# Gene expansions contributing to human brain evolution

Daniela C. Soto[1,2*‡], José M. Uribe-Salazar[1,2*], Gulhan Kaya[1,2], Ricardo Valdarrago[3], Aarthi Sekar[1,2], Nicholas K. Haghani[1,2], Keiko Hino[4], Gabriana N. La[1,2], Natasha Ann F. Mariano[1,2,5], Cole Ingamells[1,2], Aidan E. Baraban[1,2], Tychele N. Turner[6], Eric D. Green[7], Sergi Simó[4], Gerald Quon[2,3], Aida M. Andrés[8], Megan Y. Dennis[1,2†]

[1]Department of Biochemistry & Molecular Medicine, MIND Institute, University of California, Davis, CA 95616, USA

[2]Genome Center, University of California, Davis, CA 95616, USA

[3]Department of Molecular and Cellular Biology, University of California, Davis, CA 95616, USA

[4]Department of Cell Biology & Human Anatomy, University of California, Davis, CA 95616, USA

[5]Postbaccalaureate Research Education Program, University of California, Davis, CA 95616, USA

[6]Department of Genetics, Washington University School of Medicine, St Louis, MS, 63110, USA

[7]National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, 20892, USA

[8]UCL Genetics Institute, Department of Genetics, Evolution and Environment, University College, London, WC1E 6BT, UK

[*]These authors contributed equally to this work.

[‡]Current institution: Department of Psychiatry and Biobehavioral Sciences, David Geffen School of Medicine, University of California Los Angeles, Los Angeles, CA, 90095, USA

[†]Corresponding author:
Megan Y. Dennis, Ph.D.
University of California, Davis, School of Medicine
One Shields Avenue
Genome Center, 4303 GBSF
Davis, CA 95616
Email: mydennis@ucdavis.edu

## Abstract

Genomic drivers of human-specific neurological traits remain largely undiscovered. Duplicated genes expanded uniquely in the human lineage likely contributed to brain evolution, including the increased complexity of synaptic connections between neurons and the dramatic expansion of the neocortex. Discovering duplicate genes is challenging because the similarity of paralogs makes them prone to sequence-assembly errors. To mitigate this issue, we analyzed a complete telomere-to-telomere human genome sequence (T2T-CHM13) and identified 213 duplicated gene families likely containing human-specific paralogs (>98% identity). Positing that genes important in universal human brain features should exist with at least one copy in all modern humans and exhibit expression in the brain, we narrowed in on 362 paralogs with at least one copy across thousands of ancestrally diverse genomes and present in human brain transcriptomes. Of these, 38 paralogs co-express in gene modules enriched for autism-associated genes and potentially contribute to human language and cognition. We narrowed in on 13 duplicate gene families with human-specific paralogs that are fixed among modern humans and show convincing brain expression patterns. Using long-read DNA sequencing revealed hidden variation across 200 modern humans of diverse ancestries, uncovering signatures of selection not previously identified, including possible balancing selection of *CD8B*. To understand the roles of duplicated genes in brain development, we generated zebrafish CRISPR "knockout" models of nine orthologs and transiently introduced mRNA-encoding paralogs, effectively "humanizing" the larvae. Morphometric, behavioral, and single-cell RNA-seq screening highlighted, for the first time, a possible role for *GPR89B* in dosage-mediated brain expansion and *FRMPD2B* function in altered synaptic signaling, both hallmark features of the human brain. Our holistic approach provides important insights into human brain evolution as well as a resource to the community for studying additional gene expansion drivers of human brain evolution.

## Abstract (short)

Duplicated genes expanded in the human lineage likely contributed to brain evolution, yet challenges exist in their discovery due to sequence-assembly errors. We used a complete telomere-to-telomere genome sequence to identify 213 human-specific gene families. From these, 362 paralogs were found in all modern human genomes tested and brain transcriptomes, making them top candidates contributing to human-universal brain features. Choosing a subset of paralogs, we used long-read DNA sequencing of hundreds of modern humans to reveal previously hidden signatures of selection. To understand their roles in brain development, we generated zebrafish CRISPR "knockout" models of nine orthologs and introduced mRNA-encoding paralogs, effectively "humanizing" larvae. Our findings implicate two new genes in possibly contributing to hallmark features of the human brain: *GPR89B* in dosage-mediated brain expansion and *FRMPD2B* in altered synapse signaling. Our holistic approach provides new insights and a comprehensive resource for studying gene expansion drivers of human brain evolution.

# Introduction

Significant phenotypic features distinguish modern humans from closely related great apes [1–4]. Arguably, one of the most compelling innovations relates to changes in neuroanatomy, including an expanded neocortex and increased complexity of neuronal connections, which allowed the development of novel cognitive features such as reading and language [5]. While previous work implicated human-specific single-nucleotide variants (SNVs) that impact genes leading to altered brain features, including *FOXP2* [6,7] and human-accelerated regions (HARs) [8], a majority of top gene candidates are the result of segmental duplications (SDs; genomic regions >1 kbp in length that share high sequence identity [>90%]) [9–11]. SDs can give rise to new gene paralogs with the same function, altered functions, or that antagonize conserved, ancestral paralogs and contribute more to genetic divergence across species than SNVs [12]. Previous comparisons of great ape genomes have identified >30 human-specific gene families and hundreds of paralogs enriched for genes important in neurodevelopment and residing at genomic hotspots associated with neuropsychiatric disorders [13–15]. Of these, a handful of genes have been found to function in brain development using model systems, including *SRGAP2* [16,17], *NOTCH2NL* [18–20], *ARHGAP11B* [21–23], *TBC1D3* [24] *CROCCP2* [25], and *LRRC37B* [26]. Most studies have leveraged mice to study gene functions although recent studies have expanded to cortical organoids, ferrets, and primates [27]. Despite their clear importance in contributing to neural features, most duplicate genes remain functionally uncharacterized due to the arduous nature of using such models.

SDs have largely eluded analyses due to difficulties in accurate genome assembly [28] and in discovering variants across nearly identical paralogs [29–33]. As such, many human-duplicated genes are likely left to be discovered. The telomere-to-telomere (T2T) human reference genome T2T-CHM13 [34], representing a gapless sequence of all autosomes and Chromosome X, has enabled a more complete picture of SDs [35] by incorporating 238 Mbp missing from the previous human reference genome (GRCh38). In particular, this new assembly corrects >8 Mbp of collapsed duplications [36], including previously missing paralogs of human-specific duplicated gene families [13] *GPRIN2* [35] and *DUSP22* [36]. Here, using this new T2T genome, we identified thousands of recent gene duplications among hominids. By comparing genomic data between great apes and across thousands of modern humans, we narrowed in on a set of paralogs unique within and fixed across modern humans. Transcriptomic datasets from the human brain identified genes most likely to contribute to brain development and function, providing a catalog of the candidate human-specific gene families contributing to brain evolution for further functional testing in model systems. Finally, we prioritized a set of duplicate gene families to characterize in more detail using long-read sequencing and systematic analysis in zebrafish to connect gene functions to brain development.

# Results

## Genetic analysis of human-duplicated genes

### *Identification of human gene duplications in T2T-CHM13*

Understanding that highly identical SDs are enriched for human-specific duplications, we narrowed in on 97.8 Mbp of autosomal sequences sharing >98% identity with other genomic regions (or SD98) in the human T2T-CHM13 [35,37] (Figure 1A). These loci represent genes duplicated only in human lineage [13,14] as well as expansions of duplicated gene families present in other great apes. Paralogs in this latter category have experienced recent changes along the *Homo* lineage in expression (e.g., *LRRC37B* [26]) or sequence

1   content (e.g., *NOTCH2NL*, via interlocus gene conversion [18]) resulting in new novel functions. Of the
2   5,154 SD98 genes (Table S1), we focused on 698 protein-encoding genes and 1,095 unprocessed
3   pseudogenes (that could be mis-annotations of true protein-encoding genes [38]). This list includes well-
4   known duplicated genes important in neurodevelopment (*SRGAP2C*, *ARHGAP11B*), disease (*SMN1* and
5   *SMN2* [39], *KANSL1* [40]), and adaptation (e.g., amylase genes [41–43]). Sequence read depth [35] in modern
6   humans (Simons Genome Diversity Project [SGDP], n=269 [44]) verified that all paralogs had >2 gene-
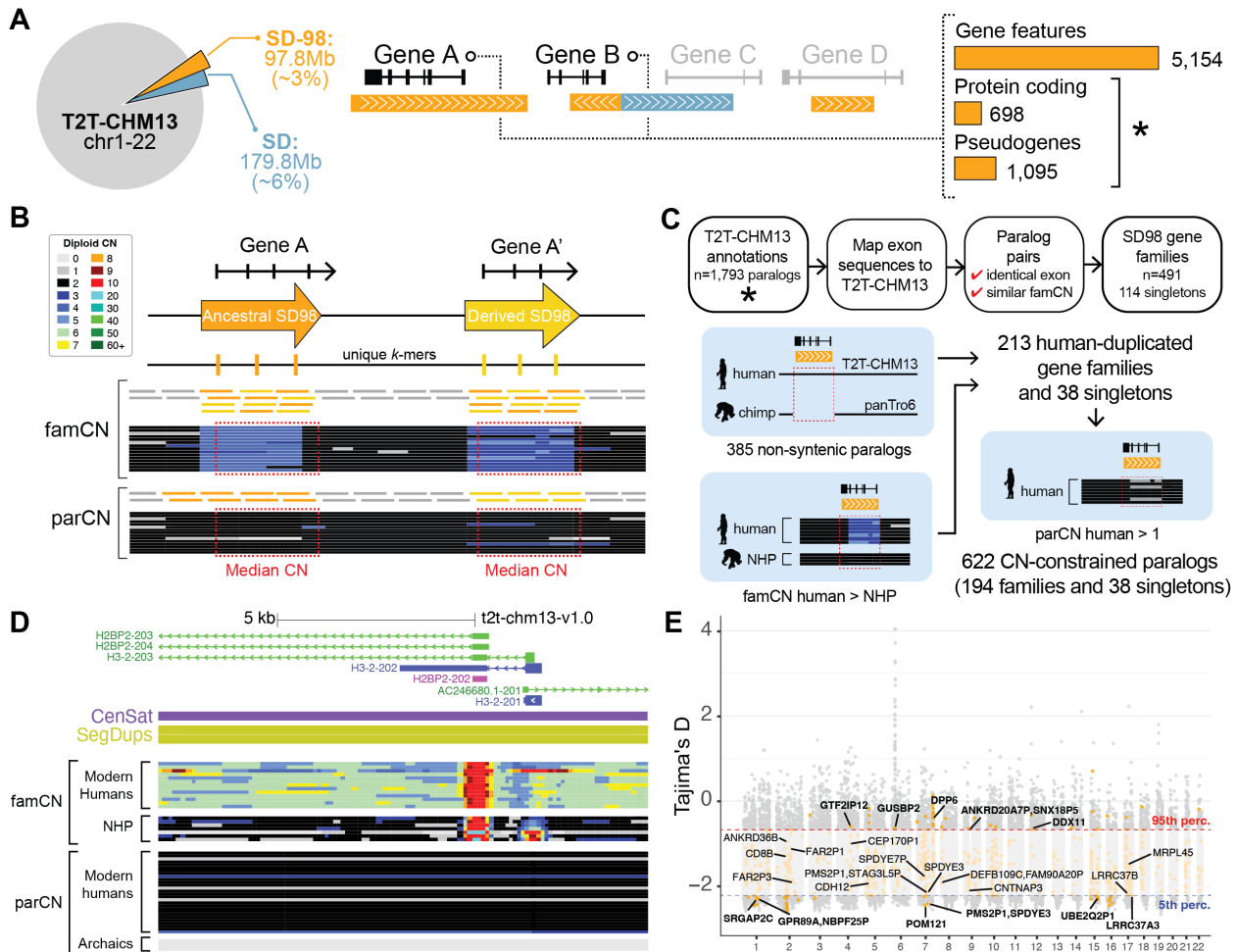7   family diploid copy number (famCN; Methods) (Table S2, Figure 1B).
8
9   Based on sequence and famCN similarity, we clustered the 1,679 paralogs into 491 multigene families
10  (Figure 1C, Figure S1), with most families having 2–3 members (n=271) (Figure S2). Three extreme
11  high-copy gene families had >50 paralogs, including macrosatellite-associated *DUX4* and *DUB/USP17* as
12  well as primate-specific *FAM90A* [45]. The remaining 114 paralogs were defined as singletons (Table S2),
13  with some failing to cluster due to high and variable copy numbers (CNs) (e.g., *CROCC* and *CROCCP2*)
14  or only a small portion of the gene duplicated (e.g., *AIDA* and *LUZP2*). Within 163 multigene families
15  and 13 singletons, we identified 385 human-specific paralogs within non-syntenic regions present in
16  human but not chimpanzee reference assemblies [35] (Figure 1C, Table S2). Several previously known
17  human-specific genes were notably absent from this list (e.g., *NPY4R2*, *ROCKP1*, and *SERF1B* [13])
18  because genome alignments across SDs can be imprecise. We next identified human-expanded gene
19  families as those with higher famCN in humans (SDGP, n=269) versus nonhuman great apes (n=4) (97
20  gene families and 27 singletons; Figure 1C, Table S3), excluding high-copy genes that were difficult to
21  accurately detect CN differences (famCN>10). In total, we conservatively predict 213 gene families and
22  38 singletons comprising at least one human-specific duplicate paralog (Table S3).
23
24  ### *Variation of duplicated genes in modern humans*

25  Positing that all humans should carry a functional version of a gene if important for a species-universal
26  trait, we used *k*-mer-based paralog-specific copy number (parCN) estimates [46] to identify 622 genes (194
27  duplicate families, 38 singletons) with at least one copy in >98% of humans ("CN constrained",
28  parCN≥0.5; 1000 Genomes Project, 1KGP; n=2,504) (Tables S1 and S4). Of these, 125 paralogs were
29  "fixed" in humans (parCN~2) and likely represent *Homo sapiens*-specific genes. We found 13 CN
30  constrained genes that were largely absent (parCN<0.5) from four archaic human genomes [47–49]. One of
31  these genes, *H3-2/H2BP2*, is a member of a core *H2B* histone family involved in the structure of
32  eukaryotic chromatin [50], homologous with another human-specific *H2BP1* and the ancestral *H2BC18*
33  paralog (Figure 1D). Another *Homo sapiens*-specific gene, *FCGR1CP*, encodes an immunoglobulin
34  gamma Fc Gamma Receptor, a family of proteins vital in regulating immune response [51]. Moving
35  forward, we consider only duplicate gene families comprising CN-constrained genes.
36
37  We identified 13 protein-encoding genes as loss-of-function intolerant (pLI ≥0.9 or LOEUF ≤0.35) using
38  SNV data from hundreds of thousands of humans from gnomAD [52] (Table S1, Figure S3), showing that
39  deleterious mutations of these genes are depleted in human populations (e.g., likely not compatible with
40  life). These conserved genes are all ancestral paralogs, including *NOTCH2*, *HERC2*, and *CORO1A*. The
41  gnomAD (v3) metrics rely on variants identified in protein-encoding genes using the human reference
42  genome hg19, which has known errors across SDs [53] and misannotated pseudogenes. As such, all
43  unprocessed pseudogenes and 32% of protein-encoding SD98 genes lacked gnomAD pLI and LOEUF
44  scores. To circumvent these issues, we assessed SNV genetic diversity by Tajima's D [54] using the T2T-

1    CHM13 reference and the 1KGP cohort [36,55]. Focusing on short-read accessible regions (Figure S4, Note
2    S1), we identified 15 CN-constrained human-duplicated genes with extreme negative D values (<5th
3    percentile of the genome-wide empirical distribution) considered signatures of positive or purifying
4    selection (Figures 1E and S5). These included human-specific paralog *SRGAP2C* previously implicated in
5    cortical neuronal migration and synaptogenesis [16,17] as well as the uncharacterized *LRRC37A3* and the
6    hominid-specific *LRRC37B*, recently found to function in cortical pyramidal neurons by impacting
7    synaptic excitability [26]. We also identified nine genes exhibiting extremely positive D values (>95th
8    percentile) as putative signatures of balancing selection, including T-cell antigen *CD8B*. Collectively,
9    variants discovered using the new T2T-CHM13 genome enabled the identification of new and interesting
10   human-duplicated genes potentially contributing to traits and diseases not previously assayed in genome-
11   wide selection screens.



13   **Figure 1. Human gene duplications in T2T-CHM13.** (**A**) Diagram of segmental duplications (SDs) with >90%
14   identity (blue) and >98% identity (orange) in T2T-CHM13 and selection of genes within SDs with >98% identity
15   (SD98 genes). Total counts are shown on the right, with protein-encoding genes and pseudogenes used for further
16   analysis indicated with an asterisk. (**B**) Schematic representation of copy number (CN) estimation methods,
17   including gene-family CN (famCN) and paralog-specific CN (parCN). Illustrated horizontal lines represent short-
18   read pileups mapping to unique (gray) and duplicated regions (orange and yellow). Read-depth diploid CN estimates
19   are shown as heatmaps with values explained in the legend (left). The CN-genotyping window is shown as red
20   dashed boxes. (**C**) Pipeline for clustering and stratification of SD98 genes. Gene families were classified as carrying
21   human duplicates based on synteny with the chimpanzee reference genome (panTro6) and famCN comparisons
22   between human and nonhuman primates (NHPs) (left). CN-constrained (fixed or nearly fixed) genes were flagged

1  based on parCN values across human populations (right). (**D**) UCSC Genome Browser snapshot of the *H3-2*/*H2BP2*
2  locus, including gene models, centromeric satellites (CenSat), SDs (SegDup), and famCN and parCN predictions
3  across modern humans, NHPs, and archaic genomes. (**E**) Distribution of Tajima's D values (y-axis) from individuals
4  of African ancestry from the 1KGP across 25-kbp windows genome wide (gray) and in the SD98 region (orange)
5  across human autosomal chromosomes (x-axis). All human-duplicated gene names with outlier D values in African,
6  European, East Asian, South Asian, and American populations are included.
7

## Human-duplicated genes implicated in brain development

9

### *Connecting genetic variation of duplicated genes with neural traits*

11  To narrow in on human-duplicate gene families contributing to neurocognitive features, we identified 187
12  genes with putative associations with brain-related phenotypes from the genome-wide association study
13  (GWAS) catalog and UK Biobank [56] (Tables S1 and S4). Three variants (rs12725078, rs17537178 and
14  rs4797876) associated with sulcal depth impact *SRGAP2*, *PTPN20*, and *ROCK1P1*, respectively. The
15  ancestral *CORO1A*, implicated in autism [57], is associated with brain morphology. Many implicated genes
16  reside at genomic hotspots (n=58), such as *GPR89* paralogs at chromosome 1q21.1 with recurrent ~2
17  Mbp deletions/duplications impacting brain size [58]. While interesting, GWAS hits are significantly
18  depleted across SD98 regions (Note S2), in part due to the common use of single-nucleotide
19  polymorphism (SNP) arrays that lack coverage across SDs [59]. As such, we assayed variation in an autism
20  cohort (Simons Simplex Collection [SSC]; n=2,459 quad families). Eighteen genes show significant
21  parCN differences in probands versus unaffected siblings (Wilcoxon signed-rank test, *q*-value<0.05)
22  (Figure S6), with all but one residing at chromosome 15q25.2 (OMIM: 614294), a region known to
23  undergo recurrent deletions/duplications [60]. The remaining gene, pseudogene AC233280.19, is associated
24  with the chromosome 3q29 genomic disorder (OMIM: 609425). *De novo* copy number variants (CNVs)
25  impact 22 human-duplicated genes in autistic probands (Table S6, Figure S7); this contrasts with six
26  events impacting five paralogs in unaffected siblings (Fisher's exact test, *p*-value = $4.5 \times 10^{-4}$). Most
27  impacted genes reside at known autism-associated genomic hotspots (n=15). The other seven, which were
28  not mutated in unaffected siblings, included protein-encoding genes *CD8B2*, *FCGR1B*, *HYDIN* and
29  *LIMS1*, representing possible contributors to autism spectrum disorder (ASD).
30

