

Supplementary figure legends

Supplementary Figure 1 | Polysome profiling of HCE KO versus wild type cells

a, Induction of expression of *Chrdl1*, *Dlx1*, and *Sema3a* by retinoic acid treatment of mESCs. On the x-axis is the time after the addition of retinoic acid. Y-axis plots the mRNA expression level normalized to the maximum value for each gene (mean, n=3). Error bars indicate standard deviation.

b-f, Polysome traces of wild-type versus HCE knockout cells. Y-axis plots the relative A260 units and X-axis is the time along the fractionation.

Supplementary Figure 2-6 | Deleted regions within hyperconserved 5'UTRs for polysome profiling

Genome browser tracks illustrating the position, multiple sequence alignment, PhastCons scores, LOD \geq 500 PhastCons elements, and the location of the deletion. Aggregate (maximum signal across all datasets) RNA-seq and H3K4me3 ChIP-seq tracks are also shown.

Supplementary Figure 7 | *Csde1* 5'UTR icM² library mutation rates

Plot of per-nucleotide mutation rates in the mutagenesis library used in icM² analysis of *Csde1* 5'UTR. Dots are colored by nucleotides. Note that flanking regions have near-zero mutation rate because they are the primer regions used for amplicon sequencing. See In-cell mutate-and-map section in methods.

Supplementary table legends

Supplementary Table 1 | List of hyperconserved 5'UTRs

This table lists each of the set of 589 mouse 5'UTRs defined as hyperconserved.

Supplementary Table 2 | List of significant GO terms enriched by hyperconserved 5'UTRs

This table lists 225 significant GO terms enriched by hyperconserved 5'UTR. Related to **Fig. 1e**.

Supplementary Table 3 | List of known disease-associated variants intersecting hyperconserved 5'UTRs

This table lists 5 known disease-associated variants that overlap hyperconserved 5'UTRs.

Supplementary Table 4 | List of deleted regions in the hyperconserved 5'UTRs analyzed by polysome profiling

This table lists exact positions of the deleted conserved segments in the 5 hyperconserved 5'UTRs analyzed by polysome profiling. Related to **Fig. 2**.

Supplementary Table 5 | List of hyperconserved 5'UTRs with significant non-canonical translation initiation activity

This table lists 90 hyperconserved 5'UTRs with high non-canonical translation activation. Related to **Fig. 3a**.

Supplementary Table 6 | List of regions with significant ATP-dependent differential accessibility

This table lists per-nucleotide DMS accessibility profiles across all probed hyperconserved 5'UTR regions. Related to **Fig. 4**.

Supplementary Table 7 | List of significant MPO terms enriched by hyperconserved 5'UTR genes with ATP-dependent differential accessibility regions

This table lists mammalian phenotype ontology terms enriched by hyperconserved 5'UTR genes with at least one significant ATP-dependent differential accessibility region.

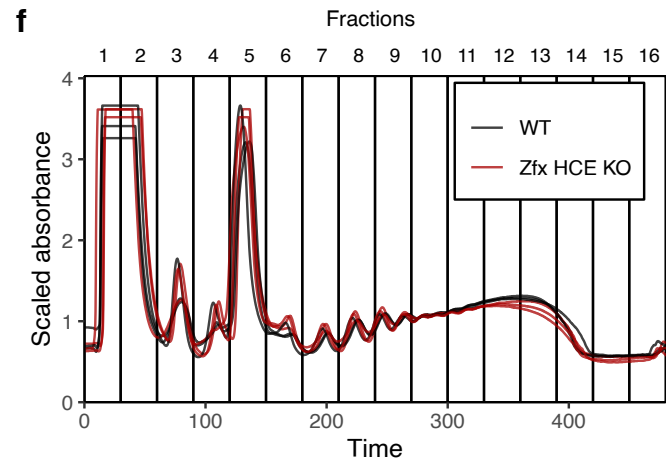
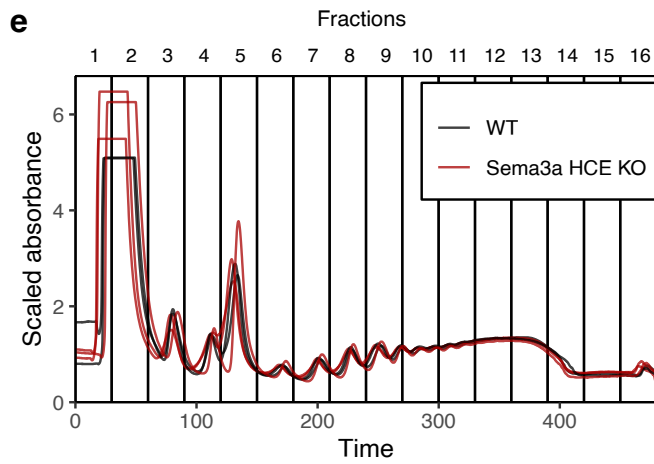
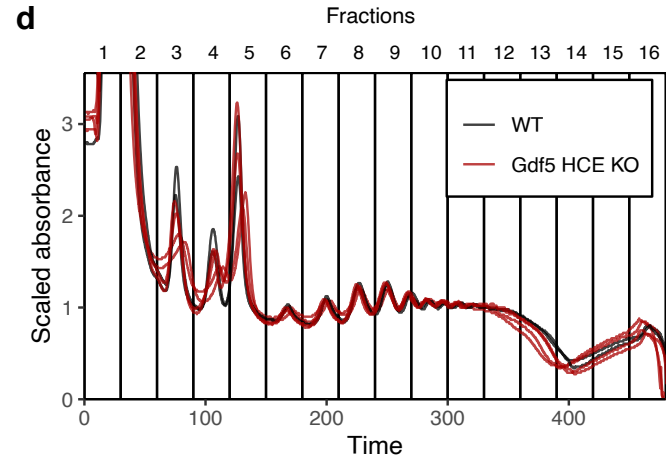
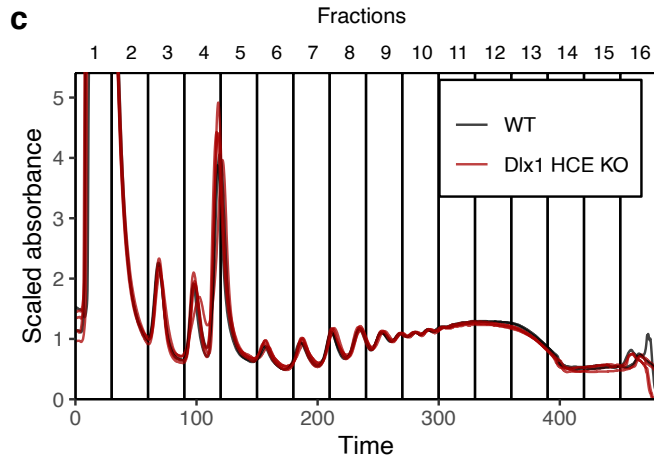
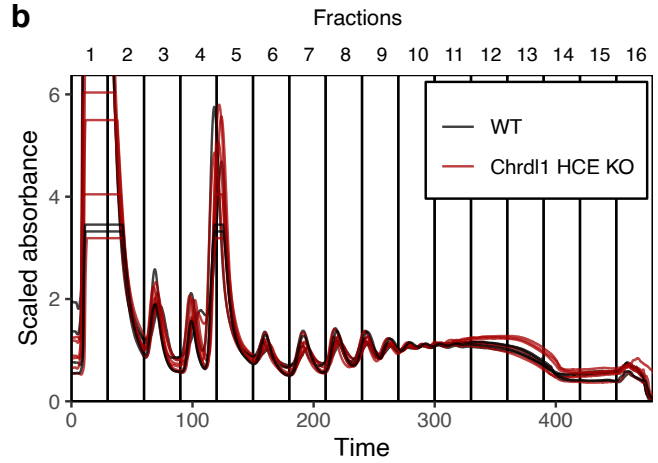
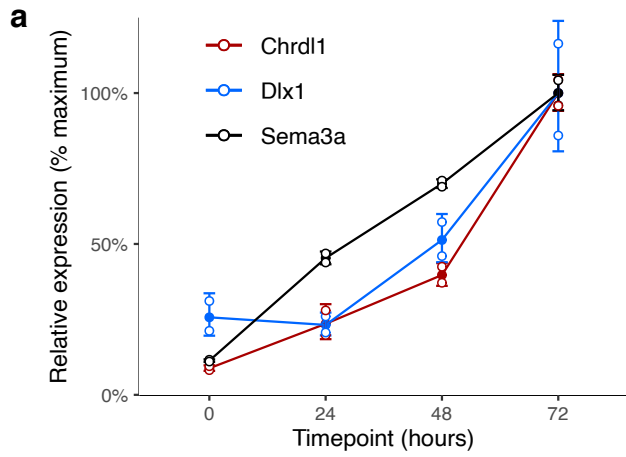
Supplementary Table 8 | List of disease associations with hyperconserved 5'UTR genes with ATP-dependent differential accessibility regions

