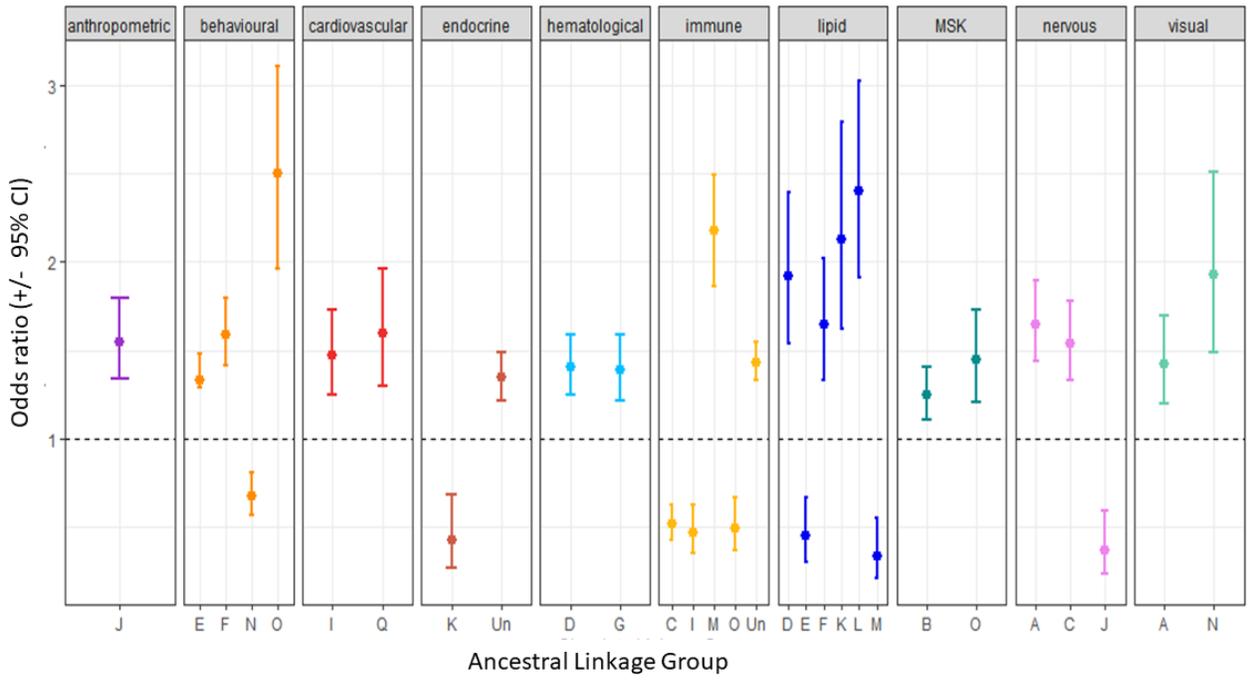


Fig. 2A



22 Fig. 2B

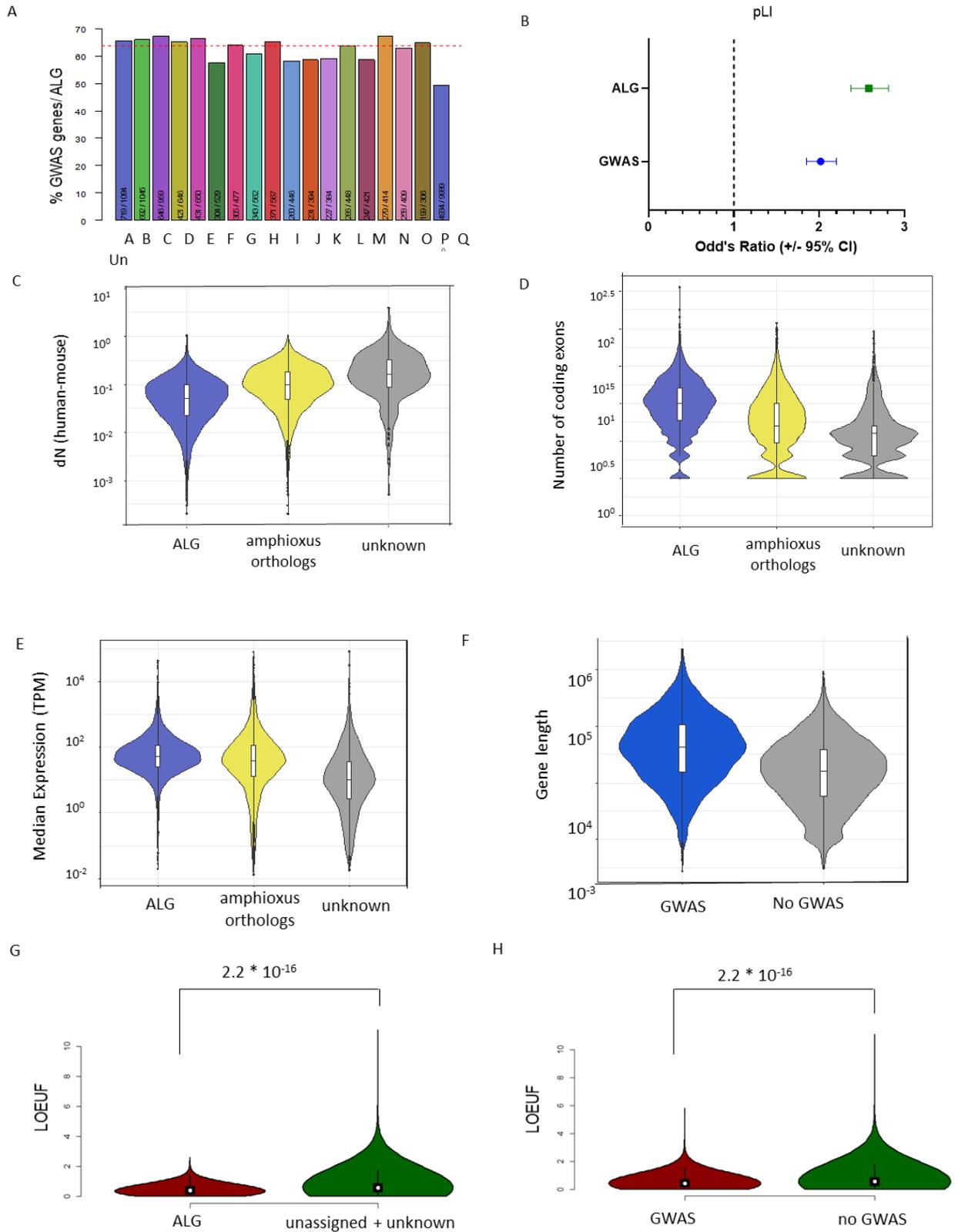
23
24
