

Genome-wide Identification of the Genetic Basis of Amyotrophic Lateral Sclerosis

Sai Zhang^{1,2,*}, Johnathan Cooper-Knock^{3,*}, Annika K. Weimer^{1,2}, Minyi Shi^{1,2},
Tobias Moll³, Calum Harvey³, Helia Ghahremani Nezhad³, John Franklin³,
Cleide dos Santos Souza³, Cheng Wang^{4,5,6,7}, Jingjing Li^{4,5,6,7}, Chen Eitan⁸,
Eran Hornstein⁸, Kevin P. Kenna⁹, Project MinE Sequencing Consortium, Jan
Veldink⁹, Laura Ferraiuolo³, Pamela J. Shaw³ and Michael P. Snyder^{1,2}

¹Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA

²Center for Genomics and Personalized Medicine, Stanford University School of Medicine, Stanford, CA, USA

³Sheffield Institute for Translational Neuroscience, University of Sheffield, Sheffield, UK

⁴Department of Neurology, School of Medicine, University of California, San Francisco, CA, USA

⁵Eli and Edythe Broad Center of Regeneration Medicine and Stem Cell Research, University of California, San Francisco, CA, USA

⁶Bakar Computational Health Sciences Institute, University of California, San Francisco, CA, USA

⁷Parker Institute for Cancer Immunotherapy, University of California, San Francisco, CA, USA

⁸Department of Molecular Genetics, Weizmann Institute of Science, Rehovot, Israel

⁹Department of Neurology, Brain Center Rudolf Magnus, University Medical Center Utrecht, Utrecht, The Netherlands

*These authors contributed equally to this work.

Correspondence should be addressed to M.P.S. (mpsnyder@stanford.edu)

ABSTRACT

Amyotrophic lateral sclerosis (ALS) is an archetypal complex disease centered on progressive death of motor neurons. Despite heritability estimates of 52%, GWAS studies have discovered only seven genome-wide significant hits, which are relevant to <10% of ALS patients. To increase the power of gene discovery, we integrated motor neuron functional genomics with ALS genetics in a hierarchical Bayesian model called RefMap. Comprehensive transcriptomic and epigenetic profiling of iPSC-derived motor neurons enabled RefMap to systematically fine-map genes and pathways associated with ALS. As a significant extension of the known genetic architecture of ALS, we identified a group of 690 candidate ALS genes, which is enriched with previously discovered risk genes. Extensive conservation, transcriptome and network analyses demonstrated the functional significance of these candidate genes in motor neurons and disease progression. In particular, we observed a genetic convergence on the distal axon, which supports the prevailing view of ALS as a distal axonopathy. Of the new ALS genes we discovered, we further characterized *KANK1* that is enriched with coding and noncoding, common and rare ALS-associated genetic variation. Modelling patient mutations in human neurons reduced *KANK1* expression and produced neurotoxicity with disruption of the distal axon. RefMap can be applied broadly to increase the discovery power in genetic association studies of human complex traits and diseases.

INTRODUCTION

ALS is an untreatable, universally fatal and relatively common neurodegenerative disease with a lifetime risk of $\sim 1/350$ in the UK. The hallmark of the disease is motor neuron loss leading to respiratory failure and death¹. 10% of ALS is autosomal dominant, and even for sporadic ALS (sALS), the heritability is estimated to be $\sim 50\%$ ^{2,3}. Genome-wide association studies (GWAS) in ALS^{4,5} have identified seven genome-wide significant loci, which have been linked to missense mutations. However, these changes occur in $<10\%$ of ALS patients, so there are likely to be a large number of missing ALS risk genes.

ALS GWAS studies to date have lost power by considering genetic variants in isolation, whereas in reality, a biological system is the product of a large number of interacting partners^{6,7}. Moreover, noncoding regulatory regions of the genome have been relatively neglected in efforts to pinpoint the genetic basis of ALS, despite their functional synergy with the coding sequence^{8,9}. Indeed, GWAS studies have suggested that a significant proportion of missing heritability in ALS is distributed throughout noncoding chromosomal regions^{4,5}. The function of noncoding DNA is often tissue, disease, or even cell-type specific¹⁰, and the understanding of the cell-type-specific biological function in complex neurological diseases has been improving^{11,12}. This therefore creates an opportunity to dramatically reduce the search space and so boost the power to discover ALS genetic risk, by focusing on genomic regions that are functional within the cell type of interest, i.e., motor neurons (MNs)¹³.

Here, we present RefMap (**R**egional **F**ine-**m**apping), a hierarchical Bayesian model to perform genome-wide identification of disease-associated genetic variation within active genomic regions. RefMap utilizes cell-type-specific epigenetic profiling to determine the prior probability of disease-association for each region. This reduces the search space by $>90\%$ given that a limited proportion of the genome is active in any specific cell type. ALS is notable for the selective vulnerability of MNs¹³. However, MNs are difficult to study in post-mortem tissues (e.g.¹⁴) because of their relative sparsity, so a different

approach is needed. We performed exhaustive transcriptomic and epigenetic profiling, including RNA-seq, ATAC-seq, histone ChIP-seq and Hi-C, for motor neurons derived from fibroblasts of neurologically normal controls. We hypothesized that the genetic variation within regulatory regions may alter the expression of their target genes, and we proposed that disease-associated variants are likely to reduce gene expression via interfering with regulation. Applying RefMap to perform genome-wide fine-mapping based on ALS GWAS data (**Fig. 1a**) identified 690 ALS-associated genes, including previous GWAS hits and even known ALS genes not previously detected in GWAS studies.

We explored the functional significance of RefMap ALS genes based on a series of orthogonal analyses. Population genetics revealed that RefMap genes consist of conserved sequences, suggesting that their functions are important and not subject to genetic redundancy. Transcriptome data from MNs, human tissues and mouse models demonstrated that RefMap genes are down-regulated in ALS patients, consistent with our aforementioned hypothesis. Network analysis of protein-protein interactions (PPIs) identified two modules enriched with RefMap genes. These modules are enriched with biological functions localized to the distal axon of MNs, suggesting that neurotoxicity may be initiated in this subcompartment, which is consistent with previous literature^{15,16}. Finally, we have further characterized a new ALS gene, i.e., *KANK1*. Common and rare genetic variants that alter *KANK1* expression were shown to be associated with ALS and neuronal toxicity. RefMap provides a promising framework to pinpoint the genetic bases of human complex traits and diseases based on GWAS data.

RESULTS

Transcriptomic and epigenetic profiling of iPSC-derived motor neurons

To identify genomic regions key to motor neuron function, we performed transcriptomic and epigenetic profiling of iPSC-derived motor neurons from neurologically normal individuals (**Supplementary Fig. 1**). The cells exhibited homogenous expression of the lower motor neuron markers, including TUJ1, Chat, SMI, MAP2 and NeuN

(Supplementary Fig. 1a). We prepared RNA-seq¹⁷, ATAC-seq¹⁸, H3K27ac, H3K4me1 and H3K4me3 ChIP-seq¹⁹, as well as Hi-C²⁰ libraries using two technical replicates and three biological replicates per assay. Sequencing data were processed and quality control (QC) was performed according to the ENCODE 4 standards²¹, and all samples exceeded ENCODE standard QC measures (**Supplementary Tables 1-4**).

ATAC-seq identifies open and functional chromatin regions, which is complementary to the profiling of transcript expression by RNA-seq. H3K27ac, H3K4me1 and H3K4me3 ChIP-seq assays pinpoint active enhancer²² regions, which are important noncoding regions for the regulation of gene expression. Hi-C profiling of three-dimensional (3D) genome structure is essential to map regulatory regions including enhancers, to their target genes. Our MN epigenetic profiling successfully reduced the search space for ALS-associated genetic variation by >90%. Specifically, total ATAC-seq peak regions across all biological replicates covered 4.9% of the genome.

To measure the consistency between distinct motor neuron profiles, we used our RNA-seq dataset to identify promoter regions for high (>90th centile) and low (<10th centile) expressed transcripts. We compared enrichment of ATAC-seq and histone ChIP-seq peak regions, and Hi-C loops in high versus low expressed promoters. Significant enrichment within highly expressed promoters was confirmed for ATAC-seq ($P=1.1e-182$, odds ratio (OR)=1.9, Fisher's exact test), H3K27ac ChIP-seq ($P=2.0e-57$, OR=2.2, Fisher's exact test), H3K4me1 ChIP-seq ($P=8.5e-57$, OR=1.9, Fisher's exact test), H3K4me3 ChIP-seq ($P=4.8e-196$, OR=2.6, Fisher's exact test), and Hi-C loops ($P=4.0e-14$, OR=1.3, Fisher's exact test) (**Fig. 1b**). Similarly, epigenetic peak regions were enriched in MN Hi-C loops: ATAC-seq ($P<1.0e-198$, OR=1.9, Fisher's exact test), H3K27ac ChIP-seq ($P<1.0e-198$, OR=2.0, Fisher's exact test), H3K4me1 ChIP-seq ($P<1.0e-198$, OR=2.0, Fisher's exact test), and H3K4me3 ChIP-seq ($P<1.0e-198$, OR=1.7, Fisher's exact test). These observations confirm that our epigenetic profiling captured functionally significant genomic variation, and that our epigenetic profiles were internally consistent.

RefMap identifies ALS risk genes

Mismatch between the relatively small number of characterized ALS risk genes and the estimate of high heritability suggests that a new approach is required to discover more ALS-associated genetic variation. Here, we designed a hierarchical Bayesian network named RefMap that exploits the epigenetic profiling of MNs to reduce the search space and so improve the statistical power to discover ALS-associated loci across the genome. Specifically, RefMap integrates the prior probability of significance derived from the epigenome of MNs, with allele effect sizes estimated from GWAS (**Figs. 1a and 1c, Methods**). Based on a linear genotype-phenotype model (**Supplementary Notes**), RefMap first disentangles effect sizes from GWAS Z-scores, which are confounded by the structure of linkage disequilibrium (LD). Effect sizes are then summarized across genomic regions in individual LD blocks. Those regions that are within active chromatin, and where the distributions of allele effect sizes are shifted from the null distribution, are prioritized by the algorithm (**Methods**).

In our study, the Z-scores were calculated based on the largest published ALS GWAS study^{4,5}, including genotyping of 12,577 sporadic ALS patients and 23,475 controls. An epigenetic signal was calculated from a linear combination of MN chromatin accessibility and histone marks specific to active enhancer regions (**Methods**). We defined LD blocks as 1Mb windows, where we assumed significant internal LD but negligible external LD²³. Within LD blocks, SNP correlations were estimated based on the European population (EUR) data from the 1000 Genomes Project²⁴. With this information, RefMap scanned the genome in 1kb windows and identified all regions that are likely to harbor ALS-associated genetic variation (**Figs. 1c and 1d, Methods, and Supplementary Table 5**).

Next, we mapped ALS-associated regions identified by RefMap to expressed transcripts in MNs (RNA-seq, TPM \geq 1), based on their regulation targets. We defined regulation targets as genes that overlap either ALS-associated regions by extension, or via their Hi-C loop anchors (**Methods**). This resulted in 690 ALS-associated genes

(**Supplementary Table 6**). Among this list, we discovered well-known ALS genes, including *C9orf72*²⁵ and *ATXN2*²⁶ (**Fig. 1d**). Indeed, RefMap genes are enriched with an independently curated list (**Supplementary Table 7**) of ALS genes including previous GWAS hits ($P=5.20e-3$, OR=2.07, Fisher's exact test) and also with clinically reportable (ClinVar²⁷) ALS genes ($P=0.03$, OR=3.06, Fisher's exact test). Interestingly, certain ALS genes, such as *UNC13A*^{28,29}, are missing from RefMap genes, but their paralogues are present, including *UNC13B*, which is consistent with a functional overlap. If we consider paralogues as equivalent to ALS genes, then the enrichment of RefMap genes with known ALS genes is further increased (curated: $P=6.12e-43$, OR=8.71; ClinVar: $P=6.40e-14$, OR=12.26; Fisher's exact test).

As a negative control, we randomly shuffled SNP Z-scores, in which case there was no overlap between RefMap outputs and known ALS genes. Additional shuffling of epigenetic features disrupted the signal further such that there were no significant RefMap outputs. This illustrates the dependence of RefMap on the two primary inputs: GWAS Z-scores and MN epigenetic features.

Conservation analysis demonstrates the functional importance of RefMap genes

A large proportion of RefMap ALS genes were identified because of ALS-associated genetic variation within noncoding regulatory regions. We hypothesized that the functional consequence of pathogenic genetic variation within regulatory regions is likely to be reduced expression of the target genes. A conservation analysis was first carried out, revealing that change in the expression of RefMap ALS genes is likely to be pathogenic based on population genetics.

Conservation refers to DNA sequences that are preserved in the population presumably because disruption would be deleterious. Conservation can be quantified by the haploinsufficiency (HI) score, which is a measure of functional similarity to known haplosufficient and haploinsufficient genes³⁰. Conservation is also related to intolerance scores, in which the rate of observed mutation of a gene in the population is compared to the expected rate in the absence of negative selection^{31,32,33}. In particular, a lower

than expected mutation rate implies intolerance to mutation. We discovered that RefMap genes are significantly haploinsufficient based on their HI score ($P=2.59e-19$, one-sided Wilcoxon rank-sum test; **Fig. 2a**), and intolerant to loss of function mutations within the Exome Aggregation Consortium (ExAC) dataset³⁴ (LoFtool score³¹: $P=2.28e-4$, one-sided Wilcoxon rank-sum test; **Fig. 2b**). They are also intolerant to other mutation types (RVIS score³²: $P=8.08e-13$, one-sided Wilcoxon rank-sum test; **Fig. 2c**), as well as within the larger gnomAD (v.2.1) dataset (o/e score³³: $P=4.08e-10$, one-sided Wilcoxon rank-sum test; **Fig. 2d**). Taken together, these results support the functional significance of RefMap ALS genes.

Transcriptome analysis supports functional significance of RefMap genes in motor neurons and in ALS

We have hypothesized that the ALS-associated genetic variation identified by RefMap is likely to be pathogenic through altered expression of the 690 RefMap genes. We have also demonstrated, based on population genetics, that the function of RefMap genes is highly sensitive to changes in expression. To explore this possibility further, we examined whether change in the expression of RefMap genes is associated with ALS, using transcriptome data from patient-derived MNs, central nervous system (CNS) tissues and an ALS animal model.

First, we inspected the expression of RefMap genes in our iPSC-derived MNs from neurologically normal individuals. RefMap genes were upregulated ($P=3.07e-17$, one-sided Wilcoxon rank-sum test; **Fig. 3a**) compared to the overall transcriptome, indicating their importance in normal MN function. No differential expression was observed for genes derived from RefMap using randomly shuffled Z-scores.

Next, we examined the expression of RefMap ALS genes in CNS tissues derived from ALS patients ($n=18$) and controls ($n=17$)³⁵. We hypothesized that RefMap genes would be downregulated in ALS patient tissues. As expected, a significant decrease in the expression of RefMap genes was observed in both frontal cortex (*C9orf72*-ALS (cALS): false discovery rate (FDR)=0.002, one-sided Wilcoxon rank-sum test) and cerebellum

(*C9orf72*-ALS: FDR=0.002; sporadic ALS: FDR=0.005) of ALS patients compared to the overall transcriptome (**Fig. 3b**). As an independent validation, we analyzed gene expression within iPSC-derived MNs from ALS patients (n=55, <https://www.answerals.org/>), and confirmed that RefMap genes were downregulated ($P=3.85e-04$, one-sided Wilcoxon rank-sum test; **Fig. 3c**) compared to neurologically normal controls (n=15).

Finally, we used longitudinal data to infer whether changes in expression of RefMap ALS genes occur upstream or downstream in the development of neuronal toxicity. To achieve this, we utilized the *SOD1*-G93A-ALS mouse model, which is the best characterized ALS model to date³⁶ and the only model featuring consistent and reproducible loss of spinal cord MNs that mirrors the human disease. We examined longitudinal gene expression averaged across spinal cord sections from *SOD1*-G93A (n=32) and *SOD1*-WT (n=24) mice³⁷. Four time points were sampled, including presymptomatic (p30), onset (p70), symptomatic (p100) and end-stage (p120). The model-estimated expression levels (β)³⁷ were adopted to quantify the gene expression difference ($\Delta\beta$) between diseased and control mice at different time points. To determine the expression changes of RefMap genes over the course of ALS pathogenesis, we first mapped RefMap genes to their mouse homologs (n=510), and then performed unsupervised clustering on gene expressions over time. We identified two different expression patterns for RefMap homologs (**Figs. 3d** and **3e**) with verified clustering quality (**Supplementary Fig. 2**). Strikingly, the largest group (286/510) of RefMap homologs were progressively downregulated through consecutive disease stages (C1; **Figs. 3d** and **3e**, **Supplementary Table 8**), consistent with our human observations. Functional enrichment analysis³⁸ of C1 genes revealed significant enrichment with functions associated with motor neuron biology (**Fig. 3f**), including 'cholinergic synapse', 'axon' and 'cytoskeleton', which is consistent with known ALS biology¹³ and with the prevailing view of ALS as a distal axonopathy^{15,16}. C2 genes do not contain significant functional enrichment (data not shown).

Systems analysis dissect ALS-associated functional modules

We have used RefMap to extend the number of ALS-associated risk genes to 690. We aimed to assess whether these genes are functionally consistent with current knowledge regarding the biology of MNs and ALS. Genes do not function in isolation and therefore, rather than examining individual genes, we mapped RefMap ALS genes to the global protein-protein interaction (PPI) network and inspected functional enrichment of ALS-associated network modules.

We first extracted high-confidence (combined score >700) PPIs from STRING v11.0³⁹, which include 17,161 proteins and 839,522 protein interactions. To eliminate the bias of hub genes⁴⁰, we performed the random walk with restart algorithm over the raw PPI network to construct a smoothed network based on those edges with weights in the top 5% (**Supplementary Table 9, Methods**). Next, this smoothed PPI network was decomposed into non-overlapping subnetworks using the Louvain algorithm⁴¹ that maximizes the modularity to detect communities from a network. This process yielded 912 different modules (**Supplementary Table 10**), in which genes within modules were densely connected with each other but sparsely connected with genes in other modules. As a negative control, we constructed 100 shuffled networks by randomly rewiring the PPI network while keeping the same number of neighbors. None of the randomized networks achieved the same modularity of our smoothed network after clustering, demonstrating the significance of our derived gene modules ($P<0.01$; **Supplementary Fig. 3a**).

RefMap ALS genes were then mapped to individual modules, and two modules were found to be significantly enriched with RefMap genes: M421 (721 genes; FDR<0.1, hypergeometric test; **Fig. 4a**) and M604 (308 genes; FDR<0.1, hypergeometric test; **Fig. 4b**) (**Supplementary Table 10**). Functionally M421 is enriched with GO/KEGG terms related to the distal axon, including synapse and axonal function within motor neurons (**Fig. 4c**). M421 is also enriched with genes related to relevant neurodegenerative diseases, including ‘amyotrophic lateral sclerosis’ and ‘Alzheimer’s

disease'. M604 is enriched with GO/KEGG terms related to the actin cytoskeleton and axonal function (**Fig. 4d**). Notably, the actin cytoskeleton is key for neuronal function and for axonal function in particular. Overall, the functional enrichment of both modules highlights an important role of the distal axon in ALS etiology (**Fig. 4e**), which is consistent with previous literature^{15,16}. Finally, both M421 and M604 were overexpressed in control iPSC-derived MNs (**Fig. 4f**), in a similar manner to the total set of RefMap genes. Interestingly, many functions ascribed to M421 and M604 overlap with the functions of the C1 cluster from our analysis of the *SOD1*-G93A mouse model (**Fig. 3f**), demonstrating a functional convergence of RefMap ALS genes.