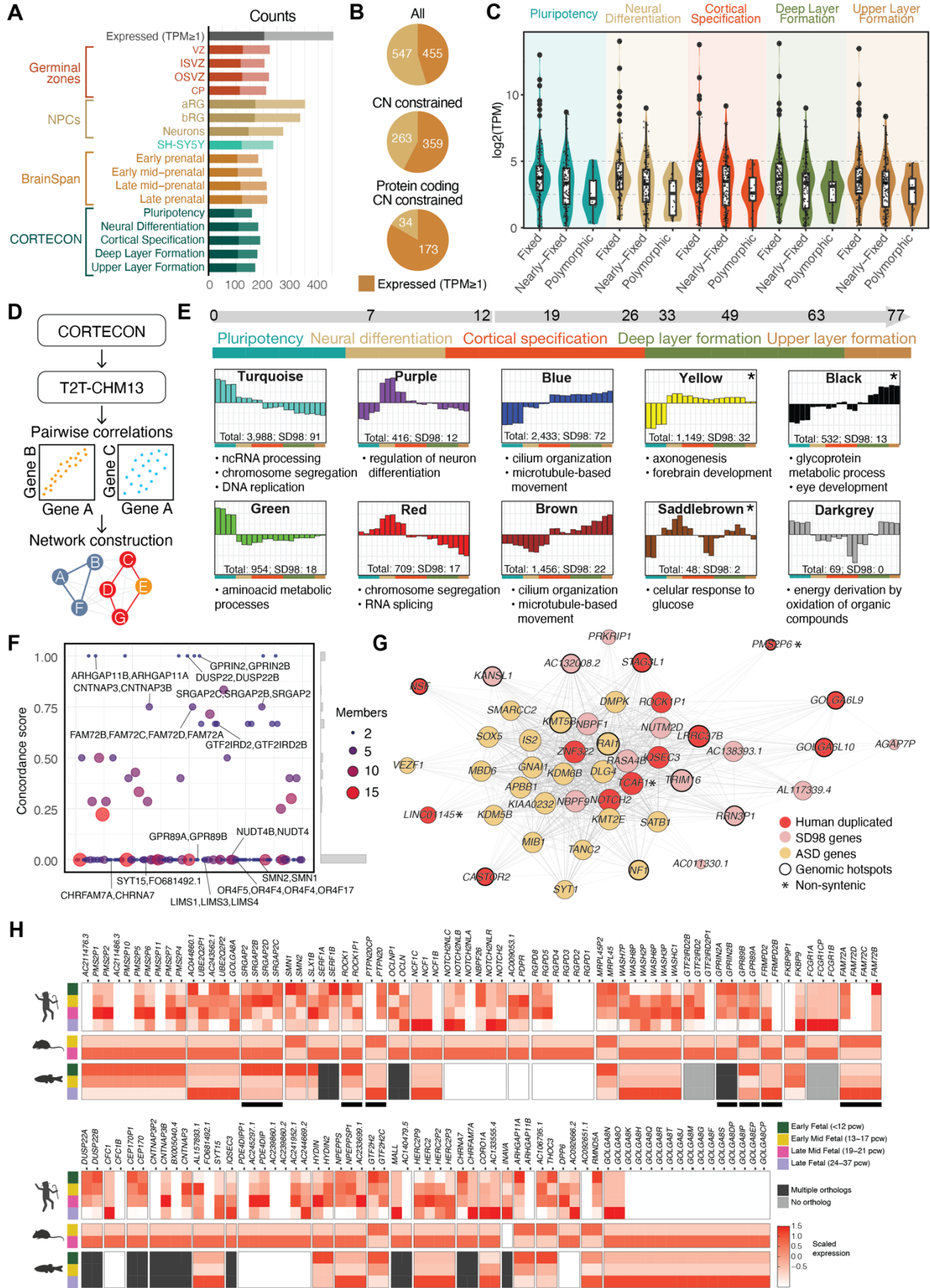### *Duplicated gene expression in the developing human brain*

32  Re-analyzing published RNA-seq datasets [21,61–64] using the new T2T-CHM13 reference, we found nearly
33  half of human-duplicated gene paralogs (455/1,002) are expressed during brain development (TPM≥1)
34  (Table S1, Figures 2A and S8), representing a depletion versus the genome-wide transcriptome
35  (21,513/23,395). This increases to 58% for CN-constrained genes (1.3-fold enrichment, *p*-value = $2.5 \times 10^{-24}$,
36  hypergeometric test) and to 84% for CN-constrained protein-encoding genes (1.4-fold enrichment, *p*-
37  value = $7.8 \times 10^{-30}$, hypergeometric test) (Figure 2B). These results suggest true functional candidates are
38  more likely to exist in the most CN-constrained protein-encoding genes (Figure 2C). In sum, 147 human-
39  duplicated families carry at least one CN-constrained and brain-expressed gene, including 39 protein-
40  encoding paralogs verified as human specific (non-syntenic). Of these, 21 genes are also expressed in the
41  postnatal brain, including *CD8B2*, which is exclusively expressed after birth.
42

43  We next used the longitudinal CORTECON dataset [63], with transcriptomes of different stages of *ex-vivo*-
44  induced neurogenesis from human embryonic stem cells, to infer developmental functions of genes using

1   weighted gene co-expression network analysis (WGCNA) [65] (Figure 2D). Expressed genes (n=15,695)
2   were clustered into 37 co-expression modules, each assigned a random color identifier. Thirty-two
3   modules comprised SD98 genes (n=399), of which 200 paralogs represented human-duplicated families
4   (55 non-syntenic) (Table S7, Figures 2E and S9). Comparing module assignment between paralogs found
5   mostly differential expression patterns, with only six duplicate gene families in complete concordance
6   (i.e., all paralogs in the same module) (Table S8, Figure 2F). This suggests that our approach largely
7   distinguishes transcriptional profiles between similar paralogs, and that expression diverges at relatively
8   short evolutionary time scales (<6 million years), as we have shown for a smaller set of genes [66].
9
10  Twenty-two of 35 modules were enriched for functional gene ontology (GO) terms ($q$-value<0.05,
11  hypergeometric test; Table S9, Figures 2E and S9). To verify module assignments, we searched for
12  duplicated genes with characterized functions. *ARHGAP11B*, which induces cortical neural progenitor
13  amplification by altering glutaminolysis in the mitochondria [23], is a member of turquoise. Genes in this
14  module are expressed highest during pluripotency and are associated with cell proliferation, including
15  DNA replication and chromosome segregation, as well as mitochondrial gene expression. The hominoid-
16  specific gene *TBC1D3*, known to promote basal progenitor amplification in the outer radial glia resulting
17  in cortical folding in mice [24] is a member of purple, a module associated with regulation of neural
18  differentiation. Human-specific *SRGAP2C*, which interacts with F-actin to produce membrane protrusions
19  required for neuronal migration [67], represents blue with co-expressed genes that peak during cortical
20  specification and upper-layer formation. This module is associated with cell motility, including motile
21  cilium organization and assembly and microtubule-based movement.
22
23  We also found autism-associated genes [57] significantly enriched in four modules (yellow, black, saddle
24  brown, and cyan), as well as the "unassigned" module (grey) ($q$-value < 0.05, hypergeometric test), and
25  included 38 paralogs from human-duplicated gene families. Remarkably, three protein-encoding paralogs
26  from the *RGPD* gene family, encoding RANBP2 Like And GRIP Domain Containing proteins, were
27  represented in these modules, including human-specific *RGPD3* (yellow) and *RGPD4* (grey) as well as
28  *RGPD8* (saddle brown). The yellow module, enriched with functions in axon guidance and
29  synaptogenesis, contains the most autism-associated genes (n=20) (Table S7, Figure 2G). Genes in this
30  module exhibit low expression during pluripotency, followed by sustained expression from neural
31  differentiation to deep layer formation, including several markers of glutamatergic neurons (e.g., *SOX5*,
32  *SLC1A6*, *OTX1*, and *TLE4*) [68]. Human-duplicate paralogs in the yellow module include *LRRC37B*,
33  important in synapse function, as well as the causal gene in the chromosome 17q21.31 microdeletion
34  syndrome, *KANSL1* [69]. We also identified compelling candidates residing at autism-associated genomic
35  hotspots (e.g., *GOLGA6L9* and *GOLGA6L10* in chromosome 15q25.2, and *CASTOR2*, *PMS2P6* and
36  *STAG3L1* in chromosome 7q11.23) (Table S1). Collectively, duplicated genes co-expressed with neural
37  and ASD-associated genes representing top candidates contributing to human brain development.

**Figure 2. Expression of human-duplicated genes during brain development.** (**A**) Counts of human-duplicated genes with transcripts per million (TPM) >1 in fetal brain datasets including germinal zones (VZ: ventricular zone, ISVZ: inner subventricular zone, OSVZ: outer subventricular zone, CP: cortical plate), neuronal progenitor cells (NPCs) (aRGs: apical radial glia, bRGs: basal radial glia), neuroblastoma cell line (SH-SY5Y), the BrainSpan dataset, and the CORTECON dataset. Counts for protein-encoding genes are represented in darker shades. (**B**) counts of expressed (TPM≥1) (dark orange) and non-expressed (light orange) human-duplicated genes across gene categories. (**C**) SD98 gene expression in $\log_2$(TPM) in the CORTECON dataset, spanning pluripotency to upper layer formation and stratified by copy number (CN) category. (**D**) Pipeline used for the weighted gene coexpression analysis (WGCNA) of CORTECON data remapped to the T2T-CHM13 reference. (**E**) Selected WGCNA co-expression modules represented with random colors. Modules were organized based on their temporal expression spanning pluripotency to upper layer formation (day 0 to 77) with overrepresented gene ontology terms shown at the bottom. (**F**) Module assignment concordance scores are shown on the vertical axis for SD98 gene families, with spacing along the horizontal axis for visual separation. The size of each point corresponds to the number of members in the respective gene family. (**G**) Network diagram of yellow module. Only genes within human-duplicated gene families (red) or SD98 (pink) and autism-associated (yellow) categories with high module membership are depicted. Genes with asterisks are non-syntenic with chimpanzee [35] and bold borders are within ±500-kbp of a genomic disorder hotspot [60]. (**H**) Scaled TPMs from post-mortem human fetal brain samples from the BrainSpan dataset, and pseudo-bulk single-cell transcriptomes from whole-brain dissected samples of mouse and zebrafish. Gene families pictured represent a subset of CN-constrained and brain-expressed human-duplicated gene families. Genes with black bars beneath them were prioritized for additional characterization.

## *Modeling functions of duplicated genes in brain development*

The next step in understanding the role of human-duplicated genes in brain development is to test their functions in model systems. Our combined analysis highlights 148 gene families with at least one CN-constrained or brain-expressed human-duplicated paralog, in addition to 30 paralogs not assigned to a family (Table S10). Of these, we found 106 with a homologous gene(s) in either mouse or zebrafish. Using matched brain-expression data from these species corresponding to human developmental stages [64,70,71] (Figure S10, as previously described [72,73]) narrowed in on 76 and 41 single-copy orthologs expressed during neurodevelopment in mice and zebrafish, respectively (Table S11), representing top candidates for functional studies. This leaves 40% of the human duplicate families with no obvious mouse/zebrafish ortholog, including fusion genes, primate-specific genes (e.g., *TBC1D3* paralogs [24,74]), or those associated with great ape ancestral "core" duplicons (e.g., *NBPF* and *NPIP*) [75]. Alternative models are required, such as *in vivo* primate or cell culture organoids, to test the functions of these genes.

## **Application of the resource: Characterizing candidate duplicated genes**

## *Genetic variation of candidate genes important in neurodevelopment*

As a proof of concept, we selected 13 priority human-specific duplicated (pHSD) gene families representing 30 paralogs from our model gene list (Table S12). Since none of the paralogs fully reside within short-read-accessible genomic regions due to their high identity (Table S1), we characterized variation using long-read sequencing. This included published draft assemblies of 47 individuals from the Human Pangenome Reference Consortium (HPRC) [76–78] and nine individuals from the Human Genome Structural Variation Consortium (HGSVC) [79] (112 total haplotypes; Figure S11). We also performed capture high-fidelity (cHiFi) sequencing on 178 individuals of diverse ancestries in the extended 1KGP cohort [55] and 22 individuals from the Human Genome Diversity Project (HGDP) [80] (Table S13, Figure S12, Note S3). Combined, we identified 46,754 variants (33,774 SNVs and 12,980 indels), or 12.7 variants/kbp, across captured regions (Table S14). Levels of variation within gene families were largely different between paralogs (Mann-Whitney U test, p≤0.05), with the exception of *FRMPD2* and *PTPN20*

1  (Figure S13). For instance, compared with the ancestral *SRGAP2* paralog, human-specific *SRGAP2B*
2  exhibited the lowest and *SRGAP2C* the highest heterozygosity levels, in line with different mutation rates
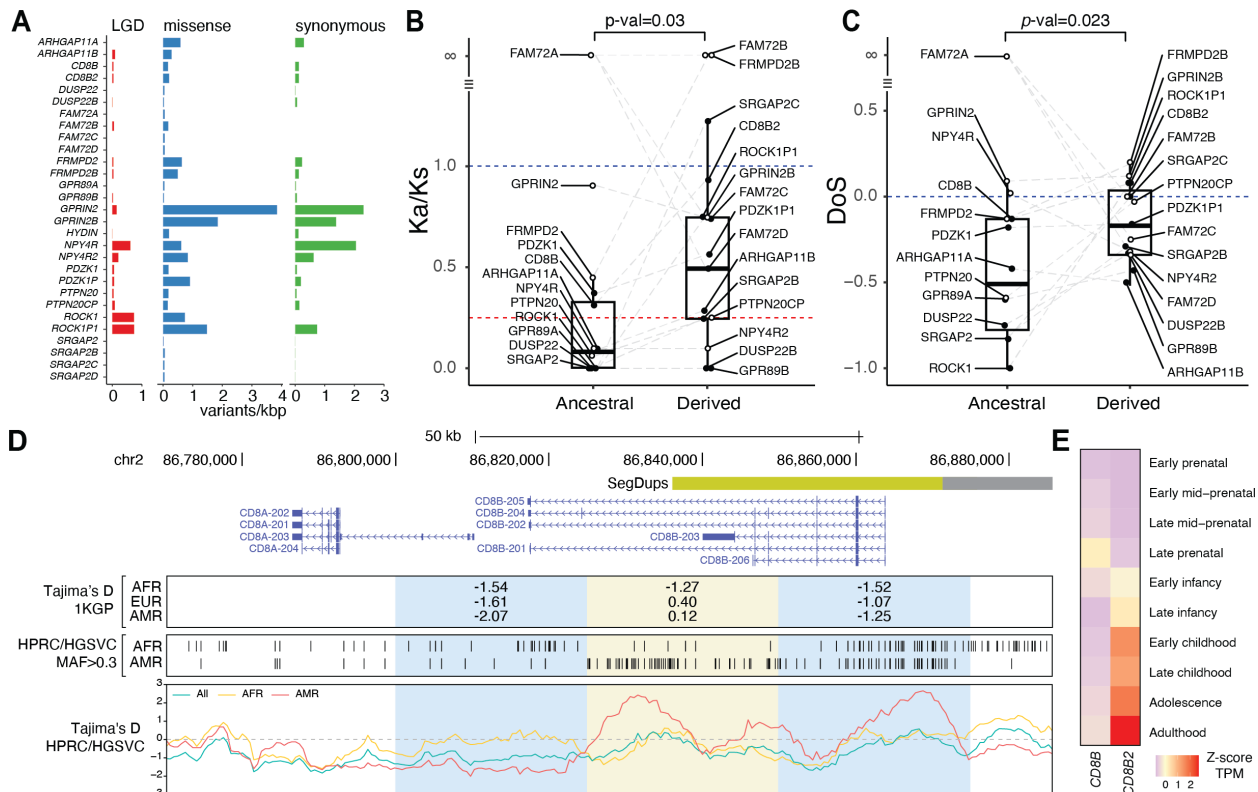3  previously observed at each loci [17].
4
5  Functional annotation [81] identified 412 gene-impacting variants (missense = 252, synonymous = 131,
6  likely gene-disruptive [LGD] = 29; Tables S15 and S16, Figure 3A), with eleven paralogs exhibiting no
7  LGD variants suggesting strong selective constraint. To infer purifying selection, an indicator of function,
8  we calculated the Ka/Ks statistic (also known as dN/dS) per gene family (Table S17). Virtually all
9  paralogs had Ka/Ks lower than 1, and seven ancestral and three derived paralogs exhibited Ka/Ks below
10 the genome-wide average (~0.25) [82]. The ancestral paralogs (Table S12) exhibited significantly lower
11 Ka/Ks values than their derived paralogs (Wilcoxon signed-rank test, *p*-value=0.03) (Figure 3B),
12 consistent with stronger purifying selection. To test for more recent selection signatures, we incorporated
13 polymorphic variation to calculate pN/pS and the direction of selection (DoS) statistic [83], which similarly
14 indicated stronger purifying selection in the ancestral versus derived paralogs (Wilcoxon signed-rank test,
15 *p*-value=0.023) (Table S17, Figure 3C). While the tests mostly agree, *NPY4R* shows discordant
16 signatures, being highly conserved according to Ka/Ks but approaching zero in DoS, in line with an
17 excess of observed LGD variants suggesting recent neutral evolution. Most paralogs within gene families
18 were under purifying selection, including *GPR89*, *CD8B*, *DUSP22*, *GPRIN2*, and *ARHGAP11* (also
19 evident from a larger phylogenetic analysis of dN/dS using a maximum likelihood approach [84]; Table
20 S18), although some show conservation in only one paralog, such as *ROCK1*. Human-specific *SRGAP2C*
21 has elevated Ka/Ks and pN/pS, together with low Tajima's D (-2.32) in African individuals from the
22 1KGP genome-wide screen (Figure 1E), suggesting *SRGAP2C* is evolving under positive selection.
23
24 We verified selection signatures of pHSDs using high-confidence variants obtained from genome
25 assemblies (n=56, HPRC/HGSVC) using nucleotide diversity π and Tajima's D. *SRGAP2C* again shows
26 negative Tajima's D (-2.14) in AFR, validating genome-wide results (Figure S14). *GPR89* gene family
27 paralogs, with low Ka/Ks, exhibit low nucleotide diversity and negative Tajima's D values across all
28 exons consistent with functional constraints (Figure S15). In contrast, *ROCK1* showed reduced nucleotide
29 diversity and more negative Tajima's D compared to *ROCK1P1*, consistent with their Ka/Ks values
30 (Figure S16). While Ka/Ks was not calculated for *FAM72* paralogs due to a lack of synonymous
31 polymorphisms, Tajima's D values similarly ranged from -2 to -1 indicating conservation of the gene
32 family members (Figure S17).
33
34 Revisiting the 1KGP genome-wide signal of balancing selection in individuals of American (Tajima's
35 D=0.12) and European ancestries (D=0.40) centered on *CD8B* (Tables S1 and S5, Figures 1E, 3D and
36 S5), we find positive Tajima's D in American (max 2.66, n=18) but not in African ancestries (max 0.62,
37 n=27) with three major peaks within the gene (Figure 3D). The ancestral *CD8B* paralog, encoding CD8
38 Subunit Beta, is highly expressed in T cells where the protein dimerizes with itself or CD8A (alpha) to
39 serve as a cell-surface glycoprotein mediating cell-cell interactions and immune response [85,86]. Leveraging
40 the assemblies, we identified two distinct haplotype clusters underlying the Tajima's D peaks, one of
41 them particularly prevalent in individuals of American ancestry (Figure S18). Expanding to the entire
42 long-read dataset (including cHiFi) shows an increase in intermediate-frequency variants, a signature of
43 balancing selection, in *CD8B* among European and American ancestries, compared with those of African
44 ancestry (Kolmogorov-Smirnov, *p*-value=$2.2 \times 10^{-16}$) (Figure S19); these variants were verified as

1   differentiating the two main haplotypes. Two of the SNPs (rs56063487 and rs6547706) are *CD8B* splice

2   eQTLs in whole blood from GTEx [87] and significantly associated with increased CD8-protein levels on

3   CD8+ T cells within a Sardinian cohort [88]. We note that *CD8B2* paralog-specific variants do not overlap

4   with the SNPs, providing confidence in these short-read-based genotype results. The haplotypes may,

5   thus, play a role in the modulation of the adaptive immune response, a frequent target of balancing

6   selection. Alternatively, the human-specific paralog *CD8B2* exhibits divergent expression in the human

7   postnatal brain rather than in T cells [38] (Figure 3E). These results provide an example of two paralogs

8   with likely divergent functions and contrasting evolutionary pressures over a relatively short evolutionary

9   time span (~5.2 million years ago [mya] [13]). Combined, we demonstrate the efficacy of long-read data to

10  uncover hidden signatures of natural selection.