This table lists known disease associations with hyperconserved 5'UTR genes with at least one significant ATP-dependent differential accessibility region.

Supplementary Table 9 | Oligonucleotides

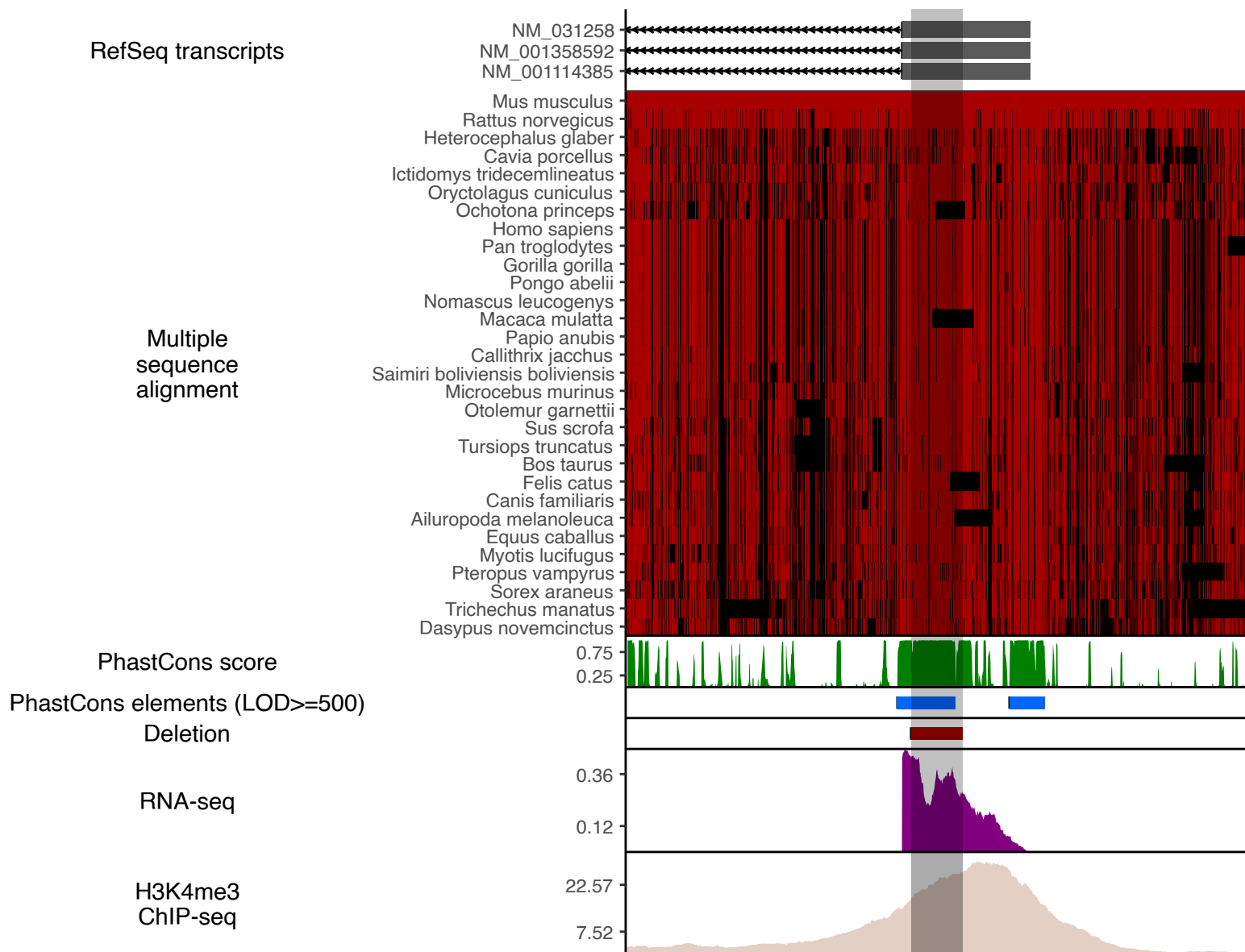
This table lists the oligonucleotides referenced in the Online Methods section.