Rare variant burden analysis is consistent with KANK1 as a novel ALS risk gene

Among all ALS-associated active regions identified by RefMap, chr9:663,001-664,000 has the highest concentration of ALS risk SNPs (22 SNPs). This region lies within intron 2 of *KANK1* and consists of independently annotated ENSEMBL regulatory features, including an enhancer element (ENSR00000873709) and a CTCF binding site (ENSR00000873710) (**Fig. 5a**). Overlap with independently annotated features supports the utility of RefMap to identify functional regulatory regions within noncoding DNA. We hypothesized that ALS-associated genetic variation within chr9:663,001-664,000 would reduce the expression of *KANK1*, leading to MN toxicity. Existing biological characterization of *KANK1* is consistent with our hypothesis: *KANK1* is expressed in motor neurons, functions in actin polymerization and deletion of this gene results in a severe developmental phenotype with MN loss⁴².

If reduced expression of *KANK1* is linked to MN toxicity, then it is reasonable to expect other loss-of-function (LoF) *KANK1* mutations to be associated with an increased risk of ALS. Thus far RefMap has utilized common genetic variants from a GWAS study⁴ so, to further investigate *KANK1*, here we performed rare variant burden tests. Rare variant analysis utilized whole-genome sequencing (WGS) data from 5,594 sporadic ALS patients and 2,238 controls⁴³. We filtered for rare, deleterious variants within *KANK1* enhancer, promoter and coding regions based on evolutionary conservation, functional

annotations and population frequency^{33,44–46} (**Methods**). Enhancer and promoter regions for *KANK1* were defined as previously described^{8,47}. Enhancer and promoter regions were independently enriched with ALS-associated rare deleterious variants ($P < 0.05$, SKAT^{48,49}; **Fig. 5b**), and nonsense coding variants were absent from controls and present in a small number ($n=4$) of ALS patients. Across all three regions, there was significant enrichment of rare deleterious variants in ALS patients compared to controls ($P=0.003$, Stouffer's method⁵⁰; **Fig. 5b**). The observation of both rare and common ALS-associated genetic variation in independent datasets utilizing independent methodology strongly suggests *KANK1* is a new ALS risk gene.

KANK1 was located within a distinct module (M826, 687 genes; **Supplementary Fig. 3b**) in our network analysis, and this module is enriched with RefMap genes ($P=5.6e-3$, hypergeometric test), but not after multiple testing correction. Functionally the *KANK1*-module is highly expressed in normal MNs (**Fig. 4f**), and is enriched for biological functions centered on the distal axon and synapse (**Supplementary Fig. 3c**), which are consistent with other RefMap-enriched modules.

Experimental validation of KANK1 in ALS development

To further investigate the role of *KANK1* in ALS, we experimentally determined the effect of ALS-associated genetic variation on gene expression and neuronal health (**Fig. 5c**). We used CRISPR/SpCas9 editing of SH-SY5Y neurons to recapitulate ALS-associated regulatory and coding mutations.

We discovered a high density of ALS-associated genetic variants within a region at chr9:663001-664000, which also contains an independently validated enhancer element (**Fig. 5a**). To replicate disruption of this sequence, we designed gRNAs to target protospacer adjacent motif (PAM) sites up- and downstream so as to delete the entire region⁵¹ (**Methods**). In addition, our rare variant analysis identified ALS-associated nonsense mutations in ALS cases but not in controls, therefore we also targeted a PAM site within *KANK1* exon 2 so as to introduce a series of indels (**Methods**). Sanger sequencing and waveform decomposition analysis⁵² in undifferentiated SH-SY5Y cells

confirmed the exon 2 editing efficiency (**Supplementary Figs. 4a** and **4b**) and the deletion of the enhancer sequence (**Supplementary Fig. 4c**). For experimental evaluation, a commercially available control gRNA targeting HPRT served as a negative control. CRISPR/SpCas9-edited SH-SY5Y cells were differentiated to a neuronal phenotype, and successful differentiation was confirmed by altered expression of PAX6 (**Supplementary Fig. 4d** and⁵³) and increased total dendritic length ($P=0.046$, paired Student's *t*-test; **Supplementary Fig. 4e** and⁵³). Differentiated cells were harvested and RNA was extracted for qPCR. We confirmed the reduced expression of *KANK1* mRNA in both exon and enhancer edited neurons (**Supplementary Fig. 4f**). Furthermore, the reduction in *KANK1* expression was associated with a trend towards reduced neuronal viability in exon edited cells, and with a significant reduction in neuronal viability in enhancer edited cells (exon: $P=0.1$, enhancer: $P=0.003$, paired Student's *t*-test; **Fig. 5d**). Finally, neurons with reduced expression of *KANK1* exhibited shorter neurites (exon: $P=0.04$, enhancer: $P=0.02$, paired Student's *t*-test; **Fig. 5e**) with reduced branch length (exon: $P=0.02$, enhancer: $P=0.01$, paired Student's *t*-test; **Fig. 5f**). In all instances, measures of neuronal toxicity are correlated with *KANK1* expression (**Supplementary Fig. 4f**), which in turn reflects editing efficiency (**Supplementary Figs. 4a-c**). These experimental observations collectively demonstrate the neuronal toxicity focused on the axon caused by ALS-associated genetic variants in *KANK1*, and further support *KANK1* as a new ALS risk gene.

DISCUSSION

Study of the genetic architectures of complex diseases has been greatly advanced by large GWAS studies. However, many of these studies have not considered cell-type-specific aspects of genomic function, which is particularly relevant for noncoding regulatory sequence¹⁰. This may explain why diseases such as ALS have been linked to relatively few risk genes despite substantial estimates of heritability^{2,3}. Fine-mapping methods have been proposed to disentangle causal SNPs from genetic associations⁵⁴⁻⁵⁹, but these approaches are not integrated with cell-type-specific

biology^{55,58}, or assume a fixed number of causal SNPs per locus^{54,56,57}, limiting their power for gene discovery. We have characterized epigenetic features within MNs, which are the key cell type for ALS pathogenesis. Integrating MN epigenetic features with ALS GWAS data in our RefMap model has discovered 690 ALS risk genes, which extends the list of candidate ALS genes by two orders of magnitude. Others have performed more limited epigenetic profiling of motor neurons⁶⁰, but our data are unique with respect to the depth and number of assessments.

Consistent with previous literature, RefMap ALS genes are functionally associated with the distal axon^{15,16}. Several known ALS risk genes are related to axonal function and axonal transport in particular⁶¹. Unlike previous literature, our work is based on a comprehensive genome-wide screening and not on a small number of rare variants. As a result, our data suggest that the distal axon may be the site of disease initiation in most ALS patients, and should be the focus of future translational research.

RefMap ALS genes include *KANK1*, which is enriched with common and rare ALS-associated genetic variation across multiple domains and datasets. *KANK1* is functionally related to a number of known ALS genes that are important for cytoskeletal function, including *PFN1*, *KIF5A* and *TUBA4A*. In particular, *PFN1*, like *KANK1*, is implicated in actin polymerization⁶². Disruption of actin polymerization has been associated with alterations in synaptic organization⁶³, including the neuromuscular junction (NMJ)⁶⁴, but also with nucleocytoplasmic transport defects⁶⁵. We have experimentally verified the link between variants identified by RefMap to ALS, and *KANK1* expression. Moreover, we have demonstrated that the reduced expression of *KANK1* in a human CNS-relevant neuron is toxic and produces axonopathy. By contrast, *KANK1* upregulation could be a new therapeutic target for ALS patients with mutations that reduce *KANK1* expression, and possibly more broadly.

In summary, our study provides a general framework for the identification of risk genes involved in complex diseases. With the expansion of genotyping data and increasing

understanding of cell-type-specific functions, it should prove valuable to the identification of the genetic underpinnings of many such diseases.

FIGURE LEGENDS

Figure 1. Genome-wide identification of ALS-associated genes.

(a) Schematic of the study design for identifying ALS risk genes by integrating ALS genetics with functional genomics from motor neurons. (1-2) We sequenced the transcriptome and epigenome of the iPSC-derived motor neurons. By integrating (3) ALS genetics with functional genomics of MNs, (4) a machine learning model called RefMap was developed to fine-map ALS-associated regions. (5) After mapping those identified regions to their target genes, 690 ALS-associated genes were pinpointed. (6) Transcriptome analysis based on iPSC-derived MNs, human tissues and mouse models, as well as (7) network analysis were performed to demonstrate the functional significance of RefMap genes in ALS. (8) CRISPR/Cas9 reproduction of newly identified ALS-associated *KANK1* mutations experimentally verified the proposed link to neuronal toxicity. The LD heatmap matrix in (4) was visualized in both R^2 (red) and D' (blue) using LDmatrix (<https://ldlink.nci.nih.gov/?tab=ldmatrix>). GRU=Gene Regulatory Unit; GO=Gene Ontology. (b) Epigenetic profiling data from motor neurons is internally consistent. Markers of genomic activity are significantly enriched in promoter regions of high-expressed genes compared to low-expressed genes. Circle area is proportional to % overlap. Hi-C data was scaled by a factor of ten for clarity. (c) Graphical representation of RefMap. Observed variables were annotated in grey, local hidden variables were in green and global latent variables were in pink (**Methods**). (d) A region (chr12:112,036,001-112,038,000) around *ATXN2* was precisely pinpointed by RefMap because of elevated SNP Z-scores and enriched epigenome peaks (ATAC-seq, H3K27ac and H3K4me3 histone ChIP-seq). The output of RefMap was labeled as Q-score. ATAC-seq and ChIP-seq signals were shown in fold change (FC) based on one replicate from sample CS14.

Figure 2. RefMap genes are intolerant to loss of function.

(a-d) Comparison of (a) haploinsufficiency score, (b) LoFtool percentile, (c) RVIS-ExAC percentile and (d) o/e score between RefMap genes and all the genes in the transcriptome. All comparisons were performed using the one-sided Wilcoxon rank-sum test. RefMap genes showed a significant increase in HI score (a) and a decrease in LoFtool percentile (b), RVIS-ExAC percentile (c) and o/e score (d). The bottom and top of the boxes indicate the first and third quartiles, respectively, where the black line in between indicates the median. The whiskers denote the minimal value within 1.5 interquartile range (IQR) of the lower quartile and the maximum value within 1.5 IQR of the upper quartile. The plus symbols represent outliers. In e, the black dashed lines indicate the lower and upper limits of the regions with regular scale. Outliers outside of the black dashed lines are visualized with compressed scale in regions surrounded by gray lines for better visualization.

Figure 3. Transcriptomics supports the functional importance of RefMap genes in motor neurons and in ALS.

(a) RefMap genes were upregulated compared to the transcriptome in iPSC-derived motor neurons from neurologically normal individuals (n=3). For a fair comparison, we only considered those genes with expressed transcripts (TPM \geq 1) in the transcriptome, following a similar procedure in mapping ALS-associated regions to their targets. (b) RefMap genes were downregulated in post-mortem CNS tissue from *C9orf72*-ALS (n=8) and sporadic ALS (n=10) patients compared to neurologically normal controls (n=17). FC=Frontal Cortex; CB=Cerebellum. (c) RefMap genes were downregulated in iPSC-derived motor neurons from ALS patients (n=55) compared to neurologically normal controls (n=15). All comparisons in a-c were performed using the one-sided Wilcoxon rank-sum test, and the Benjamini-Hochberg (BH) correction was carried out in b. In a-c, the bottom and top of the boxes indicate the first and third quartiles, respectively, where the black line in between indicates the median. The whiskers denote the minimal value within 1.5 IQR of the lower quartile and the maximum value

within 1.5 IQR of the upper quartile. The plus symbols represent outliers. In **b**, the black dashed lines indicate the lower and upper limits of the regions with regular scale. Outliers outside of the black dashed lines are visualized with compressed scale in regions surrounded by gray lines for better visualization. **(d)** Hierarchical clustering of expression changes of RefMap genes during disease progression based on the *SOD1-G93A* mouse model. RefMap genes were mapped to their mouse homologs ($n=510$). Gene expression levels were estimated using the β scores calculated in ³⁷, and were averaged across different sections of spinal cords at each time point. Time points p30, p70, p100, and p120 represent presymptomatic, onset, symptomatic and end-stage, respectively. Difference of gene expressions between *SOD1-G93A* and *SOD1-WT* mice at each time point was quantified by the difference of corresponding β ($\Delta\beta$). Before clustering, $\Delta\beta$ were standardized across genes, and one minus correlation was used as the clustering distance. **(e)** Two distinct expression patterns (C1: 286 genes; C2: 224 genes) of RefMap genes identified after clustering. The larger cluster C1 was progressively downregulated during ALS progression. Solid plot represents the mean of expressions within each cluster, and the standard error was shown as shading. **(f)** Gene ontology analysis of C1, showing that C1 is enriched with functions related to the motor neuron distal axon and synapse. GOBP=Gene Ontology Biological Process; GOCC=Gene Ontology Cellular Compartment. Dashed line represents $P=0.05$.

Figure 4. Protein-protein interaction network analyses associate RefMap genes with distal axonopathy within motor neurons.

(a-b) PPI network analysis revealed two modules that are significantly ($FDR<0.1$) enriched with RefMap genes: **(a)** M421 (721 genes) and **(b)** M604 (308 genes). Hypergeometric test was performed to quantify the enrichment followed by the BH correction. Module nodes were colored to demonstrate RefMap enrichment, where RefMap genes are in blue and other module genes are in yellow. Edge thickness is proportional to STRING confidence score (>700). **(c-d)** RefMap modules, including **(c)** M421 and **(d)** M604, are enriched for motor neuron functions localized within the distal

axon. GOBP=Gene Ontology Biological Process; GOCC=Gene Ontology Cellular Compartment. Dashed line represents $P=0.05$. (e) Representation of pathways enriched in each module (c and d) in MNs. (f) RefMap modules were highly expressed within control motor neurons, consistent with an important role in motor neuron function. All comparisons were performed using the one-sided Wilcoxon rank-sum test. The bottom and top of the boxes indicate the first and third quartiles, respectively, where the black line in between indicates the median. The whiskers denote the minimal value within 1.5 IQR of the lower quartile and the maximum value within 1.5 IQR of the upper quartile. The plus symbols represent outliers. The black dashed lines indicate the lower and upper limits of the regions with regular scale. Outliers outside of the black dashed lines are visualized with compressed scale in regions surrounded by gray lines for better visualization.

Figure 5. Reduced *KANK1* expression is associated with ALS-associated genetic variants and produces neurotoxicity.

(a) We identified a high density of ALS-associated genetic variants within a region at chr9:663001-664000, which overlaps with the regulatory regions in iPSC-derived control motor neurons as well as in the ENSEMBL regulatory build. (b) Whole genome sequencing data from sporadic ALS patients (n=5,594) and neurologically normal controls (2,238) was analyzed to determine the frequency of rare deleterious variants within *KANK1* coding and regulatory sequences. ALS-associated rare variants are shown. All variants were present in a single patient unless stated. No variant was found in a control individual. (c) To experimentally evaluate ALS-associated *KANK1* variants, we performed CRISPR/Cas9 perturbation proximate to patient mutations in enhancer and coding regions within SH-SY5Y neurons, including resection of the chr9:663001-664000 region. Edited neurons revealed (d) reduced viability, (e) reduced axonal length and (f) reduced axonal-branch length compared to HPRT-edited controls. Data shown is mean and standard deviation. Neuronal viability was quantified relative to HPRT-edited controls.

SUPPLEMENTARY INFORMATION

Supplementary Figure 1. iPSC-derived motor neurons (a) are morphologically consistent with lower motor neurons including expression of appropriate markers. **(b)** iPSC cells were derived from control fibroblasts.

Supplementary Figure 2. Clustering quality checking for clusters **(a)** C1 and **(b)** C2. The correlations of $\Delta\beta$ for individual genes and the mean of the cluster were calculated. The comparisons were performed using the Wilcoxon rank-sum test. Significantly increased between-cluster distance and significantly decreased in-cluster distance were observed, demonstrating the high quality of the clustering.

Supplementary Figure 3. Additional results from network analysis. **(a)** Distribution of modularities after Louvain for the smoothed PPI network (red) and 100 randomized networks (blue). The modularity of our smoothed network is significantly shifted from the randomized network. **(b)** M826 contains *KANK1*. M826 is enriched with RefMap genes ($P=5.6e-3$, hypergeometric test) but not after multiple testing correction. **(c)** M826 is functionally enriched for vesicle transport within the motor neuron axon. GOBP=Gene Ontology Biological Process. Dashed line represents $P=0.05$.

Supplementary Figure 4. CRISPR-editing of SH-SY5Y cells. **(a)** Sanger sequencing traces demonstrating spCas9 cut site adjacent to PAM and subsequent waveform decomposition in *KANK1* open reading frame edited cells. **(b)** Indel distribution within *KANK1* open reading frame CRISPR-edited SH-SY5Y cells. **(c)** PCR amplification of the relevant genomic segment in enhancer CRISPR-edited SH-SY5Y cells reveals that the chr9:663001-664000 region has been resected compared to HPRT-edited control cells. **(d)** Altered PAX6 expression and **(e)** increased dendrite length confirm the successful neuronal differentiation of SH-SY5Y cells. **(f)** pPCR reveals that the expression of *KANK1* mRNA was reduced in CRISPR-edited SH-SY5Y neurons.

Supplementary Tables 1

Quality control measures for RNA-seq of iPSC-derived motor neurons

Supplementary Tables 2

Quality control measures for ATAC-seq of iPSC-derived motor neurons

Supplementary Tables 3

Quality control measures for histone ChIP-seq of iPSC-derived motor neurons

Supplementary Tables 4

Quality control measures for Hi-C of iPSC-derived motor neurons

Supplementary Table 5

ALS-associated regions identified by RefMap including Q-scores

Supplementary Table 6

690 ALS-associated genes identified by RefMap

Supplementary Table 7

Manually curated ALS gene list with evidence for association including references

Supplementary Table 8

Clusters of RefMap homologs in transcriptome data from *SOD1-G93A* ALS mouse model

Supplementary Table 9

Smoothed PPI network after preserving top 5% edges predicted by random walk with restart

Supplementary Table 10

Gene modules detected from network analysis

Supplementary Table 11

Project MinE ALS Sequencing Consortium

Supplementary Notes

Mathematical and technical details of RefMap

ACKNOWLEDGEMENTS

This work used the Genome Sequencing Service Center by Stanford Center for Genomics and Personalized Medicine Sequencing Center, supported by the grant award NIH S10OD025212, and NIH/NIDDK P30DK116074. We acknowledge the Stanford Genetics Bioinformatics Service Center for providing computational infrastructure for this study. We thank W. Rheenen for the explanation of ALS GWAS data. We thank J. Adrian for the help to initiate the project. We also thank J. Zhai and X. Yang for the help with histone ChIP-seq assays, and I. Gabdank and M. Kagda for running the Hi-C pipeline. This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement n° 772376 - EScORIAL. The collaboration project is co-funded by the PPP Allowance made available by Health~Holland, Top Sector Life Sciences & Health, to stimulate public-private partnerships. This study was also supported by the ALS Foundation Netherlands, by research grants from IWT (n° 140935), the ALS Liga België, the National Lottery of Belgium, the KU Leuven Opening the Future Fund, the National Institutes of Health (CEGS 5P50HG00773504, 1P50HL083800, 1R01HL101388, 1R01-HL122939, S10OD025212, and P30DK116074, UM1HG009442 to M.P.S.), the Wellcome Trust (216596/Z/19/Z to J.C.K.) and NIHR (P.J.S.). We acknowledge support from a Kingsland fellowship (T.M.), the My Name's Doddie Foundation (J.F.), and the NIHR Sheffield Biomedical Research Centre for Translational Neuroscience. Biosample collection was supported by the MND Association and the Wellcome Trust (P.J.S.). We

are very grateful to those ALS patients and control subjects who generously donated biosamples. We acknowledge transcriptome data provided by the AnswerALS Consortium.

