

Research article

Evolution is in the details: Regulatory differences in modern human and Neanderthal

Harlan R. Barker^{a,b,*} , Seppo Parkkila^{a,b} , Martti E.E. Tolvanen^c

^a Faculty of Medicine and Health Technology, Tampere University, Tampere 33520, Finland

^b Fimlab Ltd., Tampere University Hospital, Tampere 33520, Finland

^c Department of Computing, University of Turku, Turku 20014, Finland



ARTICLE INFO

Keywords:

Transcription Factor Binding Sites
TFBS
Gene regulation
Genomic Variation
Comparative Genomics
Brain Development
Human Evolution

ABSTRACT

Transcription factor (TF) proteins play a critical role in the regulation of eukaryote gene expression via sequence-specific binding to genomic locations known as TF binding sites. We studied sites of genomic variation between modern human and Neanderthal promoters. We detected significant differences in the binding affinities of 110 TFs to the promoters of 75 target genes. The TFs were enriched for terms related to vision, motor neurons, homeobox, and brain, whereas the target genes and their direct interactors were enriched in terms related to autism, brain, connective tissue, trachea, prostate, skull morphology, and vision. Secondary analysis of single-cell data revealed that a subset of the identified TFs (CUX1, CUX2, ESRG, FOXP1, FOXP2, MEF2C, POU6F2, PRRX1 and RORA) co-occur as marker genes in L4 glutamatergic neurons. The majority of these genes have known roles in autism and/or schizophrenia and are associated with human accelerated regions (elevated divergence in humans vs. other primates). Analysis of a single-nucleus dataset of cortical tissue showed that 15 of these TFs and 16 of their target genes are differentially expressed in autism vs. control, most commonly upregulated in developing neurons. Down regulation of these genes occurred in SV2C- and somatostatin-expressing interneurons, oligodendrocytes, and oligodendrocyte precursor cells. These results support the value of gene regulation studies for the evolution of human cognitive abilities and the neuropsychiatric disorders that accompany it.

1. Introduction

1.1. North of Eden

Humans and their hominin relatives have been leaving Africa in waves for the past two million plus years. Various environmental and cultural pressures have impacted each diaspora in different ways, subsequently producing adaptations reflected in physiology, immunity, and brain size. Recently, the discovery and sequencing of DNA from the remains of Neanderthal [1–4] and Denisovan [5,6] have allowed direct comparisons of DNA from modern and ancient hominids. Among the observed genomic variations between modern humans and Neanderthals, a limited number have been identified in gene coding regions. Some of these genes are known to affect cognition and morphology (cranium, rib, dentition, and shoulder joint) [1], pigmentation and behavioral traits [2], and brain development [7]. However, as has been noted before, there is a paucity of coding variations to explain the

differences between related species; the genome of the Altai Neanderthal reveals just 96 fixed amino acid substitutions, occurring in 87 proteins [7]. Unsurprisingly, a much larger set of variants are observed in intergenic regions, owing not only to the fact that these are comparatively much larger regions but also to the expectation of less conservation in what until recently was often termed "junk DNA". While variants in coding regions can directly affect protein structure, those found in intergenic regions may affect the regulation of gene expression through alternative binding of transcription factors in promoters and enhancers and the expression of noncoding RNAs. What may surprise some is the cumulative effect of numerous small — and large — changes to gene expression arising owing to these manifold intergenic changes, which may ultimately serve as the engine of speciation. Indeed, introgressed Neanderthal DNA is more depleted in regulatory regions than in those encoding proteins [8]. Using a new computational tool that incorporates numerous transcription-relevant genomic features, we sought to reveal how the comparative differences in these regions may

* Corresponding author at: Faculty of Medicine and Health Technology, Tampere University, Tampere 33520, Finland.

E-mail address: harlan.barker@tuni.fi (H.R. Barker).

<https://doi.org/10.1016/j.csbj.2025.05.052>

Received 13 February 2025; Received in revised form 29 May 2025; Accepted 29 May 2025

Available online 30 May 2025

2001-0370/© 2025 The Author(s). Published by Elsevier B.V. on behalf of Research Network of Computational and Structural Biotechnology. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

affect regulatory differences between modern humans and Neanderthals. Specifically, our aim was to identify those gene-regulating transcription factors whose binding to DNA may vary between these species of hominids and thus drive the differences between them.

2. Methods

2.1. Analysis of modern human vs. Neanderthal genetic variation

The Neanderthal Genome Project has cataloged a total of 388,388 SNPs in modern human and Neanderthal genomes [1]. These were further reduced to 21,990 SNPs in the proximal promoters (-2500 to +2500 nucleotides relative to the TSS) of those transcripts, which are defined by Ensembl as "protein-coding" (Fig. 1C). All modern human vs. Neanderthal SNPs were mapped to the promoters of all human genes and the following were calculated: rate of SNPs per nucleotide were calculated for varying gene biotypes (Fig. 1A); the distribution of SNPs in protein-coding, RNA, T-cell receptor, and immunoglobulin genes

promoters (Fig. 1B); and total SNPs occurring on/near protein-coding and RNA genes (Fig. 1C). Using TFBSFootprinter [9–12], a tool we have developed for the prediction of TFBSs, a 50 bp region centered on each SNP was analyzed for the binding of 575 TFs for both the modern human version and the Neanderthal variant. TFBSFootprinter automatically retrieves the human sequences at a target region, and custom Python scripts were used to modify these sequences for the Neanderthal variant. All TFBS predictions that overlapped with the target SNP position were retained. The complete result set was then reduced on the basis of the combined affinity score p value via a Benjamini–Hochberg-derived critical p value corresponding to a false discovery rate (FDR) cutoff of 0.01 to address multiple testing. For each putative TFBS meeting the cutoff in either subspecies, the corresponding matched pair of PWM scores was retained. To identify significantly different scoring for TFs between subspecies, for each TF, via the compiled matched scores, the Wilcoxon rank test was performed via the SciPy stats Python library [13], and subsequent results were filtered via a Benjamini–Hochberg-derived critical p value corresponding to an FDR

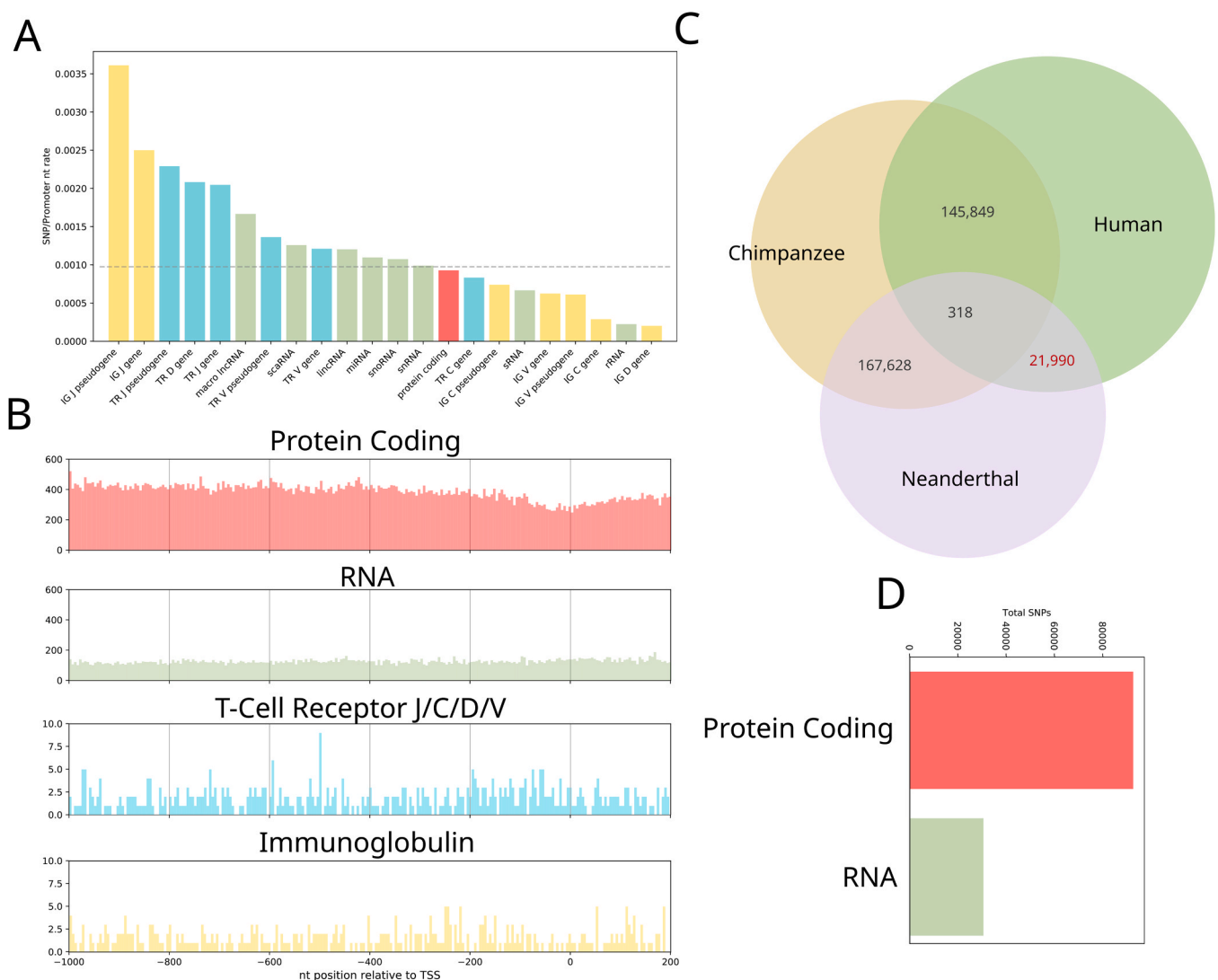


Fig. 1. Modern human and Neanderthal variants used in the analysis. (A) Incidence rate of SNPs/nucleotide in human promoters of the various transcript classes, as defined by Ensembl. The bars are colored by transcript class and correspond with those in panel B. (B) Number of observed modern human/Neanderthal SNPs at each nucleotide location in the sets of promoters of various transcript types. A visible depletion of SNPs is observed for the promoters of protein-coding transcripts corresponding with increasing proximity to the TSS (nucleotide position 0). (C) Counts of total SNPs in promoters of protein-coding transcripts between the three species cataloged in the Neanderthal Genome Project. (D) Comparison of counts of human vs. chimp SNPs in promoters of protein-coding vs RNA transcripts. IG J, immunoglobulin joining chain; TR J, T-cell receptor joining chain; TR D, T-cell receptor diversity chain; TR V, T-cell receptor variable chain; IG C, immunoglobulin constant chain; IG V, immunoglobulin variable chain; IG D, immunoglobulin diversity chain.

cutoff of 0.10 to address multiple testing.

A total of 108 TF models were identified as scoring differently across human vs. Neanderthal SNP locations. For each of these genes, we extracted RNA-Seq data from the FANTOM dataset (across all CAGE peaks associated with that TF) and retained the data for the 100 tissues with the highest aggregate expression across all of the target TF genes. In the case of heterodimer JASPAR TF models (e.g., FOS::JUN, NR1H3::RXRA, and POU5F1::SOX2), the expression of each TF component gene was used. The expression data were extracted as TPM values and normalized by \log_2 transformation. From the subsequent normalized expression data values, hierarchical clustering was performed and visualized via SciPy, Matplotlib, and Seaborn Python libraries [13–15].

The results of hierarchical clustering revealed a cluster of TF genes that were uniquely expressed in neural and immune tissues. These gene sets were used to perform PANTHER-based gene ontology enrichment analysis (using the www.geneontology.org web server) [16,17], with the default statistical settings using Fisher's exact test method and an FDR threshold of $p < 0.05$ (default setting).

2.1.1. Enrichment analysis of DB TFs and their target genes and interactors

The molecular signatures database (MSigDB) contains thousands of gene sets across several large collections, allowing enrichment analysis of, e.g., biological processes, pathways, and ontologies [18,19]. The complete set of MSigDB gene sets (version 2024) were downloaded. The set of 110 DB TFs was analyzed for enrichment for all MSigDB annotation labels via Fisher's exact test, with the full set of 529 human JASPAR TF genes as background and Benjamini–Hochberg (BH) correction for multiple testing (FDR cutoff value of 0.10). The DB TF 75 target genes were likewise analyzed for enrichment, with all genes as background.

The OmniPath DB has likewise compiled protein interactions (human, mouse and rat) from many publicly available datasets. The complete set of interactions was downloaded and filtered to determine posttranslational human protein–protein interactions. A total of 805 interactors were identified for the set of 75 DB TF target genes, and the combined list (879 genes) was used to perform enrichment analysis in MSigDB gene sets, with all genes as background and BH correction (FDR cutoff value of 0.10).

2.1.2. Analysis of the brain expression of differentially binding TFs

Expression data in the form of reads per kilobase per million reads (RPKM) were extracted for 26 unique tissues across 31 timepoints (from 8 weeks postconception to 40 years after birth) from the Allen Brain Atlas (brainspan.org). The RPKM values were then converted to TPM values and \log_2 transformed. Ages were grouped into three phases of growth for simplicity of analysis and interpretation: pcw (8–37 weeks post conception), early (4 months–4 years), and late (8 years–40 years). Corresponding to these age groups, log-transformed TPM values for each tissue were grouped and averaged and used to perform cluster map analysis to identify grouping tissues at time points with similar expression profiles.