Supplementary Figure 1



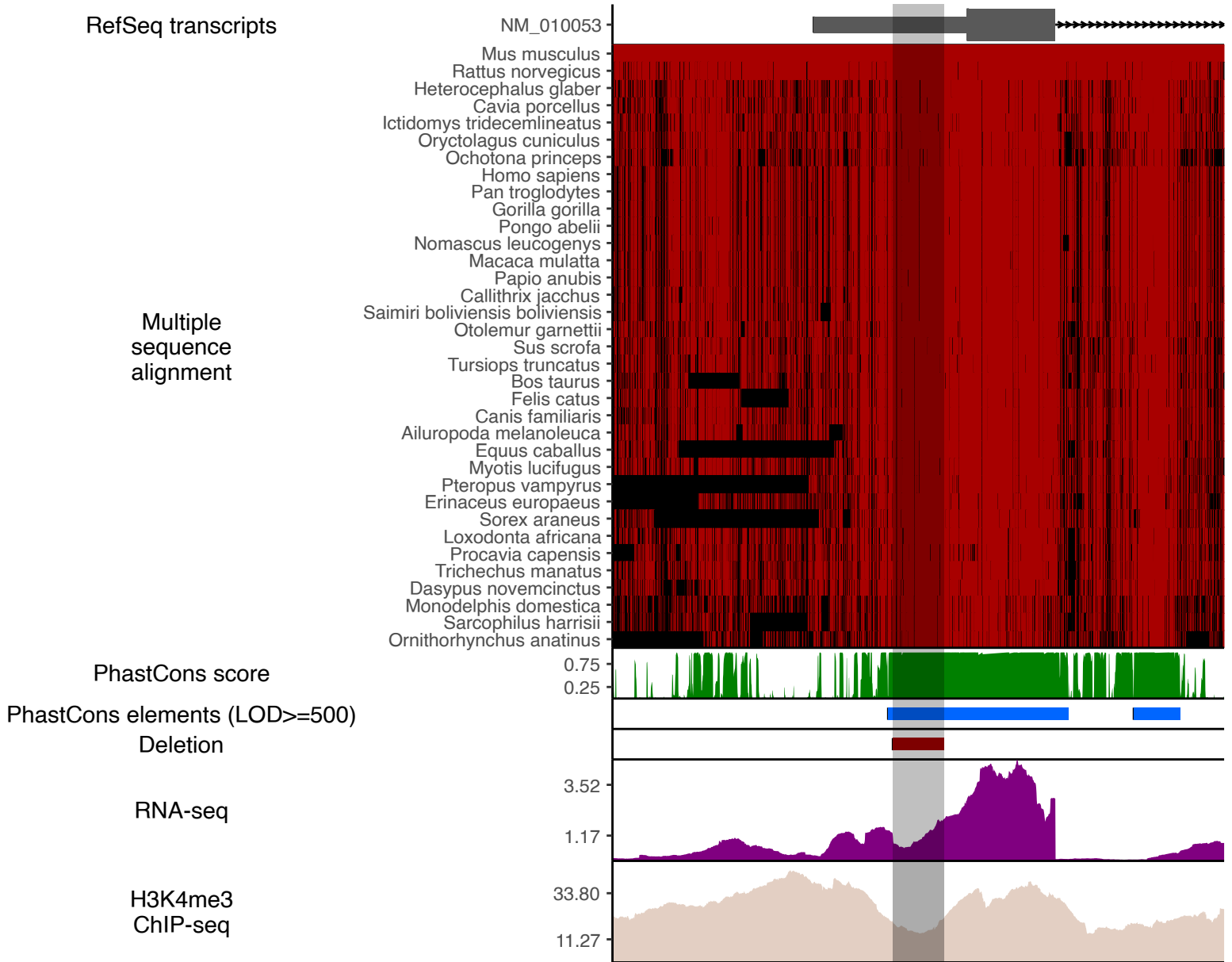
Supplementary Figure 2

ChrDI1, chrX:143392845-143395025 (Δ 143393845-143394025)