AUTHOR CONTRIBUTIONS

S.Z., J.C.K. and M.P.S. conceived and designed the study. S.Z. contributed to the design, theoretical analysis and implementation of RefMap. S.Z., J.C.K., A.K.W., M.S., T.M., C.H., H.G.N., J.F., C.S.S., L.F., P.J.S. and M.P.S. were responsible for data acquisition. S.Z., J.C.K., A.K.W., M.S., T.M., C.H., H.G.N., J.F., C.S.S., C.W., J.L., C.E., E.H., L.F., P.J.S. and M.P.S. were responsible for analysis of data. S.Z., J.C.K., A.K.W., M.S., T.M., C.H., H.G.N., J.F., C.S.S., C.W., J.L., C.E., E.H., K.P.K., J.V., L.F., P.J.S. and M.P.S. were responsible for interpretation of data. The Project MinE ALS Sequencing Consortium (**Supplementary Table 11**) was involved in data acquisition and analysis. S.Z., J.C.K. and M.P.S. prepared the manuscript with assistance from all authors. All authors meet the four ICMJE authorship criteria, and were responsible for revising the manuscript, approving the final version for publication, and for accuracy and integrity of the work.

DECLARATION OF INTERESTS

M.P.S. is a cofounder of Personalis, Qbio, Sensomics, Filtricine, Mirvie and January. He is on the scientific advisory of these companies and Genapsys. No other authors have competing interests.

METHODS

Study cohorts

iPSC-cells were derived from fibroblasts obtained from three neurologically normal controls of different ages: 55-year old male, a 52-year old female and a 6-year old male (**Supplementary Fig. 1b**). GWAS summary statistics were previously published⁴. The 6,180 patients and 2,370 controls included in this study were recruited at specialized

neuromuscular centers in the UK, Belgium, Germany, Ireland, Italy, Spain, Turkey, the United States and the Netherlands⁴³. Patients were diagnosed with possible, probable or definite ALS according to the 1994 El-Escorial criteria⁶⁶. All controls were free of neuromuscular diseases and matched for age, sex and geographical location.

Cell culture

Human induced pluripotent stem cells iPSCs were maintained in Matrigel-coated plates (Corning®, Cat.: #356230) according to the manufacturer's recommendations in complete mTeSR™-Plus™ Medium (StemCell Technologies, Cat.: #05825). The culture medium was replaced daily and confirmed mycoplasma free. Cells were passaged every four to six days as clumps using ReLeSR™ an enzyme-free reagent for dissociation (StemCell Technologies, Cat.: #05872) according to the manufacturer's recommendations. For all the experiments in this study, iPSCs were between passage 20 and 32.

iPSC-derived motor neuron differentiation

iPSCs derived from unaffected controls were differentiated to motor neurons using the modified version of the dual SMAD inhibition protocol⁶⁷. Briefly iPSC cells were transferred for Matrigel-coated plate (Corning® Matrigel® Growth Factor Reduced). On the day after plating (day 1), after the cells had reached ~100% confluence, the cells were washed once with PBS and then the medium was replaced for neural medium (50% of KnockOut™ DMEM/F-12 Cat.: #12660012, 50% of Neurobasal ThermoFisher Cat.: #12660012), 0.5× N2 supplement (ThermoFisher, Cat.: #17502001), 1x Gibco® GlutaMAX™ Supplement (ThermoFisher, Cat.: #35050061), 0.5x B-27 (ThermoFisher, Cat.: #17504001), 50 U ml⁻¹ penicillin and 50 mg ml⁻¹ streptomycin, supplemented with SMAD inhibitors (DMH-1 2 μM Tocris, Cat.: #4423; SB431542-10 μM Tocris, Cat.: #1614 and CHIR99021 3 μM, CHIR Tocris, Cat.: #4423).

The medium was changed every day for 6 days, on day 7, the medium was replaced for neural medium supplemented with DMH-1 2 μM, SB431542-10 μM and CHIR 1 μM,

All-Trans Retinoic Acid 0.1 μ M (RA, STEMCELL Technologies, Cat.: #72262), and Purmorphamine 0.5 μ M (PMN, Tocris, Cat.: #4551), the cells were kept in this medium until day 12 when is possible to see a uniform neuroepithelial sheet, the cells were split 1:6 with Accutase (StemPro® Accutase® Cell Dissociation Reagent, Gibco™ A1110501), onto matrigel substrate in the presence of 10 μ M of rock inhibitor (Y-27632 dihydrochloride, Tocris), giving rise to a sheet of neural progenitor cells (NPC). After 24 hours of incubation the medium was changed for neural medium supplemented with RA 0.5 μ M and PMN 0.1 μ M, the medium was changed every day for more 6 days. On day 19 the motor neuron progenitors were split with accutase onto to matrigel-coated plates and the medium was replaced for neural medium supplemented with RA 0.5 μ M, PMN 0.1 μ M, compound E 0.1 μ M (Cpd E, Tocris, Cat.: #6476), BDNF 10ng/mL (ThermoFisher, Cat.: #PHC7074), CNTF 10ng/mL (ThermoFisher, Cat.: #PHC7015) and IGF 10ng/mL (ThermoFisher, Cat.: #PHG0078) until day 28. On day 29, the media was replaced for Neuronal media (Neurobasal media supplemented with 1% of B27, BDNF 10ng/mL, CNTF 10ng/mL and IGF 10ng/mL). The cells were then fed alternate days with neuronal medium until day 40.

ATAC-seq

50,000 viable motor neurons were spun down at 500 RCF at 4°C for 5 min. Supernatant was discarded. 50 μ l cold ATAC Resuspension Buffer (RSB) (10 mM Tris-HCl pH 7.4, 10 mM NaCl, 3 mM MgCl₂, sterile H₂O) containing 0.1% NP40, 0.1% Tween-20, and 0.01% Digitonin was added and carefully mixed. Tubes were incubated on ice for 3 min. 1 ml of cold ATAC-RSB containing 0.1% Tween-20 was added and the tubes were inverted three times. Nuclei were spun down at 500 RCF for 10 min at 4°C. Supernatant was aspirated. Cell pellet was resuspended in 50 μ l of transposition mix (25 μ l 2x TD buffer, 2.5 μ l transposase (100 nM final), 16.5 μ l PBS, 0.5 μ l 1% digitonin, 0.5 μ l 10% Tween-20, 5 μ l H₂O) by pipetting up and down 6 times. TD buffer consists of 20 mM Tris-HCl pH 7.6, 10 mM MgCl₂, 20% DMF, sterile H₂O. pH was adjusted with acetic acid before adding DMF. The reaction was incubated at 37°C for 30 minutes in a

thermomixer while shaking at 1000 RPM. Reaction was cleaned up with a Qiagen MiElute kit. DNA was eluted in 20 μ L elution buffer. DNA was amplified using the NEBNext 2xMasterMix (M0541). Cycling conditions: 5 min at 72°C, 30 sec at 98°C, followed by 5 cycles of 10 sec at 98°C, 30 sec at 63°C and 1 min at 72°C, hold at 4°C. 5 μ l (10% of the pre-amplified mixture) were used for qPCR to determine the number of additional cycles needed (3.76 μ L H₂O, 0.5 μ L 25 μ M Primer1, 0.5 μ L 25 μ M Primer2, 0.24 μ L 25x SYBR Green, 5 μ L NEBNext MasterMix). Cycling conditions: 30 sec at 98°C, followed by 20 cycles of 10 sec at 98°C, 30 sec at 63°C and 1 min at 72°C, hold at 4°C. Amplification profiles were assessed as previously described¹⁸. The remainder of the pre-amplified DNA (45 μ L) was used to run the required number of additional cycles. The final PCR reaction was cleaned up using Qiagen MinElute kit and eluted in 20 μ l H₂O. Libraries were quantified with the KAPA Library Quantification kit (Roche) and sequenced on a NovaSeq 6000 system (Illumina). Raw data were processed with the ENCODE 4 pipeline for ATAC-seq according to ENCODE 4 standards (<https://www.encodeproject.org/atac-seq/>). All samples exceeded ENCODE 4 standards for % mapped reads, enrichment of transcription start sites, the fraction of reads that fall within peak regions (FRiP), and reproducibility between technical replicates (**Supplementary Table 1**).

Files are available at [encodeproject.org](https://www.encodeproject.org) with the following accession numbers: ENCSR065CER, ENCSR410DWV, ENCSR812ZKP, ENCSR634WYX, ENCSR459PVP, ENCSR913OWV, ENCSR704VZY, ENCSR131HOY, ENCSR516YAD, ENCSR709QRD.

Histone ChIP-seq

5 million motor neurons were crosslinked and resuspended in 10 mL of cold L1 buffer (50mM Hepes KOH, pH 7.5, 140mM NaCl, 1mM EDTA, 10% Glycerol, 0.5% NP-40, 0.25% Triton X-100, dH₂O, 1 protease inhibitor tablet (Roche Complete cat #1 697 498) per 50ml buffer). Cells were incubated on a rocking platform at 4°C for 10 minutes and spun down at 3000 rpm at 4°C for 10 minutes. Pellets were resuspended in 10 mL of L2

buffer (200mM NaCl 1mM EDTA pH 8 0.5mM EGTA 10mM Tris, pH 8, dH₂O, 1 protease inhibitor tablet (Roche Complete cat #1 697 498) per 50ml buffer, room temperature). Tubes were incubated at room temperature for 10 minutes and spun down at 3000 rpm for 10 minutes at 4°C. Nuclei were resuspended in 3 mL 1X RIPA buffer and incubated on ice for 30 minutes. Samples were sonicated with Branson 250 Sonifier to shear the chromatin. 3 mL of sheared chromatin lysate were transferred to two 2 mL tubes and spun down at 14,000 rpm at 4°C for 15 minutes. 50 µL were saved from each replicate and pooled as input (no antibody added, kept at -20°C). 2 µL histone modification antibody was added to each 3 mL lysates and incubated at 4°C on a neutator for 12-16 hours. The following antibodies were used: H3K4me1 (#5326S, lot 3, Cell Signaling Technologies), H3K4me3 (#9751S, lot 10, Cell signaling technologies), H3K27ac (#93133, lot 28518012, ActiveMotif). 80 µL of Protein A/G-agarose for each sample were washed twice with 1 mL of ice cold 1X RIPA buffer, spun down at 5000 rpm for 1 minute at 4°C and resuspended in 80µL in 1x RIPA buffer. Beads were added to tubes containing Ag-Ab complex (80 µ L 1X RIPA to wash out the beads) and incubated for 1 hour at 4°C with neutator rocking. Tubes were spun down at 1500 rpm for 3 minutes, beads were washed 3 times 15 minutes each with 10 mL of fresh, ice cold 1x RIPA buffer supplemented per 50 mL with 1 protease inhibitor tablet, 250 µL of 100 mM PMSF, 50 µL of 1M DTT, 2 ml of phosphatase inhibitor (sodium pyrophosphate 1mM, sodium orthovanadate 2mM, sodium fluoride 10mM). Afterwards, beads were washed once with ice cold 1 x PBS for 15 minutes. Beads were resuspended in 1200 µL ice cold 1x PBS, transferred to an 1.5mL Eppendorf tube and spun down at 5000 rpm for 1 minute. PBS was removed and 100 µL of Elute 1 solution (1% SDS, 1x TE, dH₂O) was added to resuspend beads and tubes were incubated at 65°C for 10 minutes with gentle mixing every 2 minutes. Beads were spun down at 5000 rpm for 1 minute at room temperature and the supernatant was kept as Elute 1. 150 µL of Elute 2 solution (0.67% SDS, 1x TE) was added to the bead pellets and incubated at 65°C for 10 minutes with gentle vortexing. After spinning down for 1 minute at 5000 rpm, the second elute was combined with the first. Input DNA was thawed and 150 µL of Elute 1 solution

was added. All samples incubated at 65°C overnight to reverse cross-linking. 250 µL 1X TE containing 100 µg RNase was added to each sample and incubated for 30 minutes at 37°C. 5 µL of 20 mg/mL Proteinase K was added to each sample and incubated at 45°C for 30 minutes. After transferring samples to 15 mL tubes, DNA was purified (PCR purification kit, Quiagen). DNA was eluted in elution buffer (50µL for input, 35µL for ChIP sample).

The following components were combined and mixed in a microfuge tube: ChIP DNA to be end-repaired (25ng) 34 µL, 5 µL 10X End-Repair Buffer, 5 µL 2.5 mM dNTP Mix , 5 µL 10 mM ATP, 1 µL End-Repair Enzyme Mix. The mixture was incubated at room temperature for 45 minutes. DNA was purified (MinElute PCR purification kit, Quiagen) and eluted in 19 µL EB. Adapter ligated DNA was run on a 2% EX-Gel and excised in the range of 450-650 bp with a clean scalpel. DNA was purified (Gel extraction kit, Quiagen) and eluted in 20 µL EB. The following components were mixed in a PCR tube: 20 µL of purified DNA, 25 µL KAPA HiFi HotStart ReadyMix (2X), 5 µL KAPA Library Amplification Primer Mix (10X). DNA was amplified with the following conditions: 45 sec at 98°C, 15x [15 sec at 98°C, 30 sec at 60°C, 30 sec at 72°C], 60 sec at 72°C, hold at 4°C. The PCR product was purified (MinElute PCR purification kit, Quiagen) and eluted in 19 µL EB. DNA was run on a 2% EX-Gel and excised in the range of 300-450 bp (or brightest smear) with a clean scalpel. DNA was purified (Gel extraction kit, Quiagen) and eluted in 12 µL EB. Library concentration was measured using Qubit and each library was run on the Bioanalyzer. Equal concentrations of different barcoded libraries were pooled and sequenced on a NovaSeq 6000 system (Illumina). Raw data were processed with the ENCODE 4 pipeline for Histone ChIP-seq according to ENCODE 4 standards (<https://www.encodeproject.org/chip-seq/histone/>). All samples exceeded ENCODE standards for % mapped reads, the fraction of reads that fall within peak regions (FRiP), and reproducibility between technical replicates (**Supplementary Table 2**)

Files are available at encodeproject.org with the following accession numbers: ENCSR754DRC, ENCSR672RKZ, ENCSR571HAY, ENCSR503HWR, ENCSR207VLY, ENCSR962OTG, ENCSR745TRI, ENCSR595HWK, ENCSR312HLG, ENCSR682BFG, ENCSR680IWU, ENCSR564EFE, ENCSR358AOC, ENCSR698HPK, ENCSR778FKK, ENCSR425FUS, ENCSR489LNU, ENCSR540KQC

Hi-C

We generated Hi-C libraries following the protocol previously described^{68,69}. In brief, 2-5 million cells were crosslinked with formaldehyde. Nuclei were permeabilized and DNA was digested with 100U of Mbol. DNA fragments were labelled with biotinylated nucleotides. Ligated DNA was purified and sheared to a length of 300-500 bp after reverse cross-linking. Ligation junctions were pulled-down with magnetic streptavidin beads. Libraries were amplified by PCR and purified. Library concentrations were measured (Qubit). Hi-C libraries were paired-end sequenced on a NovaSeq 6000 system (Illumina). Raw data were processed with the ENCODE 4 pipeline for Hi-C according to ENCODE 4 standards (<https://www.encodeproject.org/documents/75926e4b-77aa-4959-8ca7-87efc39d79/>).

RNA-seq

RNA libraries were prepared by first depleting ribosomal RNA using the Illumina Ribo-Zero rRNA depletion kit. Strand-specific libraries were then prepared using NEBext Ultra RNA prep kit. RNAseq libraries were paired-end sequences on a NovaSeq 6000 system (Illumina). Minimum 80 million reads were obtained per sample. The raw Fastq files were trimmed for the presence of Illumina adapter sequences using Cutadapt v1.2.1⁷⁰. The reads were further trimmed using Sickle v1.200 (<https://github.com/najoshi/sickle>) with a minimum window quality score of 20. Reads shorter than 15 bp after trimming were removed. If only one of a read pair passed this filter, it was included in the R0 file. Reads were aligned to hg19 transcripts (n=180,253) using Kallisto v0.46.0⁷¹.

Model design and inference of RefMap

In this study, allele Z-scores were calculated as $Z=b/se$, where b and se are effect size and standard error, respectively, and they were estimated from a mixed linear model in the ALS GWAS study^{4,5}. Given allele Z-scores and the epigenetic profiling of iPSC-derived motor neurons, we were interested in predicting causal associations of individual genomic regions with ALS risk. Suppose we have K 1Mb LD blocks with non-zero alleles, whose approximate between-block independence has been verified in previous literature²³. Also suppose each LD block contains J_k ($k=1, \dots, K$) 1kb regions and each region harbors $I_{j,k}$ ($j=1, \dots, J_k, I_{j,k}>0$) SNPs. We further denote the Z-score for the i -th SNP in the j -th region of the k -th block as $z_{i,j,k}$ ($i=1, \dots, I_{j,k}$). Under a linearity hypothesis, we can prove that \mathbf{z}_k follows a multivariate normal distribution (**Supplementary Notes**), i.e.,

$$\mathbf{z}_k | \mathbf{u}_k \sim \mathcal{N}(\Sigma_k \mathbf{u}_k, \Sigma_k), \quad k=1, \dots, K, \quad (1)$$

in which \mathbf{u}_k are the effect sizes of individual SNPs that can be expressed as

$$\mathbf{u}_k = \left[\mathbf{u}_{1:I_{1,k}, 1, k}^T, \dots, \mathbf{u}_{1:I_{j,k}, j, k}^T, \dots, \mathbf{u}_{1:I_{J_k, k}, J_k, k}^T \right]^T. \quad (2)$$

Moreover, in Eq. (1) $\Sigma_k \in \mathbb{R}^{I_k \times I_k}$ represents the in-sample LD matrix comprising of the pairwise Pearson correlation coefficients between SNPs within the k -th block, where I_k is the total number of SNPs given by $I_k = \sum_{j=1}^{J_k} I_{j,k}$. Here, since we have no access to the individual genotypes, we used EUR samples from the 1000 Genomes Project to estimate Σ_k (i.e., out-sample LD matrix).

