

ARTICLE

Received 16 Jul 2014 | Accepted 13 Feb 2015 | Published 25 Mar 2015

DOI: 10.1038/ncomms7644

A family of transposable elements co-opted into developmental enhancers in the mouse neocortex

James H. Notwell¹, Tisha Chung², Whitney Heavner^{2,3} & Gill Bejerano^{1,2,4}

The neocortex is a mammalian-specific structure that is responsible for higher functions such as cognition, emotion and perception. To gain insight into its evolution and the gene regulatory codes that pattern it, we studied the overlap of its active developmental enhancers with transposable element (TE) families and compared this overlap to uniformly shuffled enhancers. Here we show a striking enrichment of the MER130 repeat family among active enhancers in the mouse dorsal cerebral wall, which gives rise to the neocortex, at embryonic day 14.5. We show that MER130 instances preserve a common code of transcriptional regulatory logic, function as enhancers and are adjacent to critical neocortical genes. MER130, a nonautonomous interspersed TE, originates in the tetrapod or possibly Sarcopterygii ancestor, which far predates the appearance of the neocortex. Our results show that MER130 elements were recruited, likely through their common regulatory logic, as neocortical enhancers.

¹Department of Computer Science, Stanford University, 279 Campus Drive West (MC 5329), Beckman Center B-300, Stanford, California 94305-5329, USA. ²Department of Developmental Biology, Stanford University, 279 Campus Drive West (MC 5329), Beckman Center B-300, Stanford, California 94305-5329, USA. ³Department of Biology, Stanford University, 279 Campus Drive West (MC 5329), Beckman Center B-300, Stanford, California 94305-5329, USA. ⁴Department of Pediatrics, Division of Medical Genetics, Stanford University, 279 Campus Drive West (MC 5329), Beckman Center B-300, Stanford, California 94305-5329, USA. Correspondence and requests for materials should be addressed to G.B. (email: bejerano@stanford.edu).

Transcriptional enhancers control when, where and how much of each gene is produced in different cell types at different times¹. Together, groups of enhancers orchestrate the expression of many genes to form complex pathways, but a fundamental question is what produced this coordinated regulatory logic. Instances of recognizable transposable element (TE) families make up more than half the human genome², providing a large substrate of homologous sequences on which evolution can act.

Britten and Davidson³ proposed a mechanism by which the evolution of novel structures or functions could be greatly accelerated through the co-option of TEs into gene regulatory roles, activating novel groups of genes that were not co-expressed before in a new temporally and spatially specific manner. By scattering copies of the same DNA stretches across the genome near all manner of genes, TEs could greatly increase the probability that groups of enhancers would emerge to bring together a new group of genes, to be expressed together and accelerate the formation of complex new pathways and functions. Consistent with the idea of TEs seeding regulatory sequences, TEs have contributed significantly to conserved non-coding elements, via their co-option into gene regulatory roles⁴. In addition, TE families have been co-opted in different contexts such as stem cells⁵ and the placenta⁶.

p300 is a transcriptional co-activator recruited by active enhancers. It is an appealing epigenomic mark of active enhancers, with typical *in vivo* validation rates of 80% (ref. 7). In a recent paper⁸, we measured the p300 landscape of the E14.5 dorsal cerebral wall (DCW), which encompasses and gives rise to the neocortex, by harvesting tissue from mouse embryos, isolating chromatin and performing chromatin immunoprecipitation followed by high-throughput sequencing (ChIP-seq). In that study, we used multiple techniques to validate the quality of our p300 measurements, including confirming key enhancers known from the literature and enrichment analysis showing that from thousands of measured tissue-time-point combinations, our set is most correlated with genes expressed in Theiler stage 22 telencephalon (essentially E14.5 DCW). Using transgenics, we confirmed that eight of ten (80%) p300 measurements function as active neocortical enhancers, and established a primary cortical neuron transfection system where 5/8 positive transgenics showed transfection activity and 2/2 negative transgenics did not⁸.