11



**Figure 3. Genetic variation and signatures of selection of priority human-specific duplicated (pHSD) genes.**
(**A**) Number of likely gene-disruptive (LGD) (red), missense (blue), and synonymous (green) mutations identified in
pHSD genes using long-read assemblies (n=56) and PacBio capture high-fidelity (cHiFi) sequencing (n=144).
(**B**) Ka/Ks values calculated from human and chimpanzee sequences. Red dashed line indicates the average genome-
wide Ka/Ks between humans and chimpanzees. Blue line indicates neutrality in the Ka/Ks test. Differences between
the Ka/Ks of the matched ancestral and derived paralogs were tested with the Wilcoxon signed-rank test. (**C**)
Direction of selection (DoS) values derived from Ka/Ks and pN/pS estimates. Blue line indicates the threshold for
signatures of positive selection (positive values). Significant differences between ancestral and derived paralogs
were obtained with the Wilcoxon signed-rank test. Paralogs with infinite values or undetermined ancestral/derived
state (hollow dots) were excluded from Ka/Ks and DoS comparisons. (**D**) *CD8B* locus overview, including Tajima's
D values derived from 1KGP SNVs in 25-kbp windows, biallelic SNPs with a minor allele frequency greater than
0.3 identified in African (AFR, n=27) and American (AMR, n=18) individuals using continuous assemblies from the
Human Pangenome Reference Consortium (HPRC) and the Human Genome Structural Variation Consortium
(HGSVC), and Tajima's D values derived from HPRC and HGSVC SNVs using 6-kbp windows and 500-bp steps
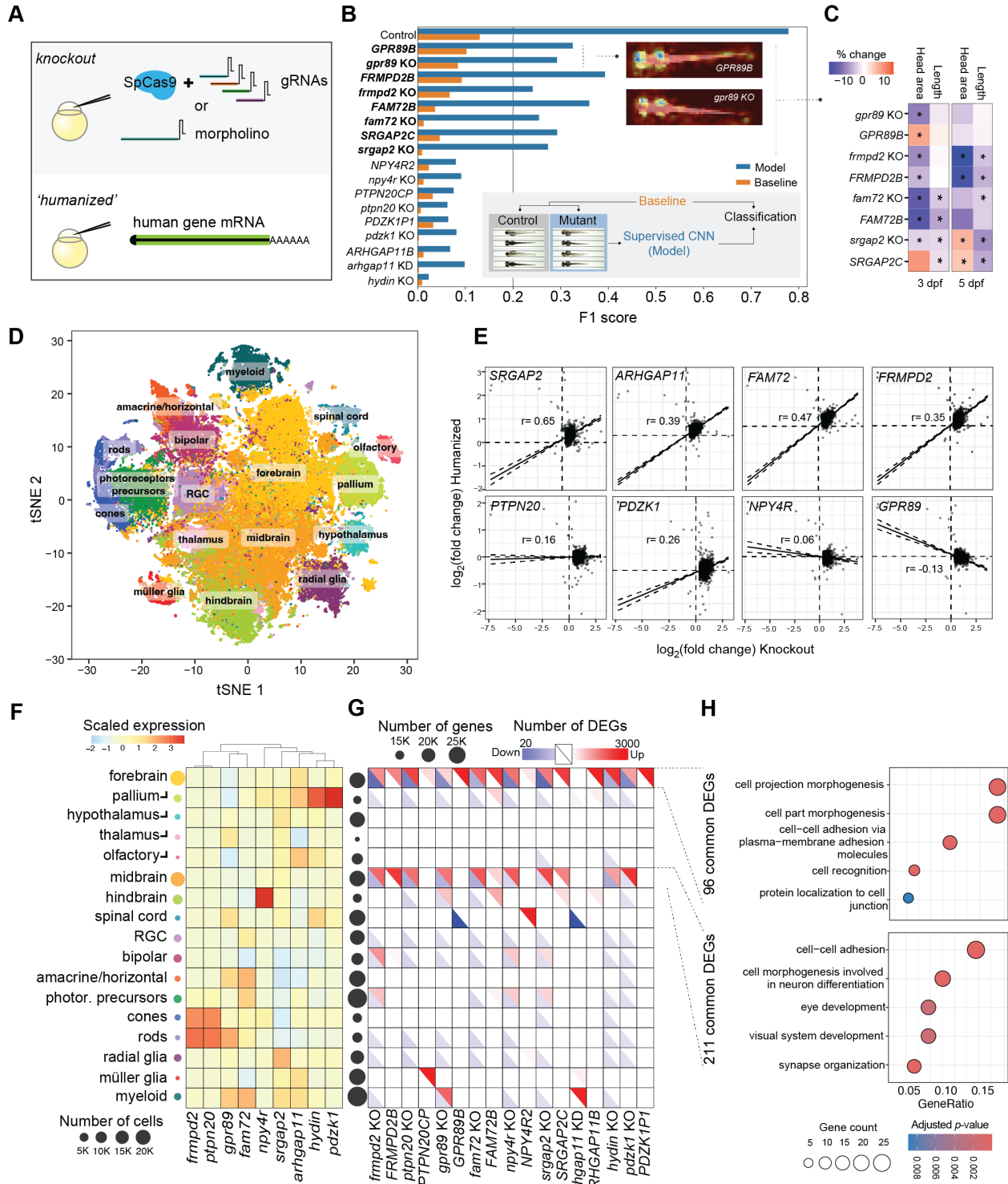for AFR, AMR, and all individuals. (**E**) Scaled transcript per million (TPM) expression of *CD8B* and *CD8B2* in
postmortem brain tissue from BrainSpan.

1    *Duplicated gene functions modeled using zebrafish*

2    We performed a high-throughput functional screen in zebrafish [89–91] of seven largely uncharacterized

3    pHSD families expressed in both human and zebrafish brain (*GPR89*, *NPY4R*, *PTPN20*, *PDZK1*, *HYDIN*,

4    *FRMPD2,* and *FAM72*; Figure 2H and Table S11). Additionally, we tested two gene families (*SRGAP2*

5    and *ARHGAP11*) previously studied in mammals [16,21–23,67,92–96]. Ancestral gene functions were assessed

6    using loss-of-function knockouts for eight zebrafish orthologs by co-injecting SpCas9 coupled with four

7    guide RNAs (gRNAs) targeting early exons resulting in ~70% ablation of alleles in $G_0$ lines [97] (termed

8    crispants). The final gene, *arhgap11*, is maternally expressed (Figure S20) prompting us to use a

9    morpholino that impedes translation. We also 'humanized' zebrafish models by introducing transiently *in*

10   *vitro* transcribed 5'-capped mRNAs encoding human-specific paralogs (Figure 4A) for all genes except

11   *HYDIN2*, due to its large size (4,797 amino acids). There were no significant morbidity differences in

12   mutants compared to controls (log-rank survival tests *p*-values > 0.05, Table S19).

13

14   To assay morphology differences in mutant zebrafish, we acquired images [98–100] for 3,146 larvae at 3 and

15   5 days post-fertilization (dpf) (average of 75±55 larvae per group, Table S20). We first used latent

16   diffusion and convolutional neural networks (CNNs) to test for significant morphological alterations

17   between mutant models and controls without predefining specific features *a priori* (Methods). Both

18   knockout and humanized models of *SRGAP2*, *GPR89*, *FRMPD2*, and *FAM72* exhibited significant

19   differences (F1 scores > 0.2, Figure 4B). Altered features were identified by quantifying body length,

20   head area, and the head-trunk angle, a classic measurement for developmental staging of zebrafish using

21   the same images. This revealed concordant phenotypes for knockout and humanized models of *SRGAP2*

22   (reduced length), and *FRMPD2* (reduced head area), and *FAM72* (both reduced body length and head

23   area) at 3 dpf (Table S21, Figure 4C). Alternatively, *GPR89* models exhibited opposing effects, with head

24   area for *gpr89* knockout larvae ~10% reduced and *GPR89B* 'humanized' larvae ~15% increased. This is

25   also evident in the feature attribution plot indicating that the CNN distinguishes both *gpr89* knockout and

26   *GPR89B* humanized larvae from controls primarily by focusing on the head (Figure 4B). At 5 dpf, the

27   alterations in *FRMPD2* and *SRGAP2* models persisted while no longer observed for *FAM72* and *GPR89*

28   (Table S21, Figure 4C). Knockout models for *gpr89* and *frmpd2* also displayed evidence of

29   developmental delay with subtle yet significant decreases in the head-trunk angle (Table S21).

30

31   We next performed single-cell RNA-sequencing (scRNA-seq) [101,102] of dissected heads of 3 dpf larvae to

32   directly characterize impacts on brain development, profiling 95,555 cells (an average of 3,822±3,227 per

33   model) (Figure 4D). Pseudo-bulk differential expression analysis using all cells in each model revealed

34   significant correlations in gene expression changes versus controls between knockout and humanized

35   models (Figure 4E). Positive correlations for *SRGAP2C*, *FAM72B*, *ARHGAP11B*, *FRMPD2B*, and

36   *PDZK1B* humanized larvae with respect to each knockout indicate loss-of-function effects. *GPR89B* gene

37   expression changes are negatively correlated with *gpr89* indicating gene dosage effects, while

38   *PTPN20CP* and *NPY4R2* show low/no relationship between models. These results are in line with our

39   morphometric findings for *SRGAP2*, *FRMPD*2, *FAM72*, and *GPR89* (Figure 4C), as well as from our

40   separate study [103] that verified the human SRGAP2C protein physically interacts with and antagonizes

41   zebrafish Srgap2.

**Figure 4. Functional evaluation of selected pHSDs using zebrafish. (A)** Functions of each pHSD gene were tested by generating knockout (co-injection of SpCas9 coupled with four gRNAs targeting early exons, or morpholino for *arhgap11*) and 'humanized' models (injection of the human-specific mRNA). **(B)** Morphological assessment using a supervised convolutional neural network (CNN) to distinguish models from matched controls (bottom inset) obtained at 3 dpf and 5 dpf. F1 score indicates the effect size of difference between models and controls and ranges from 0 to 1, where 0 indicates that no sample from that group could be distinguished from the controls. Orange bars indicate the null hypothesis that there is no difference between models and controls. A

1    threshold F1 score of 0.2 was used to define pHSD groups being robustly classified as different from their control

2    group. Pictured as a top inset are feature attribution plots for two example *GPR89B* and *gpr89* knockout larvae,

3    highlighting the region of the image used by the CNN to correctly classify and distinguish those genotypes from

4    controls. Colors range from red (region is not used for classification; zero gradient), to orange, then blue (region

5    contributes the most to classification; large magnitude gradient). **(C)** Percent change compared to the control group

6    for standard length or head area across selected pHSD models. Asterisks indicate a Benjamini Hochberg-corrected

7    *p*-value below 0.05. **(D)** t-distributed stochastic neighbor embedding (tSNE) plot highlighting the classified 17 cell

8    types from the 95,555 harvested cells across pHSD models at 3 dpf. **(E)** Fold-change comparison between knockout

9    and 'humanized' models for each pHSD across all genes (n= 29,945), versus their controls. Black lines represent the

10   Pearson correlation line and the dotted lines the 95% confidence intervals. **(F)** Endogenous z-score scaled

11   expression of each zebrafish ortholog across defined scRNA-seq cell types. Circle sizes scale with the overall

12   number of cells included in that group. **(G)** Distribution of cell-type-specific differentially expressed genes (DEGs)

13   for each pHSD model. Each square includes the downregulated genes in blue (lower diagonal) and upregulated

14   genes in red (upper diagonal). Circles next to each cell type represent the number of expressed genes. **(H)** Gene

15   ontology results for the common DEGs in forebrain (n= 96) and midbrain (n=211) across pHSD models, with circles

16   representing DEG number in the GO term and color representing the *q*-value.

17

18   We classified 17 different neuronal, retinal, and glial cell types using gene markers [71,101,104,105]. While

19   most pHSD orthologs were broadly expressed across cells, a subset showed more narrow expression in

20   specific cell types (e.g., *hydin* and *pdzk1* in the pallium, *npy4r* in the hindbrain; Figure 4F). We repeated

21   pseudo-bulk differential expression analyses across specific cell types revealing gene dysregulation in the

22   forebrain and midbrain across most pHSD models (16 out of 17, Figure 4G). Common differentially

23   expressed genes (DEGs) in the forebrain functioned in cell projection, adhesion, and recognition, while

24   DEGs in the midbrain related to neuronal differentiation and the visual system (Figure 4H). The zebrafish

25   forebrain is the closest related structure to the human cerebral cortex [106], while the midbrain primarily

26   includes the optic tectum [107], the main visual processing center. Some models also highlighted DEGs in

27   specific cell types, including Müller glia in humanized *PTPN20CP*, the spinal cord in humanized

28   *NPY4R2*, and myeloid cells in *gpr89* and *arhgap11* knockout larvae (Figure 4G). Combined, these results

29   indicate that all tested pHSD models impact the developing zebrafish brain, suggesting that they may also

30   play important roles in human brain evolution.

31

32   ***Novel human-specific genes impacting neurodevelopment***

33   <u>*GPR89B* and brain size</u>

34   Opposite phenotypes were observed for *gpr89* knockout and humanized *GPR89B* zebrafish suggesting

35   gene dosage effects. Considering both *GPR89* human paralogs are impacted by deletions and duplications

36   at the chromosome 1q21.1 genomic hotspot associated with microcephaly and macrocephaly in children

37   with neurocognitive disabilities, respectively [58], we sought to characterize mechanisms underlying larval

38   head-size phenotypes in more detail. We first verified that stable *gpr89* heterozygous and homozygous

39   knockouts exhibited reduced head size at 3 dpf, consistent with crispants. Using a neuronal reporter line

40   Tg(HuC-eGFP) [108], we generated *GPR89* mutant models finding significantly smaller and larger

41   forebrains in knockout and humanized larvae, respectively (Figures 5A). Re-examining scRNA-seq data,

42   we sub-clustered cells from the forebrain and observed endogenous expression of *gpr89* in telencephalon

43   and inner diencephalon (Figure 5B). Focusing on the telencephalon, a brain structure anatomically

44   equivalent to the mammalian forebrain with roles in higher cognitive functions such as social behavior

45   and associative learning [109,110], we performed pseudo-bulk differential expression analysis. DEGs with

46   inverse effects were enriched in negative regulation of the DNA replication and cell cycle (Figure 5C,

47   Tables S22 and S23). Several genes functioning at the G2/M checkpoint were downregulated in the
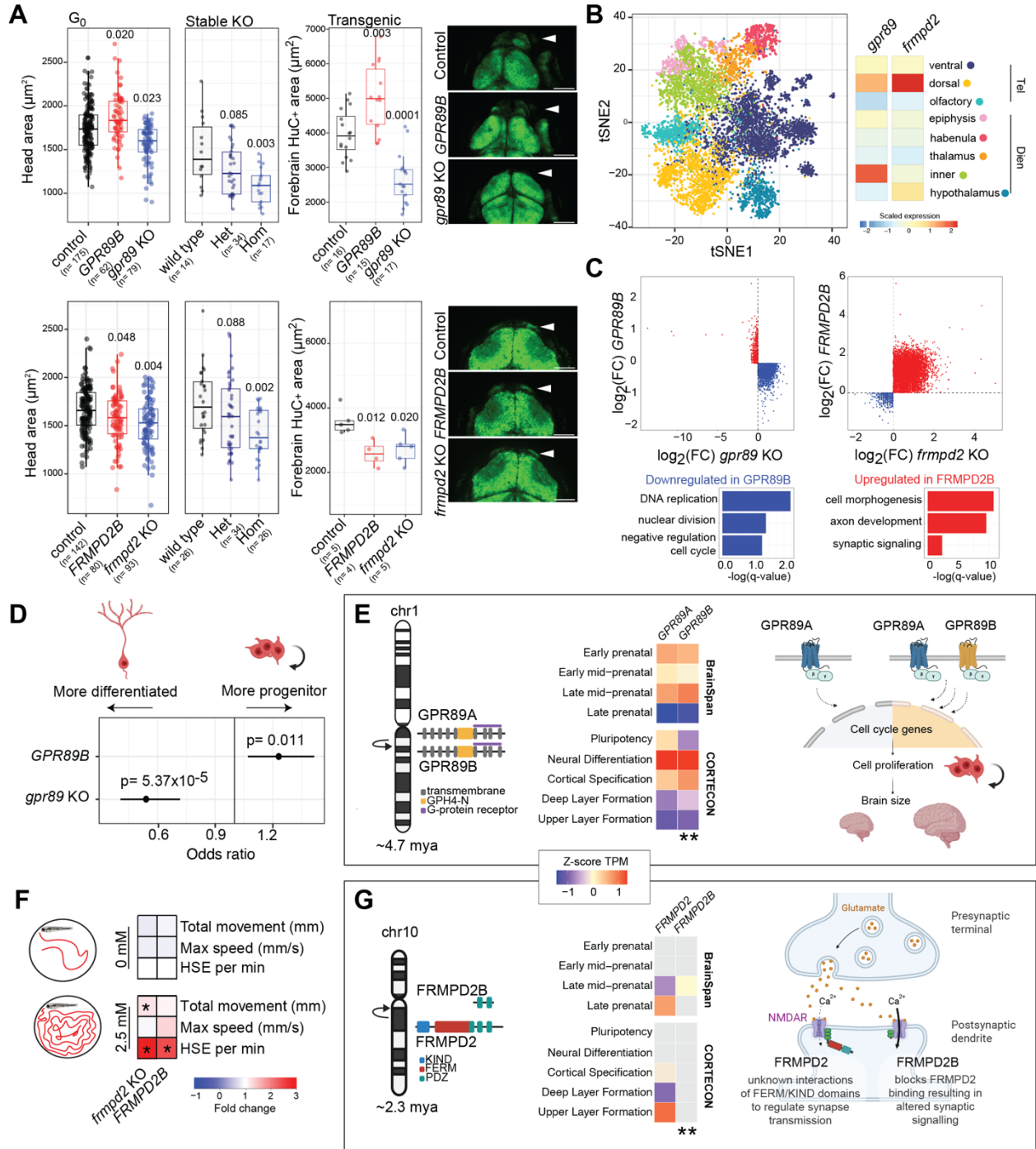
1    humanized *GPR89B* and upregulated in the knockout *gpr89* pointing to differences in cell proliferation.
2    To test this, we estimated the identity of forebrain cells based on the expression of known markers for
3    neural progenitors (*sox19a*, *sox2*, *rpl5a*, *npm1a*, *s100b*, *dla*) and differentiated neurons (*elavl3*, *elavl4*,
4    *tubb5*). This found humanized *GPR89B* cells more likely to classify as progenitors while *gpr89* knockouts
5    more likely to be differentiated (Figure 5D).

6

7    *GPR89* (G-protein receptor 89 or *GPHR*, Golgi PH regulator) encodes highly conserved transmembrane
8    proteins that participate in intracellular pH regulation in the Golgi apparatus [111]. Loss of function in
9    *Drosophila* leads to global growth deficiencies as a result of defects in the secretory pathway [112]. In
10   humans, a complete duplication of ancestral *GPR89A* ~4.7 mya produced the derived, full-length
11   *GPR89B* [13] (Figure 5E). The two paralogs maintain identical protein similarity but differential and
12   overlapping expression patterns in human brain development, with *GPR89A* highly expressed starting in
13   pluripotency (turquoise module), and *GPR89B* expression turning on slightly later during neural
14   differentiation (red module; Figures 2E and 5E). Both genes are under purifying selection (Figure S14),
15   with *GPR89A* exhibiting extreme negative Tajima's D values in individuals of AFR and AMR ancestries
16   from the 1KGP cohort (<5$^{th}$ percentile; Figure 1E, Table S1). These results provide evidence that both
17   *GPR89* paralogs function in early brain development, possibly with delayed expression of *GPR89B*
18   extending expansion of progenitor cells, a feature observed in human cerebral organoids compared with
19   those of other apes [113,114] (Figure 5E). Together with the increase in forebrain size of "humanized"
20   zebrafish, this suggests a role for *GPR89B* in contributing to the human-lineage expansion of the
21   neocortex.

22

23   *FRMPD2B* and synaptic signaling
24   While opposing traits were observed in *GPR89* models, similar phenotypes impacting head area and body
25   length suggest that the human FRMPD2B acts as a dominant negative to the endogenous Frmpd2.
26   Validating phenotypes observed in crispants (Figure 4C), we observed reduced head size in stable *frmpd2*
27   homozygous knockout larvae (Figure 5A). Additionally, both the crispant knockout *frmpd2* and
28   humanized *FRMPD2B* larvae exhibit smaller forebrains. We found that shared upregulated DEGs
29   function in cell/axon morphogenesis and growth as well as synaptic signaling in telencephalic cells
30   (Figure 5C, Tables S24 and S25). To better characterize impacts on synaptic signaling, we used motion-
31   tracking [115] to detect seizure susceptibility in mutant zebrafish. Treatment with a low dose of the GABA-
32   antagonizing drug pentylenetetrazol (PTZ) produced a significant increase in high-speed events,
33   indicative of seizures in larvae, in both *FRMPD2* mutant models (4 dpf) versus controls (Figure 5F).
34   These results suggest that Frmpd2 loss of function, through *frmpd2* knockout or antagonism via
35   FRMPD2B, disrupts excitatory synapse transmission which amplifies induced seizures, in line with the
36   known interactions of FRMPD2 with glutamate receptors [116].

37

38   *FRMPD2* (FERM and PDZ domain containing 2) encodes a scaffold protein that participates in cell-cell
39   junction and polarization [117]. Protein localization has been observed at photoreceptor synapses [118] and the
40   postsynaptic membrane in hippocampal neurons in mice [116]. A partial duplication of the ancestral
41   *FRMPD2* on human chromosome 10q11.23 created the 5'-truncated *FRMPD2B* paralog ~2.3 mya [13]. This
42   shorter *FRMPD2B*-derived paralog encodes 320 amino acids of the C-terminus, versus 1,284 amino acids
43   for the full-length ancestral, maintaining two of three PDZ domains involved in protein binding [119] while
44   lacking the KIND and FERM domains (Figure 5G).