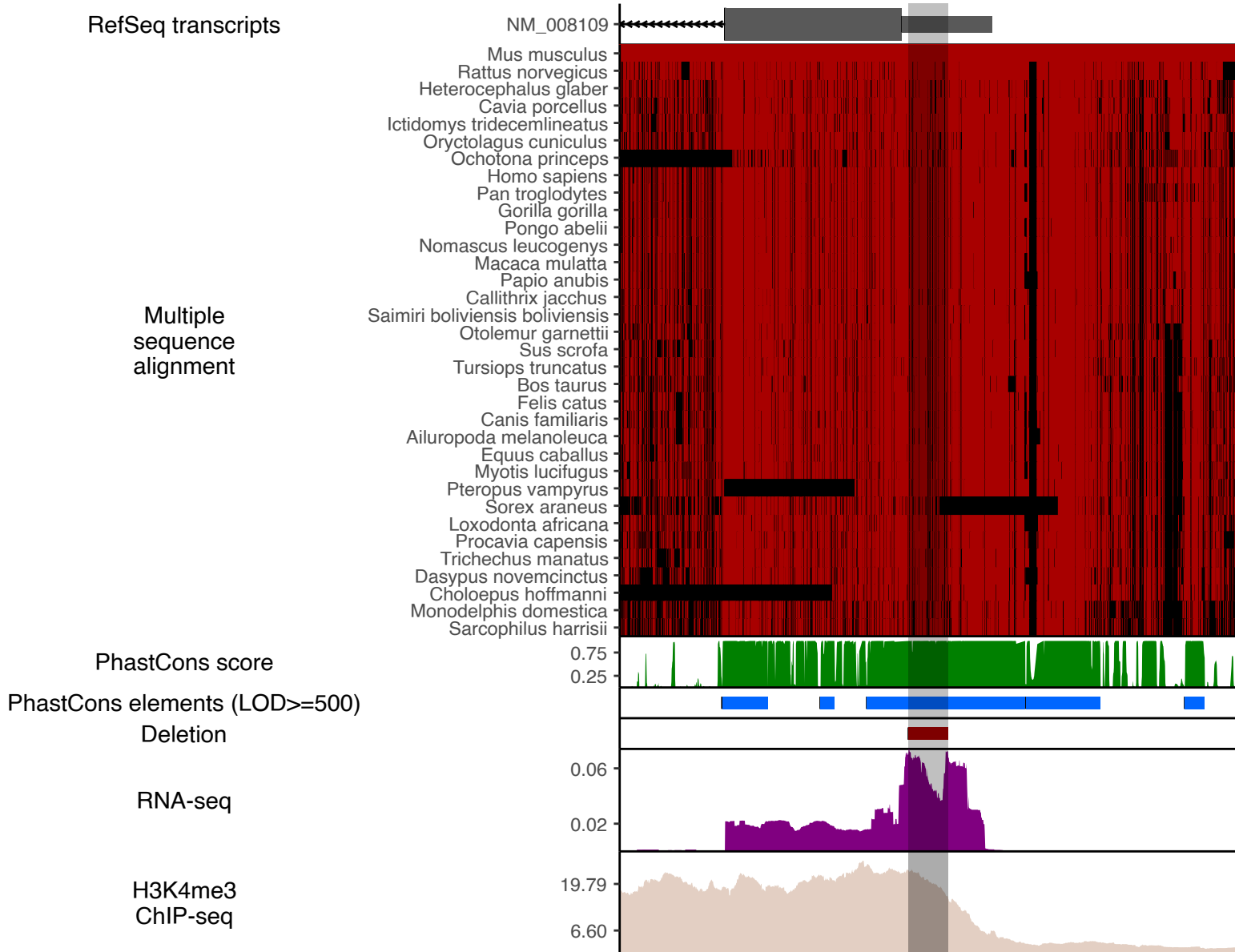
Supplementary Figure 3

Dlx1, chr2:71528727-71530908 (Δ 71529727-71529908)



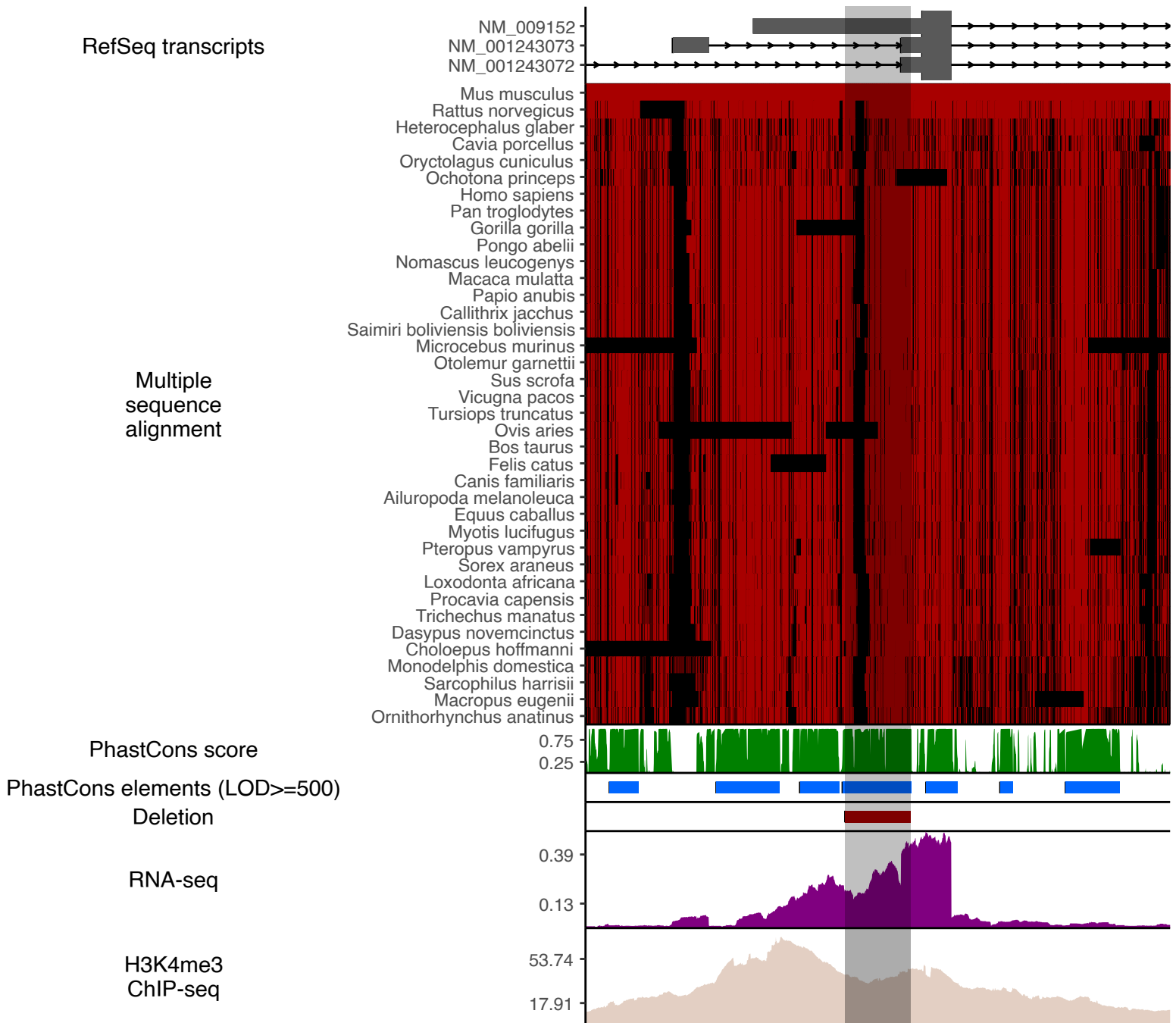
Supplementary Figure 4

Gdf5, chr2:155944078-155946216 (Δ 155945078-155945216)