Here we show that the MER130 repeat family is strongly enriched among the E14.5 enhancers in the developing neocortex, an enrichment that is specific to mouse brain tissues. MER130 is a nonautonomous interspersed repeat family that previously has been shown to be conserved in mammals and chicken⁹. In pursuit of the Britten and Davidson³ hypothesis, we show that the activity of MER130 family members is the result of a preserved core of regulatory logic that consists of multiple transcription factor binding sites. Individual MER130 instances function as enhancers, although we show that endogenous chromatin state is essential for interpreting these results. Furthermore, MER130 instances that function as enhancers lie next to important cortical genes, yet are conserved to species that far predate the emergence of the neocortex.

Results

MER130 is highly enriched among neocortex enhancers. Intrigued by the Britten and Davidson³ hypothesis and encouraged by the high specificity of the p300 mark, we asked whether any of ~1,000 known TE families are disproportionately marked by an *in vivo* p300-bound set of active enhancers (see Methods). One repeat family in one specific context stood out (Fig. 1a): MER130 in the E14.5

DCW with a 73-fold enrichment over expected (Bonferroni P value < 0.006, $n = 1,000,000$ simulations)⁸.

Interestingly, MER130 was not enriched in any of the number of other p300-bound sequences from ChIP-seq experiments in other tissues, including the forebrain at an earlier time point (Fig. 1a). The histone modification H3K27ac, deposited by p300 at active enhancers¹⁰, is easier to assay and is thus a more comprehensively surveyed mark. We next searched for MER130 activity in multiple additional tissues by performing the same enrichment test for H3K27ac ChIP-seq sets from all 24 available tissues and cell lines from the mouse Encyclopedia of DNA Elements (ENCODE)¹¹. We saw only weaker enrichments in other brain tissues: 12.9-fold enrichment for E14.5 whole brain, 12.2-fold for 8-week cerebellum, and <4-fold enrichment in non-brain tissues (Fig. 1b). Together, these data show a strong enrichment for MER130 in E14.5 DCW, with more modest enrichment in general brain tissue, compared with over 20 other tissue-time-point combinations.

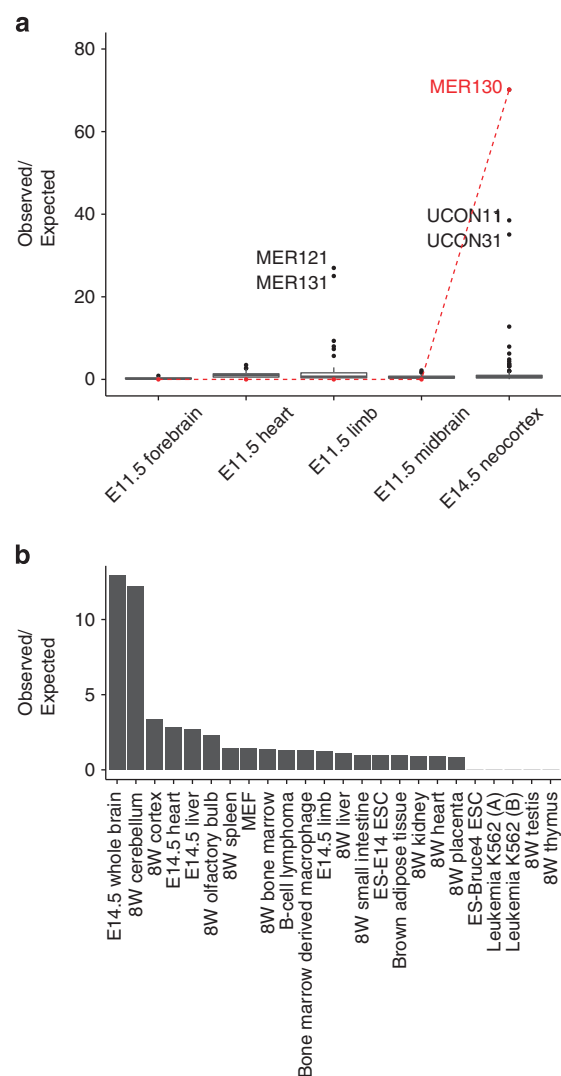


Figure 1 | MER130 is strongly enriched among neocortical enhancers.

(a) Fold enrichments compared with uniform shuffles for over 1,000 interspersed repeat families among p300 ChIP-seq data sets from developing tissues. MER130 enrichments in red. (b) Fold enrichment of observed MER130 instance overlaps compared with expected overlaps from uniform shuffles of 24 ENCODE H3K27ac sets, including different primary tissues at different developmental time points and model cell lines. No enrichment comes close to the E14.5 neocortex enrichment in a.