Here, the latent variables \mathbf{u}_k can be treated as the disentangled Z-scores from LD confounding, leaving the right place for independence assumption and facilitating downstream modelling. Indeed, we assume $u_{i,j,k}$ ($i=1, \dots, I_{j,k}$) are independent and identically distributed (i.i.d.), following a normal distribution given by

$$u_{i,j,k} | m_{j,k}, \lambda_{j,k} \sim \mathcal{N}(m_{j,k}, \lambda_{j,k}^{-1}), \quad i=1, \dots, I_{j,k}, \quad (3)$$

where the precision $\lambda_{j,k}$ follows a Gamma distribution, i.e.,

$$\lambda_{j,k} \sim \text{Gamma}(a_0, b_0) . \quad (4)$$

Moreover, to characterize the shift of the expectation in Eq. (3) from the background due to its functional effect, we model $m_{j,k}$ by a three-component Gaussian mixture model, i.e.,

$$m_{j,k} | t_{j,k}, v_{-1}, v_{+1}, \tau_0, \tau_{-1}, \tau_{+1} \sim \underbrace{\mathcal{N}(-v_{-1}, \tau_{-1}^{-1})^{t_{j,k}^{(-1)}}}_{\text{negative}} \underbrace{\mathcal{N}(0, \tau_0^{-1})^{t_{j,k}^{(0)}}}_{\text{zero}} \underbrace{\mathcal{N}(v_{+1}, \tau_{+1}^{-1})^{t_{j,k}^{(+1)}}}_{\text{positive}} , \quad (5)$$

where the precisions follow

$$\tau_{-1}, \tau_0, \tau_{+1} \sim \text{Gamma}(a_0, b_0) , \quad (6)$$

and v_{-1} and v_{+1} are non-negative variables quantifying the absolute values of effect size shifts for the negative and positive components, respectively.

To impose non-negativity over v_{-1} and v_{+1} , here we employ the rectification nonlinearity technique proposed in⁷². In particular, we assume v_{-1} and v_{+1} follow

$$v_{-1} | m_{-1}, \lambda_{-1} \sim \mathcal{R}^N(m_{-1}, \lambda_{-1}) , \quad (7)$$

$$v_{+1} | m_{+1}, \lambda_{+1} \sim \mathcal{R}^N(m_{+1}, \lambda_{+1}) , \quad (8)$$

in which the rectified Gaussian distribution is defined via a dumb variable. Specifically, we first define v_{-1} and v_{+1} by

$$v_{-1} = \max(r_{-1}, 0) , \quad (9)$$

$$v_{+1} = \max(r_{+1}, 0) , \quad (10)$$

which guarantee that v_{-1} and v_{+1} are non-negative. The dumb variable r_{-1} and r_{+1} follow Gaussian distributions given by

$$r_{-1} | m_{-1}, \lambda_{-1} \sim \mathcal{N}(m_{-1}, \lambda_{-1}^{-1}) , \quad (11)$$

$$r_{+1} | m_{+1}, \lambda_{+1} \sim \mathcal{N}(m_{+1}, \lambda_{+1}^{-1}) , \quad (12)$$

where m_{\pm} and λ_{\pm} follow the Gaussian-Gamma distributions, i.e.,

$$m_{-1}, \lambda_{-1} \sim \mathcal{N}(\mu_0, (\beta_0 \lambda_{-1})^{-1}) \text{Gamma}(a_0, b_0) , \quad (13)$$

$$m_{+1}, \lambda_{+1} \sim \mathcal{N}(\mu_0, (\beta_0 \lambda_{+1})^{-1}) \text{Gamma}(a_0, b_0) . \quad (14)$$

The indicator variables in Eq. (5) denote whether that region is ALS-associated or not. Indeed, we define the region to be disease-associated if $t_{j,k}^{(-1)} = 1$ or $t_{j,k}^{(+1)} = 1$, and to be non-associated otherwise. To simplify the analysis, we put a symmetry over $t_{j,k}^{(-1)}$ and $t_{j,k}^{(+1)}$, and define the distribution by

$$p(t_{j,k} | \pi_{j,k}) = (0.5\pi_{j,k})^{t_{j,k}^{(-1)}} (1 - \pi_{j,k})^{t_{j,k}^{(0)}} (0.5\pi_{j,k})^{t_{j,k}^{(+1)}}, j = 1, \dots, J_k, k = 1, \dots, K. \quad (15)$$

Furthermore, the probability parameter $\pi_{j,k}$ in Eq. (15) is given by

$$\pi_{j,k} = \sigma(\mathbf{w}^T \mathbf{s}_{j,k}), \quad (16)$$

where $\sigma(\cdot)$ is the sigmoid function, $\mathbf{s}_{j,k}$ is the vector of epigenetic features for the j -th region in the k -th LD block, and the weight vector \mathbf{w} follows a multivariate normal distribution, i.e.,

$$\mathbf{w} | \Lambda \sim \mathcal{N}(\mathbf{0}, \Lambda^{-1}), \quad (17)$$

and Λ follows

$$\Lambda \sim \mathcal{W}(\mathbf{W}_0, \nu_0). \quad (18)$$

In our study, the epigenetic features $\mathbf{s}_{j,k}$ were calculated as the overlapping ratios of that region with the narrow peaks of ATAC-seq and histone ChIP-seq, respectively.

Based on Eqs. (1) to (18), we are interested in calculating $p(\mathbf{T} | \mathbf{Z}, \mathbf{S})$ wherein the calculation of integrals is intractable. Here we seek for approximate inference based on the mean-field variational inference (MFVI)⁷³. To reduce false positives, we set a hard threshold for $q(t_{j,k}^{(0)})$ with respect to the ATAC-seq signal, where we set $q(t_{j,k}^{(0)} = 1) = 1$ if the corresponding region overlaps no ATAC-seq peak. This was motivated by our particular interest in active regions. More technical details, including a coordinate ascent-based inference algorithm, were provided in Supplementary Notes.

In this study, we ran the inference algorithm per chromosome to accelerate the computation. The Q^+ - and Q^- -scores were defined as $q(t^{(+1)} = 1)$ and $q(t^{(-1)} = 1)$, respectively, and we also defined the Q -score as $Q = Q^+ + Q^-$. To prioritize RefMap-scored

regions, we set a cutoff of 0.95 and defined those regions with either Q^+ - or Q^- -score larger than the cutoff as significant regions (i.e., ALS-associated regions) (**Supplementary Table 5**).

Code relevant to RefMap is available on request.

Target gene identification

After identifying ALS-associated regions based on RefMap, we mapped those active regions to their target genes for a better understanding of their functions. In particular, we performed such mapping according to two principles: (i) assign to a gene if the region overlaps the gene or the region up to 10kb either side of the gene body; (ii) assign to a gene if the region overlaps a loop anchor harboring the transcription start site (TSS) of that gene. The loops were called from the Hi-C data sequenced from the iPSC-derived MNs. Note the only transcripts with $TPM \geq 1$ were kept for downstream analysis.

Network analysis

We first downloaded the human PPIs from STRING v11, including 19,567 proteins and 11,759,455 protein interactions. To eliminate the bias caused by hub proteins, we first carried out the random walk with restart algorithm⁷⁴ over the PPI network, wherein the restart probability was set to 0.5, resulting in a smoothed network after preserving the top 5% predicted edges. To decompose the network into different subnetworks/modules, we performed the widely-used Louvain algorithm⁷⁵, a classic community detection algorithm that searches for densely connected modules by optimizing the modularity. After the algorithm converged, we obtained 912 modules with an average size of 18.39 nodes (**Supplementary Table 10**). Two modules (M421 and M604) were significantly enriched ($FDR < 0.1$) with our RefMap genes based on the hypergeometric test followed by the BH correction (**Supplementary Table 10**).

To test whether the network modularity could be observed by chance, we randomly shuffled the edges of the network while preserving the number of neighbors for each

node⁷⁶. We generated 100 such randomized networks followed by the Louvain decomposition, against which the modularity of the smoothed PPI network was tested.

Rare variant burden tests

ALS features a polygenic rare variant architecture⁴, therefore, all searches for pathogenic variants in enhancer and coding regions featured a filter for MAF within the Genome Aggregation Database (gnomAD) of <1/100 control alleles³⁴. Additional filtering varied reflecting differences in function between enhancer, promoter and coding sequence. In enhancer regions, variants were included only if evolutionary conserved based on a LINSIGHT score >0.8⁴⁴. We also utilized an independently compiled score for ALS-associated regulatory variation⁷⁷: variants were excluded with a DIVAN score <0.5. In promoter regions, we utilized two independent scores for functionality and pathogenicity: variants were included in burden testing if their CADD⁴⁵ score >25 and GWAVA⁴⁶ score >0.5. In coding regions, we filtered for variants with impact on protein function as defined by snpeff⁷⁸: variants annotated HIGH/MODERATE/LOW impact were included, but we excluded variants annotated 'synonymous' or 'TF_binding_site_variant' because these functions are independent of amino acid sequence.

The optimal unified test (SKAT-O) was used to perform burden testing in enhancer and promoter regions because it is optimized for large numbers of samples and for regions where a significant number of variants may not be causal⁴⁹. SKAT tests upweight significance of rare variants according to a beta density function of MAF in which $w_j = \text{Beta}(p_j, a_1, a_2)$, where p_j is the estimated MAF for SNP_{*j*} using all cases and controls, parameters a_1 and a_2 are prespecified, and $a_2=2500$ was chosen for all statistical tests. To adjust for confounders including population structure, burden testing used the first ten eigenvectors generated by principal components analysis of common variant profiles, sequencing platform and sex as covariates.

CRISPR/Cas9 editing of SH-SY5Y cells

Guide RNAs (gRNAs) were designed using the Crispor tool (<http://crispor.tefor.net/>) to target *KANK1* regulatory and coding regions. Design was guided by proximity to patient enhancer mutation sites, available protospacer adjacent motifs (PAM), and predicted on- and off- target efficiencies. gRNAs targeting within 30bp either side of the patient enhancer mutation site (chr9:663,001-664,000, hg19) were considered and screened for editing efficiency. One pair of guide sequences (5'-UCAUGGGAACUCUCAAUA-3' and 5'-UCAUGGGAACUCUCAAUA-3') was most efficient and chosen for subsequent experimentation. Validated, commercially available CRISPR control targeting HPRT (IDT) and *KANK1* exon-targeting (ThermoFisher Scientific, 5'-GUCUAGUUGAUACCAUAGG-3') gRNAs were also obtained. gRNA duplexes were assembled from tracrRNA and crRNA in a thermocycler according to manufacturer's instructions under RNase-free conditions. Cells were cultured to ensure 70-90% confluency on the day of transfection. 1ml antibiotic-free DMEM (Lonza) was prepared and incubated in 24-well plates at 37°C. CRISPR/Cas9 Ribonucleoproteins were formed by complexing 240ng gRNA duplex with 1250ng Alt-R V3 Cas9 Protein (IDT) in 10µL buffer R (from 10µL Neon transfection kit, ThermoFisher Scientific) - a 1:1 molar ratio - for 10 minutes. 100,000 viable cells were aliquoted per transfection and centrifuged at 400 x g for 4 minutes. Cells were washed in calcium- and magnesium-free Dulbecco's Phosphate Buffered Saline (Sigma) and centrifuged at 400 x g for 4 minutes. Cell pellets were resuspended in 10µL buffer R containing Cas9 protein and gRNA duplexes. 2µL of 10.8µM electroporation enhancer (IDT) was added and the solution mixed thoroughly to ensure a suspension of single cells. 10µL of this mixture was loaded into a Neon transfection system (ThermoFisher Scientific) and electroporated according to manufacturer's instructions (1200V, 3 pulse, 20s pulse width for SH-SY5Y cells). Cells were then transferred to pre-warmed media in 24-well plates.

Determining CRISPR editing efficiency

Genomic DNA was isolated from CRISPR-edited and control cells using a GenElute Mammalian DNA Miniprep Kit (Sigma) according to manufacturer's instructions. A ~400bp region around the expected cas9 cut site was amplified by polymerase chain reaction using VeriFi mix (PCRbio). Expected amplification was confirmed using gel electrophoresis, and the products were Sanger-sequenced. Sequencing trace files were uploaded to ICE (<https://ice.synthego.com>) and an indel efficiency calculated.

Quantitative PCR (RT-PCR)

Cells were cultured until at least 70% confluent, lysed on ice using an appropriate volume of Tri Reagent (Sigma) for 5 minutes and transferred to 1.5ml RNase-free tubes. Total RNA was extracted using a Direct-zol RNA Miniprep Kit (Zymo) according to manufacturer's instructions, and RNA concentration confirmed using a NanoDrop spectrophotometer (ThermoFisher Scientific). 2µg of total RNA was then converted to cDNA by adding 1µL 10mM dNTPs, 1µL 40µM random hexamer primer (ThermoFisher Scientific), and DNase/RNase-free water to a total volume of 14µL. This mixture was heated for 5 minutes at 70°C then placed on ice for 5 minutes. 4µL of 5x FS buffer, 2µL 0.1M DTT, and 1µL M-MLV reverse transcriptase (ThermoFisher Scientific) were then added and cDNA conversion performed in a PCR thermocycler (37°C for 50 minutes, 70°C for 10 minutes). cDNA was amplified using RT-PCR with Brilliant III SYBR Green (Agilent) as per manufacturer's instructions. Ct analysis was performed using CFX Maestro software (BioRad). GAPDH was chosen as a reference gene because expression is relatively stable in SH-SY5Y cells⁷⁹. Relative mRNA expression values were then calculated using the $2^{-\Delta\Delta CT}$ method⁸⁰.

SH-SY5Y neuronal differentiation and assessment of neurite length

Human SH-SY5Y neuroblastoma cells were seeded at densities of either 5×10^4 cells per well of a 6-well culture plate, or 2×10^3 cells per well of a 96-well culture plate in DMEM (Lonza) supplemented with 10% (v/v) FBS, 50 units/mL penicillin and 50 µg/mL of streptomycin. 24 hours after seeding the media was changed to DMEM supplemented with 5% (v/v) FBS, 50 units/mL penicillin, 50 µg/mL of streptomycin, 4mM l-glutamine

and 10 μ M retinoic acid. After 72 hours, the medium was switched to neurobasal media (ThermoFisher Scientific) containing 1% (v/v) N-2 supplement 100x, 50 units/mL penicillin, 50 μ g/mL of streptomycin, 1% l-glutamine and 50ng/mL human BDNF. Cells were cultured for an additional 3 days until fully differentiated.

To confirm neuronal differentiation and to assess for changes consistent with axonopathy, semi-automated analysis of neurite length was performed using the SimpleNeuriteTracer plugin for FIJI⁸¹. 2D images were converted to 8-bit grayscale and successive points along the midline of a neural process were selected. The software automatically identified the path between the two points. Tracing accuracy was improved using Hessian-based analysis of image curvatures. The AnalyzeSkeleton plugin⁸² was used to quantify the morphology of the traces including the length of neurites. In the case of joined neurites the shorter path length was assigned to 'branches'. To determine whether observed changes in neurite length are significant three fields of view were analyzed and differences were assessed by a t-test, where a one-tailed test was chosen based on the hypothesis that ALS-associated mutations would reduce neurite length.

Immunocytochemistry

SH-SY5Y cells were fixed with 4% paraformaldehyde for 15 minutes and washed 3x with PBS. Cells were blocked in 5% normal horse serum containing 0.1% Triton X-100 for 1 hour at RT. All primary antibodies were diluted in blocking solution (α -tubulin, 1:2000; anti-Pax6, 1:200). Cells were incubated in the primary antibody for 2 hours at RT and washed 3x in PBS before incubation in the appropriate secondary antibody (1:1000 in PBS) for 1 hour at RT. Nuclear counterstain (Hoechst 33342) was applied for 10 minutes followed by a 3x wash in PBS. Cells were imaged using an Opera Phenix High Content Screening System (PerkinElmer).

MTT assays

A colorimetric assay using 3-(4, 5-dimethylthiazol-2-yl)-2, 5-diphenyltetrazolium bromide (MTT) dye was used to assess neuronally differentiated SH-SY5Y cellular metabolic activity and hence neuronal viability. 55 μ L of 5mg/mL of MTT reagent in PBS was added per well of a 24-well culture plate and incubated at 37°C for 1 hour. 550 μ L of un-precipitated 20% SDS in 50% di-methyl formamide (DMF) + dH₂O (pH 7.4) was added per well and mixed thoroughly to lyse the cells. Cells were incubated in a dark environment on an orbital shaker for 1 hour. The colorimetric change was measured using a PHERAstar FS spectrophotometer (BMG Biotech), and absorbance readings taken at 590nm were normalized to media-only wells. Mean absorbance readings were calculated for each biological repeat and expressed as a percentage of controls.

REFERENCES

1. Hardiman, O. *et al.* Amyotrophic lateral sclerosis. *Nat Rev Dis Primers* **3**, 17071 (2017).
2. Ryan, M., Heverin, M., McLaughlin, R. L. & Hardiman, O. Lifetime Risk and Heritability of Amyotrophic Lateral Sclerosis. *JAMA Neurol.* (2019) doi:10.1001/jamaneurol.2019.2044.
3. Trabjerg, B. B. *et al.* ALS in Danish Registries: Heritability and links to psychiatric and cardiovascular disorders. *Neurol Genet* **6**, e398 (2020).
4. van Rheenen, W. *et al.* Genome-wide association analyses identify new risk variants and the genetic architecture of amyotrophic lateral sclerosis. *Nat. Genet.* **48**, 1043–1048 (2016).
5. Nicolas, A. *et al.* Genome-wide Analyses Identify KIF5A as a Novel ALS Gene. *Neuron* **97**, 1268–1283.e6 (2018).
6. Wang, L., Jia, P., Wolfinger, R. D., Chen, X. & Zhao, Z. Gene set analysis of genome-wide association studies: methodological issues and perspectives. *Genomics* **98**, 1–8 (2011).
7. Li, J., Li, X., Zhang, S. & Snyder, M. Gene-Environment Interaction in the Era of Precision Medicine. *Cell* **177**, 38–44 (2019).
8. Cooper-Knock, J. *et al.* Rare Variant Burden Analysis within Enhancers Identifies CAV1 as

- a New ALS Risk Gene. (2020) doi:10.2139/ssrn.3606796.
9. Wang, D. *et al.* Comprehensive functional genomic resource and integrative model for the human brain. *Science* **362**, (2018).
 10. Heinz, S., Romanoski, C. E., Benner, C. & Glass, C. K. The selection and function of cell type-specific enhancers. *Nat. Rev. Mol. Cell Biol.* **16**, 144–154 (2015).
 11. Bryois, J. *et al.* Genetic identification of cell types underlying brain complex traits yields insights into the etiology of Parkinson's disease. *Nat. Genet.* **52**, 482–493 (2020).
 12. Lopategui Cabezas, I., Herrera Batista, A. & Pentón Rol, G. The role of glial cells in Alzheimer disease: potential therapeutic implications. *Neurologia* **29**, 305–309 (2014).
 13. Cooper-Knock, J., Jenkins, T. & Shaw, P. J. *Clinical and Molecular Aspects of Motor Neuron Disease*. (Biota Publishing, 2013).
 14. Corces, M. R. *et al.* Single-cell epigenomic analyses implicate candidate causal variants at inherited risk loci for Alzheimer's and Parkinson's diseases. *Nat. Genet.* **52**, 1158–1168 (2020).
 15. Frey, D. *et al.* Early and selective loss of neuromuscular synapse subtypes with low sprouting competence in motoneuron diseases. *J. Neurosci.* **20**, 2534–2542 (2000).
 16. Moloney, E. B., de Winter, F. & Verhaagen, J. ALS as a distal axonopathy: molecular mechanisms affecting neuromuscular junction stability in the presymptomatic stages of the disease. *Front. Neurosci.* **8**, 252 (2014).
 17. Wang, Z., Gerstein, M. & Snyder, M. RNA-Seq: a revolutionary tool for transcriptomics. *Nat. Rev. Genet.* **10**, 57–63 (2009).
 18. Buenrostro, J. D., Wu, B., Chang, H. Y. & Greenleaf, W. J. ATAC-seq: A Method for Assaying Chromatin Accessibility Genome-Wide. *Curr. Protoc. Mol. Biol.* **109**, 21.29.1–21.29.9 (2015).