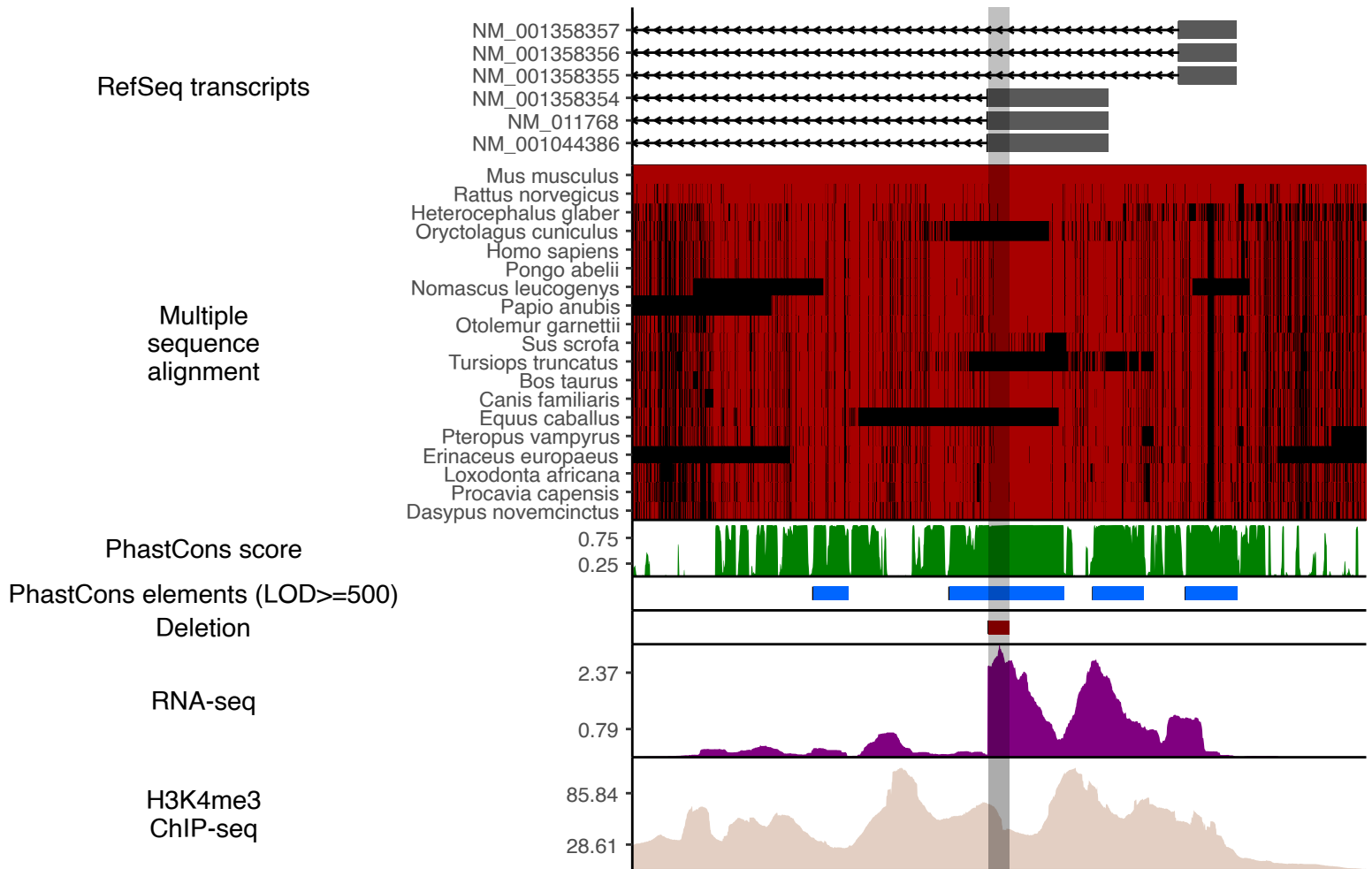
Supplementary Figure 5

Sema3a, chr5:13398661-13400912 (Δ 13399661-13399912)

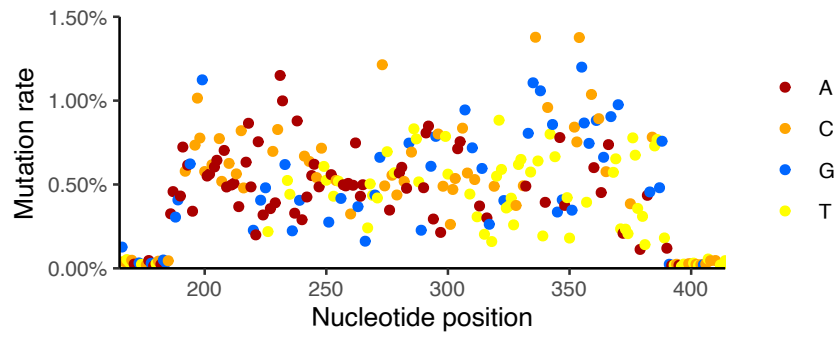


Supplementary Figure 6

Zfx, chrX:94122071-94124130 (Δ 94123071-94123130)



Supplementary Figure 7



Supplementary Notes

Data sources

The following lists the version and sources of publicly available data used in this study. RefSeq: release 84, <https://ftp.ncbi.nlm.nih.gov/refseq/>; 60-way PhastCons: UCSC mm10, <http://hgdownload.soe.ucsc.edu/goldenPath/mm10/phastCons60way/>, <http://hgdownload.soe.ucsc.edu/goldenPath/mm10/database/>; 60-way multiple sequence alignment: UCSC mm10, <http://hgdownload.soe.ucsc.edu/goldenPath/mm10/multiz60way/>; GTEX RNA-seq: V8 (GENCODE V25 annotation), <https://www.gtexportal.org/home/datasets>, ; GO term mapping: Bioconductor 3.10 org.Mm.eg.db, <http://bioconductor.org/packages/3.10/data/annotation/html/org.Mm.eg.db.html>; ENCODE: <https://www.encodeproject.org/>^{1,2,3}.

h5UTR definition

PhastCons elements represent segments of the alignment belonging to the conserved state of the phylo-HMM⁴. Each element has log odds score (LOD) = log probability under the conserved model - log probability under the non-conserved model. 60-way vertebrate PhastCons elements were downloaded from UCSC mouse genome database, and LOD \geq 500 elements were subsetted. LOD \geq 500 selects for about 92th percentile and above for all predicted PhastCons elements, and is essentially arbitrary but chosen to pick a small enough set such that lowest scores appear very extreme upon individual inspection, while providing a meaningful enough number of 5'UTRs to allow statistical confidence. Finally, for each mouse RefSeq transcript record, the total number of 5'UTR, CDS, and 3'UTR nucleotides overlapping LOD \geq 500 elements

are calculated (**Supplementary Table 1**). 5'UTRs with ≥ 250 nt overlap are labeled hyperconserved. As a crosscheck of the RefSeq annotation, transcription start sites of these h5UTRs are compared with ENCODE RAMPAGE datasets and found to overlap at least one TSS peak at a rate of 82%.