MER130 instances preserve a core regulatory logic. We next comprehensively annotated the MER130 family using nhmmer, a probabilistic alignment tool¹² (107 instances; Supplementary Table 1; see Methods). Essentially all of the MER130 instances are distal (>2.5 kb from the transcription start site), and, as opposed to some other repeat families¹³, none overlap protein coding exons. We identified 23 well-conserved MER130 instances that were strongly enriched for the p300 ChIP-seq signal but not enriched for the input control signal (see Methods). Fourteen MER130 instances also overlap E14.5 whole-brain H3K27ac ChIP-seq peaks: eight strongly marked by p300 in our set and six marked with intermediate levels of p300. No instance marked by low levels of p300 in our set was marked by H3K27Ac in E14.5 whole brain.

To achieve their specificity, enhancers contain distinct combinations of transcription factor binding sites¹⁴. When we constructed a multiple alignment from the 23 MER130 instances strongly enriched for p300 signal, it revealed a well-conserved ~100 bp core (Fig. 2b) containing 5 putative transcription factor binding sites resembling motifs from our library^{15,16}: a Neurod/Neurog motif, an Nfi dimer¹⁷ and two additional Nfi sites (Fig. 2a), representing families of factors important for brain development (see Discussion section).

p300-marked MER130 instances function as enhancers. Next, we used our primary cortical neuron transfection system to test whether the MER130 instances marked by p300 functioned as enhancers and whether the binding sites preserved in the multiple alignment were necessary for their activity (see Methods). We

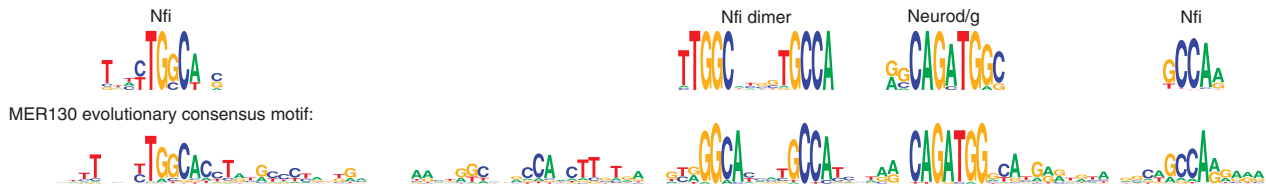
cloned each of these 23 elements upstream of a minimal promoter and luciferase reporter, and transfected these constructs into dissociated neurons isolated from E14.5 mouse cortices (see Methods). Twenty-two of twenty-three (95.7%) of the MER130 candidates hypothesized to function as neocortex enhancers produced greater than twofold expression relative to the empty vector (Fig. 3a, Supplementary Fig. 1).

Our result confirms that the MER130 instances marked by p300 function as enhancers. While virtually all MER130 instances preserve the 5 binding site core (Supplementary Fig. 2), over half (56 instances, 52%) of MER130 instances are clearly devoid of a DCW p300 mark (the remainder show weak overlap, below our peak-calling threshold; see Methods). We cloned seven instances most depleted for p300 signal and tested them in the same transfection assay. Six of seven (85.7%) of these MER130 instances drove greater than twofold expression relative to the empty vector, though at a significantly weaker level than the set marked by p300 (p300 enriched, $n = 23$, to p300 depleted, $n = 7$, t -test P value: 0.03; fold: 8.5) (Fig. 3a, Supplementary Fig. 1).