19. Creyghton, M. P. *et al.* Histone H3K27ac separates active from poised enhancers and predicts developmental state. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 21931–21936 (2010).
20. van Berkum, N. L. *et al.* Hi-C: a method to study the three-dimensional architecture of genomes. *J. Vis. Exp.* (2010) doi:10.3791/1869.
21. ENCODE Project Consortium *et al.* Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature* **583**, 699–710 (2020).
22. Pennacchio, L. A., Bickmore, W., Dean, A., Nobrega, M. A. & Bejerano, G. Enhancers: five essential questions. *Nat. Rev. Genet.* **14**, 288–295 (2013).
23. Loh, P.-R. *et al.* Contrasting genetic architectures of schizophrenia and other complex diseases using fast variance-components analysis. *Nat. Genet.* **47**, 1385–1392 (2015).
24. Consortium, T. 1000 G. P. & The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature* vol. 526 68–74 (2015).
25. DeJesus-Hernandez, M. *et al.* Expanded GGGGCC hexanucleotide repeat in noncoding region of C9ORF72 causes chromosome 9p-linked FTD and ALS. *Neuron* **72**, 245–256 (2011).
26. Elden, A. C. *et al.* Ataxin-2 intermediate-length polyglutamine expansions are associated with increased risk for ALS. *Nature* **466**, 1069–1075 (2010).
27. Landrum, M. J. *et al.* ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic Acids Res.* **46**, D1062–D1067 (2018).
28. Daoud, H. *et al.* Analysis of the UNC13A gene as a risk factor for sporadic amyotrophic lateral sclerosis. *Arch. Neurol.* **67**, 516–517 (2010).
29. Diekstra, F. P. *et al.* UNC13A is a modifier of survival in amyotrophic lateral sclerosis. *Neurobiol. Aging* **33**, 630.e3–8 (2012).
30. Huang, N., Lee, I., Marcotte, E. M. & Hurles, M. E. Characterising and predicting

- haploinsufficiency in the human genome. *PLoS Genet.* **6**, e1001154 (2010).
31. Fadista, J., Oskolkov, N., Hansson, O. & Groop, L. LoFtool: a gene intolerance score based on loss-of-function variants in 60 706 individuals. *Bioinformatics* **33**, 471–474 (2017).
 32. Petrovski, S., Wang, Q., Heinzen, E. L., Allen, A. S. & Goldstein, D. B. Genic intolerance to functional variation and the interpretation of personal genomes. *PLoS Genet.* **9**, e1003709 (2013).
 33. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans. *bioRxiv* 531210 (2020) doi:10.1101/531210.
 34. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).
 35. Prudencio, M. *et al.* Distinct brain transcriptome profiles in C9orf72-associated and sporadic ALS. *Nat. Neurosci.* **18**, 1175–1182 (2015).
 36. Philips, T. & Rothstein, J. D. Rodent Models of Amyotrophic Lateral Sclerosis. *Current Protocols in Pharmacology* vol. 69 (2015).
 37. Maniatis, S. *et al.* Spatiotemporal dynamics of molecular pathology in amyotrophic lateral sclerosis. *Science* **364**, 89–93 (2019).
 38. Kuleshov, M. V. *et al.* Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res.* **44**, W90–7 (2016).
 39. Szklarczyk, D. *et al.* STRING v11: protein-protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids Res.* **47**, D607–D613 (2019).
 40. Krishnan, A. *et al.* Genome-wide prediction and functional characterization of the genetic basis of autism spectrum disorder. *Nat. Neurosci.* **19**, 1454–1462 (2016).
 41. Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities

- in large networks. *Journal of Statistical Mechanics: Theory and Experiment* vol. 2008 P10008 (2008).
42. Lerer, I. *et al.* Deletion of the ANKRD15 gene at 9p24.3 causes parent-of-origin-dependent inheritance of familial cerebral palsy. *Hum. Mol. Genet.* **14**, 3911–3920 (2005).
 43. Project MinE ALS Sequencing Consortium. Project MinE: study design and pilot analyses of a large-scale whole-genome sequencing study in amyotrophic lateral sclerosis. *Eur. J. Hum. Genet.* **26**, 1537–1546 (2018).
 44. Huang, Y.-F., Gulko, B. & Siepel, A. Fast, scalable prediction of deleterious noncoding variants from functional and population genomic data. *Nat. Genet.* **49**, 618–624 (2017).
 45. Rentzsch, P., Witten, D., Cooper, G. M., Shendure, J. & Kircher, M. CADD: predicting the deleteriousness of variants throughout the human genome. *Nucleic Acids Res.* **47**, D886–D894 (2019).
 46. Ritchie, G. R., Dunham, I., Zeggini, E. & Flicek, P. Functional annotation of noncoding sequence variants. *Nat. Methods* **11**, 294–296 (2014).
 47. Fishilevich, S. *et al.* GeneHancer: genome-wide integration of enhancers and target genes in GeneCards. *Database* **2017**, (2017).
 48. Wu, M. C. *et al.* Rare-variant association testing for sequencing data with the sequence kernel association test. *Am. J. Hum. Genet.* **89**, 82–93 (2011).
 49. Lee, S. *et al.* Optimal unified approach for rare-variant association testing with application to small-sample case-control whole-exome sequencing studies. *Am. J. Hum. Genet.* **91**, 224–237 (2012).
 50. Whitlock, M. C. Combining probability from independent tests: the weighted Z-method is superior to Fisher's approach. *J. Evol. Biol.* **18**, 1368–1373 (2005).
 51. Zheng, Q. *et al.* Precise gene deletion and replacement using the CRISPR/Cas9 system in

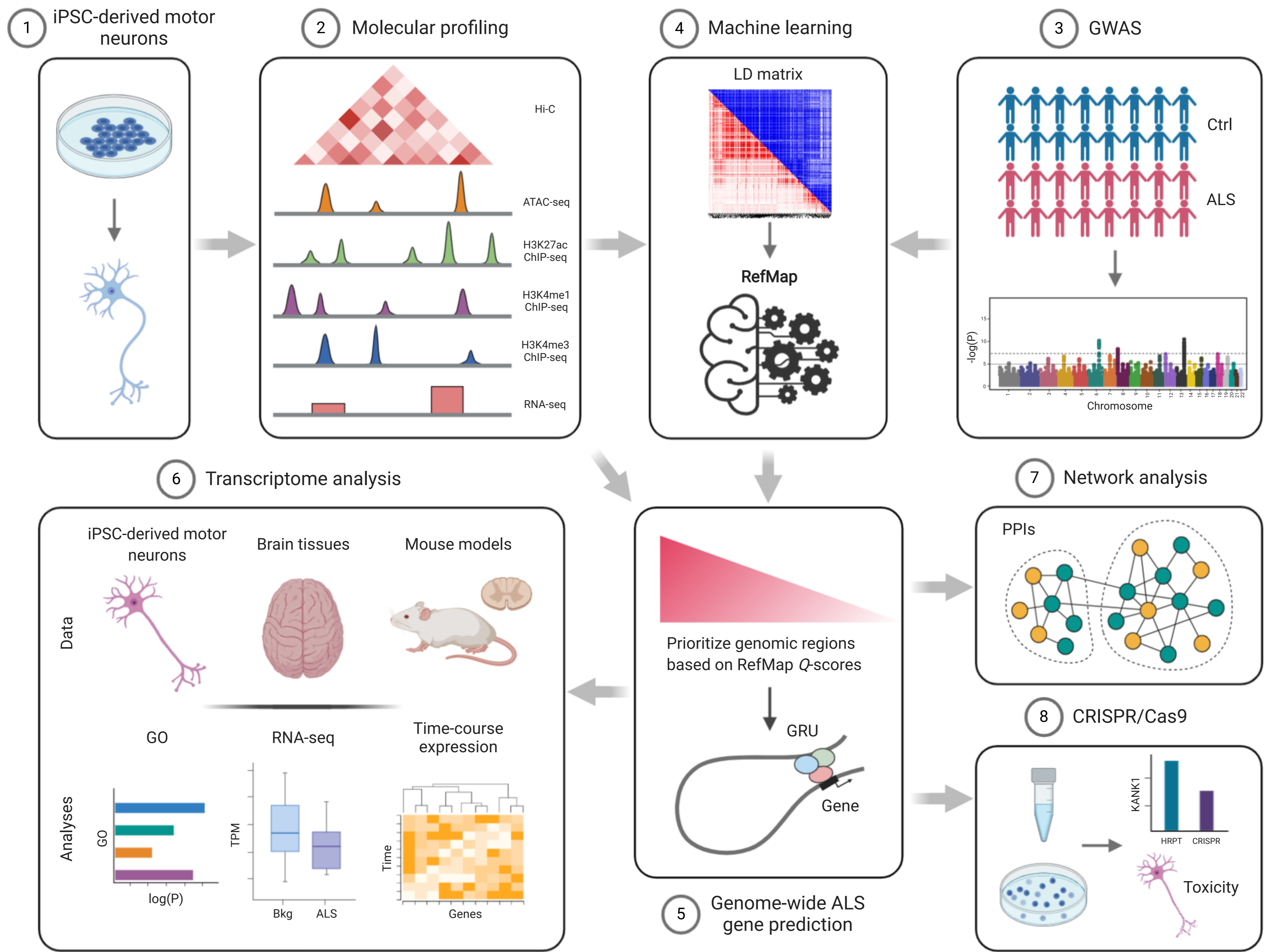
- human cells. *Biotechniques* **57**, 115–124 (2014).
52. Hsiao, T. *et al.* Inference of CRISPR Edits from Sanger Trace Data. doi:10.1101/251082.
 53. Forster, J. I. *et al.* Characterization of Differentiated SH-SY5Y as Neuronal Screening Model Reveals Increased Oxidative Vulnerability. *J. Biomol. Screen.* **21**, 496–509 (2016).
 54. Pickrell, J. K. Joint analysis of functional genomic data and genome-wide association studies of 18 human traits. *Am. J. Hum. Genet.* **94**, 559–573 (2014).
 55. Hormozdiari, F., Kostem, E., Kang, E. Y., Pasaniuc, B. & Eskin, E. Identifying causal variants at loci with multiple signals of association. *Genetics* **198**, 497–508 (2014).
 56. Kichaev, G. *et al.* Integrating functional data to prioritize causal variants in statistical fine-mapping studies. *PLoS Genet.* **10**, e1004722 (2014).
 57. Chen, W., McDonnell, S. K., Thibodeau, S. N., Tillmans, L. S. & Schaid, D. J. Incorporating Functional Annotations for Fine-Mapping Causal Variants in a Bayesian Framework Using Summary Statistics. *Genetics* vol. 204 933–958 (2016).
 58. Benner, C. *et al.* FINEMAP: efficient variable selection using summary data from genome-wide association studies. *Bioinformatics* **32**, 1493–1501 (2016).
 59. Schaid, D. J., Chen, W. & Larson, N. B. From genome-wide associations to candidate causal variants by statistical fine-mapping. *Nat. Rev. Genet.* **19**, 491–504 (2018).
 60. Song, M. *et al.* Mapping cis-regulatory chromatin contacts in neural cells links neuropsychiatric disorder risk variants to target genes. *Nat. Genet.* **51**, 1252–1262 (2019).
 61. De Vos, K. J. & Hafezparast, M. Neurobiology of axonal transport defects in motor neuron diseases: Opportunities for translational research? *Neurobiol. Dis.* **105**, 283–299 (2017).
 62. Boopathy, S. *et al.* Structural basis for mutation-induced destabilization of profilin 1 in ALS. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 7984–7989 (2015).
 63. Dillon, C. & Goda, Y. The actin cytoskeleton: integrating form and function at the synapse.

- Annu. Rev. Neurosci.* **28**, 25–55 (2005).
64. Mallik, B. & Kumar, V. Regulation of actin-Spectrin cytoskeleton by ICA69 at the *Drosophila* neuromuscular junction. *null* **11**, e1381806 (2018).
 65. Giampetruzzi, A. *et al.* Modulation of actin polymerization affects nucleocytoplasmic transport in multiple forms of amyotrophic lateral sclerosis. *Nat. Commun.* **10**, 3827 (2019).
 66. Brooks, B. R. El Escorial World Federation of Neurology criteria for the diagnosis of amyotrophic lateral sclerosis. Subcommittee on Motor Neuron Diseases/Amyotrophic Lateral Sclerosis of the World Federation of Neurology Research Group on Neuromuscular Diseases and the El Escorial ‘Clinical limits of amyotrophic lateral sclerosis’ workshop contributors. *J. Neurol. Sci.* **124 Suppl**, 96–107 (1994).
 67. Du, Z.-W. *et al.* Generation and expansion of highly pure motor neuron progenitors from human pluripotent stem cells. *Nat. Commun.* **6**, 6626 (2015).
 68. Rao, S. S. P. *et al.* Cohesin Loss Eliminates All Loop Domains. *Cell* **171**, 305–320.e24 (2017).
 69. Rao, S. S. P. *et al.* A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell* **159**, 1665–1680 (2014).
 70. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal* **17**, 10–12 (2011).
 71. Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic RNA-seq quantification. *Nat. Biotechnol.* **34**, 525–527 (2016).
 72. Harva, M. & Kabán, A. Variational learning for rectified factor analysis. *Signal Processing* vol. 87 509–527 (2007).
 73. Blei, D. M., Kucukelbir, A. & McAuliffe, J. D. Variational Inference: A Review for Statisticians. *Journal of the American Statistical Association* vol. 112 859–877 (2017).

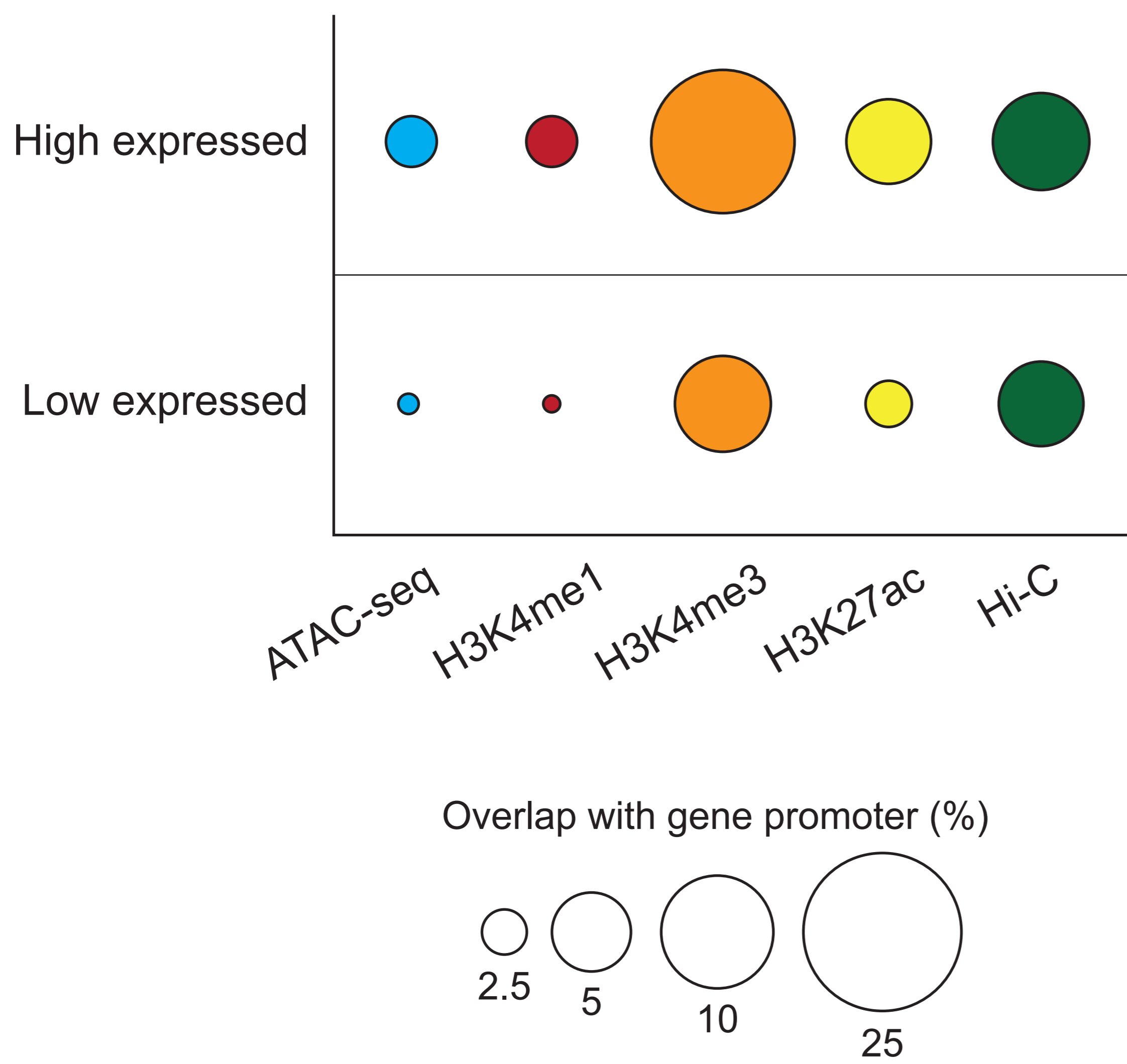
74. Wang, S., Cho, H., Zhai, C., Berger, B. & Peng, J. Exploiting ontology graph for predicting sparsely annotated gene function. *Bioinformatics* **31**, i357–64 (2015).
75. Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities in large networks. *J. Stat. Mech.* **2008**, P10008 (2008).
76. Milo, R. *et al.* Network motifs: simple building blocks of complex networks. *Science* **298**, 824–827 (2002).
77. Chen, L., Jin, P. & Qin, Z. S. DIVAN: accurate identification of non-coding disease-specific risk variants using multi-omics profiles. *Genome Biol.* **17**, 252 (2016).
78. Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly* **6**, 80–92 (2012).
79. Hoerndli, F. J., Toigo, M., Schild, A., Götz, J. & Day, P. J. Reference genes identified in SH-SY5Y cells using custom-made gene arrays with validation by quantitative polymerase chain reaction. *Anal. Biochem.* **335**, 30–41 (2004).
80. Schmittgen, T. D. & Livak, K. J. Analyzing real-time PCR data by the comparative C(T) method. *Nat. Protoc.* **3**, 1101–1108 (2008).
81. Longair, M. H., Baker, D. A. & Armstrong, J. D. Simple Neurite Tracer: open source software for reconstruction, visualization and analysis of neuronal processes. *Bioinformatics* **27**, 2453–2454 (2011).
82. Arganda-Carreras, I., Fernández-González, R., Muñoz-Barrutia, A. & Ortiz-De-Solorzano, C. 3D reconstruction of histological sections: Application to mammary gland tissue. *Microscopy Research and Technique* vol. 73 1019–1029 (2010).