**Figure 5. Neurodevelopmental impact of *GPR89* and *FRMPD2*. (A)** Head and brain area assessments at 3 dpf for $G_0$ crispants and stable knockout lines for *GPR89* (top) and *FRMPD2* (bottom) models. Results for head area of *GPR89* crispants (ANOVA *p*-values: controls vs. *GPR89B*= 0.020, controls vs. *gpr89* knockouts= 0.023) and stable knockout lines (Wilcoxon signed-rank tests *p*-values: controls vs. Het= 0.085, controls vs. Hom= 0.003), as well as forebrain area of crispants using a transgenic line with fluorescently tagged neurons (ANOVA *p*-values: controls vs. *GPR89B*= 0.003, controls vs. *gpr89* knockouts= 0.0001). Results for head area of *FRMPD2* crispants (ANCOVA *p*-values: controls vs. *FRMPD2B*= 0.048, controls vs. *gpr89* knockouts= 0.004) and stable knockout lines (Wilcoxon signed-rank tests *p*-values: controls vs. Het= 0.088, controls vs. Hom= 0.002), as well as forebrain area of crispants using a transgenic line with fluorescently tagged neurons (ANOVA *p*-values: controls vs. *FRMPD2B*= 0.012, controls vs. *frmpd2* knockouts= 0.020). Representative images of each model in the neuronal transgenic line are

1  included with scale bars representing 100 μm. **(B)** t-distributed stochastic neighbor embedding (tSNE) plot showing
2  the identified subregions classified from the forebrain (n=10,040 cells) and relative scaled endogenous expression of
3  *gpr89* and *frmpd2* across cell types. **(C)** Log$_2$ fold change (FC) of gene expression versus controls in cells from the
4  telencephalon between knockout and humanized models in *GPR89* and *FRMPD2*. Red and blue colors correspond
5  to DEGs discordant (*GPR89*) or concordant (*FRMPD2*) between the knockout and humanized models and their top
6  representative gene ontology enrichment analyses results. **(D)** Forest plot with the results from the logistic
7  regression for presence of progenitor versus differentiated states in forebrain cells across *GPR89* models.
8  **(E)** Diagram of the duplication event of *GPR89* giving rise to *GPR89A* and *GPR89B*, encoding two identical
9  proteins with different expression patterns in both neurodevelopmental timing (BrainSpan) and brain regions
10  (CORTECON) (**Wilcoxon signed-rank test, *p*-value < 0.005). A model of GPR89B gain-of-function in neuronal
11  proliferation amplification is depicted on the right. **(F)** Behavioral results from 1 h motion-tracking evaluations in 4
12  dpf larvae exposed (2.5 mM) or not (0 mM) to pentylenetetrazol (PTZ). Metrics compared included total movement
13  (mm), maximum speed (mm/s), and frequency of high-speed events (≥28 mm/s). Colors represent the fold change
14  relative to the control group and the asterisk indicates a significant Dunn's test (BH-adjusted *p*<0.05). **(G)** Diagram
15  of the duplication event of *FRMPD2* ~2.3 mya that gave rise to the 5'-truncated *FRMPD2B*, which exhibits different
16  temporal (BrainSpan) and spatial (CORTECON) expression patterns (**Wilcoxon signed-rank test, *p*-value <
17  0.005). A model of FRMPD2B antagonistic functions resulting in altered synaptic signaling is depicted on the right.
18
19  Our data shows ancestral *FRMPD2* expressed in the human prenatal cortex during upper layer formation,
20  while *FRMPD2B* is evident only postnatally [64] (Figure 5G). The paralogs also show divergent
21  evolutionary signatures, with the full-length *FRMPD2* strongly conserved and the truncated *FRMPD2B*
22  exhibiting possible positive selection (Figure 3B,C). Results in zebrafish show that loss of Frmpd2
23  function results in microcephaly and enhanced excitatory synaptic signaling. Combined, we propose a
24  model in which truncated human-specific FRMPD2B counteracts the function of full-length FRMPD2
25  leading to altered synaptic features in humans, possibly through interactions of its PDZ2 domain with
26  GluN2A of NMDA receptors at the postsynaptic terminal [116]. Its postnatal expression would avoid the
27  detrimental effects of inhibiting FRMPD2 during early fetal development (i.e., microcephaly). We note
28  that recurrent deletions and duplications in chromosome 10q11.21q11.23 impact both paralogs in children
29  with intellectual disability, autism, and epilepsy [120]. Ultimately, *FRMPD2B* could plausibly contribute to
30  the upregulation of glutamate signaling and increased synaptic plasticity observed in human brains
31  compared with other primates that is fundamental to learning and memory [121].
32

## Discussion

34  Our results provide the scientific community with a prioritized set of hundreds of genes to perform
35  functional analyses with the goal to identify drivers of human brain evolution. Using a complete T2T-
36  CHM13 reference genome, we present the most comprehensive detection of human duplicate genes to
37  date with 213 families and 1,002 total paralogs. Compared to a previous assessment of human-specific
38  duplicated genes [13], this represents an approximately fivefold increase in identified genes, in part because
39  we also included human-expanded gene families and genes with as little as one duplicated exon. We note
40  that these numbers are likely an underestimate, as we excluded 193 high-copy gene families (famCN>10),
41  as well families that have undergone independent gene expansions or incomplete lineage sorting with
42  other great apes. One compelling example is *FOXO3*, encoding the transcription factor forkhead box O-3,
43  implicated in human longevity [122], with all three paralogs CN-constrained and brain expressed (Table S1).
44  Since this gene also exists as duplicated in other great apes at similar CN, we excluded it from our list of
45  human gene expansions. This is, in part, because there is still uncertainty regarding which paralog(s) are
46  human specific for many of the gene families. SDs are often accompanied by secondary structural
47  rearrangements that hamper synteny comparisons across species [57,123]. Moving forward, the availability of

1    nonhuman primate T2T genomes will improve orthology and synteny comparisons between species [124–126]

2    revealing additional human-specific paralogs. As a resource for the community, we have made available

3    the results of our genome-wide analyses across the complete 1,793 SD98 genes (Tables S1-S11).

4

5    Collectively, 148 gene families (362 paralogs, 108 annotated as non-syntenic with the chimpanzee

6    reference) represent top candidates for contributing to human-unique neural features based on at least one

7    gene member exhibiting functional constraint across modern humans (1KGP) and brain expression (Table

8    S1). In this study, we chose zebrafish to demonstrate the efficacy of our gene list. Despite notable

9    differences with humans, such as the absence of a neocortex [127], conservation in major brain features

10    make zebrafish well-suited to characterize gene functions in neurological traits, including cranial

11    malformations [128], neuronal imbalances [129], and synaptogenesis [130]. Coupled with CRISPR mutagenesis

12    [89,90], zebrafish have been used as a higher-throughput model for human neurodevelopmental conditions

13    such as epilepsy [115], schizophrenia [131], and autism [73]. While a whole-genome teleost duplication resulted

14    in ~20% of genes with multiple zebrafish paralogs that confounds functional analysis of human gene

15    duplications [132], the nine prioritized gene families tested here were selected in part because each had only

16    one zebrafish ortholog. We characterized gene functions by knocking out the conserved ortholog and

17    introducing the human-specific paralog into developing embryos. Transient availability of the human

18    transgene by injection of *in-vitro*-transcribed mRNA limited our analysis to early developmental traits (up

19    to 5 dpf in zebrafish), approximately equivalent to human mid- to late-fetal stages in brain development

20    (Figure S10). In the future, it will be important to characterize phenotypes in adolescent and adult

21    zebrafish by generating stable transgenic humanized lines.

22

23    From our analysis, knockout and humanized models of four genes (*GPR89*, *FRMPD2*, *FAM72*, and

24    *SRGAP2*) resulted in altered morphological features, primarily to head size (often used as a proxy for

25    brain size), and all models exhibited molecular differences in single-cell transcriptomic data, most evident

26    in the fore- and midbrains of larvae (Figure 4G). Two duplicate gene families, *SRGAP2* and *ARHGAP11*,

27    have been extensively studied in diverse model systems (reviewed recently [9]). Our zebrafish model of

28    *SRGAP2*, encoding SLIT-ROBO Rho GTPase-activating protein 2, were consistent with published

29    findings in mouse where the 3'-truncated human-specific SRGAP2C inhibits the function of the

30    endogenous full-length Srgap2 [16]. Further, the shared upregulated genes identified in the forebrains of

31    *SRGAP2* mutant larvae point to alterations in axonogenesis and cell migration, matching studies in mice

32    [11,16,17,67,93,133,134] (Table S26). Alternatively, *ARHGAP11B*, encoding Rho GTPase Activating Protein 11,

33    implicated in the expansion of the neocortex through increased neurogenesis [21,23], exhibited no detectable

34    changes in head/brain size when introduced in zebrafish embryos. Upregulated DEGs were only detected

35    in the forebrains of *ARHGAP11B*-injected mutants and were enriched in cellular biosynthetic processes

36    (mRNA splicing and translation; Table S27). Given that ARHGAP11B impacts the abundance of basal

37    progenitors, a cell type unique to the mammalian neocortex [135], zebrafish may not be suitable to

38    characterize human-specific functions of this gene.

39

40    Beyond modeling gene functions, our study also highlighted the considerable amount of genetic variation

41    hiding within SD regions. Even with the resolved gaps and errors across SDs in T2T-CHM13, short-read

42    sequencing is still insufficient to identify variation. Due to high sequence identity, only 10% of SD98

43    regions are "accessible" to short reads [36] resulting in <10% sensitivity to detect variants (Note S1) and a

44    depletion of GWAS hits (Note S2). Using existing assemblies (HPRC and HGSVC) and cHiFi

1   sequencing of individuals of diverse ancestry uncovered some of this hidden variation within our 13
2   pHSD gene families. We note that, for some of the most highly identical duplicated genes (*CFC1*), our
3   cHiFi reads (~3 kbp) were still too short to accurately map to respective paralogs (data not included).
4   Nevertheless, long reads revealed that most pHSD paralogs exhibit evolutionary constraints and provided
5   support for balancing selection of *CD8B*, not previously identified in published genome-wide screens
6   [136,137]. Historically, signatures of balancing selection, which include an excess of mid-frequency alleles [138],
7   have been difficult to detect within SDs due to assembly errors [36]. In these cases, paralog-specific variants
8   are mistaken for SNPs when reads from both paralogs map to a single collapsed locus resulting in false
9   mid-frequency alleles. Scientific consortia like *All of Us* are generating long-read datasets at scale [139],
10  ushering in a new era where genomic associations and evolutionary selection may finally be uncovered
11  within human duplications to identify novel drivers of human traits and disease.
12
13  Similarly, genome sequencing of patients and their families has discovered hundreds of compelling
14  neuropsychiatric disease candidate genes impacted by rare and *de novo* variants, but the genetic risk
15  underlying conditions such as autism is still not completely elucidated [140]. SD genes may represent a
16  hidden contributor to disease etiology. Our analysis identified 82 SD98 genes (38 human duplicate
17  paralogs) co-expressed in modules enriched for ASD genes (Figure 2E), including several within disease-
18  associated genomic hotspots. Distinct SD mutational mechanisms, including ~60% higher mutation rate
19  compared to unique regions [141] and interlocus gene conversion that can occur between paralogs [142,143],
20  make duplicated genes particularly compelling to screen for *de novo* mutations contributing to idiopathic
21  conditions. For example, nonfunctional paralogs with truncating mutations can "overwrite" conserved
22  functional paralogs leading to detrimental consequences, as is the case of *SMN1* and *SMN2* in spinal
23  muscular atrophy [39]. Human-duplicated gene families include ancestral paralogs *CORO1A*, *TLK2*, and
24  *EIF4E*, with significant genetic associations with ASD [57]. We propose that interlocus gene conversion
25  between their likely nonfunctional duplicate counterparts is an understudied contributor to
26  neurodevelopmental conditions in humans. Our comprehensive list of gene families will enable future
27  work to progress in this research area.
28
29  Our study focuses on duplicate genes functioning in brain development, but primates exhibit other
30  prominent differences across musculoskeletal and craniofacial features that have diverged early in human
31  evolution [4]. Since such traits are largely universal across modern humans, our list of CN-constrained
32  genes represent top candidates though re-analysis of transcriptomes from non-brain cells/tissues is
33  required. Meanwhile, duplicate genes, such as those encoding defensins [144–147], mucins [148,149], and
34  amylases [41–43], can also play a role in metabolism and immune response that exhibit population
35  diversification due to the vast variability in diet, environment, and exposures to pathogens across modern
36  humans [27]. Our use of a single complete human T2T-CHM13 haplotype of largely European ancestry [34]
37  could miss some of these CN polymorphic genes. As additional T2T genomes are released [28], it will be
38  important to continue curating our list of duplications. Nevertheless, genes CN stratified by human
39  ancestry can be identified using metrics such as $V_{st}$ [150], as has been highlighted in other studies (reviewed
40  here [9] and most recently in a preprint [151]). Facilitating such analyses for our gene set, we provide a
41  publicly available resource to query parCN median estimates across individuals from 1KGP for our
42  complete set of SD98 paralogs (https://dcsoto.shinyapps.io/shinyc).
43

1  One notable limitation of our study is its reliance on existing gene annotations. We attempted for the first
2  time to group human duplicate paralogs into larger multigene families based on shared sequences between
3  annotated genes in SD98 regions. Due to the complexities of SDs, which can result in gene fusions and
4  altered gene structures, some genes were left unassigned to a family (n=114 singletons from SD98 genes).
5  Other noncoding transcripts and lncRNAs were excluded altogether, including a human-specific paralog
6  of *IQSEC3*, a gene implicated in GABAergic synapse maintenance [152]. Additionally, the functional
7  consequences of variants identified in 656 unprocessed pseudogenes are difficult to interpret.
8  Improvements are on the horizon, with ongoing work with long-read transcriptomes that will continue to
9  refine annotations [153] and advancements in protein-prediction models [154] and proteomic approaches [155] that
10  will confirm whether or not these genes encode proteins.
11
12  In summary, we identified and featured two genes with strong evidence of contribution to human brain
13  evolution: *GPR89B*, with a possible role in expansion of the neocortex, and *FRMPD2B*, with implications
14  in altered synaptic excitatory signaling. Taking advantage of long-read sequencing in tandem with the
15  new T2T-CHM13 reference genome, we interrogated challenging regions of the genome and
16  demonstrated a method using zebrafish to explore the functions of human-duplicated genes. Among our
17  list of hundreds, we propose that there are additional gene drivers that contribute to unique features of the
18  human brain. In the future, additional genetic analyses across modern and archaic humans and
19  experiments utilizing diverse model systems will reveal hidden roles of these genes in human traits and
20  disease.
21

## Methods

### Identification of SD98 genes

24  Duplicated regions were extracted from previously annotated SDs [156] using T2T-CHM13 (v1.0)
25  coordinates and subsequently merged using BEDTools merge [157]. SD98 regions were defined as an SD
26  with ≥98% sequence identity to another locus in the T2T-CHM13 genome using the fractMatch
27  parameter. Gene coordinates were obtained from T2T-CHM13 (v1.0) CAT/Liftoff annotations (v4) [34].
28  SD98 genes were defined as gene annotations that contain at least one exon fully contained within an
29  SD98 region, calculated with BEDTools intersect using -f 1 parameter [157]. Overall numbers of distinct
30  gene features overlapping SD98 were counted using the gene ID unique identifiers. We noticed that, in a
31  few cases, two transcript isoforms of the same gene were assigned to different gene IDs. To identify these
32  redundant transcripts, we self-intersected SD98 transcripts, selected those with different gene ID that also
33  shared >90% positional overlap, and performed manual curation of the obtained gene list, removing
34  redundant and read-through fusion transcripts.

### Gene family clustering

36  SD98 genes were grouped into gene families based on shared exons (Figure S1). Starting from T2T-
37  CHM13 (v1.0) annotations, DNA sequences of all SD98 regions were extracted using BEDTools getfasta
38  and mapped back to the reference genome using minimap2 (v2.17) with the following parameters: -c --
39  end-bonus 5 --eqx -N 50 -p 0.5 -t 64. For each SD98 exon, the BEDTools intersect with -f 0.99 parameter
40  was used to select mappings covering >99% of the exon sequence, removing self-mappings. This list was
41  refined using the previously published [35]whole-genome shotgun sequence detection (WSSD) [10] CNs
42  (famCN) of humans from the SGDP (n=269), which provides estimates of the overall CN of a gene

1  family using read depth of multi-mapping reads with nonoverlapping sliding-windows. After comparing
2  the median famCN values of SD98 genes with shared exons, groupings where the mean absolute
3  deviation of the CN was less than one were selected. The list was filtered to focus on gene families
4  containing at least one protein-coding or unprocessed pseudogene. SD98 genes associated with other gene
5  features, including lncRNAs and processed pseudogenes, were also assigned a gene family ID. On the
6  other hand, if a gene was not associated with any other gene feature, they were classified as "unassigned"
7  or "singletons". SD98 gene families were intersected with previously published DupMasker annotations
8  using BEDTools intersect, which indicate ancestral evolutionary units of duplication [35].

9  **Identification of human-duplicated genes families**
10  Human-specific and -expanded gene families were identified using CN comparisons between humans and
11  nonhuman great apes with previously published WSSD [10] (famCN) CNs from humans (SGDP n=269) and
12  four nonhuman great apes, including one representative of chimpanzee (Clint), bonobo (Mhudiblu),
13  gorilla (Kamilah), and orangutan (Susie) [35], mapped to T2T-CHM13 (v1.0). The median famCN per
14  SD98 gene was calculated using a custom Python script. For each SD98 gene, putative gene family
15  duplications and expansion were predicted, excluding genes with median famCN>10 across humans from
16  this analysis. Genes were considered expanded if the median famCN across humans was greater than the
17  maximum famCN across great apes. Human duplications and expansions were distinguished based on
18  whether the maximum famCN value across great apes was less than 2.5 (non-duplicated in great apes) or
19  greater than 2.5 (duplicated in great apes), respectively. Non-syntenic paralogs between humans and
20  chimpanzees were obtained using previously published syntenic data between human (T2T-CHM13v1.0)
21  and chimpanzee (PanTro6) references [35] intersected with SD98 genes using BEDTools intersect. For each
22  paralog, family status was designated as "Human-duplicated gene family" if it was assigned to a gene
23  family containing at least one expanded or duplicated member according to famCN and/or at least one
24  non-syntenic member based on human/chimpanzee synteny. Otherwise, family status was considered
25  "Undetermined".

26  **Paralog-specific copy number genotyping**
27  parCN estimates were obtained using QuicK-mer2 [46] for 1KGP 30× high-coverage Illumina individuals [55]
28  and four archaic genomes (including Altai Neanderthal [PRJEB1265] [47], Vindija Neanderthal
29  [PRJEB21157] [48], Mezmaiskaya Neanderthal [PRJEB1757] [47,48], and Denisova [PRJEB3092] [49]), using
30  T2T-CHM13 (v1.0) as reference [34]. The resulting BED files containing parCN estimates were converted
31  into bed9 format using a custom Python script for visualization in the UCSC Genome Browser. parCN
32  values were genotyped across SD98 regions overlapping protein-encoding and unprocessed pseudogenes
33  by calculating the mean parCN across the region of interest for each sample using a custom Python script.
34  parCN dotplots generated using the R package ggplot2 are available for SD98 genes as an interactive
35  Shiny web application in https://dcsoto.shinyapps.io/shinycn.

36  **Metrics of selective constraint**
37  Loss-of-function intolerance of SD98 genes was assayed using previously published gnomAD (v2.1.1)
38  probability of loss-of-function intolerance scores (pLI) [158] and loss-of-function observed/expected upper
39  fraction (LOEUF) [52]. We considered genes as intolerant to loss of function if either their pLI scores were
40  greater than 0.9 or their LOEUF scores were less than 0.35.

**Genome-wide Tajima's D analysis**

Additionally, Tajima's D [54] values were calculated using previously published SNPs obtained from high-coverage short-read sequencing data from unrelated 1KGP individuals (n=2,504) [55], remapped to T2T-CHM13 (v1.0) [36]. Windows were defined as SD98 if at least 10% of the bases corresponded to SD98 regions. To define short-read accessible windows, 25-kbp windows were intersected with a published short-read combined accessibility mask [36]. Considering that no SD98 windows were fully accessible (Figure S4), Tajima's D was calculated for each superpopulation using VCFtools [159] across 25-kbp windows of at least 50% accessibility and with five or more SNPs. Because previous studies have highlighted potential discrepancies of evolutionary constraints experienced between duplicated and non-duplicated genomic loci [160], outlier D values were calculated for each continental superpopulation as the $5^{th}$ and $95^{th}$ percentiles within SD98 windows only, thereby avoiding comparisons between duplicated and unique regions. Outlier threshold values for each population were defined as follows: AFR, -2.21 and -0.67; EUR, -2.37 and 0.08; EAS, -2.48 and -0.10; SAS, -2.40 and -0.28; and AMR, -2.40 and -0.41.

**Association with neural traits**

Gene-disease associations were obtained from the GWAS catalog v1.0 [161]. SNPs significantly associated with brain measurements ($p$-value < 0.05) were selected, and the GWAS "mapped genes" were intersected with the SD98 gene list using gene symbols. Similarly, previously published associations between CNVs and neural traits in the UKBB were obtained [56]. Coordinates of CNVs significantly associated with brain measurements ($p$-value<0.05) were lifted over from hg19 to hg38 and from hg38 to T2T-CHM13 (v1.0) using UCSC liftOver tool [162]. Liftover chains were obtained from the UCSC Genome Browser and T2T-CHM13 GitHub page (https://github.com/marbl/CHM13, previous assembly releases of T2T-CHM13), respectively. CNVs were intersected with SD98 gene coordinates using BEDTools intersect [157].

ParCN values from SD98 genes for families with autistic children from the SSC (n = 2,459 families, n = 9,068 individuals) mapped to the T2T-CH13v1.1 reference genome were obtained, following the same steps as described to genotype parCN across 1KGP individuals. Overall, CN differences between autistic probands and unaffected siblings were compared by rounding median CN per individual to the nearest integer, and significance was assessed using the Wilcoxon signed-rank test, correcting for multiple testing with the false discovery rate method. To identify *de novo* deletions or duplications in autistic probands and unaffected siblings, parCN values within ±0.2 of an integer were conservatively selected and rounded to the nearest integer for all family members. Intermediate values, which could potentially confound the analysis, were removed. *De novo* events were classified as cases where both parents exhibited a parCN=2, while the child showed a parCN=3 (duplication) or parCN=1 (deletion).