2.1.3. Analysis of brain samples via scRNA-Seq

The Allen Brain Atlas has performed single-nucleus RNA-Seq analysis of 49,417 nuclei derived from 8 brain cortex regions within the middle temporal gyrus (MTG), anterior cingulate gyrus (CgG), primary visual cortex (VIC), primary motor cortex (M1C), primary somatosensory cortex (S1C) and primary auditory cortex (A1C) [20]. These data were downloaded as a count matrix along with a table of the associated metadata (<https://portal.brain-map.org/atlas-and-data/rnaseq/human-multiple-cortical-areas-smart-seq>) and loaded into the Python library SCANPY [21]. Using a modified workflow described previously [22], samples were filtered by Gaussian fit of the read count ($300,000 < x < 3,500,000$), expressed gene count ($2000 < x$), and number of cells in which a gene was expressed (>50), resulting in a final count of 46,959 cells and 42,185 genes for further analysis.

Using SCANPY, counts were normalized by cell ('pp.normalize_total';

target_sum=1,000,000), log transformed ('pp.log1p'), highly variable genes identified ('pp.highly_variable_genes'; flavor='seurat_v3'; n_top_genes=5000; layer='counts'), principal component analysis ('pp.pca'; n_comps = 15; svd_solver='arpack'), and k-nearest neighbors ('pp.neighbors'; n_neighbors=15). The expression relationships between cells were graphically visualized with the Python implementation [23] of the ForceAtlas2 [24] graph layout algorithm, as called in SCANPY ('tl.draw_graph' and 'pl.draw_graph').

Annotation data regarding brain region, cortical layer, and GABAergic/glutamatergic/nonneuronal cell type features were extracted from Allen Brain Atlas sample data for mapping onto derived cell type clusters. The top 100 marker genes for each cell type+cortical layer (CT/CL) cluster (e.g., excitatory L4) were identified as those with higher expression unique to each cluster by the Welch t test in SCANPY ('tl.rank_genes_groups'). The expression of the DB TF genes which were identified as marker genes in a CT/CL cluster was mapped onto cluster figures.

Gene ontology analysis of target cluster marker genes was performed via the Protein Analysis Through Evolutionary Relationships (PANTHER) tool at the Geneontology.org web server [17]. Ontological terms associated with overabundance among cluster marker gene lists were established via Fisher's exact test, and the results were filtered by an FDR < 0.05 (default setting); analyses were performed for biological process, molecular function, and cellular component terms. Disease term gene ontology analysis was performed via Enrichr [25,26], which is based on the ontology compiled by DisGeNET [27].

To explore possible connections between DBTFs and DBTF target genes in autism a dataset comprised of 104,559 single nucleus RNA-Seq (snRNA-Seq) data from autism and control samples of cortical tissue [28] was downloaded and analyzed. Using SCANPY, counts were normalized by cell ('pp.normalize_total'; target_sum=1000,000), log transformed ('pp.log1p'), highly variable genes identified ('pp.highly_variable_genes'; flavor='seurat_v3'; n_top_genes=2000; layer='counts'), principal component analysis ('pp.pca'; n_comps = 15; svd_solver='arpack'), and k-nearest neighbors ('pp.neighbors'; n_neighbors=15). Using metadata categories of 'diagnosis', 'cluster', and 'region' differential gene expression was analyzed between control and ASD samples using Wilcoxon rank sum test in SCANPY ('tl.rank_genes_groups'). Results were subsequently filtered by adjusted p-value (≤ 0.05), fold-change (≤ 0.5 ; ≥ 1.5), and presence in either DBTFs or DBTF target genes lists.

3. Results

3.1. High-scoring TFBSs differ between modern humans and Neanderthals

In the analysis of 21,990 SNPs identified in the comparison of the modern human and Neanderthal promoteromes—the collection of all human/Neanderthal proximal promoters of protein-coding genes—a total of 108 TF models, representing 110 unique differentially binding (DB) TF proteins—showed a significant difference in scoring (Wilcoxon rank statistical test) between the two hominin species (Table 1). A significant majority of these are homeobox genes (78/110). Based on Human Protein Atlas annotations, 76 of the DB TFs are associated with sexual reproduction, 72 and have brain-specific expression. More extensive comparisons are available as [Supplementary Table 1](#).

3.2. DB TFs enriched for vision, brain, and developmental terms

Enrichment analysis of the DB TF genes compared with the JASPAR TF genes was performed via the MSigDB annotation database [18,19] via Fisher's exact test, results are presented as odds ratio (OR) with FDR, with Benjamini–Hochberg correction for multiple testing and an FDR of 0.1. The DB TFs were over-enriched for 'neuromuscular process controlling balance' (OR, 28; FDR, 7.33×10^{-2}), 'Sensory perception of

Table 1

TFs with differential binding affinities in modern humans vs. Neanderthal proximal and semi-distal promoter regions. Dev., development-associated gene; HAR, human accelerated region.

Gene	Description	HAR	Dev.
ALX3	ALX Homeobox 3		x
ALX4	ALX Homeobox 4	x	x
ARID5A	AT-Rich Interaction Domain 5A		
BARHL1	BarH Like Homeobox 1		x
BARHL2	BarH Like Homeobox 2	x	x
BARX1	BARX Homeobox 1		x
BSX	Brain Specific Homeobox		x
CRX	Cone-Rod Homeobox		x
CUX1	Cut Like Homeobox 1	x	x
CUX2	Cut Like Homeobox 2		x
DLX1	Distal-Less Homeobox 1		x
DLX2	Distal-Less Homeobox 2	x	x
DLX4	Distal-Less Homeobox 4		x
DMBX1	Diencephalon/Mesencephalon Homeobox 1		x
DMRT3	Doublesex And Mab-3 Related Transcription Factor 3	x	
EMX1	Empty Spiracles Homeobox 1		x
EN1	Engrailed Homeobox 1	x	x
ESRRG	Estrogen Related Receptor Gamma	x	
ESX1	ESX Homeobox 1		x
EVS1	Even-Skipped Homeobox 1		x
EVS2	Even-Skipped Homeobox 2	x	x
FOS	Fos Proto-Oncogene, AP-1 Transcription Factor Subunit		
FOXA1	Forkhead Box A1	x	x
FOXD3	Forkhead Box D3		x
FOXP1	Forkhead Box P1	x	x
FOXP2	Forkhead Box P2	x	x
GBX2	Gastrulation Brain Homeobox 2	x	x
GFI1	Growth Factor Independent 1 Transcriptional Repressor		
GRHL2	Grainyhead Like Transcription Factor 2		
GSC	Gooseoid Homeobox		x
GSC2	Gooseoid Homeobox 2		x
GSX2	GS Homeobox 2		x
HAND1	Heart And Neural Crest Derivatives Expressed 1		
HMX1	H6 Family Homeobox 1		x
HMX2	H6 Family Homeobox 2		x
HMX3	H6 Family Homeobox 3		x
HNF1A	HNF1 Homeobox A		x
HNF1B	HNF1 Homeobox B		x
HOXA5	Homeobox A5	x	x
HOXB2	Homeobox B2		x
HOXB3	Homeobox B3		x
HOXB5	Homeobox B5		x
HOXD8	Homeobox D8		x
IRF3	Interferon Regulatory Factor 3		
ISL2	ISL LIM Homeobox 2		x
ISX	Intestine Specific Homeobox	x	x
JUN	Jun Proto-Oncogene, AP-1 Transcription Factor Subunit		
LBX2	Ladybird Homeobox 2		x
LEF1	Lymphoid Enhancer Binding Factor 1		
LHX3	LIM Homeobox 3		x
LHX4	LIM Homeobox 4		x
LHX9	LIM Homeobox 9		x
LMX1A	LIM Homeobox Transcription Factor 1 Alpha		x
LMX1B	LIM Homeobox Transcription Factor 1 Beta		x
MAFB	MAF BZIP Transcription Factor B	x	
MAFG	MAF BZIP Transcription Factor G		
MEF2C	Myocyte Enhancer Factor 2 C	x	
MEOX1	Mesenchyme Homeobox 1		x
MEOX2	Mesenchyme Homeobox 2	x	x
MIXL1	Mix Paired-Like Homeobox		x
MNX1	Motor Neuron And Pancreas Homeobox 1		x
NEUROG1	Neurogenin 1		
NFYA	Nuclear Transcription Factor Y Subunit Alpha		
NKX2-5	NK2 Homeobox 5		x
NKX2-8	NK2 Homeobox 8		x
NKX6-2	NK6 Homeobox 2		x

Table 1 (continued)

Gene	Description	HAR	Dev.
NOBOX	NOBOX Oogenesis Homeobox		x
NR1H3	Nuclear Receptor Subfamily 1 Group H Member 3		
NR2E1	Nuclear Receptor Subfamily 2 Group E Member 1		
NRL	Neural Retina Leucine Zipper		
ONECUT1	One Cut Homeobox 1	x	x
ONECUT3	One Cut Homeobox 3		x
OTX1	Orthodenticle Homeobox 1		x
OTX2	Orthodenticle Homeobox 2		x
PAX3	Paired Box 3	x	x
PAX7	Paired Box 7		x
PHOX2A	Paired Like Homeobox 2A		x
PHOX2B	Paired Like Homeobox 2B	x	x
PITX3	Paired Like Homeodomain 3		x
POU1F1	POU Class 1 Homeobox 1	x	x
POU2F1	POU Class 2 Homeobox 1		x
POU2F3	POU Class 2 Homeobox 3		x
POU3F1	POU Class 3 Homeobox 1		x
POU3F2	POU Class 3 Homeobox 2	x	x
POU3F3	POU Class 3 Homeobox 3	x	x
POU4F1	POU Class 4 Homeobox 1		x
POU4F2	POU Class 4 Homeobox 2		x
POU4F3	POU Class 4 Homeobox 3	x	x
POU6F1	POU Class 6 Homeobox 1		x
POU6F2	POU Class 6 Homeobox 2	x	x
PROP1	PROP Paired-Like Homeobox 1		x
PRRX1	Paired Related Homeobox 1	x	x
PRRX2	Paired Related Homeobox 2		x
RARA	Retinoic Acid Receptor Alpha		
RAX	Retina And Anterior Neural Fold Homeobox		x
RAX2	Retina And Anterior Neural Fold Homeobox 2		x
RORA	RAR Related Orphan Receptor A	x	
RXRA	Retinoid X Receptor Alpha		
SHOX	Short Stature Homeobox		x
SOX10	SRY-Box Transcription Factor 10		
SOX15	SRY-Box Transcription Factor 15		
TBP	TATA-Box Binding		
TCF3	Transcription Factor 3		
UNCX	UNC Homeobox	x	x
VAX1	Ventral Anterior Homeobox 1		x
VAX2	Ventral Anterior Homeobox 2		x
VSX1	Visual System Homeobox 1		x
VSX2	Visual System Homeobox 2		x

light stimulus' (OR, 8.28; FDR, 7.87×10^{-2}), 'Sensory perception' (OR, 4.93; FDR, 6.71×10^{-2}), 'Nervous system process' (OR, 4.38; FDR, 6.71×10^{-2}), 'Cell projection morphogenesis' (OR, 4.30; FDR, 3.78×10^{-3}), 'Cell morphogenesis involved in neuron differentiation' (OR, 4.06; FDR, 2.45×10^{-2}), 'Head development' (OR, 3.13; FDR, 3.78×10^{-3}), 'Central nervous system development' (OR, 2.98; FDR, 3.78×10^{-3}), with the full results available in [Supplementary Table 2](#).

3.3. DB TF target genes

A total of 75 genes were identified, hereafter referred to as 'DB TF target genes', whose promoters contain MH vs. Neanderthal SNPs where differential binding occurs by our determined set of statistically significant DB TFs (Table 2). In all of the DB TF target genes, at least one MH vs. Neanderthal SNP position has been identified as an eQTL affecting the expression of the target gene in whose promoter it occurs, as defined by data from the GTEx database (data not shown).

The target genes with the greatest number of DB events (143) were FARS2 and LYRM4, whose bi-directional promoter contains 2 SNPs, for which 105 (97.22 %) of the DB TF binding models showed differential binding. The remaining top ten DB TF target genes with the next highest number of DB occurrences were: TEX19, NARF, SLC6A11, LSM5, FAM172A, WNT7A, EXO1, and KRT8. Among these genes, all but LSM5 (mitochondria-specific) have been associated with brain (FARS2, SLC6A11, FAM172A, WNT7A) or reproductive (FARS2, LYRM4, TEX19, NARF, WNT7A, EXO1, and KRT8) tissue- or cell-specific expression in the Protein Atlas database ([Supplementary Table 3](#)). Among all 75 DB TF

Table 2

DB TF target genes with differential binding affinity by DB TFs in modern human vs. Neanderthal proximal and semi-distal promoter regions.