Figure 1

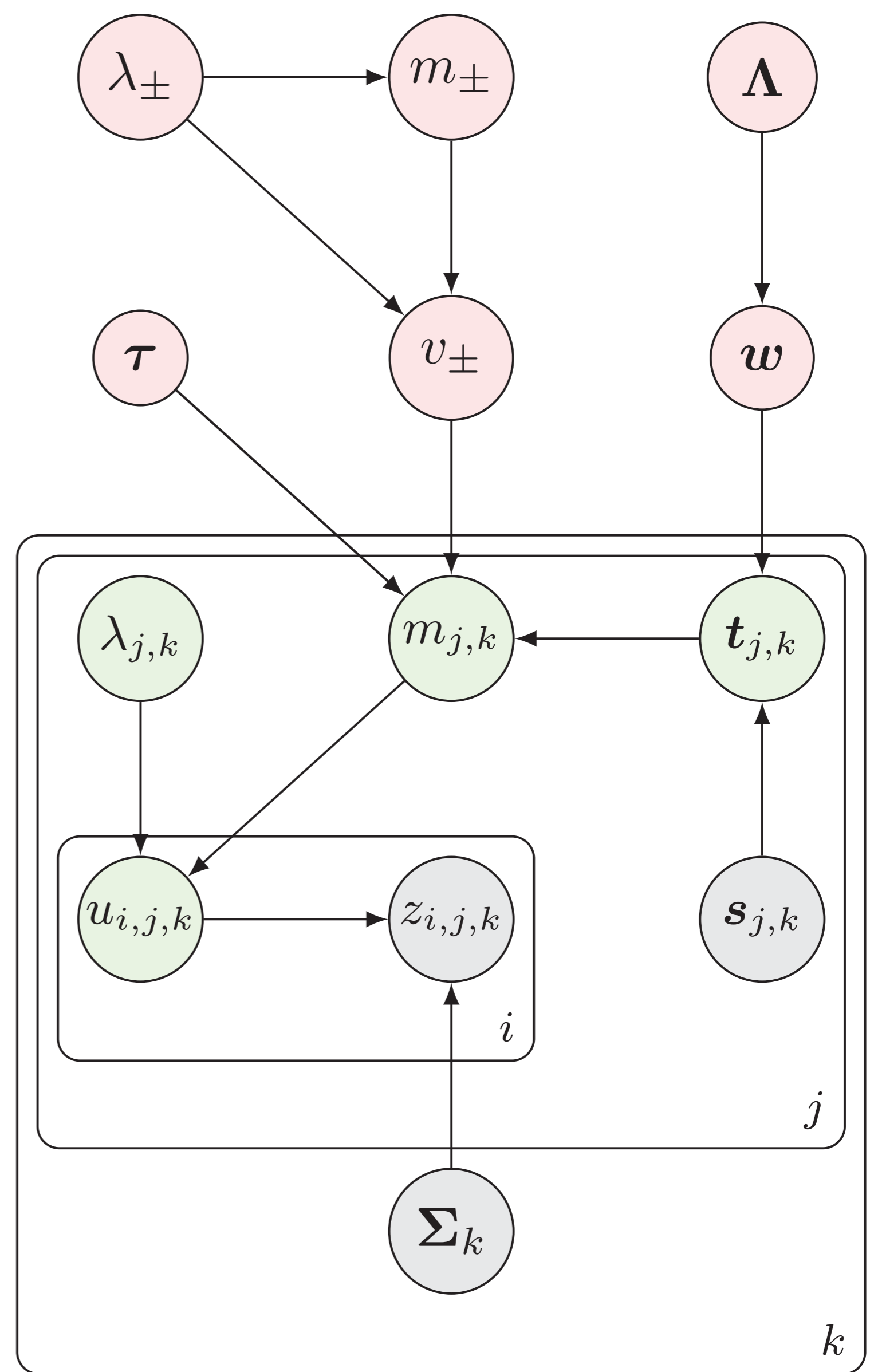
a



b



c



d

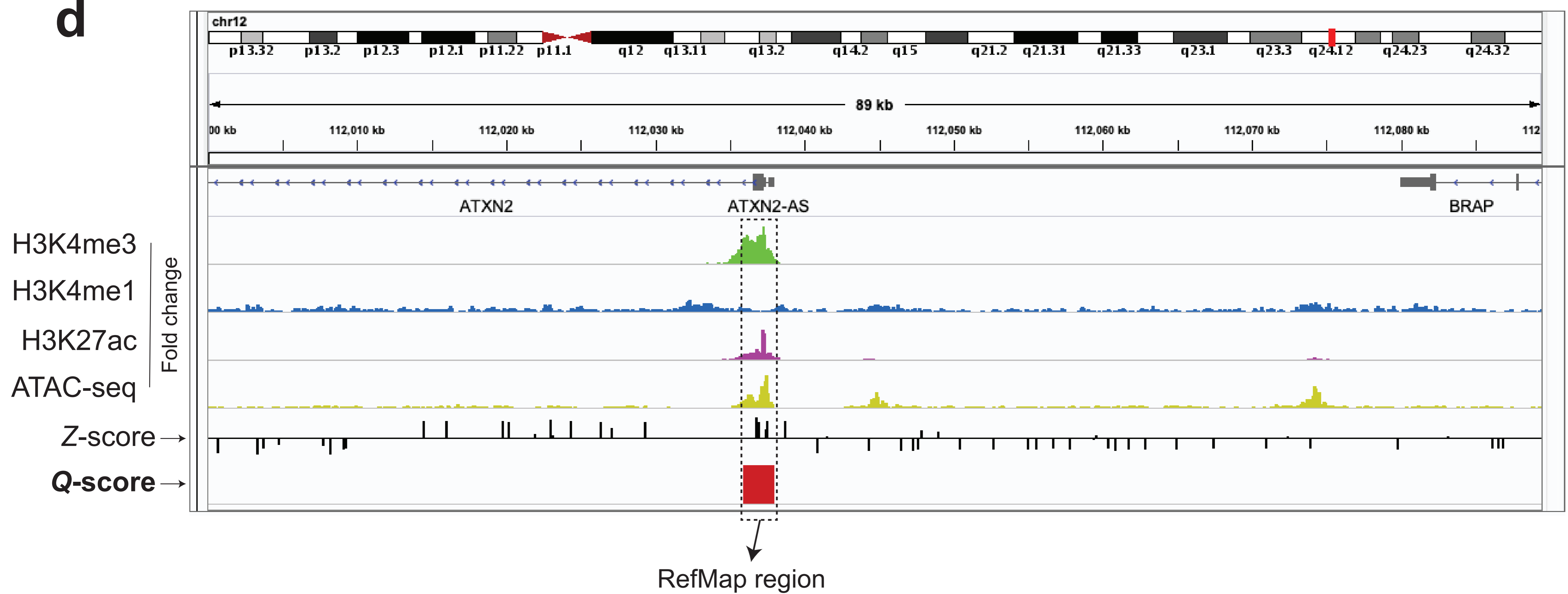
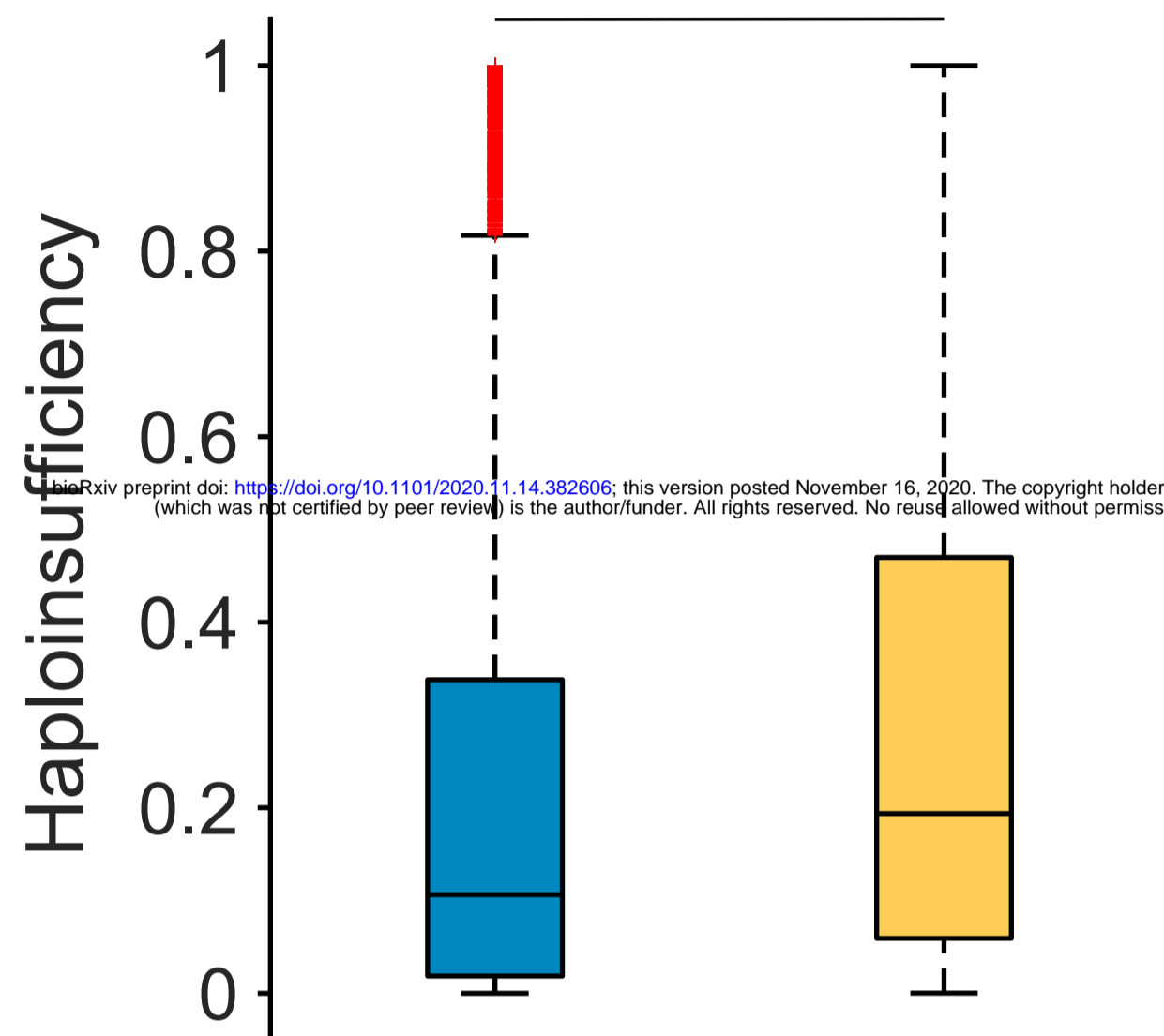