25



26 **Fig. 3**

27

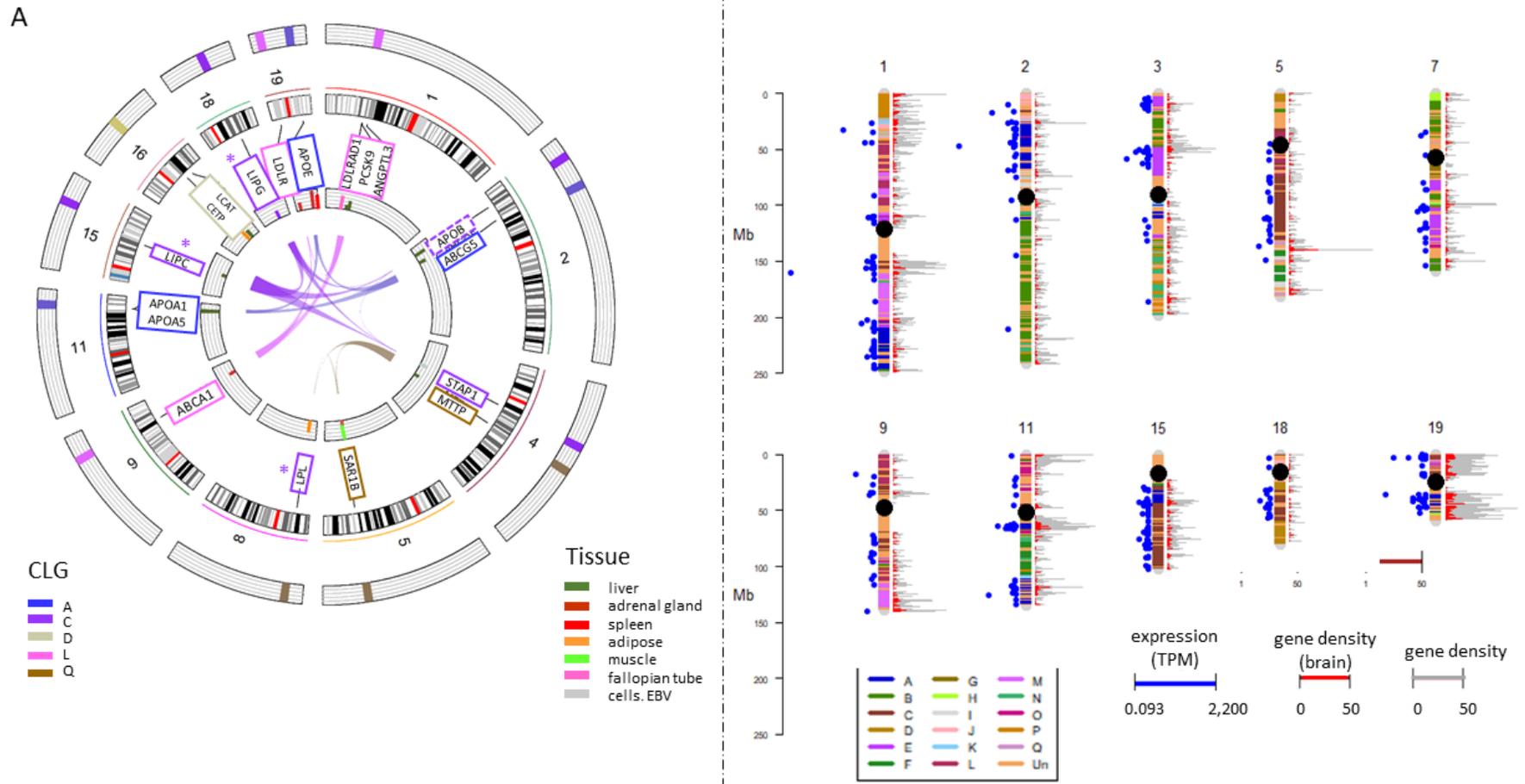
28



29

30

31 **Fig. 4**



32

33 **Fig. 5**