Gene	Gene description	Single-cell RNA cell specific	Count
FARS2	Phenylalanyl-tRNA synthetase 2, mitochondrial	astrocytes, early spermatids, excitatory neurons, inhibitory neurons, late spermatids, microglia, oligodendrocyte precursor cells, oligodendrocytes	143
LYRM4	LYR Motif-Containing Protein 4		143
TEX19	Testis expressed 19	spermatocytes, spermatogonia	60
NARF	Nuclear prelamin A recognition factor	erythroid cells, late spermatids	56
SLC6A11	Solute carrier family 6 member 11	astrocytes, excitatory neurons, inhibitory neurons, muller glia cells, oligodendrocyte precursor cells	54
LSM5	LSM5 homolog, U6 small nuclear RNA and mRNA degradation associated		49
FAM172A	Family with sequence similarity 172 member A	astrocytes, excitatory neurons, inhibitory neurons, microglia, oligodendrocyte precursor cells, oligodendrocytes	47
WNT7A	Wnt family member 7A	cytotrophoblasts, extravillous trophoblasts, pancreatic endocrine cells, syncytiotrophoblasts	43
EXO1	Exonuclease 1	early spermatids, erythroid cells, extravillous trophoblasts, plasma cells, spermatocytes, spermatogonia	42
KRT8	Keratin 8	distal enterocytes, ductal cells, exocrine glandular cells, extravillous trophoblasts, pancreatic endocrine cells, paneth cells, proximal enterocytes, syncytiotrophoblasts	35
ZNF37A	Zinc finger protein 37A		32
TIGD6	Tigger transposable element derived 6	spermatocytes	27
CSTF3	Cleavage stimulation factor subunit 3		21
GDPGP1	GDP-D-glucose phosphorylase 1		18
HLA-DQB1	Major histocompatibility complex, class II, DQ beta 1	b-cells, langerhans cells, macrophages, monocytes	16
LSM10	LSM10, U7 small nuclear RNA associated		15
ZNF132	Zinc finger protein 132	cardiomyocytes, cone photoreceptor cells, early spermatids	14
CCDC197	Coiled-coil domain containing 197		13
AK4	Adenylate kinase 4	hepatocytes, proximal tubular cells	13
DDX43	DEAD-box helicase 43	spermatocytes, spermatogonia	11
AKAP13	A-kinase anchoring protein 13		9
ZC3H3	Zinc finger CCCH-type containing 3		9
ZCCHC9	Zinc finger CCHC-type containing 9		9
MFAP3L	Microfibril associated protein 3 like	late spermatids	8
IFITM2	Interferon induced transmembrane protein 2	adipocytes, nk-cells	8

Table 2 (continued)

Gene	Gene description	Single-cell RNA cell specific	Count
GPA33	Glycoprotein A33	distal enterocytes, intestinal goblet cells, paneth cells, proximal enterocytes, undifferentiated cells	8
MMEL1	Membrane metalloendopeptidase like 1	early spermatids, late spermatids, paneth cells	8
PIGP	Phosphatidylinositol glycan anchor biosynthesis class P	spermatocytes	7
COA6	Cytochrome c oxidase assembly factor 6		7
PCNX4	Pecanex 4	dendritic cells, endometrial ciliated cells	7
PHGDH	Phosphoglycerate dehydrogenase	distal tubular cells, exocrine glandular cells	6
NPTXR	Neuronal pentraxin receptor	cardiomyocytes, excitatory neurons, inhibitory neurons, oligodendrocyte precursor cells	6
IMMT	Inner membrane mitochondrial protein	cardiomyocytes	6
CDC37L1	Cell division cycle 37 like 1		5
SMC2	Structural maintenance of chromosomes 2	endometrial ciliated cells	5
AGR2	Anterior gradient 2, protein disulphide isomerase family member	club cells, distal enterocytes, gastric mucus-secreting cells, intestinal goblet cells, paneth cells, undifferentiated cells, urothelial cells	4
SLC8A3	Solute carrier family 8 member A3	bipolar cells, excitatory neurons, inhibitory neurons, oligodendrocyte precursor cells	4
AKT1S1	AKT1 substrate 1	syncytiotrophoblasts	4
HNRNPUL1	Heterogeneous nuclear ribonucleoprotein U like 1		4
RNASEH2A	Ribonuclease H2 subunit A	cytotrophoblasts, extravillous trophoblasts	3
PRR18	Proline rich 18	oligodendrocytes	2
KCNQ3	Potassium voltage-gated channel subfamily Q member 3	excitatory neurons, inhibitory neurons, microglia	2
ITGA1	Integrin subunit alpha 1	adipocytes, endothelial cells, hepatocytes, sertoli cells, smooth muscle cells	2
GRIN2A	Glutamate ionotropic receptor NMDA type subunit 2A	excitatory neurons, inhibitory neurons	2
TACC3	Transforming acidic coiled-coil containing protein 3	early spermatids, hofbauer cells, late spermatids, spermatocytes	2
EIF2A	Eukaryotic translation initiation factor 2A		2
MPHOSPH9	M-phase phosphoprotein 9	cone photoreceptor cells, excitatory neurons, inhibitory neurons	2
F7	Coagulation factor VII	hepatocytes	2
ZNF211	Zinc finger protein 211		2
KLHL24	Kelch like family member 24	oligodendrocytes	2
DKKL1	Dickkopf like acrosomal protein 1	early spermatids, late spermatids	2
STPG2	Sperm tail PG-rich repeat containing 2	astrocytes, excitatory neurons, inhibitory neurons, microglia, oligodendrocyte precursor cells, oligodendrocytes	1
PLD3	Phospholipase D family member 3	hofbauer cells	1
MZT1	Mitotic spindle organizing protein 1	erythroid cells	1
EVI5	Ecotropic viral integration site 5	oligodendrocytes	1

(continued on next page)

Table 2 (continued)

Gene	Gene description	Single-cell RNA cell specific	Count
DGCR2	DiGeorge syndrome critical region gene 2		1
GABRA5	Gamma-aminobutyric acid type A receptor subunit alpha5	excitatory neurons, inhibitory neurons	1
CYP4F12	Cytochrome P450 family 4 subfamily F member 12	distal enterocytes, gastric mucus-secreting cells, paneth cells, proximal enterocytes	1
RNF215	Ring finger protein 215	bipolar cells, cone photoreceptor cells, rod photoreceptor cells	1
ADARB2	Adenosine deaminase RNA specific B2 (inactive)	inhibitory neurons, oligodendrocyte precursor cells	1
MYL12A	Myosin light chain 12 A	cardiomyocytes	1
RUVBL2	RuvB like AAA ATPase 2	early spermatids, late spermatids, respiratory ciliated cells, spermatocytes	1
TRIM2	Tripartite motif containing 2	oligodendrocytes	1
CADPS	Calcium dependent secretion activator	bipolar cells, cone photoreceptor cells, excitatory neurons, inhibitory neurons, oligodendrocyte precursor cells	1
PCK1	Phosphoenolpyruvate carboxykinase 1	distal enterocytes, distal tubular cells, hepatocytes, proximal enterocytes, proximal tubular cells	1
NGDN	Neuroguidin		1
ABCA8	ATP binding cassette subfamily A member 8	fibroblasts, granulosa cells, leydig cells, microglia, oligodendrocytes, skeletal myocytes	1
POGZ	Pogo transposable element derived with ZNF domain	oligodendrocytes	1
PDIA6	Protein disulfide isomerase family A member 6	extravillous trophoblasts, plasma cells	1
ACSL1	Acyl-CoA synthetase long chain family member 1	early spermatids, hepatocytes, late spermatids	1
SUPV3L1	Suv3 like RNA helicase		1
DCXR	Dicarbonyl and L-xylulose reductase	distal tubular cells, hepatocytes, proximal tubular cells	1
DPYSL5	Dihydropyrimidinase like 5	cone photoreceptor cells, late spermatids, oligodendrocytes	1
OMA1	OMA1 zinc metalloproteinase	oligodendrocytes	1
NME1	NME/NM23 nucleoside diphosphate kinase 1		1

target genes, 50 have reproductive-specific expression, 27 have brain-specific expression, and 20 have mitochondria-specific expression or localization, according to annotation data from the Human Protein Atlas.

3.4. DB TF target genes are enriched in mitochondria and neural tissue

Enrichment analysis of the DB TF target genes was performed via the MSigDB annotation database [18,19] via Fisher's exact test, with Benjamini–Hochberg correction for multiple testing and an FDR of 0.10. The top 10 terms with the strongest enrichment (odds ratio, OR) among DB TFs target genes were: 'Cristae formation' (OR, 65.47; FDR, 6.93×10^{-2}), 'Kinetochore organization' (OR, 62.02; FDR, 7.02×10^{-2}), 'Hepatic encephalopathy' (OR, 62.02; FDR, 7.02×10^{-2}), 'Generalized clonic seizure' (OR, 59.74; FDR, 2.70×10^{-2}), 'GTP metabolic process'

(OR, 53.56; FDR, 7.66×10^{-2}), 'Mitochondrial calcium ion homeostasis' (OR, 53.56; FDR, 7.66×10^{-2}), 'mRNA 3 end processing' (OR, 49.77; FDR, 3.44×10^{-2}), 'Broad neck' (OR, 43.64; FDR, 8.9×10^{-2}), 'Bilateral tonic clonic seizure with generalized onset' (OR, 40.72; FDR, 3.46×10^{-2}), and 'Positive regulation of excitatory postsynaptic potential' (OR, 39.27; FDR, 9.89×10^{-2}). Full enrichment analysis results are presented in [Supplementary Table 4](#).

3.5. DB TF target genes and their direct interactors are enriched in miRNA and developmental phenotypes

The OmniPath DB [29,30] was queried to identify all proteins that have posttranslational interactions with DB TF target proteins. This identified 805 genes that were combined with the 75 DB TF target genes for gene ontology enrichment analysis via the MSigDB annotation database [18,19] via Fisher's exact test, with Benjamini–Hochberg correction for multiple testing and an FDR of 0.1. The top 10 non-redundant terms with the strongest enrichment (odds ratio, OR) among DB TFs target genes were: 'Positive regulation of miRNA transcription' (OR, 87.53; FDR, 1.44×10^{-46}), 'Negative regulation of cellular senescence' (OR, 40.46; FDR, 4.88×10^{-14}), 'Esophageal stricture' (OR, 40.34; FDR, 4.91×10^{-11}), 'Lung cell differentiation' (OR, 37.79; FDR, 8.98×10^{-15}), 'Protein K11 linked ubiquitination' (OR, 32.45; FDR, 4.87×10^{-16}), 'Integrated stress response signaling' (OR, 31.70; FDR, 4.87×10^{-16}), 'Neonatal insulin dependent diabetes mellitus' (OR, 30.22; FDR, 2.58×10^{-9}), 'Descending aortic dissection' (OR, 30.22; FDR, 2.58×10^{-9}), 'Lung epithelium development' (OR, 29.73; FDR, 2.28×10^{-18}), and 'Secondary palate development' (OR, 29.61; FDR, 4.46×10^{-11}). Full enrichment analysis results are presented in [Supplementary Table 5](#).

Fisher's exact test revealed enrichment of target genes and their direct interactors among autism spectrum disorder (ASD) genes from the SFARI gene knowledgebase (Banerjee-Basu and Packer 2010) (sfari.org; vers. 3.0), (OR 2.68; p value 1.16×10^{-16}).

3.6. DB TFs are coexpressed in neural and immune tissues

Data derived from the bulk RNA-Seq experiments (312 different tissues; 686 samples) of the FANTOM5 project were used to generate a cluster map of DB TF expression within the 100 tissues with the highest aggregate expression of the DB TFs, presented in [Fig. 2A](#) (y-axis DB TFs; x-axis, samples). On the basis of their tissue expression in the FANTOM dataset, the DB TFs form two distinct top-level clusters: the larger cluster contains 89 genes, the majority of which are known to have a role in development (78/89) and specific expression in brain or retina tissue (72/89), whereas the smaller cluster of 17 genes all show specific expression in blood and immune cells. Data on tissue group specificity (brain and retina; male/female) from the Human Protein Atlas (bulk and single-cell RNA-Seq), status as a development-related gene (homeobox, forkhead box, and SRY-related HMG-box), and above the 90th quantile expression in immune cell types as defined by the Database of Immune Cell eQTLs (DICE), were analyzed and included in [Fig. 2C](#).