GTEX data processing and analysis

Cross-tissue quantitative proteomics data were obtained from GTEx proteomics analysis of 32 human tissues from 14 individuals that match GTEx transcriptomics samples as reported⁵.

RNA and protein abundance/variance: TMM scaling was applied on GTEx RNA-seq count matrix⁶. Log of the TMM-scaled expression values were quantile normalized and per-tissue median value across individuals are taken as the RNA expression level. Mean and variance of RNA expression levels are calculated from the per-tissue median values. Similarly for proteomics data, log of summed peptide abundances were quantile normalized and median values across individuals per tissue were taken. Mean and variance of protein expression levels are calculated from the per-tissue median values. Pearson's correlation coefficient is calculated between per-tissue medians of RNA expression and per-tissue medians of protein expression.

We began with a table of correlation values, RNA mean expression, RNA variance, protein mean expression and protein variance for each human gene. For a gene to be labeled hyperconserved in this table, it must have a homologous RefSeq annotated 5'UTR between humans and mice that overlap $\text{LOD} \geq 500$ PhastCons elements by at least 250nt in mouse and 150nt in humans. In other words, the hyperconserved 5'UTR must be known to both human and mouse RefSeq transcript annotations. For a gene to be labeled non-conserved in this table, no $\text{LOD} \geq 500$ PhastCons elements must overlap any 5'UTR annotation for the gene in both human and mouse transcript annotations. To be considered for analysis, each gene's expression must

be detected in at least 10 tissues for both protein and RNA. The total number of genes available for analysis is 8598. Size-matching or variance-matching non-conserved controls are selected as follows. The table is first ranked by maximum 5'UTR length annotated for the gene. For each h5UTR gene, the gene closest in 5'UTR length or RNA variance rank to the h5UTR but has a non-conserved label is selected without replacement.

Cell culture

mESC: E14 mESCs were cultured on 0.1% gelatin-coated dishes using the following media recipe: Knockout DMEM (ThermoFisher Scientific, 10829018), 15% Embryomax FBS (MilliporeSigma, ES-009-B), 2 mM non-essential amino acids (MilliporeSigma, TMS-001-C), 2 mM L-Glutamine (MilliporeSigma, TMS-002-C), 0.1 mM 2-mercaptoethanol (ThermoFisher Scientific, 21985023), and 10^3 U/mL mLIF (MilliporeSigma, ESG1107; Gemini 400-495).

10T1/2: C3H10T1/2 cells were grown using the following media recipe: DMEM (ThermoFisher Scientific, 11965118), 10% FBS (ThermoFisher Scientific 26140079), 1X penicillin-streptomycin (MilliporeSigma, TMS-AB2-C).

NSC: NSC were differentiated from E14 mES cells⁷. After differentiation, NSC were grown using the following media recipe: 50% DMEM-F12 (ThermoFisher Scientific, 11320082), 50% Neurobasal medium (ThermoFisher Scientific, 21103049), 1x modified N2 (ThermoFisher Scientific, 17502048), 1x B27 supplement (ThermoFisher Scientific, 17504044), 0.0007% BSA Fraction V (ThermoFisher Scientific, 15260037), 0.01 ug/mL murine EGF (ThermoFisher Scientific, PMG8044), 0.01 ug/mL murine FGF-basic (ThermoFisher Scientific, PMG0031), 1X penicillin-streptomycin (MilliporeSigma, TMS-AB2-C). Accutase (MilliporeSigma, A6964 was used for dissociation of NSCs.

Neuron: Neurons were differentiated from mESC derived NSCs⁷. $1-2 \times 10^3$ NSCs were plated onto poly-ornithine/laminin coated 96-well plates using the following differentiation media recipe: 25% DMEM-F12 (ThermoFisher Scientific, 11320082), 75% Neurobasal medium (ThermoFisher Scientific, 21103049), 1x modified N2 (ThermoFisher Scientific, 17502048), 1x B27 supplement (ThermoFisher Scientific, 17504044), 0.005% BSA Fraction V (ThermoFisher Scientific, 15260037), 0.01 ug/mL murine FGF-basic (ThermoFisher Scientific, PMG0031). Half of the media was replaced every 2 days. After 7 days, the media was exchanged for further neuronal differentiation and maturation to: 25% DMEM-F12 (ThermoFisher Scientific, 11320082), 75% Neurobasal medium (ThermoFisher Scientific, 21103049), 0.25x modified N2 (ThermoFisher Scientific, 17502048), 0.25x B27 supplement (ThermoFisher Scientific, 17504044), 0.0007% BSA Fraction V (ThermoFisher Scientific, 15260037).

Limb mesenchyme culture: Stage E11.5 mouse embryos were isolated and washed in PBS. The limbs were dissected in dissection media: DMEM-F12 (ThermoFisher Scientific, 11320082), 10% FBS (ThermoFisher Scientific, 26140079), 1X penicillin-streptomycin (MilliporeSigma, TMS-AB2-C). Dissected limbs were trypsinized in 1% (w/v) trypsin (ThermoFisher Scientific, 27250018) in HBSS (ThermoFisher Scientific, 14175095) at 37°C, 5% CO₂ incubator for 30 min. After trypsinization, tissue was resuspended in dissection media, passed through a cell strainer and 2.5×10^4 - 3×10^4 cells per 96 well plate was cultured in dissection media overnight before transfection.

Embryoid body: E14 mESCs were trypsinized and plated on low-attachment dish at 5×10^6 per 10cm in mESC media without LIF: Knockout DMEM (ThermoFisher Scientific, 10829018), 15% Embryomax FBS (MilliporeSigma, ES-009-B), 2 mM non-essential amino acids (MilliporeSigma, TMS-001-C), 2 mM L-Glutamine (MilliporeSigma, TMS-002-C), 0.1 mM 2-mercaptoethanol (ThermoFisher Scientific, 21985023).