Previously published genomic hotspots [60] were obtained in hg19 coordinates and lifted over to hg38 and from hg38 to T2T-CHM13 (v1.0) using the UCSC liftOver tool and associated chain files (described above). Three regions failed the liftover process due to differences in reference genome sequences. An extra 500 kbp were added upstream and downstream of each reported genomic hotspot to account for breakpoint errors. SD98 genes, including those exhibiting putative *de novo* events in the SSC dataset, were intersected with expanded genomic hotspots coordinates using BEDTools intersect.

**Gene expression and network analysis**

Previously published brain transcriptomic datasets, including post-mortem tissue and cell lines, were obtained. These datasets included neocortical germinal zones [61], neural stem and progenitor cells [21], a neuroblastoma cell line SHSY5Y [62], and two longitudinal studies of *in vitro* induced neurogenesis from human embryonic stem cells [63] (CORTECON), and post-mortem brain (BrainSpan) [64]—the latter of which was separated into prenatal and postnatal samples. Raw reads were pseudo-mapped to T2T-CHM13 (v2.0) CAT/Liftoff transcriptome and the CHM13v2.0 assembly as decoy sequence using Salmon v1.8.0 [163] with the flags "--validateMappings --gcBias". Transcripts per million (TPM) values and raw counts were summed to the gene level using tximport [164]. An SD98 gene was considered expressed during development if TPM values were greater than one in at least one of these samples, excluding postnatal BrainSpan data. Conversely, an SD98 gene was considered expressed postnatally if TPM values were greater than one in at least one postnatal stage of BrainSpan.

Gene co-expression analysis was performed using the WGCNA R package [65]. Briefly, samples were analyzed using principal components and hierarchical clustering to assess outliers, removing two samples (SRR1238515 and SRR1238516). Features with consistently low counts across remaining samples (counts <10 in 90% samples) were removed from this analysis. Raw counts for each sample were normalized using variance stabilizing transformation before performing a signed network construction with function blockwiseModules, with parameters soft power = 24, deepSplit = 4, detectCutHeight = 0.995, minModuleSize = 30, and MergeCutHeight = 0.15. GO terms enrichment analysis was performed using the R package clusterProfiler ego function [165]. Enrichment of gene categories were performed using the hypergeometric test in R for autism genes [57], expanded genomic hotspots [60], and cell markers [68].

Visualization of the yellow network was constructed by selecting genes with module membership greater than 0.5, generating an adjacency matrix with remaining genes, and then reconstructing a signed network with soft threshold = 18. Edges with Pearson correlation <0.1 were removed. The network visualization was built with the igraph R package (https://r.igraph.org/), using layout_with_fr for vertex placement. Vertex size was proportional to the degree and edges width was proportional to the Pearson's correlation coefficient. Some vertices were manually adjusted to improve aesthetics of the plot.

**Mouse and zebrafish orthologs**

Mouse-human orthologs were obtained from the Mouse Genome Informatics (MGI) complete list of human and mouse homologs and ENSEMBL BioMart, intersected with SD98 genes using gene symbols, and manual curation. Zebrafish-human orthologs were obtained from combined ENSEMBL BioMart annotations, MGI complete list of vertebrate homology classes, and manual curation. MGI files were downloaded from their website (https://www.informatics.jax.org/homology.shtml) and BioMart analyses were performed using the R package biomaRt. Comparison of developmental brain expression of SD98 orthologs in model organisms was performed using previously published expression data for mouse (PRJNA637987) [70] and zebrafish (GSE158142) [71], calculating Z-score normalized TPM values. Matching of developmental stages across human, mouse, and zebrafish was done as previously described [72]. In brief, genes with one-to-one orthologs with human genes were identified (mouse n= 19,949; zebrafish n= 16,910) and the principle component analysis rotations of the human BrainSpan data used to predict PC coordinates for the mouse and zebrafish data in human principle component space.

**Capture HiFi (cHiFi) sequencing**

We performed cHiFi sequencing of 172 individuals from the 1KGP, two trios from Genome in a Bottle [166], and 22 HGDP individuals with available linked-read data via the 10X Genomics platform [167], totaling 200 samples and 18 family trios (Table S13). DNA samples for 1KGP and Genome in a Bottle were obtained from the Coriell Institute (Camden, NJ, USA) and HGDP samples were obtained from the CEPH Biobank at the Fondation Jean Dausset-CEPH (Paris, France). PacBio cHiFi sequencing was performed using the RenSeq protocol [168]. Briefly, genomic DNA (~4 μg) was sheared to approximately 3 kbp with the Covaris E220 sonicator using Covaris blue miniTUBEs, followed by purification and size selection with AMPure XP beads. End repair and adapter ligation were performed using the NEBNext Ultra DNA Library Prep Kit. Barcodes to distinguish each sample were added via PCR using Kapa HiFi Polymerase (Roche, CA, USA). After the first PCR (fewer than 9 cycles), the libraries were purified and size-selected. For target enrichment, 80-mer RNA baits were designed and tiled at 2× coverage across targeted SD regions and unique exonic regions (Table S28). pHSD regions of interest were targeted and enriched for using a custom myBaits kit (Arbor Biosciences, MI, USA) following manufacturer's recommended protocol. Eight pooled barcoded libraries were hybridized overnight to the baits, and the captured DNA was bound to Dynabeads MyOne Streptavidin C1 beads. A second PCR was performed post-hybridization to generate sufficient material for sequencing. A PCR cycle test was conducted prior to the second amplification to limit PCR duplication bias.

The final libraries were size-selected using the Blue Pippin system to enrich for fragments >2 kbp and sequenced on the PacBio Sequel II platform (Maryland Genomics, University of Maryland). Briefly, Sequel II libraries were constructed using SMRTbell Express Template Prep Kit 2.0 (Pacific Biosciences, Menlo Park, CA) according to manufacturer's instructions. In brief, DNA samples were treated with DNA-damage repair enzymes followed by end-repair enzymes before being ligated to overhang sequencing adaptors. Libraries were then purified with SPRI beads (Beckman Coulter, Indianapolis, IN) and quantified on the Femto Pulse instrument (Agilent Technologies, Santa Clara, CA). Prior to sequencing, libraries were bound to Sequel II polymerase, then sequenced with Sequel II Sequencing kit and SMRT cell 8M on the Sequel II instrument (Pacific Biosciences, Menlo Park, CA). The approach yielded ~3-kbp reads with an average coverage of 27× across regions of interest, considering reads with MAPQ greater than 10 (Table S29).

**Long-read genetic variation discovery and analysis**

cHiFi reads were processed using the standard PacBio SMRT sequencing software tools available in the Conda repository pbbioconda. Circular consensus was obtained from subreads using CCS command with the following parameters --minPasses 3 and --minPredictedAccuracy 0.9. PacBio adapters and sample barcodes were removed using lima software and duplicates were removed with pbmarkdup. Resulting cHiFi reads were aligned to T2T-CHM13v1.0 reference using pbmm2 align, a wrapper of minimap2, with the CCS preset and default parameters. Read groups were added with Picard AddOrReplaceReadGroups and variants were called on each sample using GATK HaplotypeCaller [169], using ploidy = 2 and minimum mapping quality thresholds for genotyping of 0, 2, 5, 10 and 20, resulting in gVCF files per sample for joint genotyping. Joint genotyping was performed with GATK CombineGVCFs and GenotypeGVCFs tools using the pedigree file for accurate calculation of inbreeding coefficients. Genotyping was performed using minimum confidence threshold of 0, 10, 20 and 30. As the technical profile of variants in SDs differs from Variant Quality Score Recalibration training sets, a hard-filtering approach was utilized,

1 including genotyping quality threshold of 0, 20, 50, 70, and depth of 0, 4, 8, 12, and 16. Based on a
2 comprehensive benchmark, including comparison of cHiFi with HPRC/HGSCV and trio mendelian
3 concordance (Note S3), the following minimum thresholds were selected: mapping quality of 20,
4 confidence of 30, genotype quality of 20, and depth of 8. Only unrelated samples from the 1KGP were
5 selected for downstream population analyses (n=144).

6

7 Fully phased haplotypes from 47 individuals from the HPRC Year 1 freeze (https://github.com/human-
8 pangenomics/HPP_Year1_Data_Freeze_v1.0) and 15 from the HGSVC [79] were downloaded. Each
9 haplotype was mapped to T2T-CHM13v1.0 reference genome using minimap2 with parameters -a --eqx -
10 -cs -x asm5 --secondary=no -s 25000 -K 8G, and unmapped contigs and non-primary alignments were
11 discarded. For each region of interest, the longest alignment spanning the locus was selected and
12 additional alignments were removed. This process ensured that one single contiguous contig was used for
13 variant detection. Variants were called with htsbox pileup with parameters -q 0 -evcf and converted into
14 diploid calls using dipcall-aux.js vcfpair. For each region of interest, individual sample calls were merged
15 into a multi-sample VCF file using BCFtools merge, only including individuals whose two haplotypes
16 fully spanned the region of interest. Redundant samples between the HPRC and HGSVC (HG00733,
17 HG02818, HG03486, NA19240, NA24385) were removed, prioritizing HPRC assemblies. Finally, the
18 HPRC/HGSVC dataset was merged with cHiFi variants from 144 unrelated samples into a combined
19 dataset for downstream analyses using BCFtools merge. Functional consequences of the combined dataset
20 were assessed using the ENSEMBL Variant Effect Prediction (VEP) tool.

21

22 Haplotype networks for *CD8B* were constructed using HPRC/HGSVC continuous haplotypes extracted
23 with BEDtools getfasta and aligned with Muscle using Mega Software [170]. Networks were generated
24 using a minimum spanning tree with the software PopArt [171].

25 **Tests for signatures of natural selection**
26 Ka/Ks ratios (also known as dN/dS) were calculated for pHSD paralogs, performing pairwise comparison
27 between human and chimpanzee sequences, based on T2T-CHM13v1.0 and panTro6 reference genomes,
28 respectively. Alignments between human and chimpanzee canonical transcripts sequences were manually
29 curated and used as input for seqinr package for Ka/Ks estimation. pN/pS ratios were calculated using as
30 input variant sites estimated by seqinr package as well as polymorphic variation from the combined cHiFi
31 and HPRC/HGSCV dataset, considering only biallelic SNPs from unrelated samples (n=144).
32 Synonymous and nonsynonymous mutations were defined based on previously calculated VEP
33 consequences. Ka/Ks and pN/pS values were jointly analyzed using the Direction of Selection (DoS)
34 statistic, a derivation of McDonald–Kreitman's neutrality index, defined as DoS = Dn/(Dn + Ds) - Pn/(Pn
35 + Ps) [83]. Significant differences in Ka/Ks or DoS between ancestral and derived paralogs were assessed
36 using Wilcoxon signed-rank test, pairing each derived paralog to its ancestral counterpart. dN/dS was
37 determined, in parallel, across gene families using codeml as part of the Phylogenetic Analysis by
38 Maximum Likelihood (PAML [84]) from generated multiple-species alignments for each gene family
39 (MAFFT [172]), using T2T-CHM13 for human paralog sequences and orthologous sequences from
40 respective genomes for chimpanzee (panTro6), gorilla (gorGor6), orangutan (ponAbe3), rhesus
41 (rheMac10), mouse (mm39), and rat (rn7). Ancestral and derived states for pHSD genes were assigned
42 based on previously published predicted states [13]. Conservatively, the evolutionary status of four gene
43 families was considered as "unknown" and excluded from calculations of statistical differences

1 (*FRMPD2*/*FRMPD2B*, *PTPN20*/*PTPN20CP*, *GPRIN2*/*GPRIN2B*, and *NPY4R*/*NPY4R2*). Paralogs with
2 infinite values were also excluded from the analysis.
3
4 Nucleotide diversity (π) and Tajima's D statistics were calculated across selected pHSD loci using
5 biallelic SNPs derived from continuous haplotypes from HPRC and HGSVC assemblies, utilizing the
6 PopGenome R package and its functions F_ST.stats and neutrality.stats, respectively. For the gene bodies
7 of *GPR89*, *ROCK1*, *FAM72*, and *CD8B*, π and Tajima's D values were calculated using 15-kbp windows
8 with 1-kbp steps. For GPR89 paralogs, π was calculated across extended surrounding duplicated regions
9 using 20-kbp windows and 1-kbp steps. For *CD8B* paralogs, Tajima's D was calculated in surrounding
10 regions using 6-kbp windows and 500-bp steps.

11 **Generation of zebrafish lines**
12 Wild-type NHGRI-1[173] and Tg[HuC-GFP] [108] adult zebrafish were kept in a temperature (28±0.5°C) and
13 light (10h dark / 14h light) controlled environment following standard protocols [174] with flowing water
14 filtered via UV (Aquaneering, San Diego, CA). As described previously [97,100], feeding included brine
15 shrimp (Artemia Brine Shrimp 90% hatch, Aquaneering, San Diego, CA) and flakes (Select Diet,
16 Aquaneering, San Diego, CA). To obtain embryos for the different assays, males and females were placed
17 in a 1L breeding tank in a 1:1 or 1:2 ratio and eggs from at least five crosses collected and kept in Petri
18 dishes with E3 media (0.03% Instant Ocean salt in deionized water) in an incubator at 28°C until used.
19 All animal use was approved by the Institutional Animal Care and Use Committee from the Office of
20 Animal Welfare Assurance, University of California, Davis.
21
22 Creation of CRISPR lines to knockout genes of interest was done as previously described [97,100,175]. Briefly,
23 crRNAs were annealed with tracrRNA (Alt-R system, Integrated DNA Technologies, Newark, NJ) in a
24 100 µM final concentration to make the sgRNA duplex, which was then coupled with SpCas9 (20 µM,
25 New England BioLabs, Ipswich, MA) to prepare injection mixes. All oligonucleotide sequences can be
26 found in Table S30. Microinjection of one-cell stage zebrafish embryos was performed using an air
27 injector (Pneumatic MPP1-2 Pressure Injector) to release ~1 nl of injection mix into each embryo.
28 Injection mixes to knockout-specific genes included ribonucleoproteins with four different sgRNAs
29 targeting early exons in equimolar concentrations. In parallel, stable CRISPR knockout lines were made
30 using a single sgRNA (Table S30) and adults carrying candidate knockout alleles for each gene of interest
31 were further outcrossed to remove potential off-target edits and then incrossed to generate a batch of wild-
32 type, heterozygous, and homozygous larvae, following standard protocols [176]. Knockout alleles in stable
33 lines corresponded to a 5 bp deletion in *frmpd2* and an 8 bp deletion in *gpr89* (allele sequences can be
34 found in Figure S21). For *arhgap11* knockdown, morpholinos blocking translation (GeneTools,
35 Philomath, OR) were reconstituted to 2 mM and ~1 nl of a 2 ng/nl mix was microinjected into one-cel-
36 stage embryos. Assessments of potential off-target sites for all sgRNAs used in this study were performed
37 with the CIRCLE-seq protocol [177,178] and top potential off-target sites were evaluated via Sanger
38 sequencing as previously described [97]. No editing was observed in potential off-target sites for any
39 sgRNA used in this study, suggesting that phenotypes observed are due to the targeted knockout.
40
41 "Humanized" zebrafish larvae were generated by temporal expression of transcribed mRNAs. Expression
42 vectors containing human transcripts were used to generate mRNA, including pEF-DEST51 (*SRGAP2C*
43 and *ARHGAP11B*), pGCS1 (*GPR89B*, *PDZK1P1*, and *PTPN20CP*), pCR4 (*NPY4R*), and pCMV-

1  SPORT6 (*FAM72B* and *FRMPD2B*). The cDNA inserts of two genes were synthesized (Twist
2  Biosciences, San Francisco, CA) based on transcript evidence from IsoSeq data from the ENCODE [179]
3  project (*PDZK1P1*: ENCFF158KCA, ENCFF939EUU; *PTPN20CP*: ENCFF305AFY). All plasmids were
4  sequenced through either Azenta or Plasmidsaurus and are included as Data S1. Following plasmid
5  linearization using restriction enzyme digest and DNA purification, 5'-capped *in vitro* mRNA was
6  generated using the MEGAshortscript transcription kit (Thermo Fisher, Waltham, MA) following the
7  manufacturer's protocol with a 3.5 h 56°C incubation with T7 or SP6 RNA polymerase, depending on the
8  plasmid. The resulting transcripts were purified with the MEGAclear transcription clean-up kit (Thermo
9  Fisher, Waltham, MA), measured quantity with the Qubit, and visualized on a 2% agarose gel to ensure
10  intact transcript. All mRNA injection mixes included mRNA at a 100 ng/ul concentration and ~1 nl of the
11  mix microinjected into one-cell stage embryos, as described above.

## Morphometric assessments

13  High-throughput imaging of the zebrafish larvae was performed using the VAST BioImager system
14  (Union Biometrica, Holliston, MA) as previously described [96,97]. Mutant and control larvae at 3 or 5 dpf
15  were placed into 96-well plates where they were then acquired by a robotic arm, placing the larvae in a
16  rotating 600 μm capillary coupled with a camera, allowing for the automatic acquisition of images from
17  four sides. Images were then processed and analyzed using the TableCreator tool in FishInspector v1.7 [95]
18  to measure the head area and body length of 3,146 larvae—discarding images with general issues (e.g.,
19  dead or truncated larvae). To validate changes in head area, a neuronal reporter transgenic zebrafish line
20  Tg[HuC-GFP][108] was used to create CRISPR-knockouts or humanized larvae that were then kept in an
21  incubator at 28°C until imaged at 3 dpf using tricaine as anesthesia (0.0125%) and low-melting agarose.
22  Imaging was performed in the Dragonfly spinning disk confocal microscope system with an iXon camera
23  (Andor Technology, Belfast, United Kingdom). Z-stacks of 10 μm slices for each larva were collected
24  and processed using Fiji [180] to generate hyperstacks with maximum intensity projections. Forebrain areas
25  were measured in a blinded manner by a different trained investigator by manually delimiting the
26  forebrain region. Any image with tilted larvae or unclear definition of the different brain regions were not
27  included.

## Supervised classification of mutants and controls

29  As an alternative to performing statistical tests to identify changes in predefined morphological
30  measurements between mutants and controls, we employed a CNN to broadly identify differences
31  between mutants and controls without the need to measure predefined features. Due to the use of multiple
32  96-well plates for each mutant, we observed significant batch effects in the resulting images, where larvae
33  images from the same plate were significantly more similar to each other than to genotypically matched
34  larvae from different plates. Therefore, before training our CNN-based classifier, we trained a latent
35  diffusion model (LDM) to minimize the plate batch effect before input into the CNN. The broad goal of
36  the LDM is to use the larvae with control genotypes present on each plate to learn the plate-specific batch
37  effects, or 'style'. We then select a single plate as a reference and use the LDM to transform all images to
38  the reference plate style, therefore making them comparable. The LDM removes batch effects by
39  computing four transformations of the original image $x$. First, the original image $x$ is transformed into a
40  latent representation $z_0$ through the use of an variational autoencoder function $\varepsilon$ that summarizes the input
41  image $x$ but does not remove any batch effect:

$$z_0 = \varepsilon(x)$$

44  We then pass the encoded image $z_0$ through an LDM forward process, in which Gaussian noise is
45  gradually added to the latent representation over $T$ time steps $z_t$, calculated as:

$$z_t = \sqrt{\bar{a}_t}z_0 + \sqrt{1 - \bar{a}_t}\varepsilon_t \ \varepsilon, \dots, \varepsilon_{t-2}, \varepsilon_{t-1} \sim N(\mathbf{0}, \mathbf{I})$$

Where $\bar{a}_t$ is the cumulative product of the noise scheduler, and $\varepsilon_t$ is the Gaussian noise sampled from a standard normal distribution $N(0,I)$ at time step $t$. This ultimately transforms the initial image embedding $z_0$ into the embedding $z_T$. This embedding $z_T$ represents the larvae as an embedded image, free of association with any batch effect.

In the third step, we apply a reverse process of the LDM, we successively transform the image $z_T$ into a new $z_0$, but 'add back in' the effect of a reference plate batch by introducing a condition variable $c$ which comprises the desired batch and mutant ID. This conditional reverse process can be expressed as:

$$p_\theta(z_t-1|z_t,c) = \mathbf{N}(z_t-1; \mu_\theta(z_t,t,c), \Sigma_\theta(z_t,t,c))$$

Where $\mu_\theta(z_t,t,c)$ represents the predicted mean of our denoised latent image through the weighted autoencoder with weights $\theta$. Likewise, on top of predicting the mean noise we can predict the variance presented by $\Sigma_\theta(z_t,t,c)$. Finally, we pass our model through the decoder of a variational autoencoder to reconstruct the original image $x$ into a new image, $x'$, that represents the original image $x$ but in the new reference plate style, suitable for input into the CNN classifier.

Our LDM is trained by minimizing the mean squared error between the true $\varepsilon$ and predict noise $\varepsilon_\theta(z_t,t)$ expressed as:

$$L_{LDM} := E_{\varepsilon,t \sim N(0,1),t}\left[||\varepsilon - \varepsilon_\theta(z_t,t)||^2\right]$$

This ensures that the model is accurately predicting the noise applied during the forward process, and by conditioning on the batch and mutant ID, we ensure we can reconstruct into our desired batch with respect to the original mutant.

Furthermore, we trained the model using 350 diffusion steps with a linear noise scheduler [181,182]. After training the model, we applied the model to transform all images to one reference plate, which is selected as the one with the highest number of controls. This transformation process minimizes the batch effect by generating images that appear as if they were collected from the same plate.