Within the larger cluster of 89 developmental- and brain-specific genes, there is a subcluster of 12 genes whose expression is most limited to FANTOM 5 neural tissues: NKX6-2, CUX2, NR2E1, POU3F2, POU6F2, EMX1, VAX1, POU3F3, DLX1, DLX2, OTX1, and POU3F1. All of these genes are identified in the Protein Atlas as being enriched or enhanced in the brain or retina. Compared with all protein-coding genes, GO analysis of these 12 genes with the PANTHER-based gene ontology.org web server [16,17] revealed the top 10 enriched biological process terms: 'cerebral cortex GABAergic interneuron fate commitment' (GO:0021893; 100x enrichment), 'positive regulation of amacrine cell differentiation' (GO:1902871; 100x), 'forebrain ventricular zone progenitor cell division' (GO:0021869; 100x), 'negative regulation of photoreceptor cell differentiation' (GO:0046533; 100x), 'regulation of photoreceptor cell differentiation' (GO:0046532; 100x), 'regulation of

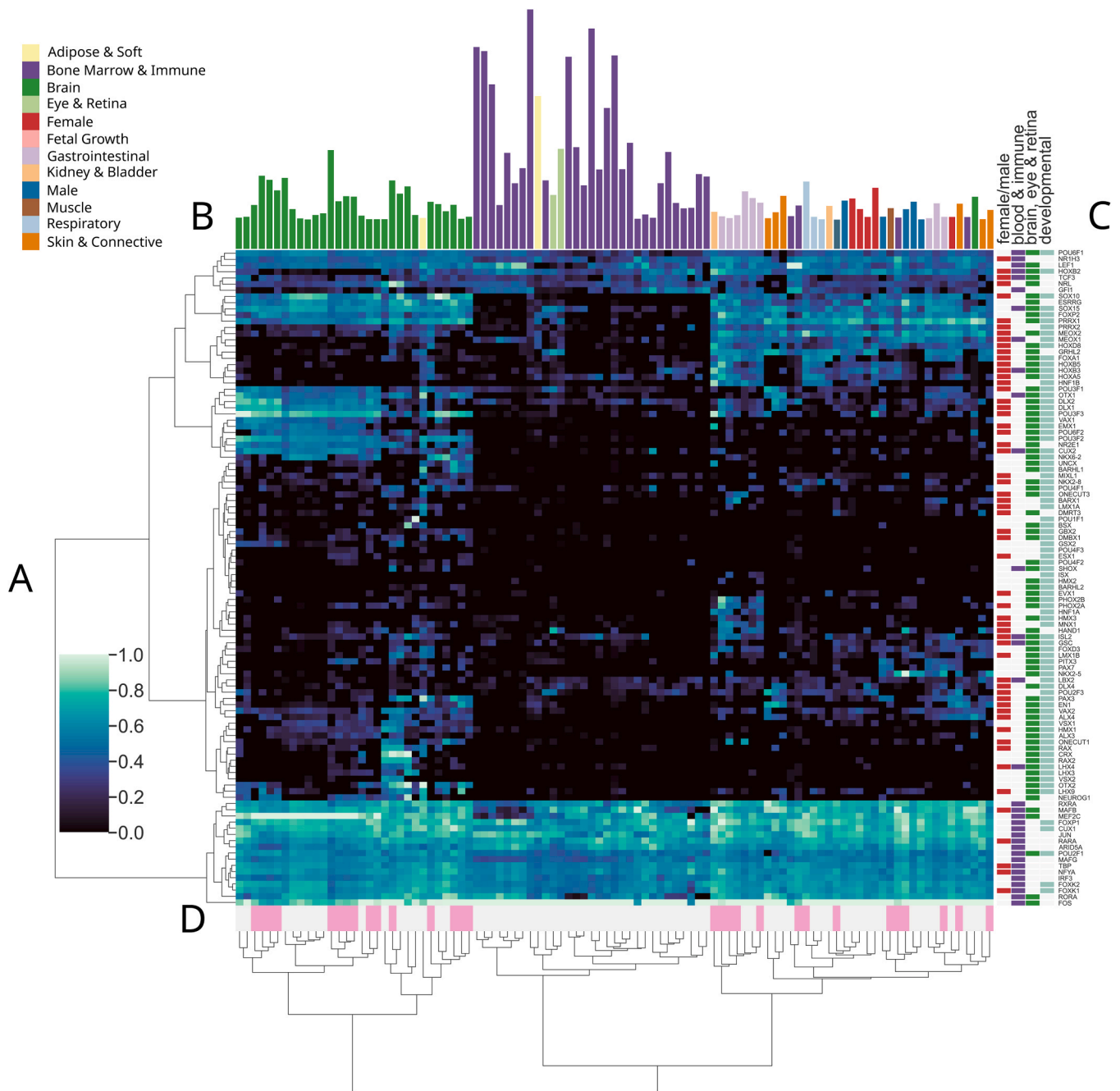


Fig. 2. Cluster map of the expression of transcription factors displaying variable binding in modern human vs. Neanderthal. (A) Expression data derived from the FANTOM 5 dataset were used to cluster 106 DB transcription factors (y-axis) in 100 tissues (x-axis). Expression data were not available for the DUXA and PROX1 genes. (B) Aggregate expression for each tissue across all genes is presented as bars, and bars/tissues are colored by tissue group. Distinct clusters of brain and eye/retina tissues as well as blood and immune cells/tissues are observed. A bar plot with the tissue names is presented in [Supplementary Figure 1](#). (C) TF genes are labeled with colored bars corresponding to their identification as homeobox, SOX, or FOX genes (developmental), above the 90th quantile expression in an immune cell type as defined by the DICE database (blood & immune) [31], or to have enriched or enhanced expression by tissue as defined by proteinatlas.org (female/male; brain, eye & retina) [32]. (D) Life stage for each tissue is described, where white and pink represent adult and fetal/newborn tissues, respectively. The expression values (TPM) were normalized via log transformation, and each tissue column was adjusted on a zero-to-one scale.

amacrine cell differentiation’ (GO:1902869; 100x), ‘neuroblast differentiation’ (GO:0014016; 100x), ‘negative regulation of oligodendrocyte differentiation’ (GO:0048715; 100x), ‘cerebral cortex GABAergic interneuron differentiation’ (GO:0021892; 100x), and ‘negative regulation of glial cell differentiation’ (GO:0045686; 100x).

Additionally, within the brain tissue cluster (Fig. 2B), we observed subclusters of fetal and newborn tissues (Fig. 2D). Overall, DB TFs presented the highest expression in the medial frontal gyrus (newborn), medial temporal gyrus (adult), parietal lobe (fetal), occipital lobe (fetal),

eye (fetal), pineal gland (adult), temporal lobe (fetal), retina (adult), occipital cortex (newborn), parietal lobe (newborn), dura mater (adult), medial temporal gyrus (newborn), and spinal cord (fetal) ([Supplementary Figure 1](#)). Similarly, the cluster with high expression in blood and immune cells/tissues is comprised of 17 genes: FOS, RORA, FOXP1, FOXP2, IRF3, NFYA, TBP, MAFG, POU2F1, ARID5A, RARA, JUN, CUX1, FOXP1, MEF2C, MAFB, and RXRA. All of these genes have an expression ≥ 90 th quantile (compared with all protein-coding genes) in at least one immune cell type, as cataloged in the Database of Immune

Cell eQTLs (DICE) [31]. GO analysis of the 17 genes (geneontology.org web server) revealed enrichment of immune-related biological process terms: ‘CD4-positive, alpha-beta T-cell differentiation involved in immune response’ (GO:0002294; 76x enrichment), ‘T-helper cell differentiation’ (GO:0042093; 76x), and ‘alpha-beta T-cell differentiation involved in immune response’ (GO:0002293; 73x), among others.

Mapping all modern human vs Neanderthal SNPs to the promoters of protein-coding genes and classifying by gene biotype, we observed the highest rate of SNPs per nucleotide in IG J (joining chain immunoglobulin) pseudogenes and genes, and the lowest rate for IG D (diversity chain immunoglobulin) genes (Fig. 1A). Additionally, we observe a small reduction in the number of SNPs near the TSS of protein-coding genes (Fig. 1B).

3.7. DB TFs have developmental time-point-related expression profiles in the brain

Cluster analysis of the Allen Brain Atlas bulk RNA-Seq expression data for DB TF genes revealed distinct clusters of time-point-specific expression. Those genes comprising the cluster with the highest aggregate expression in the dataset (module 1 genes: NFYA, TCF3, MEF2C, SOX10, CUX2, TBP, FOXP1, IRF3, SOX15, RARA, FOXP2, NKX6–2, FOXP1, FOS, JUN, POU6F1, POU3F2, MAFG, POU3F3, MAFB, and CUX1) show a dynamic and cyclic pattern of expression during the prebirth timepoints, which then stabilizes in infancy and trends downward in early childhood (Fig. 3A). Despite differences in physiology, different tissues generally cluster together by time point rather than tissue type: pcw (8–37 weeks post conception), early (4 months–4 years), and late (8 years–40 years) (Fig. 3B).

3.8. Single-nucleus RNA-Seq of cortical cells reveals differential expression of DB TFs

An analysis of the Allen Brain Atlas single-nucleus RNA-Seq data from 49,417 cortical brain cells [20] was performed to identify which of the MH-Neanderthal DB TF genes may play a functional role in specific annotated brain regions, layers, and cell types. The analyses in the original paper defined metadata indicating cell class (excitatory,

inhibitory, and nonneuronal), layer (L1–L6), and predicted type (e.g., astrocyte, microglia, endothelial, etc.), as well as brain regions (Fig. 4A). Cell groups defined by the intersection of class, layer, and region were analyzed for differentially expressed genes (DEGs), which are the genes that have the most distinct expression in the target group(s) compared to the other groups. Several DB TFs were DEGs in specific categories, occurring in nearly all cases in glutamatergic neurons, usually in layers L4–L6, and across nearly all brain regions. The categories in which four or more DB TFs were DEGs are presented in Table 3 and included cut-like homeobox 1 (CUX1), cut-like homeobox 2 (CUX2), estrogen-related receptor gamma (ESRRG), forkhead box protein 1 (FOXP1), forkhead box protein 2 (FOXP2), Myocyte Enhancer Factor 2 C (MEF2C), POU Class 6 Homeobox 2 (POU6F2), paired related homeobox 1 (PRRX1), and RAR-related orphan receptor alpha (RORA). In all categories, POU6F2 was a DEG, and in all but two categories, either FOXP2 or FOXP1 (or both) were DEGs. FOXP2 was present as a DEG in excitatory L4 cells in A1C (primary visual cortex), MTG (medial temporal gyrus), S1_{lm} (lower limb somatosensory cortex), and S1_{ul} (upper limb somatosensory cortex).

3.9. Gene ontology analysis of L4 excitatory neurons

In the Allen Brain Atlas cortical cell snRNA-Seq data, multiple DB TFs were consistently up-regulated DEGs in both layer 4 and excitatory cell groups, and their intersection. Gene ontology analysis of the top 100 marker genes of L4 excitatory cells via the DisGeNET ontology within the Enrichr tool [26] revealed the top enriched terms ‘autism spectrum disorders’, ‘autistic disorder’, ‘narcolepsy’, ‘neurodevelopmental disorders’, ‘intelligence’, ‘epilepsy’, ‘schizophrenia’, ‘refractive errors’, ‘cognition’, and ‘apraxia, developmental verbal’ (Fig. 4D). The top 10 enriched biological process terms for L4 excitatory cells were: ‘corticospinal neuron axon guidance’ (GO:0021966; >100x enrichment), ‘corticospinal tract morphogenesis’ (GO:0021957; 91x), ‘ventricular cardiac muscle cell differentiation’ (GO:0055012; 42x), ‘synaptic membrane adhesion’ (GO:0099560; 35x), ‘central nervous system projection neuron axonogenesis’ (GO:0021952; 30x), ‘synaptic transmission’, ‘GABAergic’ (GO:0051932; 29x), ‘positive regulation of excitatory postsynaptic potential’ (GO:2000463; 28x), ‘negative regulation of

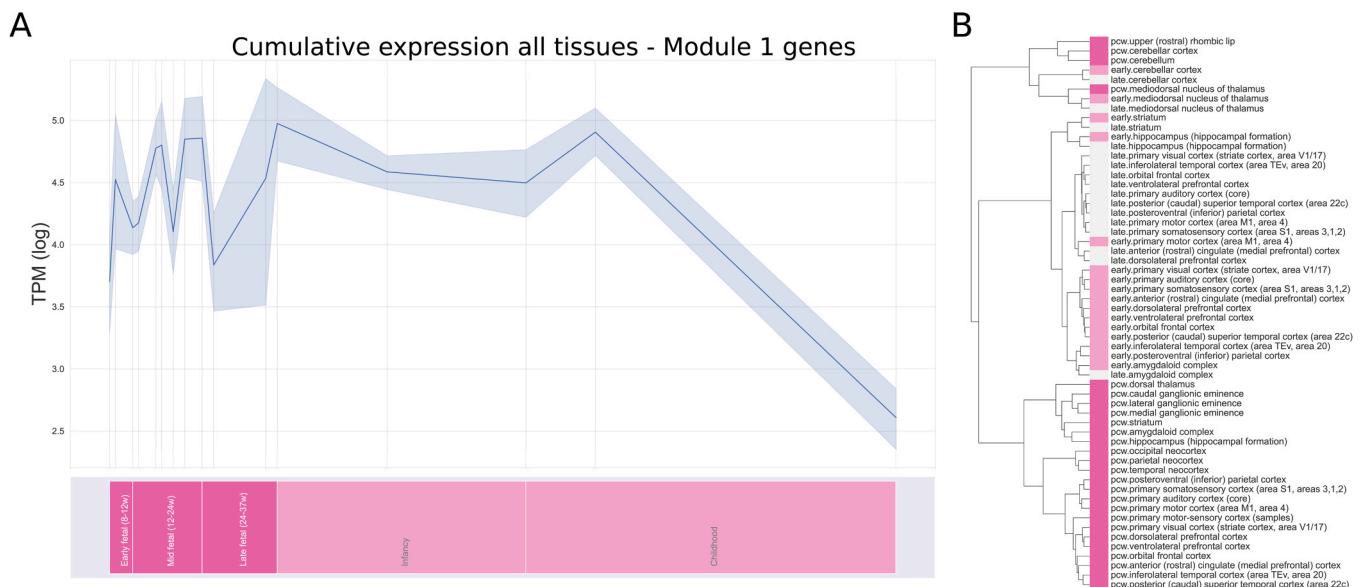


Fig. 3. Expression of modern human vs Neanderthal DB TFs in developing brain tissue. The expression data from the Allen Brain Atlas [33] were extracted as RPKM values for 26 unique tissues across 31 timepoints, converted to TPM values, log transformed, and used to construct a cluster map. A) Aggregate expression of the cluster of DB TF genes with the highest total expression across all tissues (module 1) from 8–144 weeks post conception. B) Clustering of brain tissues defined by development time point. The colors represent development stages for both figures: 8–37 weeks post conception (pcw, magenta), 4 months–4 years (early, pink), and 8–40 years (late, white).