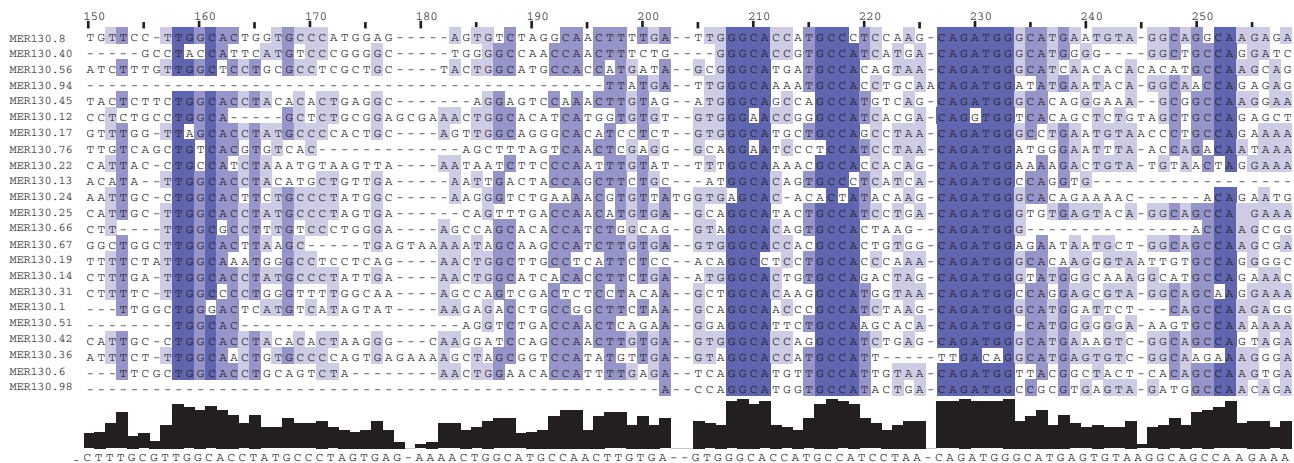
To better understand the *in vivo* state of these sequences, we examined the chromatin state from mouse E14.5 whole-brain DNase-seq data. Remarkably, the MER130 instances that were not marked by p300, yet functioned as enhancers in dissociated cortical neurons, were significantly (p300 enriched, $n = 23$, to p300 depleted, $n = 7$, t -test P value: $3.9e - 09$; fold: 51.7) depleted for DNase I cleavage when compared with the p300-marked MER130 instances that functioned in transfection assays (Fig. 3b).

Of the 107 MER130 instances we annotate in the mouse genome, 103 (96%) are conserved to human. We explored the

a Transcription factor biophysical binding preferences:



b Multiple alignment:



c MER130.31 transfection mutations:



Figure 2 | MER130 contains a preserved core of transcriptional regulatory logic. (a) Biophysical transcription factor binding preferences (top) and matching consensus motif (bottom) derived from multiple alignment of MER130 instances marked by p300. (b) Multiple alignment of 23 MER130 instances strongly enriched for p300 signal. (c) Bases mutated in transfections (see Fig. 4).

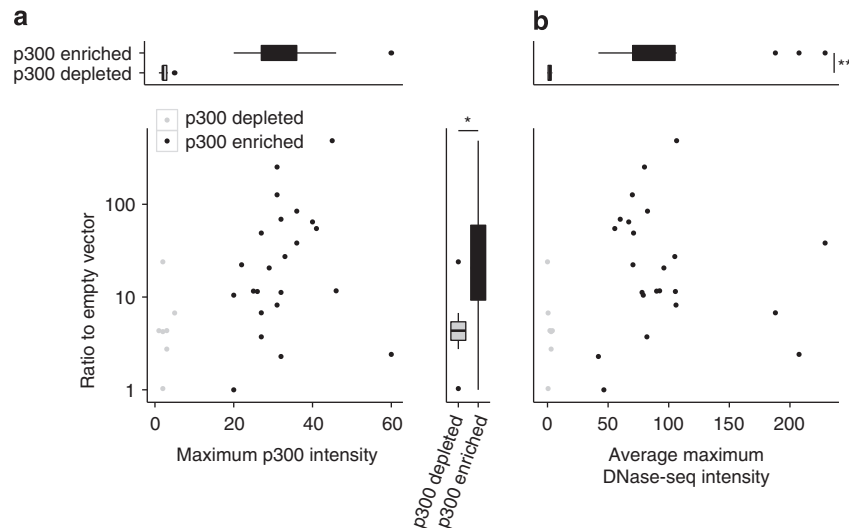


Figure 3 | MER130 members function as enhancers but may not occupy open chromatin. (a) Maximum neocortex p300 ChIP-seq intensity (x-axis) plotted against average fold activity relative to the empty vector when transfected into dissociated cortical neurons (y-axis) for each MER130 instance tested. (b) Average DNase-seq intensity across two mouse E14.5 whole-brain replicates (x-axis) plotted against same y-axis as (a) for the same set of MER130 instances. In both panels, MER130 elements enriched for the p300 signal in black ($n = 23$) and depleted for the p300 signal in grey ($n = 7$). t -test. * P value < 0.05 ; ** P value < 0.01 .