Retinoic acid treated mESC: E14 mESCs were seeded at 1×10^6 cells per 10cm dish in mESC media without LIF + 10uM retinoic acid: Knockout DMEM (ThermoFisher Scientific, 10829018), 15% Embryomax FBS (MilliporeSigma, ES-009-B), 2 mM non-essential amino acids (MilliporeSigma, TMS-001-C), 2 mM L-Glutamine (MilliporeSigma, TMS-002-C), 0.1 mM 2-mercaptoethanol (ThermoFisher Scientific, 21985023), 10uM retinoic acid (MilliporeSigma, R2625). Media was changed every 24 hours and cells were harvested for polysome profiling after 3 days. For the initial time course to check expression of targeted h5UTR genes, 1.5×10^5 cells were seeded per one 6-well. Media was changed every 24 hours until lysis with Trizol (ThermoFisher Scientific, 15596026) and RNA extraction from aqueous phase on silica column (Zymo, R1013). 100ng RNA was reverse transcribed using iScript reverse transcriptase (Biorad, 1708890) in a 10uL reaction. qPCR was performed using Ssoadvanced Universal SYBR Green Supermix (Biorad, 1725270) with 2uL of 1:4 diluted reverse transcription reaction and primer pairs targeting mouse Actb and targeted h5UTRs. $2^{\Delta\Delta Cq}$ values relative to Actb and maximally expressed time point for each h5UTR is plotted.

Term enrichment analysis

GO term enrichment analysis is performed using R package topGO (version 2.38.1)⁸. GO term-gene mappings are obtained from Bioconductor annotation package org.Mm.eg.db. Only terms with at least 10 genes annotated and at least 1 gene mapping to h5UTRs or size-matched non-conserved gene sets are tested. First, two-tailed Fisher's exact test p-value is calculated for each term, and those with Benjamini-Hochberg FDR estimate ≤ 0.05 are retained. topGO weight01 p-value ≤ 0.05 , observed/expected ratio ≥ 3 , and minimum number of genes mapping ≥ 3 is used to further filter the term list. For the final set of GO terms, semantic similarity adjacency matrix is calculated using Lin method in R package GOSemSim (version 2.12.1)⁹.

Adjacency matrix is ranked on the column. The rank matrix is compared with the transposed rank matrix and higher rank is taken. Network with terms as nodes are constructed and terms with maximum rank (total number of terms - 1) is connected. Clustering is performed with cluster_edge_betweenness algorithm in R package iGraph. (version 1.2.5). Mammalian phenotype ontology term enrichment analysis is performed using MouseMine¹⁰. MPO terms with Benjamini-Hochberg FDR estimate ≤ 0.05 and minimum number of genes mapping ≥ 5 are reported as significant.

Primer design and pooling for multiplexed mutational profiling

A total of 384 primer pairs were designed across these h5UTRs using ThermoAlign (version 1.0.0) and following parameters (others are left at default): primer_size_range=16-40, GC_range=25-75, Tm_range=64-72, primer_conc=300, Na=0, K=75, Tris=10, Mg=2, dNTPs=1.2, amplicon_size_min=250, amplicon_size_max=250¹¹. Primers were split into 4 pools of 96 pairs using PrimerPooler (version 1.41)¹². A final set of 380 pairs were synthesized by Eurofins Genomics (See **Supplementary Table 9** for primer sequences).

1D accessibility data analysis

After bcl conversion and demultiplexing with Illumina bcl2fastq, adaptor sequences are further trimmed using cutadapt (version 1.18)¹³. Overlapping paired-end reads are merged using BBMerge (version 38.22)¹⁴. Trimmed and merged reads are aligned to indexed reference of amplicon sequences using Bowtie2 (version 2.3.4.3) with the following parameters: -D 20 -R 3 -N 0 -L 20 -i S,1,0.50 --mp 3,1 --rdg 5,1 --dpad 30¹⁵. Alignments are filtered to have at least 4 mutations with minimum PHRED estimate of 30 (quality cutoff for substitutions only) and

minimum length 150nt. Co-occurrence of modifications at more than 4 positions is considered highly unlikely and assumed to be PCR jackpots. Duplicates are grouped using Hamming distance and directional network algorithm for clustering as implemented for UMIs in `umi_tools` (version 0.5.4)¹⁶. The consensus reads are then parsed and counted for different classes of mutations using functions in `ShapeMapper 2` (version 2.1.5) library¹⁷. This mutational profiling data processing pipeline results in a matrix of mutation counts, where rows are each nucleotide position of the amplicons and columns are the different samples.

For each amplicon, the counts are first normalized by TMM⁶. As a quality filter, we require that the average pairwise Pearson correlations between replicate samples for both conditions must be greater than 0.65 for the amplicon to be further analyzed. For statistical significance of differential mutation rates between ATP depletion and untreated conditions at each nucleotide, we use `voom-limma` (version 3.42.2) which models mean-variance bias and calculates moderated T-statistics^{18,19}. To analyze per-window accessibility pattern differences, for each sliding window of 11nt, we calculate Anderson-Darling statistic between per-nucleotide T-statistic values of the window versus per-nucleotide T-statistic values of the whole amplicon. False discovery rates are estimated by Benjamini-Hochberg procedure. Overlapping significant windows above the chosen cutoff are merged.

2D accessibility data analysis and structure models

After `bcl` conversion and demultiplexing with Illumina `bcl2fastq`, adaptor sequences are further trimmed using `cutadapt` (version 1.18)¹³. Overlapping paired-end reads are merged using `BBMerge` (version 38.22)¹⁴. Trimmed and merged reads are aligned to indexed reference of amplicon sequences using `Bowtie2` (version 2.3.4.3) with the following parameters: `-D 20 -R 3 -N 0 -L 20 -i S,1,0.50 --mp 3,1 --rdg 5,1 --dpad 30`¹⁵. Alignments are filtered to have at least 4 mutations with minimum PHRED estimate of 30 (quality cutoff for substitutions only) and

minimum length 150nt. The consensus reads are then parsed for different classes of mutations using functions in ShapeMapper 2 (version 2.1.5) library¹⁷. Only substitution mutations are taken and pairwise mutation co-occurrence counts are recorded into a symmetric 2D matrix for each mutation spectrum, where rows and columns are positions of the nucleotide along the amplicon. For a detailed explanation of the mechanisms and analysis of M2 signal generated by mutational profiling, see M2-seq paper²⁰.