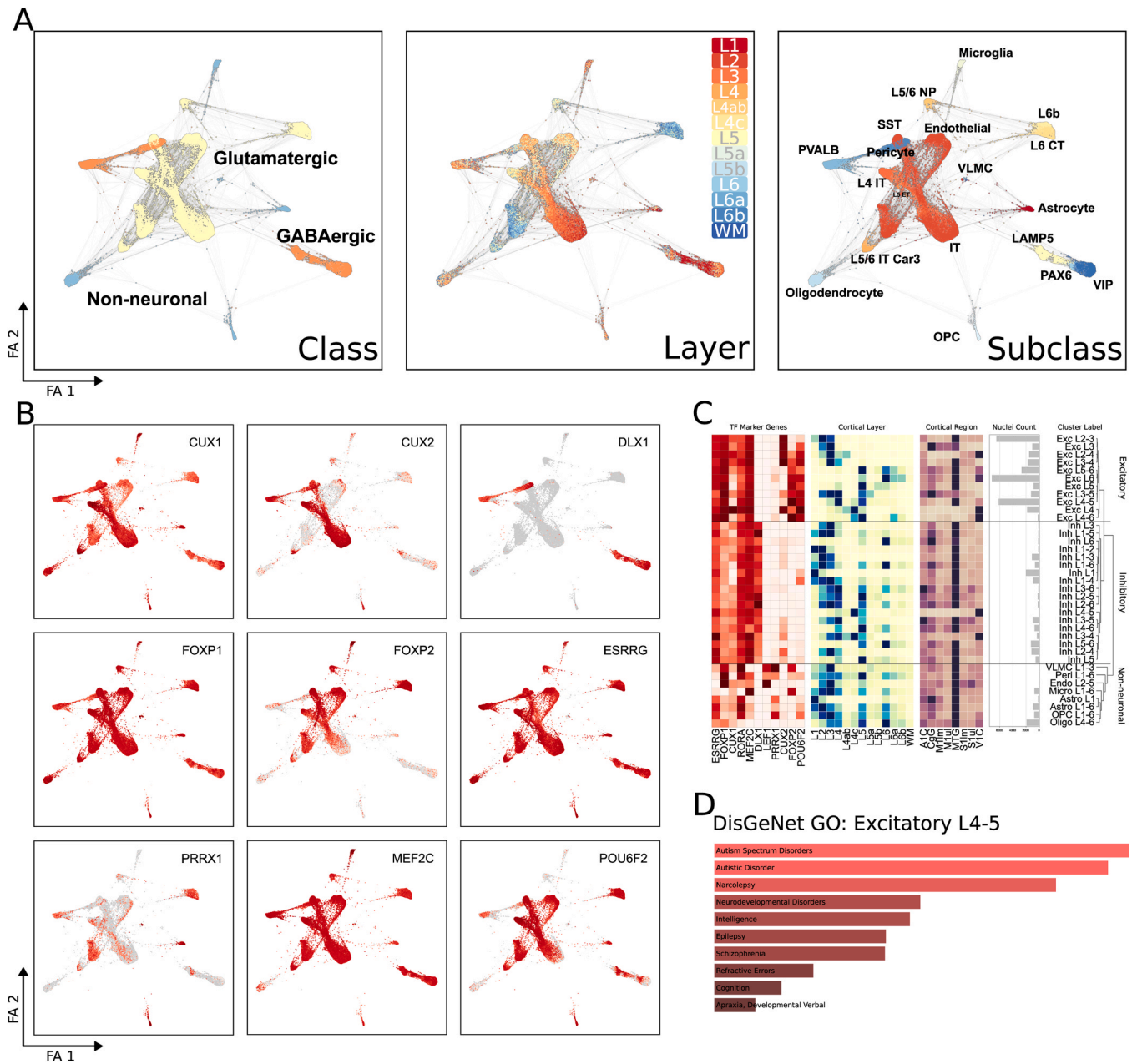


Fig. 4. Analysis of single-nucleus RNA-Seq data from brain cortical regions. (A) Force-directed graph (FA) visualization of Allen Brain Atlas single-nucleus RNA-Seq data from 49,417 brain cortex cells; the brain excitatory/inhibitory status, cortical layer, and predicted cell type (subclass) data are mapped from annotations in the Allen Brain Atlas sample data [20]. (B) Expression of DB TFs, which are DEGs mapped onto FA visualization. (C) On the x-axis: DB TF DEGs, cortical layer, and cortical region; on the y-axis: cell class/type and cortical layer, with the dendrogram indicating expression relationships. The color intensity indicates the normalized expression (0–1 scale). Nuclei count indicates the number of nuclei for that cell class/type at the indicated layer. (D) DisGeNET disease ontology terms enriched for the top 100 marker genes for L4–5 excitatory cells. ET, extratelencephalic/pyramidal tract; IT, intratelencephalic; NP, near-projecting; OPC, oligodendrocyte precursor; VLMC, vascular leptomenigeal cell; SST, LAMP5, PAX6, VIP, PVALB, labels on the basis of expression of the corresponding gene.

smooth muscle cell migration' (GO:0014912: 27x), 'positive regulation of dendrite morphogenesis' (GO:0050775: 25x), and 'modulation of excitatory postsynaptic potential' (GO:0098815: 25x).

3.10. Gene ontology analysis of L4 excitatory neurons

In the GO analysis of the CC category, the top 10 enriched terms for L4 excitatory cells were: 'anchored component of presynaptic membrane' (GO:0099026; 61x enrichment), 'presynaptic cytosol' (GO:0099523; 47x), 'NMDA selective glutamate receptor complex' (GO:0017146; 47x), 'intrinsic component of presynaptic active zone membrane' (GO:0098945; 45x), 'node of Ranvier' (GO:0033268; 42x),

'integral component of presynaptic active zone membrane' (GO:0099059; 40x), 'GABA-A receptor complex' (GO:1902711; 34x), 'GABA receptor complex' (GO:1902710; 30x), 'presynaptic active zone membrane' (GO:0048787; 27x), and 'intrinsic component of presynaptic membrane' (GO:0098889; 22x).

In the GO analysis of molecular function, the top 10 enriched terms for L4 excitatory cells were: 'transmembrane receptor protein tyrosine phosphatase activity' (GO:0005001; 50x enrichment), 'transmembrane receptor protein phosphatase activity' (GO:0019198; 50x), 'GABA-gated chloride ion channel activity' (GO:0022851; 49x), 'inhibitory extracellular ligand-gated ion channel activity' (GO:0005237; 42x), 'ligand-gated anion channel activity' (GO:0099095; 35x), 'GABA-A receptor

Table 3

DB TFs that co-occur as DEGs in specific regions, layers, and cell types according to the Allen Brain Atlas scRNA-Seq data from the cerebral cortex.

Layer	Class	Region	DEG DB TF genes
L4	Glutamatergic	A1C	FOXP2, FOXP1, POU6F2, MEF2C
L4	Glutamatergic	MTG	FOXP2, PRRX1, POU6F2, MEF2C
L4	Glutamatergic	S1lm	FOXP2, PRRX1, POU6F2, ESRRG
L4	Glutamatergic	S1ul	FOXP2, FOXP1, RORA, POU6F2, MEF2C
L4	Glutamatergic		FOXP1, CUX2, RORA, POU6F2, CUX1, MEF2C
L4		IT	FOXP1, CUX2, RORA, POU6F2, CUX1, MEF2C
L4ab	Glutamatergic	V1C	CUX2, RORA, POU6F2, CUX1, MEF2C
L4ab			CUX2, RORA, POU6F2, CUX1
L5	Glutamatergic	S1lm	FOXP2, PRRX1, POU6F2, ESRRG
L5	Glutamatergic	V1C	FOXP2, FOXP1, RORA, POU6F2, MEF2C
L6a			FOXP2, PRRX1, POU6F2, ESRRG

activity' (GO:0004890; 34x), 'syntaxin-1 binding' (GO:0017075; 32x), 'GABA receptor activity' (GO:0016917; 29x), 'transmitter-gated channel activity' (GO:0022835; 17x), and 'transmitter-gated ion channel activity' (GO:0022824; 17x).

3.11. DB TFs and DB TF target genes in autism

Analysis of snRNA-Seq data from autism and control samples of cortical tissue showed differential expression, adjusted p-value (≤ 0.05) and fold-change (≤ 0.5 ; ≥ 1.5), of 15 DB TFs and 16 DB TF target genes (Table 4). Broadly, we observed upregulation (higher expression in ASD) of both DB TFs and their target genes in neurons and astrocytes, and downregulation in interneurons and oligodendrocytes. The three most upregulated genes were DB TF target genes (MFAP3L, NPTXR, and GABRA5) with higher expression in maturing neurons. The gene occurring most commonly as differentially expressed in various comparisons was FOS (10 comparisons) followed by NPTXR (6 comparisons) and GABRA5 (4 comparisons).

3.12. DB TFs and human accelerated regions

Doan et al. identified and cataloged a list of genes in close proximity to human accelerated regions (HARs; elevated divergence in humans vs. other primates) as well as those that directly interact with HARs through chromatin conformation changes, e.g., enhancers [34]. Among the set of 110 DB TF genes, 28 were HAR-associated (Table 1). Fisher's exact test revealed that the overrepresentation of HAR associations in DB TFs versus all JASPAR TFs was not significant (OR 0.98; p value 0.57). However, compared with all protein-coding genes, the JASPAR TFs as a group were enriched for HAR associations (OR 3.27; p value 6.73×10^{-26}).

4. Discussion

4.1. Altered regulation centers on sexual reproduction, neurological development, and mitochondrial function

In our study we identified the transcription factors with statistically different binding affinities at MH vs. Neanderthal promoter SNPs. We performed enrichment analyses on these TFs and the genes whose promoters contain differential binding sites. The 110 DB TFs are comprised of several distinct features, specifically they are predominantly homeobox (78 of 110 genes) and have sexual reproduction (76/110) and brain-specific expression (72/110). They are enriched for biological process and phenotype terms related to vision (e.g., sensory perception of light stimulus), brain (e.g., cell morphogenesis involved in neuron differentiation), and development (e.g., head development). The 75 genes in whose promoters the DB TFs bind in a manner differential between modern human and Neanderthal, the DB TF target genes, also share several characteristics. They likewise have specific expression in

reproductive (50/75) and brain (27/75) tissues, and also mitochondria-associated functions (20/75). They are enriched in biological process and phenotype terms related to mitochondria function (e.g., cristae formation) and brain (e.g., positive regulation of excitatory postsynaptic potential). The DB TFs and their target genes appear to have clear function related to reproduction and neurological development which

4.2. Modern human brain regulatory changes

Hierarchical clustering of TF expression revealed that DB TF genes associated with neural and immune functions formed distinct groups. Within the tissue axis, in the neural cluster, adult and fetal/newborn neural tissues cluster separately in several instances (Fig. 2). From these discoveries in a broader set of data, it became important to focus more closely on the expression of these DB TFs in developing neural tissues specifically. Analysis of RNA-Seq data from whole-brain subtissues further revealed differences in the expression of these genes across time and likely in a developmentally dependent manner. Specifically, the expression of DB TFs in all tissues for the two groups occurring within 8 weeks post conception to 4 years post birth was largely segregated from all tissues for the group defined by 8 years to 40 years of age (Fig. 3). Additionally, several modules of DB TFs identified by cluster analysis appear to show a cyclical pattern of expression during fetal development, which then stabilizes post birth (Fig. 3A). Additional work on the developmental expression patterns of these DB TFs is warranted.

4.3. DB TFs co-occur as DEGs according to single-cell RNA-Seq of glutamatergic cells in the L4 human cortex

Recently, the proliferation of single-cell transcriptomics has begun to clarify the direct role of transcription factors in cell type determination and identity [35]. To determine in which cell types DB TFs are expressed in the adult brain, single-nucleus data for 49,417 cortical cells were examined. A subset of the DB TFs was annotated as commonly co-occurring DEGs: CUX1, CUX2, ESRRG, FOXP1, FOXP2, MEF2C, POU6F2, PRRX1, and RORA. DEGs, also known as marker genes, are highly relevant as they are often the genes most relevant to the characteristics of the cell type/state. Interestingly, we observed that these DB TF DEGs were primarily co-occurring in glutamatergic cells of the L4 cortex. Among this set of nine genes, six (CUX1, CUX2, FOXP1, FOXP2, MEF2C, and RORA) are cataloged in the SFARI Gene knowledgebase [36] (sfari.org; vers. 3.0) as ASD candidate genes. Similarly, six (CUX1, FOXP1, FOXP2, MEF2C, POU6F2, and RORA) are present in the genome-wide association studies (GWAS) catalog [37] (www.ebi.ac.uk/gwas) for the 'schizophrenia' (MONDO_0005090) term. In support of this connection, a recent large-scale study on schizophrenia genetics in 74,776 individuals with schizophrenia (and 101,023 control individuals) revealed enrichment of genes that are highly expressed in glutamatergic neurons of both the human and mouse cerebral cortex [38]. Similarly, a recent large-scale study on ASD genetics in 11,986 individuals with ASD (and 23,598 control individuals) identified 102 ASD risk genes and revealed that fetal-development excitatory neurons accounted for the greatest number of these genes among all the identified cell types [39].