Having minimized batch effects on the larvae images, we then trained a CNN image classifier to determine the extent to which each mutant genotype differs from matched controls on the basis of the raw morphometric images alone. Higher classification accuracy, as measured by F1 score, indicates a larger effect size of mutant genotype on morphology. Our CNN framework involved fine-tuning a pretrained Alexnet classifier on the transformed larvae images [183]. More specifically, we trained 17 different Alexnet classifiers, one per mutant genotype, to perform binary classification to distinguish one specific mutant genotype from controls. The models were trained and evaluated in an x-fold cross validation framework, with F1 scores averaged over all folds. To generate feature attribution heat maps highlighting the morphological regions used to distinguish each mutant genotype from controls (Figure 4B), we used the GradCAM (Gradient-weighted Class Activation Mapping) approach [165]. We selected a *GPR89B* and *gpr89KO* sample representative of the pattern exhibited across mutants from this family.

1 **sciRNA-seq**

2 We performed cellular assessments using the single-cell combinatorial indexing RNA sequencing

3 (sciRNA-seq) protocol [102]. Zebrafish larvae from CRISPR knockout or mRNA-injected lines were

4 generated as described above and kept in an incubator at 28°C until 3 dpf when they were euthanized in

5 cold tricaine (0.025%) and their heads immediately dissected, pooling 15 heads together per sample.

6 Dissociation of the dissected heads was performed following two washes in 1 ml of cold 1x PBS on ice

7 with a 15 min incubation in dissociation mix (480 µl of 0.25% trypsin-EDTA and 20 µl of collagenase P

8 at 100 mg/ml), gently pipetting each sample every 5 min with a cut-open P1000 tip for complete

9 dissociation. Once all tissue was visibly dissociated, 800 µl of DMEM with 10% FBS was added to each

10 sample and centrifuged for 5 min at 700g at 4°C, resuspended in cold 1x PBS and centrifuged again at

11 700g for 5 min at 4°C. Cells were then resuspended in 800 µl of DMEM with 10% FBS and filtered

12 through a 40 µm cell strainer (Flowmi, Sigma Aldrich, St. Louis, MO) using low-bind DNA tubes

13 (Eppendorf, Hamburg, Germany). Cells were counted using a Countess II (Thermo Fisher, Waltham,

14 MA) and all samples with viability >65% used further. Immediately after viability confirmation, cells

15 were fixed as previously described [184] with a 10 min incubation in 1.33% formaldehyde in 1x PBS on ice

16 followed by permeabilization with 5% Triton-X for 3 min on ice, and neutralization with 10% Tris-HCl

17 (1M, pH 8). Cells were then filtered through a 40 µm cell strainer again, 15 µl of DMSO added to each

18 sample, and then slowly freezed in a Mr Frosty (Thermo Fisher, Waltham, MA) freezing container filled

19 with isopropanol at -80°C overnight.

20

21 Library preparation was performed following the sciRNA-seq protocol as described [102], including three

22 rounds of combinatorial indexing of the cells (all primer sequences correspond to Plate 1 of the original

23 protocol and can be found in www.github.com/JunyueC/sci-RNA-seq3_pipeline). The first round

24 involved reverse transcription with barcoded oligo-dT primers to introduce the initial index. Cells were

25 then pooled and redistributed into new wells for the second round, where a second index was added via

26 ligation. The third round included second-strand synthesis, tagmentation with Tn5 transposase, and PCR

27 amplification to incorporate the final index. Libraries were evaluated for quality control in a BioAnalyzer

28 and Qubit to check integrity and concentration, and then sequenced in three NovaSeq 6000 lanes

29 (Novogene, Sacramento, CA). Raw fastq files were processed following the available sci-RNA-seq3

30 pipeline [185] (www.github.com/JunyueC/sci-RNA-seq3_pipeline). This pipeline includes attachment of the

31 unique molecular identifier (UMI) sequence to each read2 based on the identified RT and ligation

32 barcodes from read1 (edit distance ≤1), and trimming with TrimGalore v0.4.1

33 (https://zenodo.org/records/7598955), using cutadapt [186] and fastqc [187]. Reads were then mapped to the

34 improved zebrafish transcriptome [188] with STAR [189] using the --outSAMstrandField intronMotif option.

35 Duplicates (reads with the same UMI) were removed with the available custom-made python scripts

36 found in the Cao lab GitHub repository. Lastly, filtered SAM files were splitted by their UMI sequences

37 (corresponding to individual cells) and gene-cell count matrices constructed by mapping reads to the

38 zebrafish v4 GTF file [188].

39

40 Gene-cell count matrices were loaded into R to generate Seurat v4 [190] objects and cells with transcript

41 counts below 150 or above two standard deviations over the mean, mitochondrial or ribosomal gene

42 counts >5%, or potential doublets (with a ~4% doublet expectation based on previous reports [185,191] and

43 estimated using DoubletFinder [192]) were removed (Figure S22). Cells from different libraries were

44 normalized using SCTransform [193] with the glmGamPoi method and regressing by the percentage of

mitochondrial and ribosomal counts. Then, normalized counts across sequencing libraries were integrated with Harmony [194] with a PCA reduction using batch as a grouping variable. Hierarchical clustering was performed by calculating the euclidean distances across all cells using the Harmony cell embeddings and clustering with the hclust function using the ward.D2 method. The hierarchical tree was cut at a K of 50, gene markers for each cluster estimated using the FindAllMarkers function (logfc.threshold=0.10, test.use="MAST", min.pct=0.15, min.diff.pct=0.10), and classification into cell types using available zebrafish brain scRNA-seq atlases [71,105] and the Zebrafish Information Network (ZFIN [195]) website. Focusing on neuronal, glial, and eye-related clusters left a total of 95,555 cells for further analysis (Tables S31 and S32). General correlations across samples (knockout vs. "humanized" models for each gene of interest) were done with a balanced number of cells for each pair and pseudo-bulking gene counts by sample and cluster, so counts across cells were summed together for each sample, allowing for biological replicates to be maintained. Then, pseudo-counts were processed with DESeq2 [196] with the Wald test option to obtain fold-change values for each gene compared to their respective control (SpCas9-scrambled gRNA injected for crispants, GFP-mRNA-injected for "humanized", and control-morpholino-injected for *arhgap11*-knockdown). Then, cell-type-specific differential gene expression tests were performed similarly but with previous subsetting of the matrix for each cell type. For *FRMPD2* and *GPR89* models, forebrain cells were further re-clustered to obtain more detailed cell types; gene expression across samples correlated as described above using a pseudo-bulk approach with the telencephalic cells. Progenitor and differentiated cell classification was performed using known neural progenitor (*sox19a*, *sox2*, *rpl5a*, *npm1a*, *s100b*, *dla*) or mature neuron (*elavl3*, *elavl4*, *tubb5*) markers and the PercentageFeatureSet function to estimate the weight of these genes per cell. Enrichment of DEGs in gene ontology terms was estimated with clusterProfiler [165] using only the expressed genes as the background list for the tests.

## Seizure susceptibility

To assess changes to chemically induced seizure susceptibility, we employed an optimized published protocol [115]. Briefly, larvae were collected and kept in an incubator at 28°C until 4 dpf, when they were distributed in a 96-well plate and placed in a Zebrabox system chamber (ViewPoint, Montreal, Canada) that has a camera with an acquisition speed of 30 frames per second. Treatments included 0 or 2.5 mM of pentylenetetrazol (PTZ, #P6500, Sigma-Aldrich, St. Louis, MO) in a total volume of 200 µl per well. Once placed in the Zebrabox chamber, larvae were left for 10 min unbothered before starting a 15 min recording (acquisition in 1 s bins) to then extract the frequency of high-speed events (>28 mm/s) using a published MATLAB script [115] to compare against batch-sibling controls.

## Data Availability

NCBI GenBank numbers of deposited data pending: cHiFi sequencing, scRNA-seq for zebrafish. Code and data: https://github.com/mydennislab/public_data/ (zenodo pending).

## Acknowledgements

# References

11  1.  Carroll, S.B. (2003). Genetics and the making of Homo sapiens. Nature *422*, 849–857.

12  2.  Varki, A., and Altheide, T.K. (2005). Comparing the human and chimpanzee genomes: Searching for
13      needles in a haystack. Preprint, https://doi.org/10.1101/gr.3737405
14      https://doi.org/10.1101/gr.3737405.

15  3.  Pääbo, S. (2014). The Human Condition—A Molecular Approach. Preprint,
16      https://doi.org/10.1016/j.cell.2013.12.036 https://doi.org/10.1016/j.cell.2013.12.036.

17  4.  Pollen, A.A., Kilik, U., Lowe, C.B., and Camp, J.G. (2023). Human-specific genetics: new tools to
18      explore the molecular and cellular basis of human evolution. Nat. Rev. Genet. *24*, 687–711.

19  5.  Sousa, A.M.M., Meyer, K.A., Santpere, G., Gulden, F.O., and Sestan, N. (2017). Evolution of the
20      Human Nervous System Function, Structure, and Development. Cell *170*, 226–247.

21  6.  Enard, W., Gehre, S., Hammerschmidt, K., Holter, S.M., Blass, T., Somel, M., Bruckner, M.K.,
22      Schreiweis, C., Winter, C., Sohr, R., et al. (2009). A humanized version of Foxp2 affects cortico-
23      basal ganglia circuits in mice. Cell *137*, 961–971.

24  7.  Enard, W., Przeworski, M., Fisher, S.E., Lai, C.S., Wiebe, V., Kitano, T., Monaco, A.P., and Paabo,
25      S. (2002). Molecular evolution of FOXP2, a gene involved in speech and language. Nature *418*,
26      869–872.

27  8.  Pollard, K.S., Salama, S.R., Lambert, N., Lambot, M.A., Coppens, S., Pedersen, J.S., Katzman, S.,
28      King, B., Onodera, C., Siepel, A., et al. (2006). An RNA gene expressed during cortical
29      development evolved rapidly in humans. Nature *443*, 167–172.

30  9.  Soto, D.C., Uribe-Salazar, J.M., Shew, C.J., Sekar, A., McGinty, S.P., and Dennis, M.Y. (2023).
31      Genomic structural variation: A complex but important driver of human evolution. Am J Biol
32      Anthropol. https://doi.org/10.1002/ajpa.24713.

33  10. Bailey, J.A., Gu, Z., Clark, R.A., Reinert, K., Samonte, R.V., Schwartz, S., Adams, M.D., Myers,
34      E.W., Li, P.W., and Eichler, E.E. (2002). Recent segmental duplications in the human genome.
35      Science *297*, 1003–1007.

36  11. Dennis, M.Y., and Eichler, E.E. (2016). Human adaptation and evolution by segmental duplication.
37      Curr. Opin. Genet. Dev. *41*, 44–52.

38  12. Porubsky, D., and Eichler, E.E. (2024). A 25-year odyssey of genomic technology advances and

1    structural variant discovery. Cell *187*, 1024–1037.

2    13.  Dennis, M.Y., Harshman, L., Nelson, B.J., Penn, O., Cantsilieris, S., Huddleston, J., Antonacci, F.,
3         Penewit, K., Denman, L., Raja, A., et al. (2017). The evolution and population diversity of human-
4         specific segmental duplications. Nat Ecol Evol *1*, 69.

5    14.  Sudmant, P.H., Kitzman, J.O., Antonacci, F., Alkan, C., Malig, M., Tsalenko, A., Sampas, N.,
6         Bruhn, L., Shendure, J., 1000 Genomes Project, et al. (2010). Diversity of human copy number
7         variation and multicopy genes. Science *330*, 641–646.

8    15.  Sudmant, P.H., Huddleston, J., Catacchio, C.R., Malig, M., Hillier, L.W., Baker, C., Mohajeri, K.,
9         Kondova, I., Bontrop, R.E., Persengiev, S., et al. (2013). Evolution and diversity of copy number
10        variation in the great ape lineage. Genome Res. *23*, 1373–1382.

11   16.  Charrier, C., Joshi, K., Coutinho-Budd, J., Kim, J.-E., Lambert, N., de Marchena, J., Jin, W.-L.,
12        Vanderhaeghen, P., Ghosh, A., Sassa, T., et al. (2012). Inhibition of SRGAP2 function by its human-
13        specific paralogs induces neoteny during spine maturation. Cell *149*, 923–935.

14   17.  Dennis, M.Y., Nuttle, X., Sudmant, P.H., Antonacci, F., Graves, T.A., Nefedov, M., Rosenfeld, J.A.,
15        Sajjadian, S., Malig, M., Kotkiewicz, H., et al. (2012). Evolution of human-specific neural SRGAP2
16        genes by incomplete segmental duplication. Cell *149*, 912–922.

17   18.  Fiddes, I.T., Lodewijk, G.A., Mooring, M., Bosworth, C.M., Ewing, A.D., Mantalas, G.L., Novak,
18        A.M., van den Bout, A., Bishara, A., Rosenkrantz, J.L., et al. (2018). Human-Specific NOTCH2NL
19        Genes Affect Notch Signaling and Cortical Neurogenesis. Cell *173*, 1356–1369.e22.

20   19.  Suzuki, I.K., Gacquer, D., Van Heurck, R., Kumar, D., Wojno, M., Bilheu, A., Herpoel, A.,
21        Lambert, N., Cheron, J., Polleux, F., et al. (2018). Human-Specific NOTCH2NL Genes Expand
22        Cortical Neurogenesis through Delta/Notch Regulation. Cell *173*, 1370–1384.e16.

23   20.  Florio, M., Heide, M., Pinson, A., Brandl, H., Albert, M., Winkler, S., Wimberger, P., Huttner,
24        W.B., and Hiller, M. (2018). Evolution and cell-type specificity of human-specific genes
25        preferentially expressed in progenitors of fetal neocortex. Elife *7*.
26        https://doi.org/10.7554/eLife.32332.

27   21.  Florio, M., Albert, M., Taverna, E., Namba, T., Brandl, H., Lewitus, E., Haffner, C., Sykes, A.,
28        Wong, F.K., Peters, J., et al. (2015). Human-specific gene *ARHGAP11B* promotes basal progenitor
29        amplification and neocortex expansion. Science *347*, 1465–1470.

30   22.  Florio, M., Namba, T., Pääbo, S., Hiller, M., and Huttner, W.B. (2016). A single splice site mutation
31        in human-specific ARHGAP11B causes basal progenitor amplification. Sci Adv *2*, e1601941.

32   23.  Namba, T., Dóczi, J., Pinson, A., Xing, L., Kalebic, N., Wilsch-Bräuninger, M., Long, K.R., Vaid,
33        S., Lauer, J., Bogdanova, A., et al. (2020). Human-Specific ARHGAP11B Acts in Mitochondria to
34        Expand Neocortical Progenitors by Glutaminolysis. Neuron *105*, 867–881.e9.

35   24.  Ju, X.-C., Hou, Q.-Q., Sheng, A.-L., Wu, K.-Y., Zhou, Y., Jin, Y., Wen, T., Yang, Z., Wang, X., and
36        Luo, Z.-G. (2016). The hominoid-specific gene TBC1D3 promotes generation of basal neural
37        progenitors and induces cortical folding in mice. Elife *5*. https://doi.org/10.7554/eLife.18197.

38   25.  Van Heurck, R., Bonnefont, J., Wojno, M., Suzuki, I.K., Velez-Bravo, F.D., Erkol, E., Nguyen,
39        D.T., Herpoel, A., Bilheu, A., Beckers, S., et al. (2023). CROCCP2 acts as a human-specific

1   modifier of cilia dynamics and mTOR signaling to promote expansion of cortical progenitors.
2   Neuron *111*, 65–80.e6.

3 26. Libé-Philippot, B., Lejeune, A., Wierda, K., Louros, N., Erkol, E., Vlaeminck, I., Beckers, S.,
4   Gaspariunaite, V., Bilheu, A., Konstantoulea, K., et al. (2023). LRRC37B is a human modifier of
5   voltage-gated sodium channels and axon excitability in cortical neurons. Cell *186*, 5766–5783.e25.

6 27. Karageorgiou, C., Gokcumen, O., and Dennis, M.Y. (2024). Deciphering the role of structural
7   variation in human evolution: a functional perspective. Curr. Opin. Genet. Dev. *88*, 102240.

8 28. Taylor, D.J., Eizenga, J.M., Li, Q., Das, A., Jenike, K.M., Kenny, E.E., Miga, K.H., Monlong, J.,
9   McCoy, R.C., Paten, B., et al. (2024). Beyond the Human Genome Project: The Age of Complete
10   Human Genome Sequences and Pangenome References. Annu. Rev. Genomics Hum. Genet.
11   https://doi.org/10.1146/annurev-genom-021623-081639.

12 29. Ebbert, M.T.W., Jensen, T.D., Jansen-West, K., Sens, J.P., Reddy, J.S., Ridge, P.G., Kauwe, J.S.K.,
13   Belzil, V., Pregent, L., Carrasquillo, M.M., et al. (2019). Systematic analysis of dark and
14   camouflaged genes reveals disease-relevant genes hiding in plain sight. Genome Biol. *20*, 97.

15 30. Amemiya, H.M., Kundaje, A., and Boyle, A.P. (2019). The ENCODE Blacklist: Identification of
16   Problematic Regions of the Genome. Sci. Rep. *9*, 9354.

17 31. Cabanski, C.R., Wilkerson, M.D., Soloway, M., Parker, J.S., Liu, J., Prins, J.F., Marron, J.S., Perou,
18   C.M., and Hayes, D.N. (2013). BlackOPs: increasing confidence in variant detection through
19   mappability filtering. Nucleic Acids Res. *41*, e178.

20 32. Derrien, T., Estellé, J., Marco Sola, S., Knowles, D.G., Raineri, E., Guigó, R., and Ribeca, P. (2012).
21   Fast computation and applications of genome mappability. PLoS One *7*, e30377.

22 33. Lee, H., and Schatz, M.C. (2012). Genomic dark matter: the reliability of short read mapping
23   illustrated by the genome mappability score. Bioinformatics *28*, 2097–2105.

24 34. Nurk, S., Koren, S., Rhie, A., Rautiainen, M., Bzikadze, A.V., Mikheenko, A., Vollger, M.R.,
25   Altemose, N., Uralsky, L., Gershman, A., et al. (2022). The complete sequence of a human genome.
26   Science *376*, 44–53.

27 35. Vollger, M.R., Guitart, X., Dishuck, P.C., Mercuri, L., Harvey, W.T., Gershman, A., Diekhans, M.,
28   Sulovari, A., Munson, K.M., Lewis, A.P., et al. (2022). Segmental duplications and their variation in
29   a complete human genome. Science *376*, eabj6965.

30 36. Aganezov, S., Yan, S.M., Soto, D.C., Kirsche, M., Zarate, S., Avdeyev, P., Taylor, D.J., Shafin, K.,
31   Shumate, A., Xiao, C., et al. (2022). A complete reference genome improves analysis of human
32   genetic variation. Science *376*, eabl3533.

33 37. Numanagic, I., Gökkaya, A.S., Zhang, L., Berger, B., Alkan, C., and Hach, F. (2018). Fast
34   characterization of segmental duplications in genome assemblies. Bioinformatics *34*, i706–i714.

35 38. Dougherty, M.L., Underwood, J.G., Nelson, B.J., Tseng, E., Munson, K.M., Penn, O., Nowakowski,
36   T.J., Pollen, A.A., and Eichler, E.E. (2018). Transcriptional fates of human-specific segmental
37   duplications in brain. Genome Res. *28*, 1566–1576.

38 39. Larson, J.L., Silver, A.J., Chan, D., Borroto, C., Spurrier, B., and Silver, L.M. (2015). Validation of
39   a high resolution NGS method for detecting spinal muscular atrophy carriers among phase 3

1  participants in the 1000 Genomes Project. BMC Med. Genet. *16*, 100.

2  40.  Moreno-Igoa, M., Hernández-Charro, B., Bengoa-Alonso, A., Pérez-Juana-del-Casal, A., Romero-
3        Ibarra, C., Nieva-Echebarria, B., and Ramos-Arroyo, M.A. (2015). KANSL1 gene disruption
4        associated with the full clinical spectrum of 17q21.31 microdeletion syndrome. BMC Med. Genet.
5        *16*, 68.

6  41.  Perry, G.H., Dominy, N.J., Claw, K.G., Lee, A.S., Fiegler, H., Redon, R., Werner, J., Villanea, F.A.,
7        Mountain, J.L., Misra, R., et al. (2007). Diet and the evolution of human amylase gene copy number
8        variation. Nat. Genet. *39*, 1256–1260.

9  42.  Bolognini, D., Halgren, A., Lou, R.N., Raveane, A., Rocha, J.L., Guarracino, A., Soranzo, N., Chin,
10       C.-S., Garrison, E., and Sudmant, P.H. (2024). Recurrent evolution and selection shape structural
11       diversity at the amylase locus. Nature, 1–9.

12 43.  Yilmaz, F., Karageorgiou, C., Kim, K., Pajic, P., Scheer, K., Human Genome Structural Variation
13       Consortium, Beck, C.R., Torregrossa, A.-M., Lee, C., and Gokcumen, O. (2024). Paleolithic Gene
14       Duplications Primed Adaptive Evolution of Human Amylase Locus Upon Agriculture. bioRxiv.
15       https://doi.org/10.1101/2023.11.27.568916.