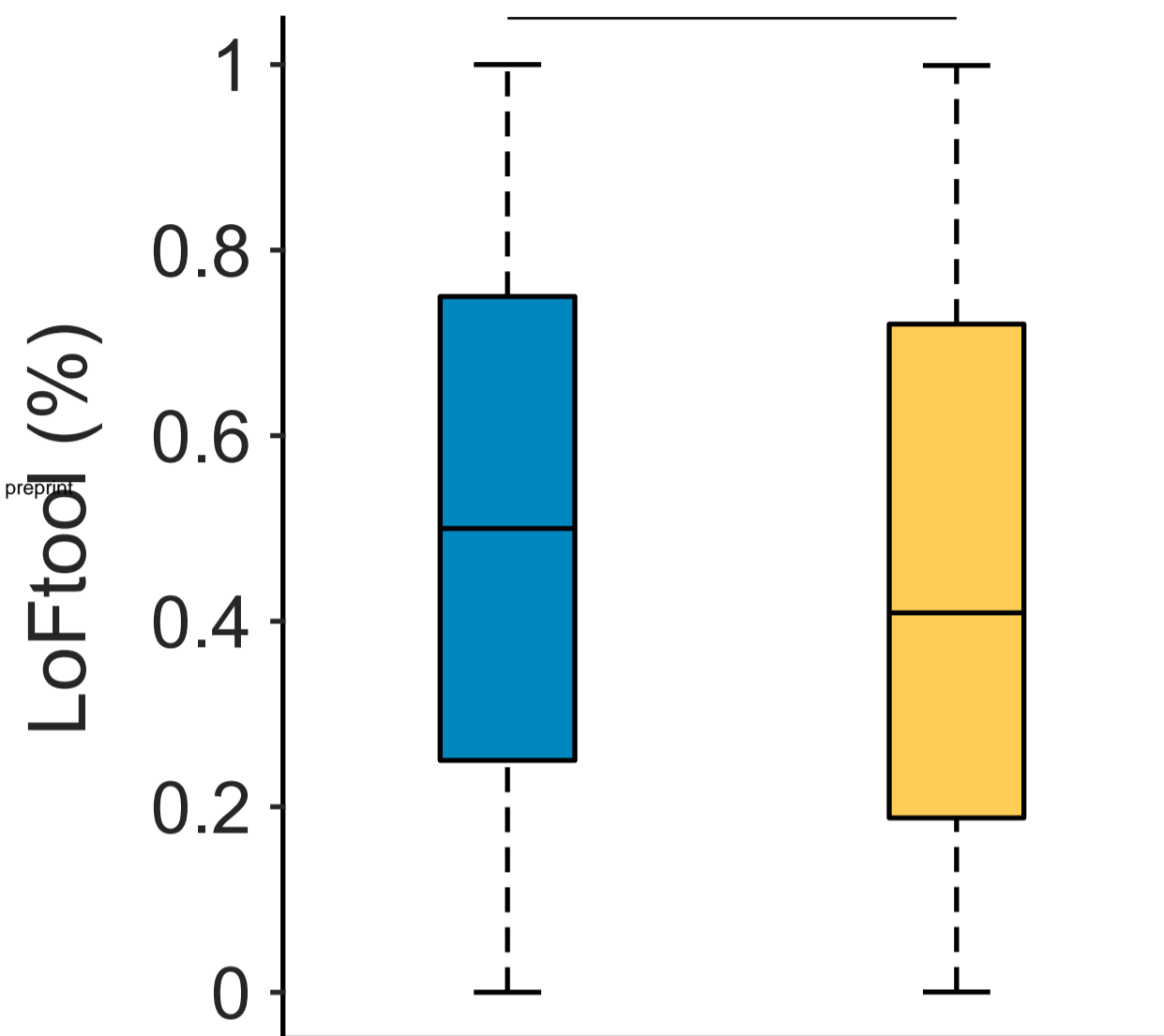
Figure 2

a

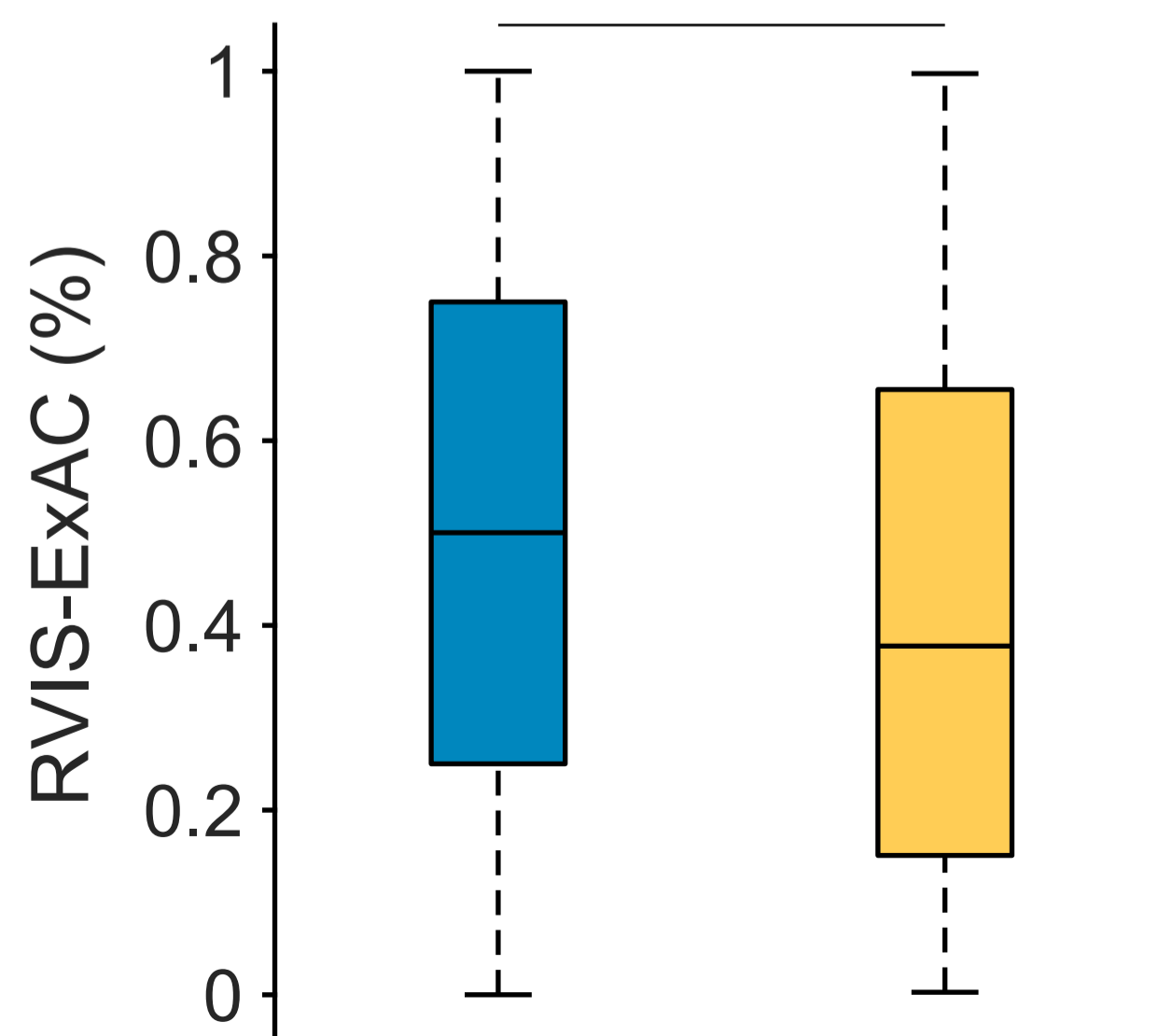
$$P = 2.59 \times 10^{-19}$$

**b**

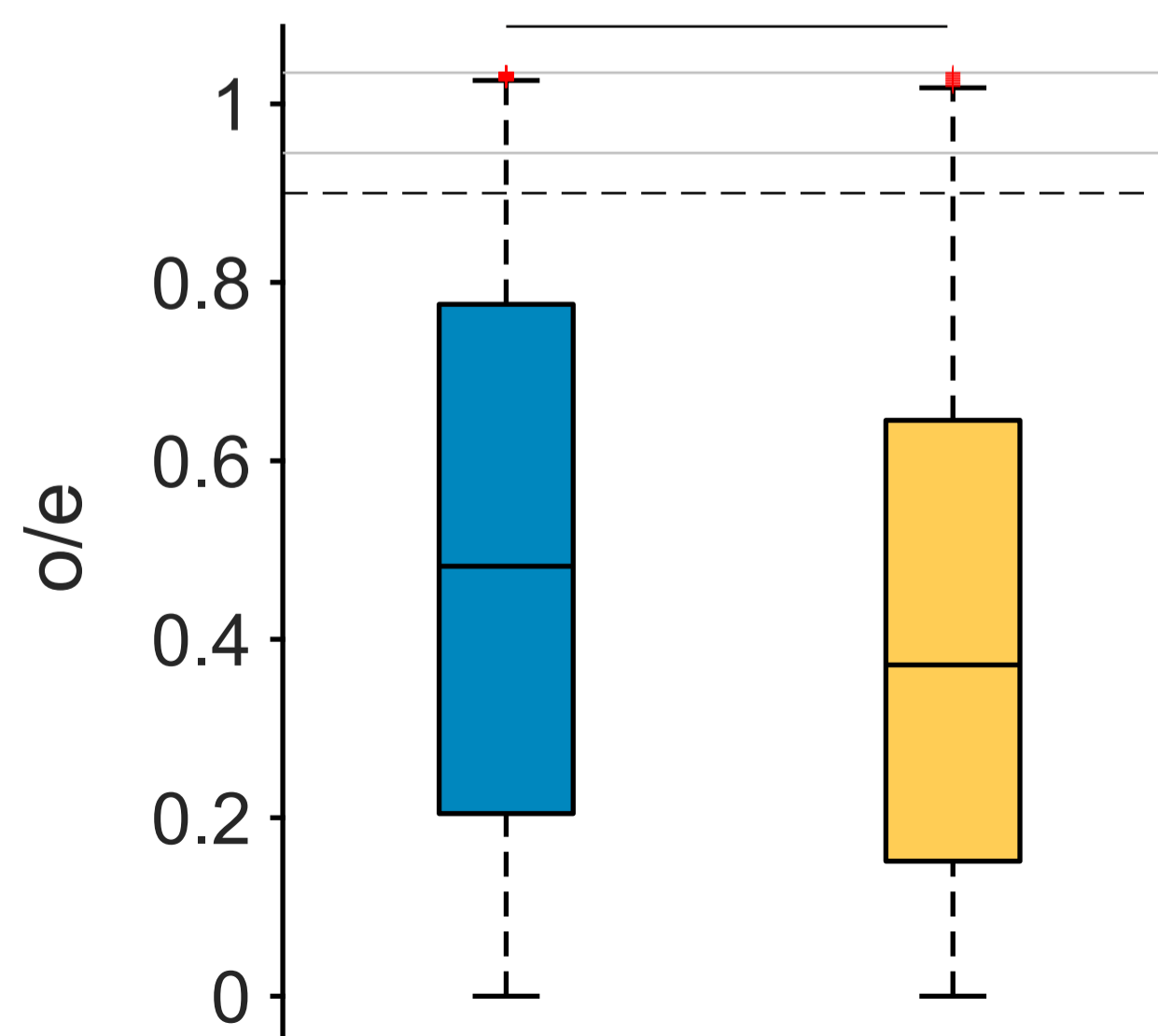
$$P = 2.28 \times 10^{-4}$$

**c**

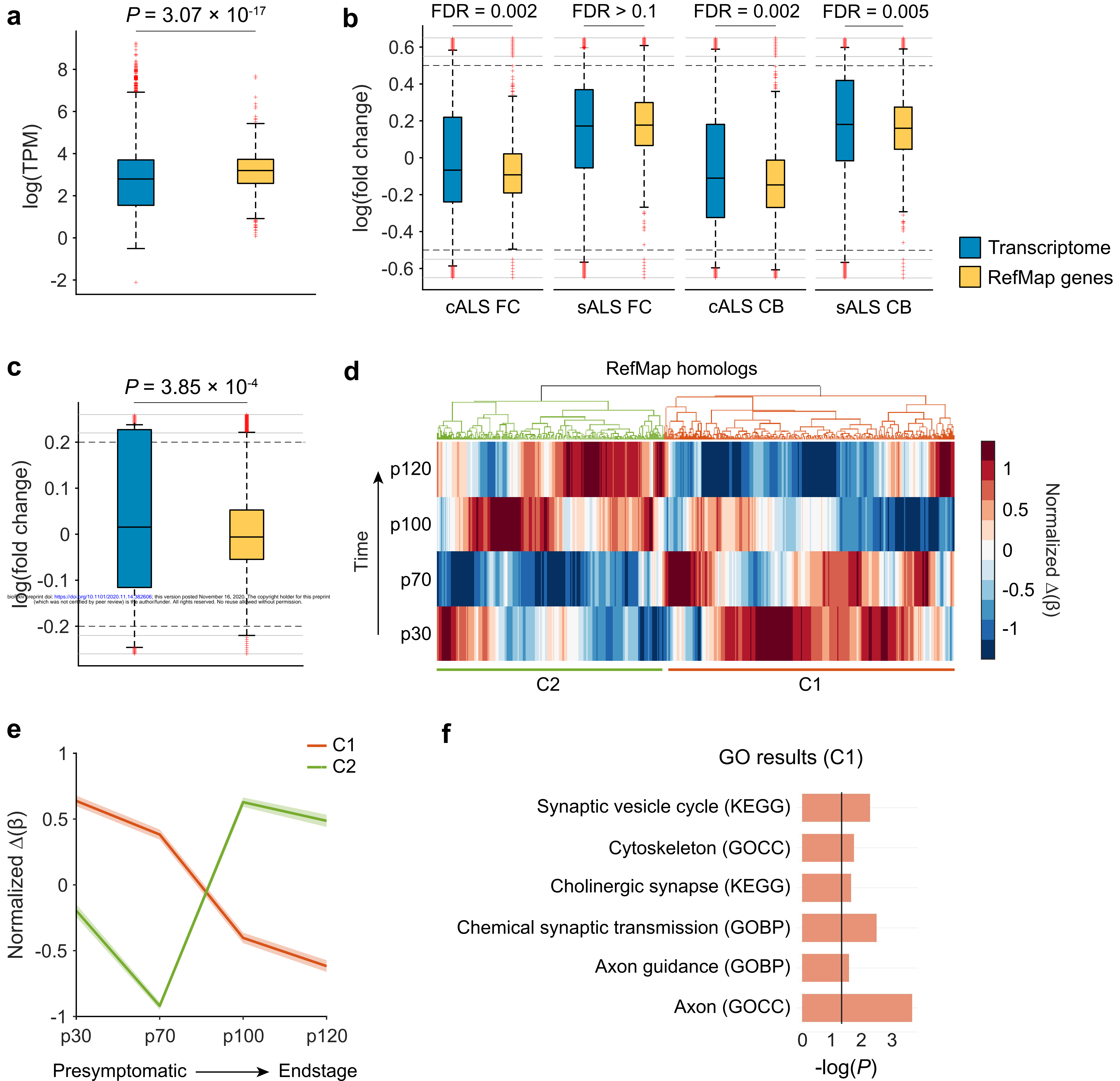
$$P = 8.08 \times 10^{-13}$$

**d**

$$P = 4.08 \times 10^{-10}$$



■ Transcriptome
■ RefMap genes

Figure 3

bioRxiv preprint doi: <https://doi.org/10.1101/2020.11.14.382606>; this version posted November 16, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

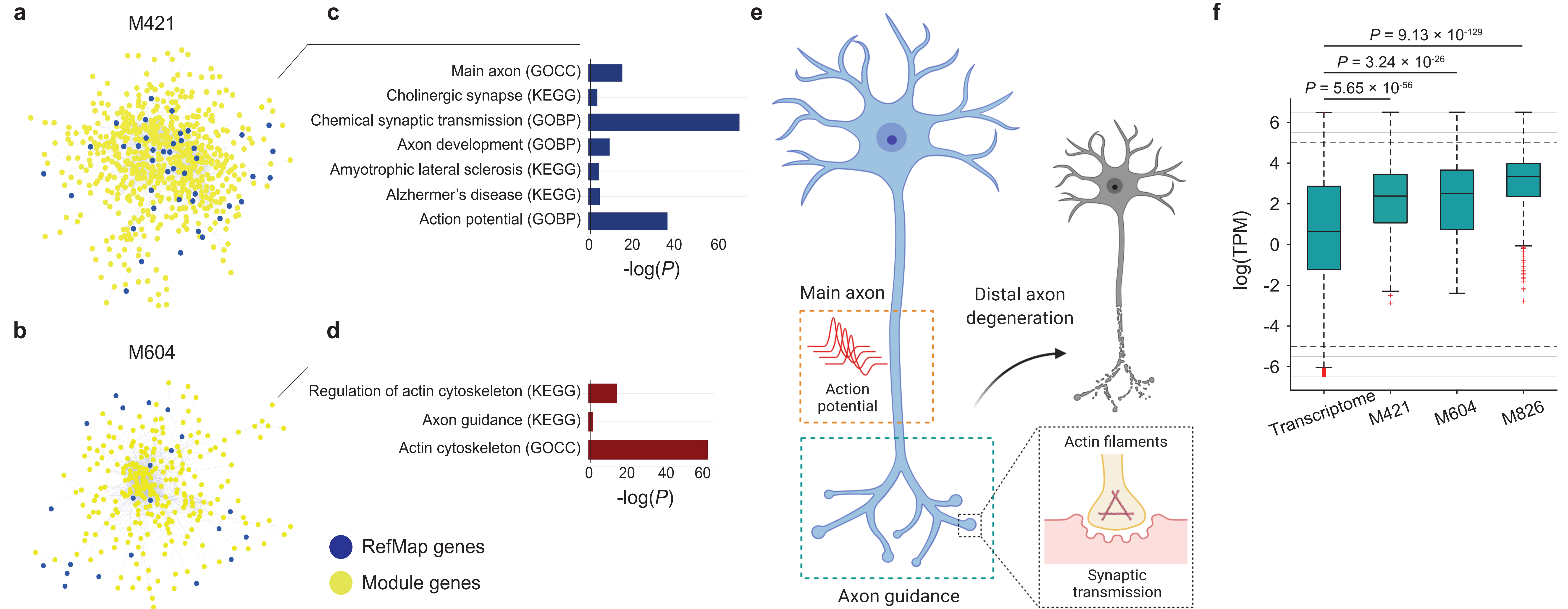
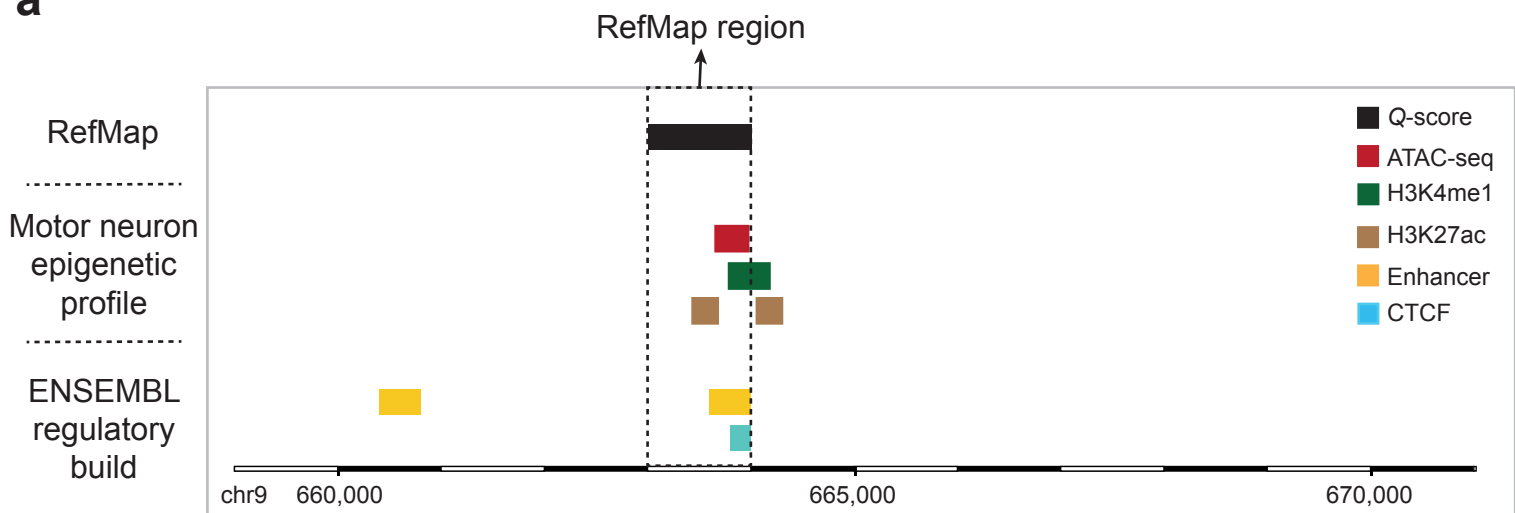
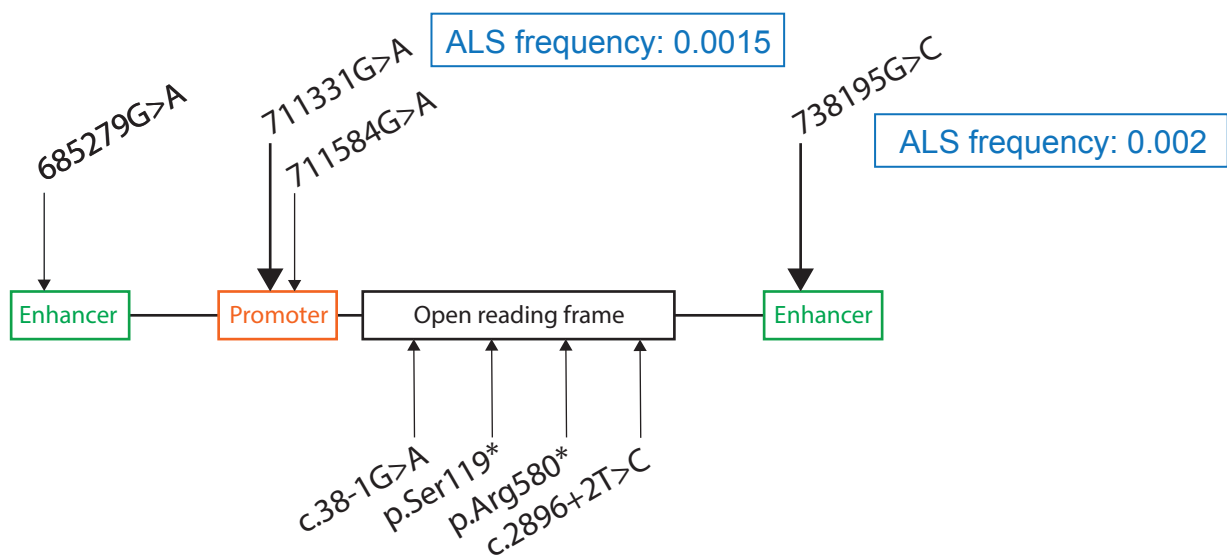
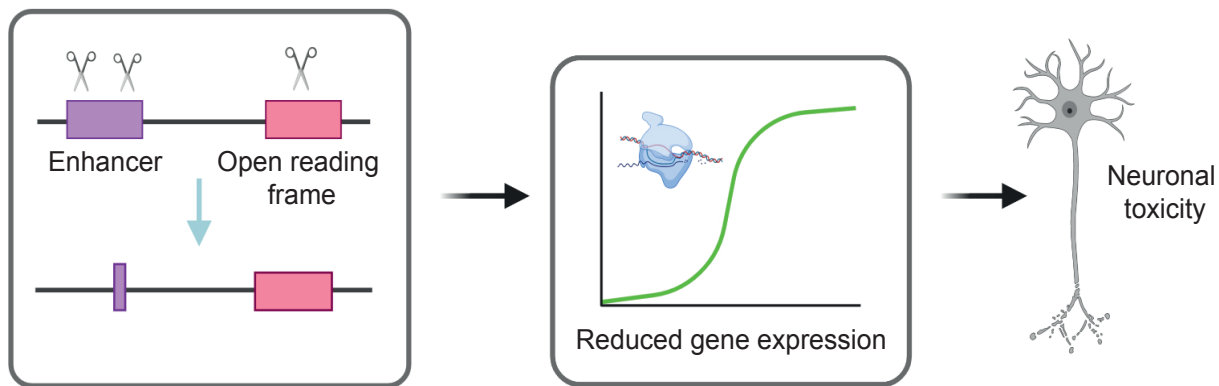
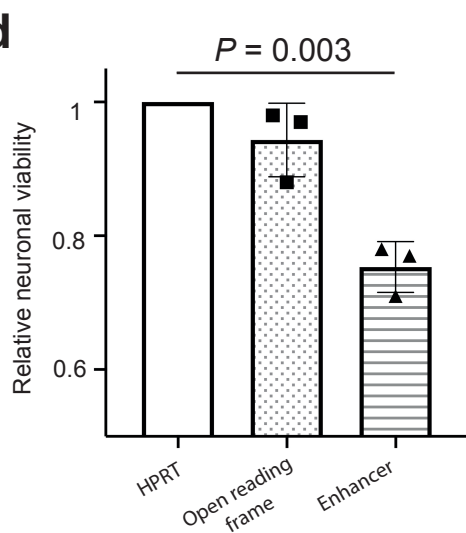
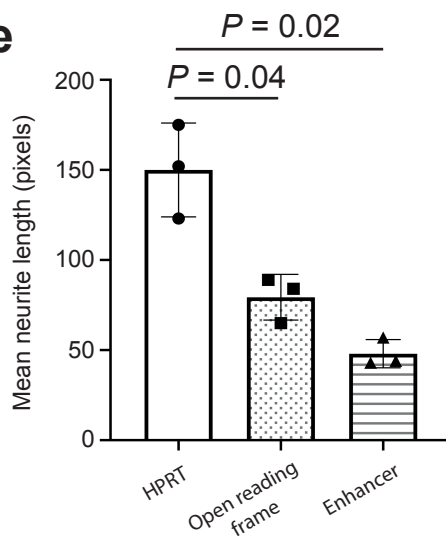
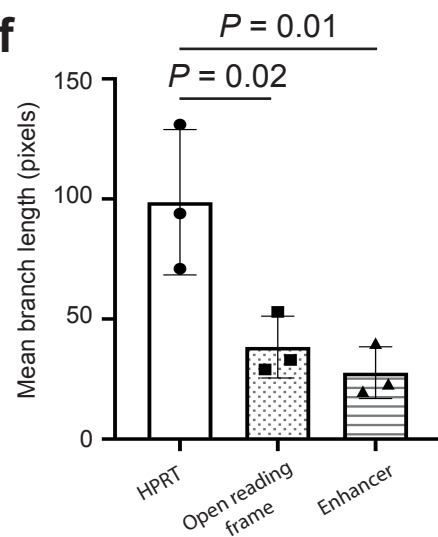
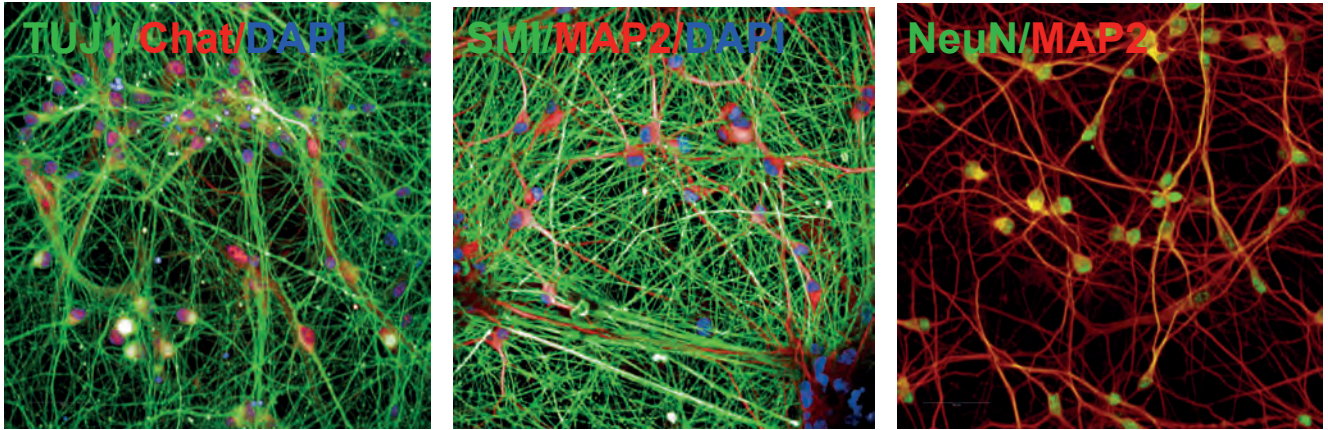
Figure 4