4.4. Co-occurring DB TF DEGs associate with HARs and neuropsychiatric disease

Human-accelerated regions (HARs) are locations in the genome where the rate of evolutionary change has accelerated since divergence from chimpanzee. A study of 2737 HARs revealed that they are enriched for TFBSs generally and for TFs associated with neural development specifically (Doan et al., 2016). While the DB TFs were not enriched for HAR association compared with the JASPAR TFs (OR 0.98; p value 0.57), the JASPAR TFs are compared to all genes (OR 3.27; p value 6.73×10^{-26}). Using the HARs cataloged by [34], we identified 28 DB

Table 4

DB TFs and their target genes are differentially expressed in ASD vs. control.

Gene	Fold change	P adj.	Category	Observation	DB TF	DB TF target gene
MFAP3L	2.68	1.38E-03	cluster	Neurons - maturing		x
NPTXR	2.66	2.53E-12	cluster	Neurons - maturing		x
GABRA5	2.58	6.56E-06	cluster	Neurons - maturing		x
FOS	2.44	6.74E-08	cluster	Astrocytes - fibrous	x	
FOXP2	2.36	2.33E-10	cluster	Microglia	x	
JUN	2.36	1.63E-03	cluster	Neurons - maturing	x	
FOS	2.34	1.56E-14	cluster	L5/6-cortico-cortical projection neurons	x	
POU2F1	2.26	1.03E-07	cluster	Neurons - maturing	x	
FOS	2.25	8.28E-08	cluster	Astrocytes - protoplasmic	x	
LEF1	2.12	2.21E-08	cluster	Endothelial	x	
FOS	2.11	4.32E-05	cluster	Parvalbumin interneurons	x	
NPTXR	2.03	4.38E-18	cluster	NRGN-expressing neurons II		x
DGCR2	2.01	8.27E-07	cluster	NRGN-expressing neurons I		x
NPTXR	2.00	3.04E-20	cluster	NRGN-expressing neurons I		x
JUN	1.99	5.15E-04	cluster	Parvalbumin interneurons	x	
CUX1	1.87	1.09E-05	cluster	Neurons - maturing	x	
FAM172A	1.83	2.58E-05	cluster	Neurons - maturing		x
NPTXR	1.81	1.81E-69	region	Anterior cingulate cortex		x
POGZ	1.77	2.58E-04	cluster	Neurons - maturing		x
GABRA5	1.77	1.94E-36	region	Anterior cingulate cortex		x
NPTXR	1.73	2.19E-45	cluster	L5/6-cortico-cortical projection neurons		x
RORA	1.71	1.42E-158	diagnosis	ASD	x	
PLD3	1.71	5.30E-04	cluster	Neurons - maturing		x
RORA	1.69	9.46E-106	region	Prefrontal cortex	x	
KLHL24	1.65	5.45E-03	cluster	Neurons - maturing		x
JUN	1.65	1.47E-02	cluster	Astrocytes - protoplasmic	x	
RORA	1.65	3.91E-47	region	Anterior cingulate cortex	x	
GABRA5	1.64	2.12E-04	cluster	NRGN-expressing neurons I		x
POU3F2	1.64	5.11E-05	region	Anterior cingulate cortex	x	
FOS	1.63	3.05E-13	region	Prefrontal cortex	x	
FOXP1	1.62	7.08E-08	cluster	Endothelial	x	
PRRX1	1.62	4.71E-02	cluster	Endothelial	x	
GABRA5	1.62	4.55E-30	cluster	L5/6-cortico-cortical projection neurons		x
GRIN2A	1.62	6.25E-08	cluster	NRGN-expressing neurons II		x
FOXP1	1.61	7.66E-04	cluster	NRGN-expressing neurons II	x	
MAFB	1.59	9.15E-06	cluster	Parvalbumin interneurons	x	
FOS	1.59	1.28E-02	cluster	L5/6	x	
FOS	1.57	7.70E-03	cluster	Somatostatin interneurons	x	
FOXP1	1.56	6.93E-04	cluster	NRGN-expressing neurons I	x	
GRIN2A	1.56	1.26E-02	cluster	Neurons - maturing		x
MFAP3L	1.56	2.12E-02	cluster	Astrocytes - fibrous		x
IFITM2	1.54	5.75E-03	cluster	Endothelial		x
NPTXR	1.53	2.99E-24	cluster	L4		x
FAM172A	1.53	2.64E-03	cluster	Endothelial		x
MFAP3L	1.53	6.75E-13	region	Anterior cingulate cortex		x
GRIN2A	1.52	3.70E-07	cluster	NRGN-expressing neurons I		x
ITGA1	1.51	9.55E-04	cluster	Endothelial		x
PRR18	0.50	1.68E-02	region	Prefrontal cortex		x
LSM5	0.50	4.20E-12	cluster	L5/6-cortico-cortical projection neurons		x
FOS	0.49	6.90E-04	cluster	VIP interneurons	x	
NKX6-2	0.48	2.20E-05	region	Prefrontal cortex	x	
LSM10	0.48	1.48E-16	cluster	L2/3		x
LSM10	0.48	3.37E-06	cluster	L4		x
ESRRG	0.47	1.65E-24	cluster	L4	x	
NARF	0.46	3.47E-05	cluster	SV2C interneurons		x
LSM10	0.46	4.84E-03	cluster	Somatostatin interneurons		x
COA6	0.45	2.29E-02	cluster	NRGN-expressing neurons I		x
CUX2	0.45	5.55E-12	cluster	Somatostatin interneurons	x	
CUX2	0.45	5.90E-04	cluster	SV2C interneurons	x	
JUN	0.41	5.92E-05	cluster	Oligodendrocyte precursor cells	x	
ARID5A	0.41	1.10E-02	cluster	Endothelial	x	
FOS	0.26	5.99E-03	cluster	Oligodendrocytes	x	
FOS	0.24	1.18E-09	cluster	Oligodendrocyte precursor cells	x	

TFs that are either in close proximity to a HAR or directly interact with one through chromatin conformation changes. Additionally, all but one (CUX2) of the co-occurring DEGs DB TFs (CUX1, ESRRG, FOXP1, FOXP2, MEF2C, POU6F2, PRRX1, and RORA) identified in the L4 excitatory neurons of the cerebral cortex (Allen Brain Atlas scRNA-Seq data) contain HAR regions or HAR interactivity. Several of these contain or associate with HAR variants that have been noted to cause neurological disease. POU6F2 possesses an intronic HAR in which a mutant allele (GRCh38:chr7:39,033,595) is associated with ASD (Doan

et al., 2016). The promoter of the CUX1 gene has been determined via ChIA-Pet analysis of chromatin to interact with a HAR located ~200 kb away; unrelated individuals with intellectual disability (ID) (IQ<40) and ASD have been identified with a homozygous mutation (GRCh38: chr7:101,606,361) in this HAR. A luciferase reporter assay revealed that when the mutant version of this HAR interacts with the promoter of the CUX1 gene, its expression is increased threefold, while cultured differentiated neurons with increased CUX1 expression exhibit increased synaptic spine density. Similarly, another HAR (GRCh38:chr5:88,480,

873) has been shown to interact with the promoter of MEF2C, and a mutation in this HAR creates a putative MEF2A binding site, reducing expression by ~50 %. Mutations in the MEF2C gene are associated with autism [40,41], mental retardation [40–42], schizophrenia [43,44], epilepsy [40,41], and speech abnormalities [40]. Additionally, binding motifs for several of the DB TFs identified in our study were found to be enriched in HAR regions; POU6F1 and POU2F1 in ultra-conserved HARs, and HNF1A in all HARs [34]. While we did not observe enrichment of DB TFs with HARs, further analysis of those DB TFs which do associate with HARs is warranted.

4.5. FOXP2 is a DB TF and differentially expressed in ASD microglia

Importantly, one of the identified DB TFs in our data is the FOXP2 TF gene, which was the first noted "speech gene" [45]. For the early post-natal period in mice, FOXP2 is a negative regulator of MEF2C, and the likely result is the promotion of synaptogenesis in the cortical striatum [46]. In addition to this interaction with MEF2C, there is a POU3F2 (which, according to our results is also a DB TF) binding site within the FOXP2 gene, which has been shown to affect FOXP2 expression and is associated with a selective sweep that has occurred since the divergence of humans and Neanderthal [47,48]. FOXP2 has recently been shown to have human-specific expression in microglia [49], and in our analysis of an autism snRNA-Seq dataset we observed upregulated expression in ASD microglia (2.36-fold higher). Dysregulated gene expression in microglia has been associated with ASD [50,51]. Because of its significant role in human evolution, further analysis is warranted regarding the regulation of FOXP2 in other relevant areas (e.g., introns, 3' UTRs, enhancers, etc.) and for genome-wide binding sites of FOXP2 itself.

4.6. DB TFs and their target genes show differential expression in autism

Our analysis of an snRNA-Seq data set with both control and autism samples revealed differential expression of 15 DB TFs and 16 DB TF target genes (Table 4). Interestingly, compared with control, in autism samples we observed upregulation of genes from both DB TFs and their target genes in neurons and astrocytes. Many of the genes upregulated in autism samples (three DB TFs and eight DB TF target genes) had higher expression in maturing neurons. This includes the three genes with highest overall upregulation, all of which are DB TF target genes, MFAP3L, NPTXR, and GABRA5. In addition, several DB TFs (FOXP1) and their target genes (NPTXR, DGCR2, GABRA5, GRIN2A, COA6) were differentially expressed in neurons expressing NGRN, a gene associated with neuron development [52] and intellectual ability in schizophrenia [53]. Greater focus on neurogenesis as a cause of autism is currently unfolding, with 88 % of high-risk genes associated with neuron development [54].

In the cases where DB TFs and their target genes were downregulated in autism, seven of the nine largest changes occurred in SV2C- and somatostatin-expressing interneurons (CUX2, LSM10, and NARF), and oligodendrocytes and oligodendrocyte precursor cells (FOS and JUN). Conversely, only upregulation was observed in parvalbumin-expressing interneurons (FOS, JUN, and MAFB).

Among all genes, FOS showed altered expression in the greatest number of observations, with expression being higher in ASD fibrous astrocytes, L5/6-cortico-cortical projection neurons, protoplasmic astrocytes, parvalbumin interneurons, prefrontal cortex, L5/6 region, and somatostatin interneurons; while lower in oligodendrocyte precursor cells, oligodendrocytes, and VIP interneurons. This is notable as FOS is an immediate early gene, involved in rapid response to external stimuli [55], and has been extensively documented to contribute to neuronal function and implicated in psychiatric disorders [56].

4.7. DB TF target genes are expressed in oligodendrocytes, testes/sperm, the maternal–fetal interface, and mitochondria

There were 75 genes identified as DB TF target genes. All but 6 of these genes are specifically expressed in reproductive (50), brain and retina (27), or mitochondrial (25) cells or tissues. Among the 30 DB TF target genes with reproductive cell- and tissue-specific expression, the majority (19) are expressed in male tissues (testis, early and late spermatids, spermatocytes, and spermatogonia). A total of 12 genes were expressed in female-specific cells and tissues and, in many cases, specifically in those involved in the maternal–fetal interface (cervix, granulosa cells, cytotrophoblasts, extravillous trophoblasts, syncytiotrophoblasts, and endometrial ciliated cells). Among the 28 DB TF target genes with brain- and retina-specific expression, 16 are expressed in oligodendrocytes, 12 in inhibitory neurons, and 11 in excitatory neurons.

The DB TF target gene with the greatest number of differential binding events in its promoter region is FARS2, a member of the mitochondrial aminoacyl-tRNA synthetases (mt-aARSs). scRNA-Seq data from the Protein Atlas indicate that FARS2 expression is increased in excitatory neurons, oligodendrocyte precursor cells, inhibitory neurons, oligodendrocytes, astrocytes, microglia, late spermatids, and early spermatids. mt-aARSs charge their cognate mt-tRNA with the appropriate amino acid, and mutations in these genes cause a variety of diseases, most predominantly those that affect the central nervous system [57], perhaps due to delayed myelination, demyelination, or both [58]. In addition to their canonical roles, mt-aARSs are hypothesized to function in monitoring amino acid levels, as sensors for the mitochondrial environment, and in transcriptional regulation [57], and through the addition of new protein domains have been associated with neural development and immune response, among other processes, as reviewed by [59]. The FARS2 gene produces the mt-PheRS protein (phenylalanyl-tRNA synthetase), which is mitochondria-locating and is responsible for attaching phenylalanine to its corresponding mt-tRNA for mitochondrial protein translation [60]. Intragenic variants in the FARS2 gene have been linked to two primary clinical manifestations, early-onset epileptic mitochondrial encephalopathy and spastic paraplegia, and for patients in both groups, symptoms can also include intellectual disability or developmental delay [60]. Deletions within FARS2 and reduced expression levels have also been associated with schizophrenia [61].

Importantly, FARS2 shares a bidirectional promoter with another mitochondrial gene, LYRM4 (LYR motif-containing protein 4), which has a TSS just ~400 bp away. Together with ISCU, NFS1, and FXN, the ISD11 protein of the LYRM4 gene is involved in the formation of iron–sulfur (Fe–S) clusters [62,63], which are essential cofactors in many basic biological processes, such as the formation of respiratory chain complexes I, II, and III [64] and subsequent oxidative phosphorylation [65]. Mutations in LYRM4 cause deficits in oxidation phosphorylation reactions [66], which are critical in neuronal development and schizophrenia, as reviewed previously [67]. In support of this connection, polymorphisms occurring in the FARS2-LYRM4 bidirectional promoter region have been shown to be associated with cognitive deficit and schizophrenia [68]. Similarly, a 900 kb microdeletion in this region (encompassing RPP40, PPP1R3G, LYRM4, and part of the FARS2 and CDYL genes) has been shown to produce gyral pattern anomalies, intellectual disability and speech and language disorders [69].