16 44.  Mallick, S., Li, H., Lipson, M., Mathieson, I., Gymrek, M., Racimo, F., Zhao, M., Chennagiri, N.,
17       Nordenfelt, S., Tandon, A., et al. (2016). The Simons Genome Diversity Project: 300 genomes from
18       142 diverse populations. Nature *538*, 201–206.

19 45.  Bosch, N., Cáceres, M., Cardone, M.F., Carreras, A., Ballana, E., Rocchi, M., Armengol, L., and
20       Estivill, X. (2007). Characterization and evolution of the novel gene family FAM90A in primates
21       originated by multiple duplication and rearrangement events. Hum. Mol. Genet. *16*, 2572–2582.

22 46.  Shen, F., and Kidd, J.M. (2020). Rapid, Paralog-Sensitive CNV Analysis of 2457 Human Genomes
23       Using QuicK-mer2. Genes *11*, 141.

24 47.  Prüfer, K., Racimo, F., Patterson, N., Jay, F., Sankararaman, S., Sawyer, S., Heinze, A., Renaud, G.,
25       Sudmant, P.H., de Filippo, C., et al. (2014). The complete genome sequence of a Neanderthal from
26       the Altai Mountains. Nature *505*, 43–49.

27 48.  Prüfer, K., de Filippo, C., Grote, S., Mafessoni, F., Korlević, P., Hajdinjak, M., Vernot, B., Skov, L.,
28       Hsieh, P., Peyrégne, S., et al. (2017). A high-coverage Neandertal genome from Vindija Cave in
29       Croatia. Science *358*, 655–658.

30 49.  Meyer, M., Kircher, M., Gansauge, M.-T., Li, H., Racimo, F., Mallick, S., Schraiber, J.G., Jay, F.,
31       Prüfer, K., de Filippo, C., et al. (2012). A high-coverage genome sequence from an archaic
32       Denisovan individual. Science *338*, 222–226.

33 50.  Strahl, B.D., and Allis, C.D. (2000). The language of covalent histone modifications. Nature *403*,
34       41–45.

35 51.  Nimmerjahn, F., and Ravetch, J.V. (2008). Fcγ receptors as regulators of immune responses. Nat.
36       Rev. Immunol. *8*, 34–47.

37 52.  Karczewski, K.J., Francioli, L.C., Tiao, G., Cummings, B.B., Alföldi, J., Wang, Q., Collins, R.L.,
38       Laricchia, K.M., Ganna, A., Birnbaum, D.P., et al. (2020). The mutational constraint spectrum
39       quantified from variation in 141,456 humans. Nature *581*, 434–443.

53. Schneider, V.A., Graves-Lindsay, T., Howe, K., Bouk, N., Chen, H.-C., Kitts, P.A., Murphy, T.D., Pruitt, K.D., Thibaud-Nissen, F., Albracht, D., et al. (2017). Evaluation of GRCh38 and de novo haploid genome assemblies demonstrates the enduring quality of the reference assembly. Genome Res. *27*, 849–864.

54. Tajima, F. (1989). Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. Genetics *123*, 585–595.

55. Byrska-Bishop, M., Evani, U.S., Zhao, X., Basile, A.O., Abel, H.J., Regier, A.A., Corvelo, A., Clarke, W.E., Musunuri, R., Nagulapalli, K., et al. (2022). High-coverage whole-genome sequencing of the expanded 1000 Genomes Project cohort including 602 trios. Cell *185*, 3426–3440.e19.

56. Aguirre, M., Rivas, M.A., and Priest, J. (2019). Phenome-wide Burden of Copy-Number Variation in the UK Biobank. Am. J. Hum. Genet. *105*, 373–383.

57. Trost, B., Thiruvahindrapuram, B., Chan, A.J.S., Engchuan, W., Higginbotham, E.J., Howe, J.L., Loureiro, L.O., Reuter, M.S., Roshandel, D., Whitney, J., et al. (2022). Genomic architecture of autism from comprehensive whole-genome sequence annotation. Cell *185*, 4409–4427.e18.

58. Brunetti-Pierri, N., Berg, J.S., Scaglia, F., Belmont, J., Bacino, C.A., Sahoo, T., Lalani, S.R., Graham, B., Lee, B., Shinawi, M., et al. (2008). Recurrent reciprocal 1q21.1 deletions and duplications associated with microcephaly or macrocephaly and developmental and behavioral abnormalities. Nat. Genet. *40*, 1466–1471.

59. Winchester, L., Yau, C., and Ragoussis, J. (2009). Comparing CNV detection methods for SNP arrays. Brief. Funct. Genomics *8*, 353–366.

60. Satterstrom, F.K., Kosmicki, J.A., Wang, J., Breen, M.S., De Rubeis, S., An, J.-Y., Peng, M., Collins, R., Grove, J., Klei, L., et al. (2020). Large-Scale Exome Sequencing Study Implicates Both Developmental and Functional Changes in the Neurobiology of Autism. Cell *180*, 568–584.e23.

61. Fietz, S.A., Lachmann, R., Brandl, H., Kircher, M., Samusik, N., Schröder, R., Lakshmanaperumal, N., Henry, I., Vogt, J., Riehn, A., et al. (2012). Transcriptomes of germinal zones of human and mouse fetal neocortex suggest a role of extracellular matrix in progenitor self-renewal. Proc. Natl. Acad. Sci. U. S. A. *109*, 11836–11841.

62. Pezzini, F., Bettinetti, L., Di Leva, F., Bianchi, M., Zoratti, E., Carrozzo, R., Santorelli, F.M., Delledonne, M., Lalowski, M., and Simonati, A. (2017). Transcriptomic Profiling Discloses Molecular and Cellular Events Related to Neuronal Differentiation in SH-SY5Y Neuroblastoma Cells. Cell. Mol. Neurobiol. *37*, 665–682.

63. van de Leemput, J., Boles, N.C., Kiehl, T.R., Corneo, B., Lederman, P., Menon, V., Lee, C., Martinez, R.A., Levi, B.P., Thompson, C.L., et al. (2014). CORTECON: a temporal transcriptome analysis of in vitro human cerebral cortex development from human embryonic stem cells. Neuron *83*, 51–68.

64. Miller, J.A., Ding, S.-L., Sunkin, S.M., Smith, K.A., Ng, L., Szafer, A., Ebbert, A., Riley, Z.L., Royall, J.J., Aiona, K., et al. (2014). Transcriptional landscape of the prenatal human brain. Nature *508*, 199–206.

65. Langfelder, P., and Horvath, S. (2008). WGCNA: an R package for weighted correlation network analysis. BMC Bioinformatics *9*, 559.

66. Shew, C.J., Carmona-Mora, P., Soto, D.C., Mastoras, M., Roberts, E., Rosas, J., Jagannathan, D., Kaya, G., O'Geen, H., and Dennis, M.Y. (2021). Diverse Molecular Mechanisms Contribute to Differential Expression of Human Duplicated Genes. Mol. Biol. Evol. *38*, 3060–3077.

67. Guerrier, S., Coutinho-Budd, J., Sassa, T., Gresset, A., Jordan, N.V., Chen, K., Jin, W.-L., Frost, A., and Polleux, F. (2009). The F-BAR domain of srGAP2 induces membrane protrusions required for neuronal migration and morphogenesis. Cell *138*, 990–1004.

68. Parikshak, N.N., Luo, R., Zhang, A., Won, H., Lowe, J.K., Chandran, V., Horvath, S., and Geschwind, D.H. (2013). Integrative functional genomic analyses implicate specific molecular pathways and circuits in autism. Cell *155*, 1008–1021.

69. Koolen, D.A., Kramer, J.M., Neveling, K., Nillesen, W.M., Moore-Barton, H.L., Elmslie, F.V., Toutain, A., Amiel, J., Malan, V., Tsai, A.C., et al. (2012). Mutations in the chromatin modifier gene KANSL1 cause the 17q21.31 microdeletion syndrome. Nat. Genet. *44*, 639–641.

70. La Manno, G., Siletti, K., Furlan, A., Gyllborg, D., Vinsland, E., Mossi Albiach, A., Mattsson Langseth, C., Khven, I., Lederer, A.R., Dratva, L.M., et al. (2021). Molecular architecture of the developing mouse brain. Nature *596*, 92–96.

71. Raj, B., Farrell, J.A., Liu, J., El Kholtei, J., Carte, A.N., Navajas, A.J., Du, L.Y., McKenna, A., Relić, Đ., Leslie, J.M., et al. (2020). Emergence of Neuronal Diversity during Vertebrate Brain Development. Neuron *108*. https://doi.org/10.1016/j.neuron.2020.09.023.

72. Willsey, H.R., Exner, C.R.T., Xu, Y., Everitt, A., Sun, N., Wang, B., Dea, J., Schmunk, G., Zaltsman, Y., Teerikorpi, N., et al. (2021). Parallel in vivo analysis of large-effect autism genes implicates cortical neurogenesis and estrogen in risk and resilience. Neuron *109*, 1409.

73. Weinschutz Mendes, H., Neelakantan, U., Liu, Y., Fitzpatrick, S.E., Chen, T., Wu, W., Pruitt, A., Jin, D.S., Jamadagni, P., Carlson, M., et al. (2023). High-throughput functional analysis of autism genes in zebrafish identifies convergence in dopaminergic and neuroimmune pathways. Cell Rep. *42*, 112243.

74. Guitart, X., Porubsky, D., Yoo, D., Dougherty, M.L., Dishuck, P., Munson, K.M., Lewis, A.P., Hoekzema, K., Knuth, J., Chang, S., et al. (2024). Independent expansion, selection and hypervariability of the gene family in humans. Genome Res. https://doi.org/10.1101/gr.279299.124.

75. Jiang, Z., Tang, H., Ventura, M., Cardone, M.F., Marques-Bonet, T., She, X., Pevzner, P.A., and Eichler, E.E. (2007). Ancestral reconstruction of segmental duplications reveals punctuated cores of human genome evolution. Nat. Genet. *39*, 1361–1368.

76. Jarvis, E.D., Formenti, G., Rhie, A., Guarracino, A., Yang, C., Wood, J., Tracey, A., Thibaud-Nissen, F., Vollger, M.R., Porubsky, D., et al. (2022). Automated assembly of high-quality diploid human reference genomes. bioRxiv, 2022.03.06.483034. https://doi.org/10.1101/2022.03.06.483034.

77. Liao, W.-W., Asri, M., Ebler, J., Doerr, D., Haukness, M., Hickey, G., Lu, S., Lucas, J.K., Monlong, J., Abel, H.J., et al. (2022). A Draft Human Pangenome Reference. bioRxiv, 2022.07.09.499321. https://doi.org/10.1101/2022.07.09.499321.

78. Wang, T., Antonacci-Fulton, L., Howe, K., Lawson, H.A., Lucas, J.K., Phillippy, A.M., Popejoy, A.B., Asri, M., Carson, C., Chaisson, M.J.P., et al. (2022). The Human Pangenome Project: a global resource to map genomic diversity. Nature *604*, 437–446.

79. Ebler, J., Ebert, P., Clarke, W.E., Rausch, T., Audano, P.A., Houwaart, T., Mao, Y., Korbel, J.O., Eichler, E.E., Zody, M.C., et al. (2022). Pangenome-based genome inference allows efficient and accurate genotyping across a wide spectrum of variant classes. Nat. Genet. *54*, 518–525.

80. Cavalli-Sforza, L.L. (2005). The Human Genome Diversity Project: past, present and future. Nat. Rev. Genet. *6*, 333–340.

81. McLaren, W., Gil, L., Hunt, S.E., Riat, H.S., Ritchie, G.R.S., Thormann, A., Flicek, P., and Cunningham, F. (2016). The Ensembl Variant Effect Predictor. Genome Biol. *17*, 122.

82. Ellegren, H. (2005). Evolution: natural selection in the evolution of humans and chimps. Curr. Biol. *15*, R919–R922.

83. Stoletzki, N., and Eyre-Walker, A. (2011). Estimation of the neutrality index. Mol. Biol. Evol. *28*, 63–70.

84. Yang, Z. (2007). PAML 4: phylogenetic analysis by maximum likelihood. Mol. Biol. Evol. *24*, 1586–1591.

85. Connolly, J.M., Hansen, T.H., Ingold, A.L., and Potter, T.A. (1990). Recognition by CD8 on cytotoxic T lymphocytes is ablated by several substitutions in the class I alpha 3 domain: CD8 and the T-cell receptor recognize the same class I molecule. Proc. Natl. Acad. Sci. U. S. A. *87*, 2137–2141.

86. Salter, R.D., Benjamin, R.J., Wesley, P.K., Buxton, S.E., Garrett, T.P., Clayberger, C., Krensky, A.M., Norment, A.M., Littman, D.R., and Parham, P. (1990). A binding site for the T-cell co-receptor CD8 on the alpha 3 domain of HLA-A2. Nature *345*, 41–46.

87. Garrido-Martín, D., Borsari, B., Calvo, M., Reverter, F., and Guigó, R. (2021). Identification and analysis of splicing quantitative trait loci across multiple tissues in the human genome. Nat. Commun. *12*, 727.

88. Orrù, V., Steri, M., Sidore, C., Marongiu, M., Serra, V., Olla, S., Sole, G., Lai, S., Dei, M., Mulas, A., et al. (2020). Complex genetic signatures in immune cells underlie autoimmunity and inform therapy. Nat. Genet. *52*, 1036–1045.

89. Holtzman, N.G., Iovine, M.K., Liang, J.O., and Morris, J. (2016). Learning to Fish with Genetics: A Primer on the Vertebrate Model Danio rerio. Genetics *203*, 1069–1089.

90. Meyers, J.R. (2018). Zebrafish: Development of a Vertebrate Model Organism: Zebrafish : Development of a Vertebrate Model Organism. Current Protocols Essential Laboratory Techniques *16*, e19.

91. Sakai, C., Ijaz, S., and Hoffman, E.J. (2018). Zebrafish Models of Neurodevelopmental Disorders: Past, Present, and Future. Front. Mol. Neurosci. *11*, 294.

92. Schmidt, E.R.E., Kupferman, J.V., Stackmann, M., and Polleux, F. (2019). The human-specific paralogs SRGAP2B and SRGAP2C differentially modulate SRGAP2A-dependent synaptic development. Sci. Rep. *9*, 18692.

93. Schmidt, E.R.E., Zhao, H.T., Park, J.M., Dipoppa, M., Monsalve-Mercado, M.M., Dahan, J.B., Rodgers, C.C., Lejeune, A., Hillman, E.M.C., Miller, K.D., et al. (2021). A human-specific modifier of cortical connectivity and circuit function. Nature *599*, 640.

94. Heide, M., Haffner, C., Murayama, A., Kurotaki, Y., Shinohara, H., Okano, H., Sasaki, E., and Huttner, W.B. (2020). Human-specific ARHGAP11B increases size and folding of primate neocortex in the fetal marmoset. Science *369*. https://doi.org/10.1126/science.abb2401.

95. Kalebic, N., Gilardi, C., Albert, M., Namba, T., Long, K.R., Kostic, M., Langen, B., and Huttner, W.B. (2018). Human-specific ARHGAP11B induces hallmarks of neocortical expansion in developing ferret neocortex. Elife *7*. https://doi.org/10.7554/eLife.41241.

96. Meng, X., Lin, Q., Zeng, X., Jiang, J., Li, M., Luo, X., Chen, K., Wu, H., Hu, Y., Liu, C., et al. (2023). Brain developmental and cortical connectivity changes in transgenic monkeys carrying the human-specific duplicated gene SRGAP2C. Natl Sci Rev *10*, nwad281.

97. Uribe-Salazar, J.M., Kaya, G., Sekar, A., Weyenberg, K., Ingamells, C., and Dennis, M.Y. (2022). Evaluation of CRISPR gene-editing tools in zebrafish. BMC Genomics *23*, 12.

98. Teixidó, E., Kießling, T.R., Krupp, E., Quevedo, C., Muriana, A., and Scholz, S. (2019). Automated Morphological Feature Assessment for Zebrafish Embryo Developmental Toxicity Screens. Toxicol. Sci. *167*. https://doi.org/10.1093/toxsci/kfy250.

99. Pulak, R. (2016). Tools for automating the imaging of zebrafish larvae. Methods *96*. https://doi.org/10.1016/j.ymeth.2015.11.021.

100. Colón-Rodríguez, A., Uribe-Salazar, J.M., Weyenberg, K.B., Sriram, A., Quezada, A., Kaya, G., Jao, E., Radke, B., Lein, P.J., and Dennis, M.Y. (2020). Assessment of Autism Zebrafish Mutant Models Using a High-Throughput Larval Phenotyping Platform. Front. Cell Dev. Biol. *8*, 586296.

101. Saunders, L.M., Srivatsan, S.R., Duran, M., Dorrity, M.W., Ewing, B., Linbo, T.H., Shendure, J., Raible, D.W., Moens, C.B., Kimelman, D., et al. (2023). Embryo-scale reverse genetics at single-cell resolution. Nature *623*, 782–791.

102. Martin, B.K., Qiu, C., Nichols, E., Phung, M., Green-Gladden, R., Srivatsan, S., Blecher-Gonen, R., Beliveau, B.J., Trapnell, C., Cao, J., et al. (2023). Optimized single-nucleus transcriptional profiling by combinatorial indexing. Nat. Protoc. *18*, 188–207.

103. Uribe-Salazar, J.M., Kaya, G., Weyenberg, K.B., Radke, B., Hino, K.K., Soto, D.C., Shiu, J.-L., Zhang, W., Ingamells, C., Haghani, N.K., et al. (2024). Zebrafish models of human-duplicated gene SRGAP2 reveal novel functions in microglia and visual system development. bioRxiv. https://doi.org/10.1101/2024.09.11.612570.

104. d'Amora, M., and Giordani, S. (2018). The Utility of Zebrafish as a Model for Screening Developmental Neurotoxicity. Front. Neurosci. *12*. https://doi.org/10.3389/fnins.2018.00976.

105. Zhang, H., Wang, H., Shen, X., Jia, X., Yu, S., Qiu, X., Wang, Y., Du, J., Yan, J., and He, J. (2021). The landscape of regulatory genes in brain-wide neuronal phenotypes of a vertebrate brain. Elife *10*. https://doi.org/10.7554/eLife.68224.

106. Kozol, R.A., Abrams, A.J., James, D.M., Buglo, E., Yan, Q., and Dallman, J.E. (2016). Function Over Form: Modeling Groups of Inherited Neurological Conditions in Zebrafish. Front. Mol. Neurosci. *9*, 55.

107. Kimmel, C.B., Ballard, W.W., Kimmel, S.R., Ullmann, B., and Schilling, T.F. (1995). Stages of embryonic development of the zebrafish. Dev. Dyn. *203*, 253–310.

1	108. Park, H.C., Kim, C.H., Bae, Y.K., Yeo, S.Y., Kim, S.H., Hong, S.K., Shin, J., Yoo, K.W., Hibi, M.,
2	Hirano, T., et al. (2000). Analysis of upstream elements in the HuC promoter leads to the
3	establishment of transgenic zebrafish with fluorescent neurons. Dev. Biol. *227*, 279–293.

4	109. Porter, B.A., and Mueller, T. (2020). The Zebrafish Amygdaloid Complex - Functional Ground Plan,
5	Molecular Delineation, and Everted Topology. Front. Neurosci. *14*, 608.

6	110. Anneser, L., Satou, C., Hotz, H.-R., and Friedrich, R.W. (2024). Molecular organization of neuronal
7	cell types and neuromodulatory systems in the zebrafish telencephalon. Curr. Biol. *34*, 298–312.e4.

8	111. Deckstein, J., van Appeldorn, J., Tsangarides, M., Yiannakou, K., Müller, R., Stumpf, M.,
9	Sukumaran, S.K., Eichinger, L., Noegel, A.A., and Riyahi, T.Y. (2015). The Dictyostelium
10	discoideum GPHR ortholog is an endoplasmic reticulum and Golgi protein with roles during
11	development. Eukaryot. Cell *14*, 41–54.

12	112. Charroux, B., and Royet, J. (2014). Mutations in the Drosophila ortholog of the vertebrate Golgi pH
13	regulator (GPHR) protein disturb endoplasmic reticulum and Golgi organization and affect systemic
14	growth. Biol. Open *3*, 72–80.

15	113. Otani, T., Marchetto, M.C., Gage, F.H., Simons, B.D., and Livesey, F.J. (2016). 2D and 3D Stem
16	Cell Models of Primate Cortical Development Identify Species-Specific Differences in Progenitor
17	Behavior Contributing to Brain Size. Cell Stem Cell *18*, 467–480.

18	114. Benito-Kwiecinski, S., Giandomenico, S.L., Sutcliffe, M., Riis, E.S., Freire-Pritchett, P., Kelava, I.,
19	Wunderlich, S., Martin, U., Wray, G.A., McDole, K., et al. (2021). An early cell shape transition
20	drives evolutionary expansion of the human forebrain. Cell *184*, 2084–2102.e19.

21	115. Griffin, A., Carpenter, C., Liu, J., Paterno, R., Grone, B., Hamling, K., Moog, M., Dinday, M.T.,
22	Figueroa, F., Anvar, M., et al. (2021). Phenotypic analysis of catastrophic childhood epilepsy genes.
23	Commun Biol *4*, 680.

24	116. Lu, X., Zhang, Q., and Wang, T. (2019). The second PDZ domain of scaffold protein Frmpd2 binds
25	to GluN2A of NMDA receptors. Biochem. Biophys. Res. Commun. *516*, 63–67.