For each amplicon, 2D matrices for each type of substitution mutations are summed. TMM factors are used to normalize the matrix per column and logarithm of normalized values are used⁶. Based on mean-variance (per column) plot, minimum count threshold is selected. We chose to mask the nucleotide positions which average count less than the threshold. This results in filtering of nearly all G nucleotides as expected, which have lower signal (lower modification-mutation rates) than other bases and different variance distribution²¹. The diagonal of the matrix within 4nt is also masked due to extremely low co-variation. The matrix is further z-scaled per row and centered by the median of the entire matrix. For visualization, the z-score values are winsorized to 2 standard deviations. For clustering similar mutant (row) accessibility patterns, missing values in the diagonals are first imputed by k-nearest neighbor method with k=10. Classical multidimensional scaling is used to translate pairwise distances (Euclidean, i.e. PCA) between the mutant accessibilities into N dimensions. The resulting points from N=2 are heuristically grouped by inspection of scatter plot visualization. Two clusters are identified and per-column average z-scores are calculated.

To model structures for region B, we subsetting positions 215-315. Cluster average scores are used as constraints for partition function calculation in Vienna RNA (version 2.4.14)²². 250 structures are sampled for each cluster constrained partition function and used as input suboptimals to REEFIT (version 0.6.3)²³. As the input accessibility data to REEFIT, we processed the TMM normalized log count matrix as follows. First, all cell values within each

nucleotide were scaled and centered. The diagonal of the matrix within 4nt was masked; missing values are imputed by k-nearest neighbor method with k=10. The distribution of all values in the matrix are bimodal, similar to distributions observed for log data for existing M² datasets generated with SHAPE and capillary electrophoresis that are available in RMDB (<https://rmdb.stanford.edu/>). Since REEFIT estimates paired and unpaired reactivity distributions by mixture modelling of exponential distributions, we rescaled the DMS icM² matrix on exponential distribution such that the quantile/rank is preserved but the distribution of values closely approximate SHAPE-CE M² datasets downloaded from RMDB. For visualization of the landscape, we used pairwise distance metrics, structure clustering, and medoid assignment produced by REEFIT. Sum of weights for structures belonging to each of the 3 clusters represented by a medoid structure is presented in the main figure. We used bootstrapping to robustly estimate population fraction errors. Each bootstrap runs (200 times total) using the same parameters, except that the columns (positions) of the data matrix are shuffled. This gives bounds on the frequency of the structures in the population. Reported values in the manuscript are estimated mean±standard deviation of the weight.

References

1. Haeussler, M. *et al.* The UCSC Genome Browser database: 2019 update. *Nucleic Acids Res.* **47**, D853–D858 (2019).
2. Carithers, L. J. *et al.* A Novel Approach to High-Quality Postmortem Tissue Procurement: The GTEx Project. *Biopreserv. Biobank.* **13**, 311–319 (2015).
3. Huber, W. *et al.* Orchestrating high-throughput genomic analysis with Bioconductor. *Nat. Methods* **12**, 115–121 (2015).
4. Siepel, A. *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* **15**, 1034–1050 (2005).

5. Jiang, L. *et al.* A quantitative proteome map of the human body. *Cell* **183**, 269–283.e19 (2020).
6. Robinson, M. D. & Oshlack, A. A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biol.* **11**, R25 (2010).
7. Conti, L. *et al.* Niche-independent symmetrical self-renewal of a mammalian tissue stem cell. *PLoS Biol.* **3**, e283 (2005).
8. Alexa, A., Rahnenführer, J. & Lengauer, T. Improved scoring of functional groups from gene expression data by decorrelating GO graph structure. *Bioinformatics* **22**, 1600–1607 (2006).
9. Yu, G. *et al.* GOSemSim: an R package for measuring semantic similarity among GO terms and gene products. *Bioinformatics* **26**, 976–978 (2010).
10. Motenko, H., Neuhauser, S. B., O’Keefe, M. & Richardson, J. E. MouseMine: a new data warehouse for MGI. *Mamm. Genome* **26**, 325–330 (2015).
11. Francis, F., Dumas, M. D. & Wisser, R. J. ThermoAlign: a genome-aware primer design tool for tiled amplicon resequencing. *Sci. Rep.* **7**, 44437 (2017).
12. Brown, S. S. *et al.* PrimerPooler: automated primer pooling to prepare library for targeted sequencing. *Biol. Methods Protoc.* **2**, bpx006 (2017).
13. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet j.* **17**, 10 (2011).
14. Bushnell, B., Rood, J. & Singer, E. BBMerge - accurate paired shotgun read merging via overlap. *PLoS ONE* **12**, e0185056 (2017).
15. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
16. Smith, T., Heger, A. & Sudbery, I. UMI-tools: modeling sequencing errors in Unique Molecular Identifiers to improve quantification accuracy. *Genome Res.* **27**, 491–499 (2017).
17. Busan, S. & Weeks, K. M. Accurate detection of chemical modifications in RNA by mutational profiling (MaP) with ShapeMapper 2. *RNA* **24**, 143–148 (2018).

18. Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. voom: Precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biol.* **15**, R29 (2014).
19. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.* **43**, e47 (2015).
20. Cheng, C. Y., Kladwang, W., Yesselman, J. D. & Das, R. RNA structure inference through chemical mapping after accidental or intentional mutations. *Proc Natl Acad Sci USA* **114**, 9876–9881 (2017).
21. Mustoe, A. M., Lama, N. N., Irving, P. S., Olson, S. W. & Weeks, K. M. RNA base-pairing complexity in living cells visualized by correlated chemical probing. *Proc Natl Acad Sci USA* **116**, 24574–24582 (2019).
22. Lorenz, R. *et al.* ViennaRNA Package 2.0. *Algorithms Mol. Biol.* **6**, 26 (2011).
23. Cordero, P. & Das, R. Rich RNA Structure Landscapes Revealed by Mutate-and-Map Analysis. *PLoS Comput. Biol.* **11**, e1004473 (2015).