A comparison of metabolite levels in the prefrontal cortex, visual cortex, cerebellum, kidney, and muscle between humans and other primates (chimpanzee and macaque), showed 'aminoacyl-tRNA biosynthesis' was the most significantly enriched pathway for metabolites having higher levels in human; for all three brain regions, but not muscle or kidney [70]. Similarly, the 'Phenylalanine, tyrosine and tryptophan biosynthesis' pathway was enriched for all three brain regions. This observation fits with the ongoing hypothesis that there is a requirement for the coevolution of nuclear and mitochondrial genes

coding for mitochondrial proteins, known as mitonuclear compensatory coevolution, as reviewed in [71], and that this coevolution is a driver of speciation [72,73]. In the mitonuclear compensatory coevolution hypothesis, nuclear-encoded genes encoding aminoacyl tRNA synthetases (such as FARS2) and OXPHOS complex components (such as LYRM4) are expected and, in some cases, have been shown to evolve more rapidly in response to high rates of evolutionary change in mitochondrially encoded genes, with which their protein products subsequently directly interact [71]. Another consideration is how to regulate the translation of OXPHOS complex components when the corresponding genes exist in both the mitochondrial and nuclear genomes [74,75]. The high number of observed DB TFBSs in the bidirectional promoter of the nuclear FARS2 and LYRM4 genes thus potentially allows more well-regulated coordination of cross-compartment OXPHOS component gene expression and may even contribute to the effects of both mitonuclear compensation and speciation.

Taken together, the LYRM4-FARS2 locus is potentially of great interest in the divergence of modern humans from other hominins and hominids. It is reasonable to imagine that differences in neuronal development and function, and even diet, between MH and Neanderthal could require corresponding differences in mitochondrial activity, metabolism, and amino acid sensing.

4.8. Future perspectives

The divergence of modern humans and Neanderthals is estimated to have occurred between 400,000 and 800,000 years ago [76–78]. Since then, both subspecies have experienced unique evolutionary, social, and cultural paths, many of which overlap and intertwine. What appears to be both divergent and unifying in the hominin lineage is a significant change in the brain, leading to differing abilities in terms of cognition, social function, language and creativity. In the case of humans, and likely for Neanderthals as well, many of the same genes, which are at the core of our novel capabilities, are the same which are commonly the root of our maladies. Among other neuropsychiatric disorders, autism and schizophrenia appear to be maintained in humans at a greater rate than would be expected if they were strictly deleterious: the global prevalence of ASD is approximately 1 % [79], and 0.28 % for schizophrenia [80]. Likewise, as the study of these diseases continues, the direction of understanding tends toward a continuum of effects rather than a discrete on/off state. A study by Linscott et al. reported that the prevalence of psychotic experiences in the general population was 7.2 % [81]. Our results revealed that DB TFs are enriched for development and the brain, and studies on HAR regions have shown that mutations in regulatory regions alone are sufficient to cause significant neuropsychiatric effects [34]. As a result, there appear to be significant opportunities to analyze gene regulation to understand not only the evolution of human cognitive abilities but also the neuropsychiatric disorders that accompany them.

4.9. Limitations

We have not included analyses of 3' UTRs, introns, or enhancers. The single-most expressed transcript for each protein-coding gene was chosen for analysis, while MH vs. Neanderthal SNPs may have occurred more closely to TSSs of lower-expressed transcripts and thereby had higher and perhaps more significant scoring. Only transcripts for protein-coding genes were analyzed, despite the existence of numerous other transcribed sequence types. The SNP dataset we used for analysis is based on multiple Neanderthal individuals and has lower sequencing coverage than newer datasets. The inclusion of multiple individuals is useful for comparing modern humans to Neanderthals as a group but can provide only a more general comparison. Using a higher coverage dataset would allow greater assurance that the SNPs under inquiry are legitimate.

CRedit authorship contribution statement

Barker Harlan: Writing – review & editing, Writing – original draft, Visualization, Validation, Software, Project administration, Methodology, Investigation, Funding acquisition, Formal analysis, Data curation, Conceptualization. **Seppo Parkkila:** Writing – review & editing, Supervision, Funding acquisition. **Tolvanen Martti:** Writing – review & editing, Supervision, Methodology, Conceptualization.

Funding

This work was supported by the Finnish Cultural Foundation and Fimlab to HB and the Academy of Finland and the Jane & Aatos Erkko Foundation to SP.

Declaration of Competing Interest

The authors declare that they have no conflicts of interest.

Acknowledgements

The nonprofit CSC—IT Center for Science Ltd., owned by the state of Finland and Finnish higher education institutions—is acknowledged for providing computational resources for analyses.

Appendix A. Supporting information

Supplementary data associated with this article can be found in the online version at doi:10.1016/j.csbj.2025.05.052.

Data Availability

Code reproducing the analyses and figures in the manuscript are available on GitHub at https://github.com/thirtysix/TFBS_footprinting_manuscript. The results of the MH vs. Neanderthal promoter analysis are available at <https://osf.io/r2mtw/>.

References

- [1] Green RE, Krause J, Briggs AW, Maricic T, Stenzel U, Kircher M, et al. Research article a draft sequence of the neanderthal genome. *Science* 2010;328:710–23. <https://doi.org/10.1126/science.1188021>.
- [2] Castellano S, Parra G, Sanchez-Quinto FA, Racimo F, Kuhlwil M, Kircher M, et al. Patterns of coding variation in the complete exomes of three Neanderthals. *Proc Natl Acad Sci* 2014;111:6666–71. <https://doi.org/10.1073/pnas.1405138111>.
- [3] Prüfer K, Hajdinjak M, Vernot B, Skov L, Hsieh P, Peyrégue S, et al. A high-coverage neanderthal genome from vindija cave in Croatia. *Science* 2017;358:655–8. (<http://science.sciencemag.org/>) (Available).
- [4] Hajdinjak M, Fu Q, Hubner A, Petr M, Mafessoni F, Grote S, et al. Reconstructing the genetic history of late Neanderthals. *Nature* 2018;555:652–6. <https://doi.org/10.1038/nature26151>.
- [5] Reich D, Green RE, Kircher M, Krause J, Patterson N, Durand EY, et al. Genetic history of an archaic hominin group from Denisova cave in Siberia. *Nature* 2010;468:1053–60. <https://doi.org/10.1038/nature09710>.
- [6] Meyer M, Kircher M, Gansauge MT, Li H, Racimo F, Mallick S, et al. A high-coverage genome sequence from an archaic Denisovan individual. *Science* 2012;338:222–6. <https://doi.org/10.1126/science.1224344>.
- [7] Prüfer K, Racimo F, Patterson N, Jay F, Sankararaman S, Sawyer S, et al. The complete genome sequence of a Neanderthal from the Altai Mountains. *Nature* 2014;505:43–9. <https://doi.org/10.1038/nature12886>.
- [8] Petr M, Paabo S, Kelso J, Vernot B. Limits of long-term selection against Neanderthal introgression. *Proc Natl Acad Sci USA* 2019;116:1639–44. <https://doi.org/10.1073/pnas.1814338116>.
- [9] Barker H, Aaltonen M, Pan P, Vähätupa M, Kaipainen P, May U, et al. Role of carbonic anhydrases in skin wound healing. *Exp Mol Med* 2017. <https://doi.org/10.1038/emmm.2017.60>.
- [10] Karjalainen SL, Haapasalo HK, Aspatwar A, Barker H, Parkkila S, Haapasalo JA. Carbonic anhydrase related protein expression in astrocytomas and oligodendroglial tumors. *BMC Cancer* 2018;18. <https://doi.org/10.1186/s12885-018-4493-4>.
- [11] Barker H, Parkkila S. Bioinformatic characterization of angiotensin-converting enzyme 2, the entry receptor for SARS-CoV-2. *PLoS One* 2020;15:e0240647. <https://doi.org/10.1371/journal.pone.0240647>.