26	117. Stenzel, N., Fetzer, C.P., Heumann, R., and Erdmann, K.S. (2009). PDZ-domain-directed basolateral
27	targeting of the peripheral membrane protein FRMPD2 in epithelial cells. J. Cell Sci. *122*, 3374–
28	3384.

29	118. Ueno, A., Omori, Y., Sugita, Y., Watanabe, S., Chaya, T., Kozuka, T., Kon, T., Yoshida, S.,
30	Matsushita, K., Kuwahara, R., et al. (2018). Lrit1, a Retinal Transmembrane Protein, Regulates
31	Selective Synapse Formation in Cone Photoreceptor Cells and Visual Acuity. Cell Rep. *22*, 3548–
32	3561.

33	119. Lee, H.-J., and Zheng, J.J. (2010). PDZ domains and their binding partners: structure, specificity,
34	and modification. Cell Commun. Signal. *8*, 8.

35	120. Stankiewicz, P., Kulkarni, S., Dharmadhikari, A.V., Sampath, S., Bhatt, S.S., Shaikh, T.H., Xia, Z.,
36	Pursley, A.N., Cooper, M.L., Shinawi, M., et al. (2012). Recurrent deletions and reciprocal
37	duplications of 10q11.21q11.23 including CHAT and SLC18A3 are likely mediated by complex
38	low-copy repeats. Hum. Mutat. *33*, 165–179.

39	121. Muntané, G., Horvath, J.E., Hof, P.R., Ely, J.J., Hopkins, W.D., Raghanti, M.A., Lewandowski,

A.H., Wray, G.A., and Sherwood, C.C. (2015). Analysis of synaptic gene expression in the neocortex of primates reveals evolutionary changes in glutamatergic neurotransmission. Cereb. Cortex *25*, 1596–1607.

122. Willcox, B.J., Donlon, T.A., He, Q., Chen, R., Grove, J.S., Yano, K., Masaki, K.H., Willcox, D.C., Rodriguez, B., and Curb, J.D. (2008). FOXO3A genotype is strongly associated with human longevity. Proc. Natl. Acad. Sci. U. S. A. *105*, 13987–13992.

123. Antonacci, F., Dennis, M.Y., Huddleston, J., Sudmant, P.H., Steinberg, K.M., Rosenfeld, J.A., Miroballo, M., Graves, T.A., Vives, L., Malig, M., et al. (2014). Palindromic GOLGA8 core duplicons promote chromosome 15q13.3 microdeletion and evolutionary instability. Nat. Genet. *46*, 1293–1302.

124. Yoo, D., Rhie, A., Hebbar, P., Antonacci, F., Logsdon, G.A., Solar, S.J., Antipov, D., Pickett, B.D., Safonova, Y., Montinaro, F., et al. (2024). Complete sequencing of ape genomes. bioRxiv. https://doi.org/10.1101/2024.07.31.605654.

125. Makova, K.D., Pickett, B.D., Harris, R.S., Hartley, G.A., Cechova, M., Pal, K., Nurk, S., Yoo, D., Li, Q., Hebbar, P., et al. (2024). The complete sequence and comparative analysis of ape sex chromosomes. Nature *630*, 401–411.

126. L Rocha, J., Lou, R.N., and Sudmant, P.H. (2024). Structural variation in humans and our primate kin in the era of telomere-to-telomere genomes and pangenomics. Curr. Opin. Genet. Dev. *87*, 102233.

127. Tropepe, V., and Sive, H.L. (2003). Can zebrafish be used as a model to study the neurodevelopmental causes of autism? Genes Brain Behav. *2*, 268–281.

128. Golzio, C., Willer, J., Talkowski, M.E., Oh, E.C., Taniguchi, Y., Jacquemont, S., Reymond, A., Sun, M., Sawa, A., Gusella, J.F., et al. (2012). KCTD13 is a major driver of mirrored neuroanatomical phenotypes of the 16p11.2 copy number variant. Nature *485*, 363–367.

129. Hoffman, E.J., Turner, K.J., Fernandez, J.M., Cifuentes, D., Ghosh, M., Ijaz, S., Jain, R.A., Kubo, F., Bill, B.R., Baier, H., et al. (2016). Estrogens Suppress a Behavioral Phenotype in Zebrafish Mutants of the Autism Risk Gene, CNTNAP2. Neuron *89*, 725–733.

130. Shah, A.N., Davey, C.F., Whitebirch, A.C., Miller, A.C., and Moens, C.B. (2015). Rapid reverse genetic screening using CRISPR in zebrafish. Nat. Methods *12*, 535–540.

131. Thyme, S.B., Pieper, L.M., Li, E.H., Pandey, S., Wang, Y., Morris, N.S., Sha, C., Choi, J.W., Herrera, K.J., Soucy, E.R., et al. (2019). Phenotypic Landscape of Schizophrenia-Associated Genes Defines Candidates and Their Shared Functions. Cell *177*, 478–491.e20.

132. Glasauer, S.M.K., and Neuhauss, S.C.F. (2014). Whole-genome duplication in teleost fishes and its evolutionary consequences. Mol. Genet. Genomics *289*, 1045–1060.

133. Fossati, M., Pizzarelli, R., Schmidt, E.R., Kupferman, J.V., Stroebel, D., Polleux, F., and Charrier, C. (2016). SRGAP2 and Its Human-Specific Paralog Co-Regulate the Development of Excitatory and Inhibitory Synapses. Neuron *91*, 356–369.

134. Schmidt, E.R.E., Kupferman, J.V., and Stackmann, M. (2019). The human-specific paralogs SRGAP2B and SRGAP2C differentially modulate SRGAP2A-dependent synaptic development.

Scientific Reports *9*. https://doi.org/10.1101/596940.

135. Kalebic, N., and Huttner, W.B. (2020). Basal Progenitor Morphology and Neocortex Evolution. Trends Neurosci. *43*, 843–853.

136. Bitarello, B.D., de Filippo, C., Teixeira, J.C., Schmidt, J.M., Kleinert, P., Meyer, D., and Andrés, A.M. (2018). Signatures of Long-Term Balancing Selection in Human Genomes. Genome Biol. Evol. *10*, 939–955.

137. Andrés, A.M., Hubisz, M.J., Indap, A., Torgerson, D.G., Degenhardt, J.D., Boyko, A.R., Gutenkunst, R.N., White, T.J., Green, E.D., Bustamante, C.D., et al. (2009). Targets of balancing selection in the human genome. Mol. Biol. Evol. *26*, 2755–2764.

138. Bitarello, B.D., Brandt, D.Y.C., Meyer, D., and Andrés, A.M. (2023). Inferring Balancing Selection From Genome-Scale Data. Genome Biol. Evol. *15*. https://doi.org/10.1093/gbe/evad032.

139. Mahmoud, M., Huang, Y., Garimella, K., Audano, P.A., Wan, W., Prasad, N., Handsaker, R.E., Hall, S., Pionzio, A., Schatz, M.C., et al. (2024). Utility of long-read sequencing for All of Us. Nat. Commun. *15*, 837.

140. Searles Quick, V.B., Wang, B., and State, M.W. (2021). Leveraging large genomic datasets to illuminate the pathobiology of autism spectrum disorders. Neuropsychopharmacology *46*, 55–69.

141. Vollger, M.R., DeWitt, W.S., Dishuck, P.C., Harvey, W.T., Guitart, X., Goldberg, M.E., Rozanski, A.N., Lucas, J., Asri, M., The Human Pangenome Reference Consortium, et al. (2022). Increased mutation rate and interlocus gene conversion within human segmental duplications. bioRxiv, 2022.07.06.498021. https://doi.org/10.1101/2022.07.06.498021.

142. Dumont, B.L. (2015). Interlocus gene conversion explains at least 2.7% of single nucleotide variants in human segmental duplications. BMC Genomics *16*, 456.

143. Dumont, B.L., and Eichler, E.E. (2013). Signals of historical interlocus gene conversion in human segmental duplications. PLoS One *8*, e75949.

144. Hardwick, R.J., Machado, L.R., Zuccherato, L.W., Antolinos, S., Xue, Y., Shawa, N., Gilman, R.H., Cabrera, L., Berg, D.E., Tyler-Smith, C., et al. (2011). A worldwide analysis of beta-defensin copy number variation suggests recent selection of a high-expressing DEFB103 gene copy in East Asia. Hum. Mutat. *32*, 743–750.

145. Hughes, T., Hansson, L., Akkouh, I., Hajdarevic, R., Bringsli, J.S., Torsvik, A., Inderhaug, E., Steen, V.M., and Djurovic, S. (2020). Runaway multi-allelic copy number variation at the α-defensin locus in African and Asian populations. Sci. Rep. *10*, 9101.

146. Linzmeier, R.M., and Ganz, T. (2005). Human defensin gene copy number polymorphisms: comprehensive analysis of independent variation in alpha- and beta-defensin regions at 8p22-p23. Genomics *86*, 423–430.

147. Mohajeri, K., Cantsilieris, S., Huddleston, J., Nelson, B.J., Coe, B.P., Campbell, C.D., Baker, C., Harshman, L., Munson, K.M., Kronenberg, Z.N., et al. (2016). Interchromosomal core duplicons drive both evolutionary instability and disease susceptibility of the Chromosome 8p23.1 region. Genome Res. *26*, 1453–1467.

148. Sherman, R.M., Forman, J., Antonescu, V., Puiu, D., Daya, M., Rafaels, N., Boorgula, M.P.,

Chavan, S., Vergara, C., Ortega, V.E., et al. (2019). Assembly of a pan-genome from deep sequencing of 910 humans of African descent. Nat. Genet. *51*, 30–35.

149. Plender, E.G., Prodanov, T., Hsieh, P., Nizamis, E., Harvey, W.T., Sulovari, A., Munson, K.M., Kaufman, E.J., O'Neal, W.K., Valdmanis, P.N., et al. (2024). Structural and genetic diversity in the secreted mucins MUC5AC and MUC5B. Am. J. Hum. Genet. *111*, 1700–1716.

150. Redon, R., Ishikawa, S., Fitch, K.R., Feuk, L., Perry, G.H., Andrews, T.D., Fiegler, H., Shapero, M.H., Carson, A.R., Chen, W., et al. (2006). Global variation in copy number in the human genome. Nature *444*, 444–454.

151. Jeong, H., Dishuck, P.C., Yoo, D., Harvey, W.T., Munson, K.M., Lewis, A.P., Kordosky, J., Garcia, G.H., Human Genome Structural Variation Consortium (HGSVC), Yilmaz, F., et al. (2024). Structural polymorphism and diversity of human segmental duplications. bioRxiv. https://doi.org/10.1101/2024.06.04.597452.

152. Kim, S., Kim, H., Park, D., Kim, J., Hong, J., Kim, J.S., Jung, H., Kim, D., Cheong, E., Ko, J., et al. (2024). Loss of IQSEC3 Disrupts GABAergic Synapse Maintenance and Decreases Somatostatin Expression in the Hippocampus. Cell Rep. *43*, 114254.

153. Pardo-Palacios, F.J., Wang, D., Reese, F., Diekhans, M., Carbonell-Sala, S., Williams, B., Loveland, J.E., De María, M., Adams, M.S., Balderrama-Gutierrez, G., et al. (2023). Systematic assessment of long-read RNA-seq methods for transcript identification and quantification. bioRxiv. https://doi.org/10.1101/2023.07.25.550582.

154. Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Žídek, A., Potapenko, A., et al. (2021). Highly accurate protein structure prediction with AlphaFold. Nature *596*, 583–589.

155. Timp, W., and Timp, G. (2020). Beyond mass spectrometry, the next step in proteomics. Science Advances. https://doi.org/10.1126/sciadv.aax8978.

156. Vollger, M.R., Guitart, X., Dishuck, P.C., Mercuri, L., Harvey, W.T., Gershman, A., Diekhans, M., Sulovari, A., Munson, K.M., Lewis, A.M., et al. (2021). Segmental duplications and their variation in a complete human genome. bioRxiv, 2021.05.26.445678. https://doi.org/10.1101/2021.05.26.445678.

157. Quinlan, A.R., and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for comparing genomic features. Bioinformatics *26*, 841–842.

158. Lek, M., Karczewski, K.J., Minikel, E.V., Samocha, K.E., Banks, E., Fennell, T., O'Donnell-Luria, A.H., Ware, J.S., Hill, A.J., Cummings, B.B., et al. (2016). Analysis of protein-coding genetic variation in 60,706 humans. Nature *536*, 285–291.

159. Danecek, P., Auton, A., Abecasis, G., Albers, C.A., Banks, E., DePristo, M.A., Handsaker, R.E., Lunter, G., Marth, G.T., Sherry, S.T., et al. (2011). The variant call format and VCFtools. Bioinformatics *27*, 2156–2158.

160. Hartasánchez, D.A., Brasó-Vives, M., Heredia-Genestar, J.M., Pybus, M., and Navarro, A. (2018). Effect of Collapsed Duplications on Diversity Estimates: What to Expect. Genome Biol. Evol. *10*, 2899–2905.

161. Buniello, A., MacArthur, J.A.L., Cerezo, M., Harris, L.W., Hayhurst, J., Malangone, C., McMahon, A., Morales, J., Mountjoy, E., Sollis, E., et al. (2019). The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. Nucleic Acids Res. *47*, D1005–D1012.

162. Hinrichs, A.S., Karolchik, D., Baertsch, R., Barber, G.P., Bejerano, G., Clawson, H., Diekhans, M., Furey, T.S., Harte, R.A., Hsu, F., et al. (2006). The UCSC Genome Browser Database: update 2006. Nucleic Acids Res. *34*, D590–D598.

163. Patro, R., Duggal, G., Love, M.I., Irizarry, R.A., and Kingsford, C. (2017). Salmon provides fast and bias-aware quantification of transcript expression. Nat. Methods *14*, 417–419.

164. Soneson, C., Love, M.I., and Robinson, M.D. (2015). Differential analyses for RNA-seq: transcript-level estimates improve gene-level inferences. F1000Res. *4*, 1521.

165. Wu, T., Hu, E., Xu, S., Chen, M., Guo, P., Dai, Z., Feng, T., Zhou, L., Tang, W., Zhan, L., et al. (2021). clusterProfiler 4.0: A universal enrichment tool for interpreting omics data. Innovation (Camb) *2*, 100141.

166. Zook, J.M., Catoe, D., McDaniel, J., Vang, L., Spies, N., Sidow, A., Weng, Z., Liu, Y., Mason, C.E., Alexander, N., et al. (2016). Extensive sequencing of seven human genomes to characterize benchmark reference materials. Scientific Data *3*, 160025.

167. Bergström, A., McCarthy, S.A., Hui, R., Almarri, M.A., Ayub, Q., Danecek, P., Chen, Y., Felkel, S., Hallast, P., Kamm, J., et al. (2020). Insights into human genetic variation and population history from 929 diverse genomes. Science *367*. https://doi.org/10.1126/science.aay5012.

168. Witek, K., Jupe, F., Witek, A.I., Baker, D., Clark, M.D., and Jones, J.D.G. (2016). Accelerated cloning of a potato late blight-resistance gene using RenSeq and SMRT sequencing. Nat. Biotechnol. *34*, 656–660.

169. Poplin, R., Ruano-Rubio, V., DePristo, M.A., Fennell, T.J., Carneiro, M.O., Van der Auwera, G.A., Kling, D.E., Gauthier, L.D., Levy-Moonshine, A., Roazen, D., et al. (2018). Scaling accurate genetic variant discovery to tens of thousands of samples. bioRxiv, 201178. https://doi.org/10.1101/201178.

170. Stecher, G., Tamura, K., and Kumar, S. (2020). Molecular Evolutionary Genetics Analysis (MEGA) for macOS. Mol. Biol. Evol. *37*, 1237–1239.

171. Leigh, J.W., and Bryant, D. (2015). Popart: Full-feature software for haplotype network construction. Methods Ecol. Evol. *6*, 1110–1116.

172. Katoh, K., and Standley, D.M. (2013). MAFFT multiple sequence alignment software version 7: improvements in performance and usability. Mol. Biol. Evol. *30*, 772–780.

173. LaFave, M.C., Varshney, G.K., Vemulapalli, M., Mullikin, J.C., and Burgess, S.M. (2014). A defined zebrafish line for high-throughput genetics and genomics: NHGRI-1. Genetics *198*, 167–170.

174. Westerfield, M. (1995). The Zebrafish Book: A Guide for the Laboratory Use of Zebrafish (Danio Rerio).

175. Jao, L.-E., Wente, S.R., and Chen, W. (2013). Efficient multiplex biallelic zebrafish genome editing using a CRISPR nuclease system. Proc. Natl. Acad. Sci. U. S. A. *110*, 13904–13909.

176. Varshney, G.K., Carrington, B., Pei, W., Bishop, K., Chen, Z., Fan, C., Xu, L., Jones, M., LaFave, M.C., Ledin, J., et al. (2016). A high-throughput functional genomics workflow based on CRISPR/Cas9-mediated targeted mutagenesis in zebrafish. Nat. Protoc. *11*, 2357–2375.

177. Tsai, S.Q., Nguyen, N.T., Malagon-Lopez, J., Topkar, V.V., Aryee, M.J., and Joung, J.K. (2017). CIRCLE-seq: a highly sensitive in vitro screen for genome-wide CRISPR-Cas9 nuclease off-targets. Nat. Methods *14*, 607–614.

178. Lazzarotto, C.R., Nguyen, N.T., Tang, X., Malagon-Lopez, J., Guo, J.A., Aryee, M.J., Joung, J.K., and Tsai, S.Q. (2018). Defining CRISPR-Cas9 genome-wide nuclease activities with CIRCLE-seq. Nat. Protoc. *13*, 2615–2642.

179. ENCODE Project Consortium (2012). An integrated encyclopedia of DNA elements in the human genome. Nature *489*, 57–74.

180. Schindelin, J., Arganda-Carreras, I., Frise, E., Kaynig, V., Longair, M., Pietzsch, T., Preibisch, S., Rueden, C., Saalfeld, S., Schmid, B., et al. (2012). Fiji: an open-source platform for biological-image analysis. Nat. Methods *9*, 676–682.

181. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. (2021). High-Resolution Image Synthesis with Latent Diffusion Models.

182. Krizhevsky, A., Sutskever, I., and Hinton, G.E. (2017). ImageNet classification with deep convolutional neural networks. Commun. ACM *60*, 84–90.

183. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. (2017). Grad-CAM: Visual explanations from deep networks via gradient-based localization. In 2017 IEEE International Conference on Computer Vision (ICCV) (IEEE). https://doi.org/10.1109/iccv.2017.74.

184. Rosenberg, A.B., Roco, C.M., Muscat, R.A., Kuchina, A., Sample, P., Yao, Z., Graybuck, L.T., Peeler, D.J., Mukherjee, S., Chen, W., et al. (2018). Single-cell profiling of the developing mouse brain and spinal cord with split-pool barcoding. Science *360*, 176–182.

185. Cao, J., Spielmann, M., Qiu, X., Huang, X., Ibrahim, D.M., Hill, A.J., Zhang, F., Mundlos, S., Christiansen, L., Steemers, F.J., et al. (2019). The single-cell transcriptional landscape of mammalian organogenesis. Nature *566*, 496–502.

186. Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing reads. EMBnet.journal *17*, 10–12.

187. Andrews, S., and Others (2010). FastQC: a quality control tool for high throughput sequence data. Preprint at Babraham Bioinformatics, Babraham Institute, Cambridge, United Kingdom.

188. Lawson, N.D., Li, R., Shin, M., Grosse, A., Yukselen, O., Stone, O.A., Kucukural, A., and Zhu, L. (2020). An improved zebrafish transcriptome annotation for sensitive and comprehensive detection of cell type-specific genes. Elife *9*. https://doi.org/10.7554/eLife.55792.

189. Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M., and Gingeras, T.R. (2013). STAR: ultrafast universal RNA-seq aligner. Bioinformatics *29*, 15–21.

190. Chen, G., Liu, Z., and Peng, C. (2021). Multimodal and Integrative Analysis of Single-Cell or Bulk Sequencing Data (Frontiers Media SA).

191. Tran, V., Papalexi, E., Schroeder, S., Kim, G., Sapre, A., Pangallo, J., Sova, A., Matulich, P., Kenyon, L., Sayar, Z., et al. (2022). High sensitivity single cell RNA sequencing with split pool barcoding. bioRxiv, 2022.08.27.505512. https://doi.org/10.1101/2022.08.27.505512.

192. McGinnis, C.S., Murrow, L.M., and Gartner, Z.J. (2019). DoubletFinder: Doublet Detection in Single-Cell RNA Sequencing Data Using Artificial Nearest Neighbors. Cell Syst *8*, 329–337.e4.

193. Hafemeister, C., and Satija, R. (2019). Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. Genome Biol. *20*, 296.

194. Korsunsky, I., Millard, N., Fan, J., Slowikowski, K., Zhang, F., Wei, K., Baglaenko, Y., Brenner, M., Loh, P.-R., and Raychaudhuri, S. (2019). Fast, sensitive and accurate integration of single-cell data with Harmony. Nat. Methods *16*, 1289–1296.

195. Bradford, Y.M., Van Slyke, C.E., Ruzicka, L., Singer, A., Eagle, A., Fashena, D., Howe, D.G., Frazer, K., Martin, R., Paddock, H., et al. (2022). Zebrafish information network, the knowledgebase for Danio rerio research. Genetics *220*. https://doi.org/10.1093/genetics/iyac016.

196. Love, M.I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. Genome Biol. *15*, 550.