Figure 5**a****b****c****d****e****f**

Supplementary Figure 1

a



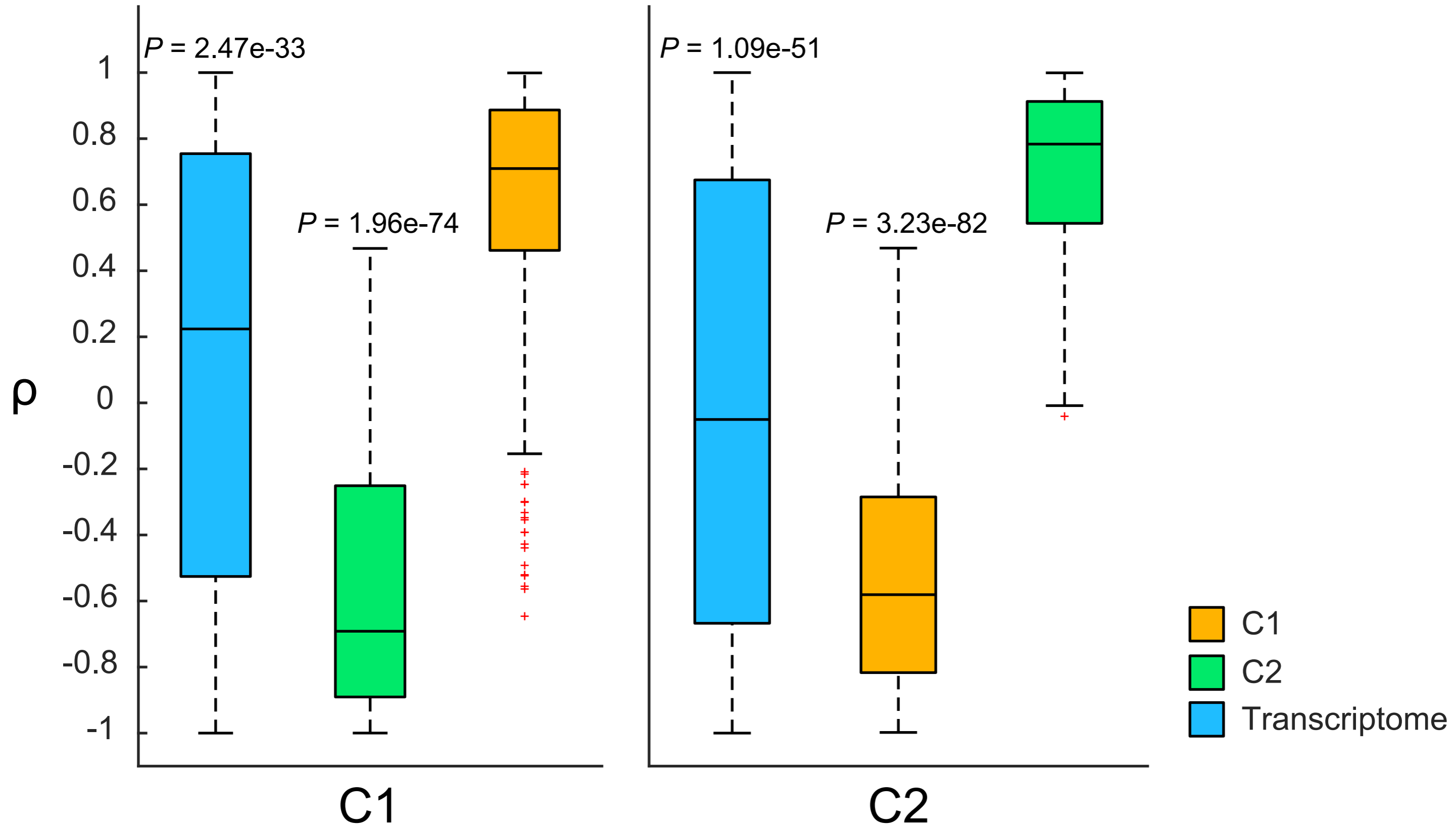
b

Cell Line	Short name	Source Tissue	Gender	Age at Sampling	Biobank
CS00iCTR-nxx	SMA	Fibroblast	Male	6	Cedars-Sinai
14iCTR-21nxx	CS14	Fibroblast	Female	52	Cedars-Sinai
GM23338	PGP	Fibroblast	Male	55	Coriell

Supplementary Figure 2

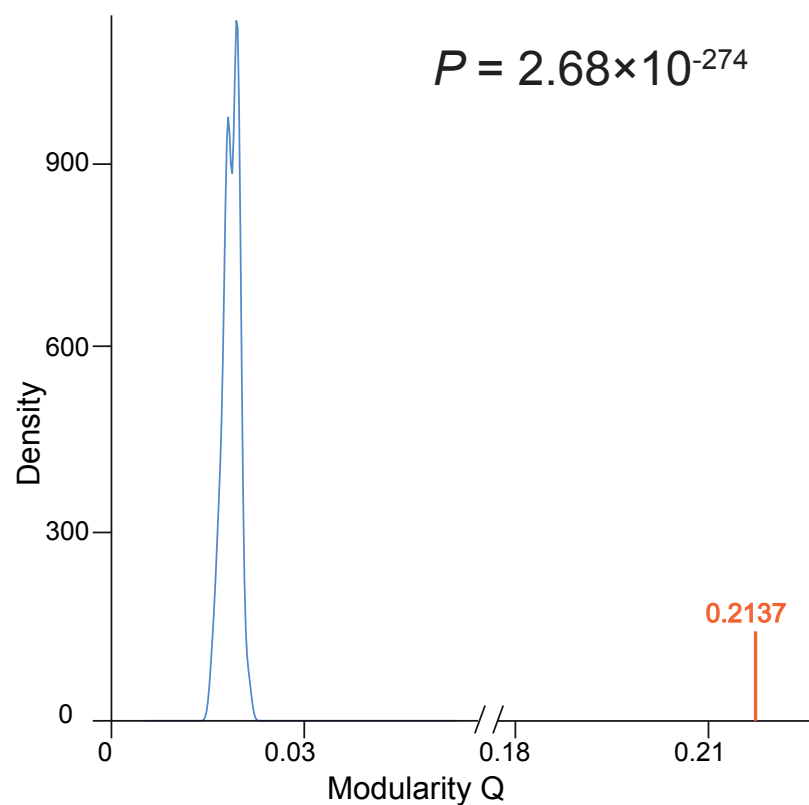
a

b

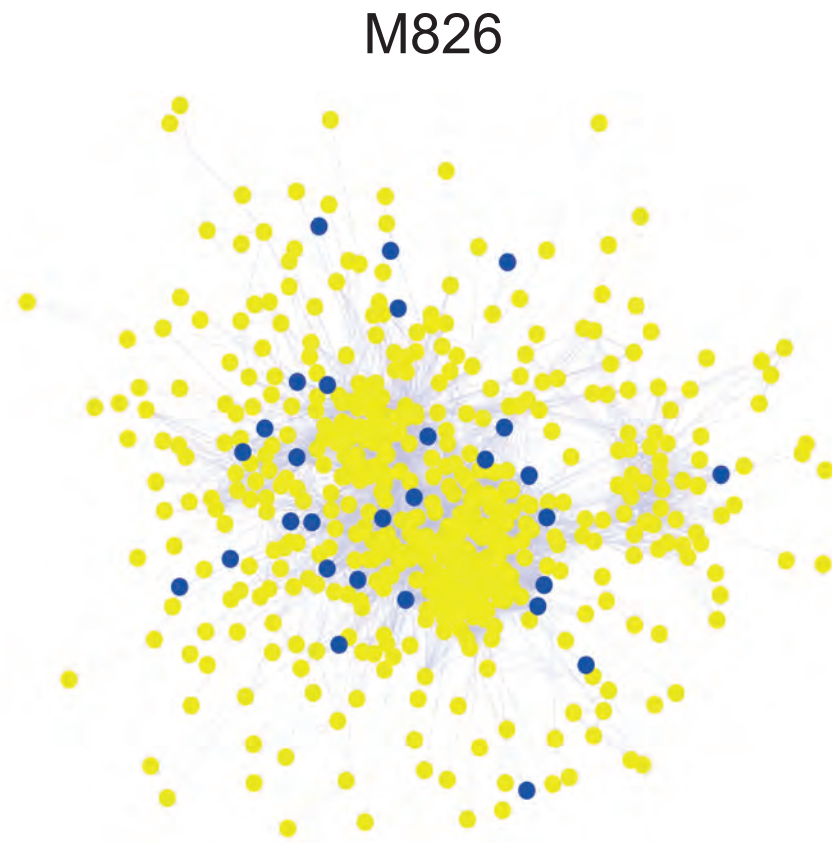


Supplementary Figure 3

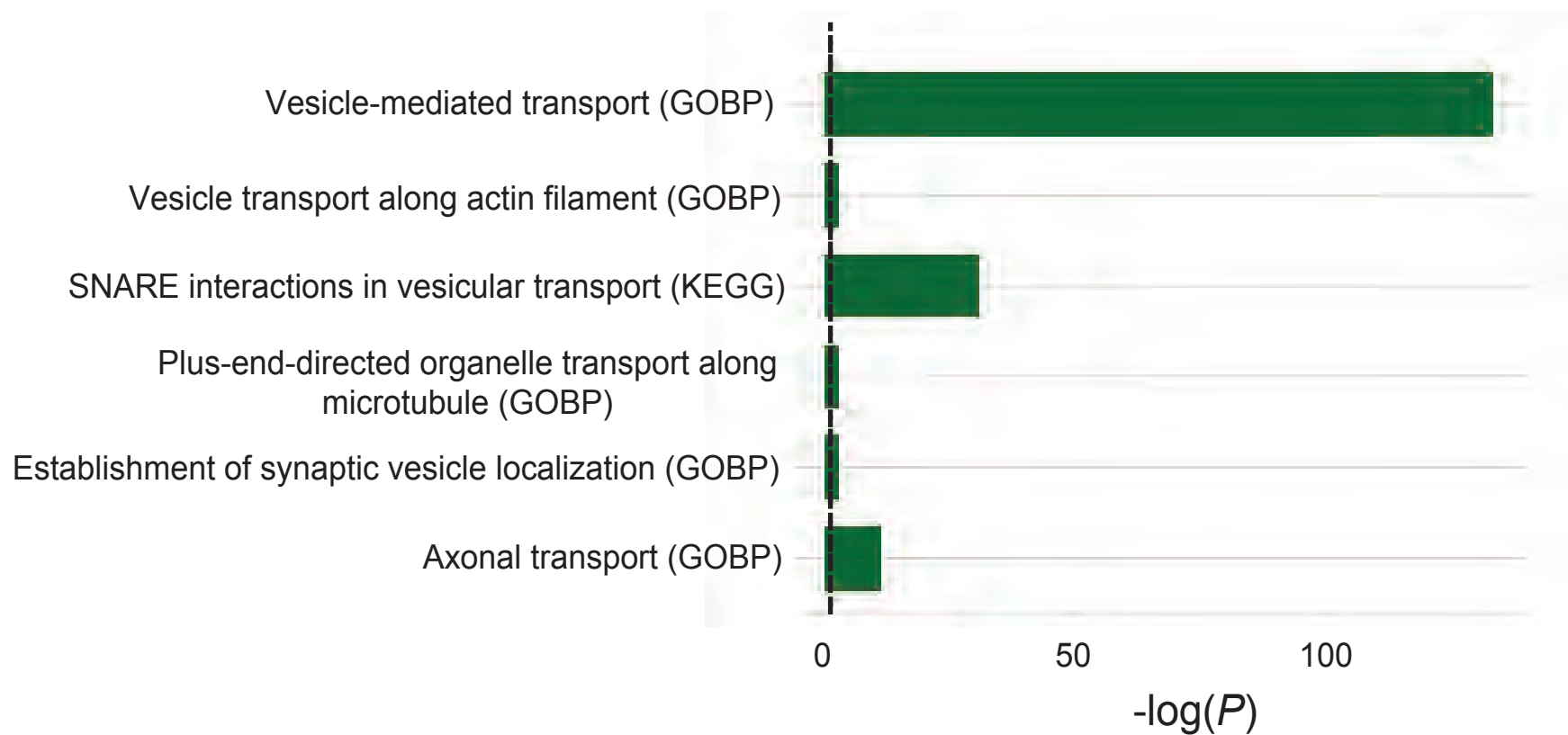
a



b

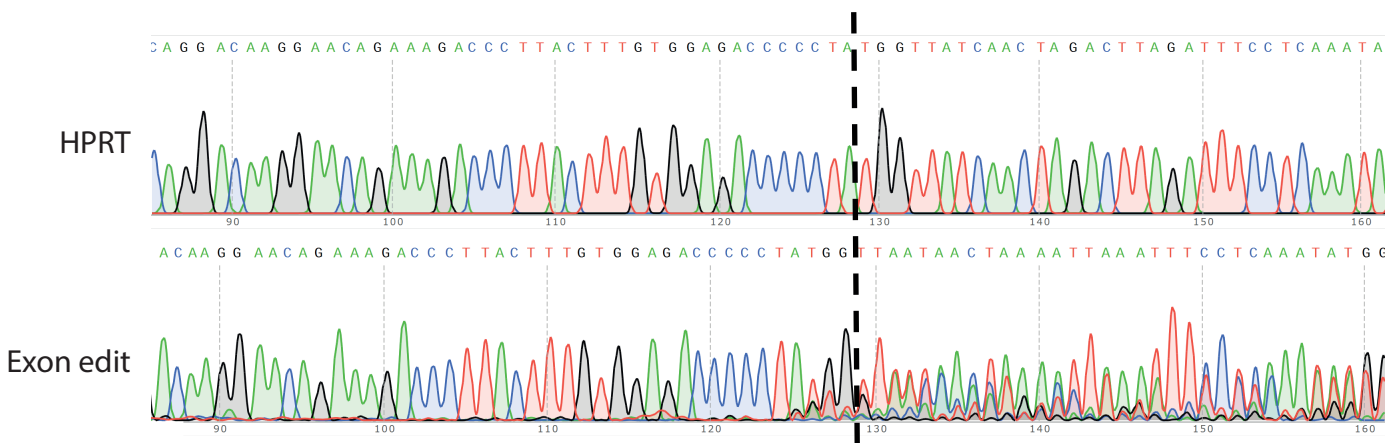


c

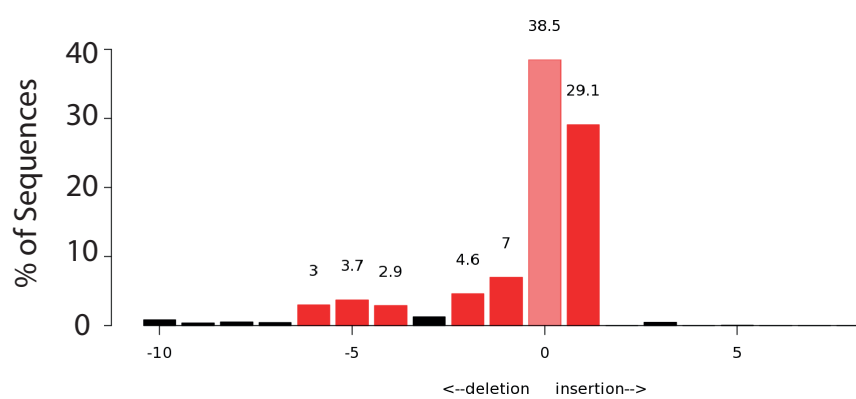


Supplementary Figure 4

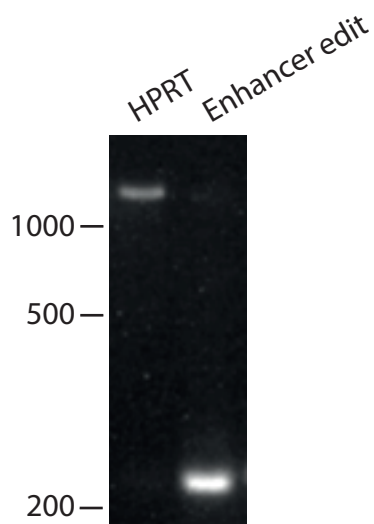
a



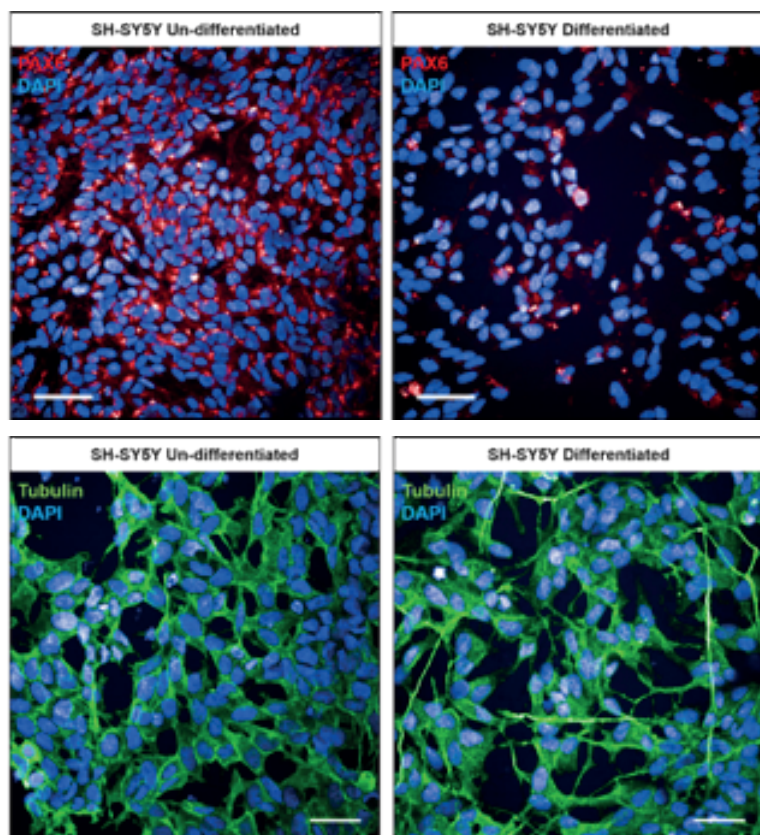
b



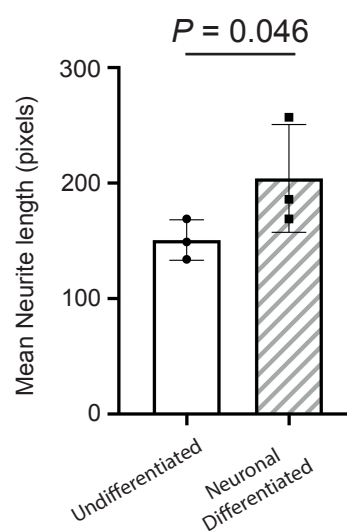
c



d



e



f