- [12] Arppo A, Barker H, Parkkila S. Bioinformatic characterization of ENPEP, the gene encoding a potential cofactor for SARS-CoV-2 infection. *PLoS One* 2024;19: e0307731. <https://doi.org/10.1371/journal.pone.0307731>.
- [13] Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, et al. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat Methods* 2020;17:261–72. <https://doi.org/10.1038/s41592-019-0686-2>.
- [14] Hunter JD. Matplotlib: A 2D Graphics Environment. *Comput Sci Eng* 2007;9:90–5. <https://doi.org/10.1109/MCSE.2007.55>.
- [15] Waskom M, Gelbart M, Botvinnik O, Ostblom J, Hobson P, Lukauskas S, et al. *mwaskom/seaborn: v0.11.1*, 2020 (December 2020). [doi:10.5281/zenodo.4379347](https://doi.org/10.5281/zenodo.4379347).
- [16] Carbon S, Ireland A, Mungall CJ, Shu S, Marshall B, Lewis S, et al. AmiGO: online access to ontology and annotation data. *Bioinformatics* 2009;25:288–9. <https://doi.org/10.1093/bioinformatics/btm615>.
- [17] Mi H, Muruganujan A, Ebert D, Huang X, Thomas PD. PANTHER version 14: more genomes, a new PANTHER GO-slim and improvements in enrichment analysis tools. *Nucleic Acids Res* 2019;47:D419–26. <https://doi.org/10.1093/nar/gky1038>.
- [18] Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci USA* 2005;102:15545–50. <https://doi.org/10.1073/pnas.0506580102>.
- [19] Liberzon A, Subramanian A, Pinchback R, Thorvaldsdóttir H, Tamayo P, Mesirov JP. Molecular signatures database (MSigDB) 3.0. *Bioinformatics* 2011;27: 1739–40. <https://doi.org/10.1093/bioinformatics/btr260>.
- [20] Hodge RD, Bakken TE, Miller JA, Smith KA, Barkan ER, Grayback LT, et al. Conserved cell types with divergent features in human versus mouse cortex. *Nature* 2019;573:61–8. <https://doi.org/10.1038/s41586-019-1506-7>.
- [21] Wolf FA, Angerer P, Theis FJ. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol* 2018;19:15. <https://doi.org/10.1186/s13059-017-1382-0>.
- [22] Luecken MD, Theis FJ. Current best practices in single-cell RNA-seq analysis: a tutorial. *Mol Syst Biol* 2019;15:e8746. <https://doi.org/10.15252/msb.20188746>.
- [23] Chippada B. forceatlas2: Fastest Gephi's ForceAtlas2 graph layout algorithm implemented for Python and NetworkX. Github; Available: (<https://github.com/bhargavchippada/forceatlas2>).
- [24] Jacomy M, Venturini T, Heymann S, Bastian M. ForceAtlas2, a continuous graph layout algorithm for handy network visualization designed for the Gephi software. *PLoS One* 2014;9:e98679. <https://doi.org/10.1371/journal.pone.0098679>.
- [25] Chen EY, Tan CM, Kou Y, Duan Q, Wang Z, Meirelles GV, et al. Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. *BMC Bioinforma* 2013;14:128. <https://doi.org/10.1186/1471-2105-14-128>.
- [26] Kuleshov MV, Jones MR, Rouillard AD, Fernandez NF, Duan Q, Wang Z, et al. Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res* 2016;44:W90–7. <https://doi.org/10.1093/nar/gkw377>.
- [27] Pinero J, Bravo A, Queralt-Rosinach N, Gutierrez-Sacristan A, Deu-Pons J, Centeno E, et al. DisGeNET: a comprehensive platform integrating information on human disease-associated genes and variants. *Nucleic Acids Res* 2017;45:D833–9. <https://doi.org/10.1093/nar/gkw943>.
- [28] Velmeshev D, Schirmer L, Jung D, Haessler M, Perez Y, Mayer S, et al. Single-cell genomics identifies cell type-specific molecular changes in autism. *Science* 2019; 364:685–9. <https://doi.org/10.1126/science.aav8130>.
- [29] Türei D, Korcsmáros T, Saez-Rodriguez J. OmniPath: guidelines and gateway for literature-curated signaling pathway resources. *Nat Methods* 2016;13:966–7. <https://doi.org/10.1038/nmeth.4077>.
- [30] Türei D, Valdeolivas A, Gul L, Palacio-Escat N, Klein M, Ivanova O, et al. Integrated intra- and intercellular signaling knowledge for multicellular omics analysis. *Mol Syst Biol* 2021;17:e9923. <https://doi.org/10.15252/msb.20209923>.
- [31] Schmiedel BJ, Singh D, Madrigal A, Valdovino-Gonzalez AG, White BM, Zapardiel-Gonzalo J, et al. Impact Of Genetic Polymorphisms On Human Immune Cell Gene Expression. *Cell* 2018;175:1701–15. <https://doi.org/10.1016/j.cell.2018.10.022>.
- [32] Uhlen M, Fagerberg L, Hallstrom BM, Lindskog C, Oksvold P, Mardinoglu A, et al. Proteomics. Tissue-based map of the human proteome. *Science* 2015;347:1260419. <https://doi.org/10.1126/science.1260419>.
- [33] Ding SL, Royall JJ, Sunkin SM, Ng L, Facer BA, Lesnar P, et al. Comprehensive cellular-resolution atlas of the adult human brain. *J Comp Neurol* 2016;524: 3127–481. <https://doi.org/10.1002/cne.24080>.
- [34] Doan RN, Bae BI, Cubelos B, Chang C, Hossain AA, Al-Saad S, et al. Mutations in human accelerated regions disrupt cognition and social behavior. *Cell* 2016;167: 341–54. <https://doi.org/10.1016/j.cell.2016.08.071>.
- [35] Arendt D, Bertucci PY, Achim K, Musser JM. Evolution of neuronal types and families. *Curr Opin Neurobiol* 2019 Jun;56:144–52. <https://doi.org/10.1016/j.conb.2019.01.022>.
- [36] Banerjee-Basu S, Packer A. SFARI Gene: an evolving database for the autism research community. *Dis Model Mech* 2010;3:133–5. <https://doi.org/10.1242/dmm.005439>.
- [37] Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, et al. The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res* 2019;47: D1005–12. <https://doi.org/10.1093/nar/gky1120>.
- [38] Trubetskoy V, Pardiñas AF, Qi T, Panagiotaropoulou G, Awasthi S, Bigdeli TB, et al. Mapping genomic loci implicates genes and synaptic biology in schizophrenia. *Nature* 2022;604:502–8. <https://doi.org/10.1038/s41586-022-04434-5>.
- [39] Satterstrom FK, Kosmicki JA, Wang J, Breen MS, De Rubeis S, An J-Y, et al. Large-scale exome sequencing study implicates both developmental and functional changes in the neurobiology of autism. *Cell* 2020;180:568–584.e23. <https://doi.org/10.1016/j.cell.2019.12.036>.
- [40] Paciorkowski AR, Traylor RN, Rosenfeld JA, Hoover JM, Harris CJ, Winter S, et al. MEF2C Haploinsufficiency features consistent hyperkinesia, variable epilepsy, and has a role in dorsal and ventral neuronal developmental pathways. *Neurogenetics* 2013;14:99–111. <https://doi.org/10.1007/s10048-013-0356-y>.
- [41] Zhou WZ, Zhang J, Li Z, Lin X, Li J, Wang S, et al. Targeted resequencing of 358 candidate genes for autism spectrum disorder in a Chinese cohort reveals diagnostic potential and genotype-phenotype correlations. *Hum Mutat* 2019;40: 801–15. <https://doi.org/10.1002/humu.23724>.
- [42] Nowakowska BA, Obersztyn E, Szymanska K, Bekiesinska-Figatowska M, Xia Z, Ricks CB, et al. Severe mental retardation, seizures, and hypotonia due to deletions of MEF2C. *Am J Med Genet B Neuropsychiatr Genet* 2010;153B:1042–51. <https://doi.org/10.1002/ajmg.b.31071>.
- [43] Schizophrenia Working Group of the Psychiatric Genomics Consortium. Biological insights from 108 schizophrenia-associated genetic loci. *Nature* 2014;511:421–7. <https://doi.org/10.1038/nature13595>.
- [44] Mitchell AC, Javidfar B, Pothula V, Ibi D, Shen EY, Peter CJ, et al. MEF2C transcription factor is associated with the genetic and epigenetic risk architecture of schizophrenia and improves cognition in mice. *Mol Psychiatry* 2018;23:123–32. <https://doi.org/10.1038/mp.2016.254>.
- [45] Lai CS, Fisher SE, Hurst JA, Vargha-Khadem F, Monaco AP. A forkhead-domain gene is mutated in a severe speech and language disorder. *Nature* 2001;413: 519–23. <https://doi.org/10.1038/35097076>.
- [46] Chen YC, Kuo HY, Bornschein U, Takahashi H, Chen SY, Lu KM, et al. Foxp2 controls synaptic wiring of corticostriatal circuits and vocal communication by opposing Mef2c. *Nat Neurosci* 2016;19:1513–22. <https://doi.org/10.1038/nn.4380>.
- [47] Maricic T, Gunther V, Georgiev O, Gehre S, Curlin M, Schreiweis C, et al. A recent evolutionary change affects a regulatory element in the human FOXP2 gene. *Mol Biol Evol* 2013;30:844–52. <https://doi.org/10.1093/molbev/mss271>.
- [48] Atkinson EG, Audesse AJ, Palacios JA, Bobo DM, Webb AE, Ramachandran S, et al. No evidence for recent selection at FOXP2 among diverse human populations. *Cell* 2018;174:1424–35. <https://doi.org/10.1016/j.cell.2018.06.048>.
- [49] Ma S, Skarica M, Li Q, Xu C, Risgaard RD, Tebbenkamp ATN, et al. Molecular and cellular evolution of the primate dorsolateral prefrontal cortex. *Science* 2022; eabo7257. <https://doi.org/10.1126/science.abo7257>.
- [50] Gupta S, Ellis SE, Ashar FN, Moes A, Bader JS, Zhan J, et al. Transcriptome analysis reveals dysregulation of innate immune response genes and neuronal activity-dependent genes in autism. *Nat Commun* 2014;5:5748. <https://doi.org/10.1038/ncomms6748>.
- [51] Gandal MJ, Haney JR, Wamsley B, Yap CX, Parhami S, Emami PS, et al. Broad transcriptomic dysregulation occurs across the cerebral cortex in ASD. *Nature* 2022;611:532–9. <https://doi.org/10.1038/s41586-022-05377-7>.
- [52] Gao Y, Dong Q, Arachchilage KH, Risgaard RD, Syed M, Sheng J, et al. Multimodal analyses reveal genes driving electrophysiological maturation of neurons in the primate prefrontal cortex. *cit ed 28 May 2025*. *Neuron* 2025. <https://doi.org/10.1016/j.neuron.2025.04.025>.
- [53] Ohi K, Hashimoto R, Yasuda Y, Fukumoto M, Yamamori H, Umeda-Yano S, et al. Influence of the NRGN gene on intellectual ability in schizophrenia. *J Hum Genet* 2013;58:700–5. <https://doi.org/10.1038/jhg.2013.82>.
- [54] Casanova EL, Casanova MF. Genetics studies indicate that neural induction and early neuronal maturation are disturbed in autism. *Front Cell Neurosci* 2014;8:397. <https://doi.org/10.3389/fncel.2014.00397>.
- [55] Greenberg ME, Greene LA, Ziff EB. Nerve growth factor and epidermal growth factor induce rapid transient changes in proto-oncogene transcription in PC12 cells. *J Biol Chem* 1985;260:14101–10. [https://doi.org/10.1016/s0021-9258\(17\)38689-1](https://doi.org/10.1016/s0021-9258(17)38689-1).
- [56] Gallo FT, Katche C, Morici JF, Medina JH, Weisstaub NV. Immediate Early Genes, memory and psychiatric disorders: focus on c-Fos, Egr1 and arc. *Front Behav Neurosci* 2018;12:79. <https://doi.org/10.3389/fnbeh.2018.00079>.
- [57] Sissler M, González-Serrano LE, Westhof E. Recent advances in mitochondrial aminoacyl-tRNA synthetases and disease. *Trends Mol Med* 2017;23:693–708. <https://doi.org/10.1016/j.molmed.2017.06.002>.
- [58] Fine AS, Nemeth CL, Kaufman ML, Fatemi A. Mitochondrial aminoacyl-tRNA synthetase disorders: an emerging group of developmental disorders of myelination. *J Neurodev Disord* 2019;11:29. <https://doi.org/10.1186/s11689-019-9292-y>.
- [59] Guo M, Yang X-L, Schimmel P. New functions of aminoacyl-tRNA synthetases beyond translation. *Nat Rev Mol Cell Biol* 2010;11:668–74. <https://doi.org/10.1038/nrm2956>.
- [60] Barcia G, Rio M, Assouline Z, Zangarelli C, Roux C-J, de Lonlay P, et al. Novel FARS2 variants in patients with early onset encephalopathy with or without epilepsy associated with long survival. *Eur J Hum Genet* 2021;29:533–8. <https://doi.org/10.1038/s41431-020-00757-x>.
- [61] Duan J, Sanders AR, Moy W, Drigalenko EI, Brown EC, Freda J, et al. Transcriptome outlier analysis implicates schizophrenia susceptibility genes and enriches putatively functional rare genetic variants. *Hum Mol Genet* 2015;24: 4674–85. <https://doi.org/10.1093/hmg/ddv199>.
- [62] Shi Y, Ghosh MC, Tong W-H, Rouault TA. Human ISD11 is essential for both iron-sulfur cluster assembly and maintenance of normal cellular iron homeostasis. *Hum Mol Genet* 2009;18:3014–25. <https://doi.org/10.1093/hmg/ddp239>.
- [63] Schmucker S, Martelli A, Colin F, Page A, Wattenhofer-Donzé M, Reutenauer L, et al. Mammalian frataxin: an essential function for cellular viability through an interaction with a preformed ISCU/NFS1/ISD11 iron-sulfur assembly complex. *PLoS One* 2011;6:e16199. <https://doi.org/10.1371/journal.pone.0016199>.

- [64] Mayr JA, Haack TB, Freisinger P, Karall D, Makowski C, Koch J, et al. Spectrum of combined respiratory chain defects. *J Inher Metab Dis* 2015;38:629–40. <https://doi.org/10.1007/s10545-015-9831-y>.
- [65] Lill R. Function and biogenesis of iron-sulphur proteins. *Nature* 2009;460:831–8. <https://doi.org/10.1038/nature08301>.
- [66] Lim SC, Friemel M, Marum JE, Tucker EJ, Bruno DL, Riley LG, et al. Mutations in LYRM4, encoding iron-sulfur cluster biogenesis factor ISD11, cause deficiency of multiple respiratory chain complexes. *Hum Mol Genet* 2013;22:4460–73. <https://doi.org/10.1093/hmg/ddt295>.
- [67] Bergman O, Ben-Shachar D. Mitochondrial oxidative phosphorylation system (OXPHOS) deficits in schizophrenia: possible interactions with cellular processes. *Can J Psychiatry* 2016;61:457–69. <https://doi.org/10.1177/0706743716648290>.
- [68] Jablensky A, Angelicheva D, Donohoe GJ, Cruickshank M, Azmanov DN, Morris DW, et al. Promoter polymorphisms in two overlapping 6p25 genes implicate mitochondrial proteins in cognitive deficit in schizophrenia. *Mol Psychiatry* 2012;17:1328–39. <https://doi.org/10.1038/mp.2011.129>.
- [69] Bozza M, Bernardini L, Novelli A, Bovedani P, Moretti E, Canapicchi R, et al. 6p25 interstitial deletion in two dizygotic twins with gyral pattern anomaly and speech and language disorder. *Eur J Paediatr Neurol* 2013;17:225–31. <https://doi.org/10.1016/j.ejpn.2012.09.008>.
- [70] Stepanova V, Moczulska KE, Vacano GN, Kurochkin I, Ju X, Riesenberger S, et al. Reduced purine biosynthesis in humans after their divergence from Neandertals. *Elife* 2021;10. <https://doi.org/10.7554/eLife.58741>.
- [71] Hill GE. Mitonuclear Compensatory Coevolution. *Trends Genet* 2020;36:403–14. <https://doi.org/10.1016/j.tig.2020.03.002>.
- [72] Hill GE. Mitonuclear coevolution as the genesis of speciation and the mitochondrial DNA barcode gap. *Ecol Evol* 2016;6:5831–42. <https://doi.org/10.1002/ece3.2338>.
- [73] Tobler M, Barts N, Greenway R. Mitochondria and the origin of species: bridging genetic and ecological perspectives on speciation processes. *Integr Comp Biol* 2019;59:900–11. <https://doi.org/10.1093/icb/icz025>.
- [74] Soto I, Couvillion M, Hansen KG, McShane E, Moran JC, Barrientos A, et al. Balanced mitochondrial and cytosolic translomes underlie the biogenesis of human respiratory complexes. *Genome Biol* 2022;23:170. <https://doi.org/10.1186/s13059-022-02732-9>.
- [75] Couvillion MT, Soto IC, Shipkovenska G, Churchman LS. Synchronized mitochondrial and cytosolic translation programs. *Nature* 2016;533:499–503. <https://doi.org/10.1038/nature18015>.
- [76] Mendez FL, Poznik GD, Castellano S, Bustamante CD. The divergence of neandertal and modern human Y chromosomes. *Am J Hum Genet* 2016;98:728–34. <https://doi.org/10.1016/j.ajhg.2016.02.023>.
- [77] Endicott P, Ho SYW, Stringer C. Using genetic evidence to evaluate four palaeoanthropological hypotheses for the timing of Neanderthal and modern human origins. *J Hum Evol* 2010;59:87–95. <https://doi.org/10.1016/j.jhevol.2010.04.005>.
- [78] Gómez-Robles A. Dental evolutionary rates and its implications for the Neanderthal-modern human divergence. *Sci Adv* 2019;5:eaaw1268. <https://doi.org/10.1126/sciadv.aaw1268>.
- [79] Zeidan J, Fombonne E, Scorah J, Ibrahim A, Durkin MS, Saxena S, et al. Global prevalence of autism: a systematic review update. *Autism Res* 2022;15:778–90. <https://doi.org/10.1002/aur.2696>.
- [80] Charlson FJ, Ferrari AJ, Santomauro DF, Diminic S, Stockings E, Scott JG, et al. Global epidemiology and burden of schizophrenia: findings from the global burden of disease study 2016. *Schizophr Bull* 2018;44:1195–203. <https://doi.org/10.1093/schbul/sby058>.
- [81] Linscott RJ, van Os J. An updated and conservative systematic review and meta-analysis of epidemiological evidence on psychotic experiences in children and adults: on the pathway from proneness to persistence to dimensional expression across mental disorders. *Psychol Med* 2013;43:1133–49. <https://doi.org/10.1017/S0033291712001626>.