chromatin state of these conserved instances in day 85 human foetal brain tissue from the Roadmap Epigenomics Project¹⁸. Day 87 is the human time point matched to cortical neurogenesis in E14.5 mice¹⁹. Strikingly, MER130 instances not marked by p300 occupied closed chromatin not only in mouse, but also in human brain (Supplementary Fig. 3).

We then tested the necessity of the preserved core for driving specific expression levels in dissociated cortical neurons. Of the p300-enriched MER130 instances robustly expressed in our transfection assay and preserving the five binding sites, we chose two at random. We mutated each of the 5 preserved binding sites in these 2 MER130 neocortex enhancers (Fig. 2c) and tested the resulting 10 constructs in our primary cortical neuron transfection system. Mutating 7/10 (70%) predicted binding sites in the two enhancers resulted in a significant reduction (wild-type to mutant fold change of 1.8–7.7) of reporter activity (Fig. 4).

MER130 instances near critical neocortical development genes.

To identify functional coherence among the genes putatively regulated by the MER130 neocortex enhancers, we analysed them with GREAT²⁰, a genomic region enrichment tool. We did not observe any statistically significant enrichments for either the full set of 107 MER130 elements or the set of 22 validated MER130 neocortex enhancers. Six of twenty-two validated MER130 cortical enhancers, however, do reside next to genes whose perturbations result in abnormal telencephalon morphology: *Robo1*, *Id4*, *Dgkb*, *Ap3b1*, *Prok2* and *Zfp423*. We also observed a second MER130 instance marked by intermediate levels of p300 adjacent to two of these six genes: *Robo1* and *Zfp423*. From the remaining MER130 instances not enriched for p300 at E14.5, we observe MER130 instances next to eight additional genes annotated with abnormal telencephalon morphology: *Cadm1*, *Enah*, *Epha4*, *Foxg1*, *Mycn*, *Npas3*, *Park2* and *Pax3* (ref. 21).

MER130 originated in or before the tetrapod ancestor. The MER130 family is likely an interspersed repeat due to its consensus length (~470 bp), composition and distribution across the genome¹³. We queried all known repeat families and databases of RNAs but found no homologous sequences that may help shed

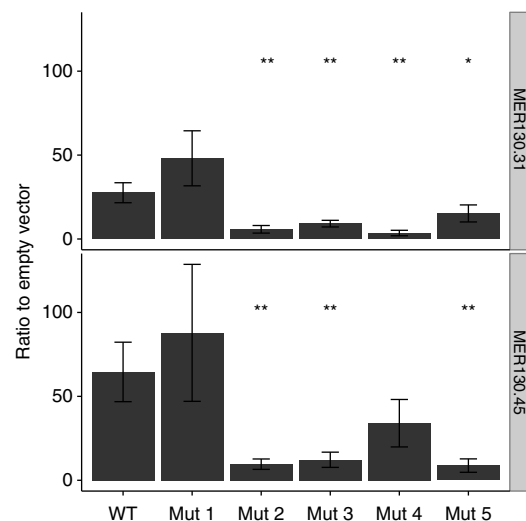


Figure 4 | MER130 core mutations modulate enhancer activity similarly across different instances. Two MER130 instances were mutated each at five putative binding sites (nomenclature as in Fig. 2c). Average fold activity relative to the empty vector when mutated MER130 instances were transfected into dissociated cortical neurons. Error bars represent s.d. $n = 3$ biological replicates \times 3 technical replicates for each condition. t -test. * P value < 0.05 ; ** P value < 0.01 . All five pairs trend in the same direction, with 7/10 mutations significantly reducing expression.

light on the identity of MER130 (see Methods). We also searched 75 vertebrate genome drafts and found no species with recent MER130 activity (Supplementary Table 2 and Methods). While MER130 was previously thought to originate in the amniote ancestor⁹, we discovered that in fact it is conserved in *Xenopus tropicalis*, suggesting this element originated no later than the tetrapod ancestor (Fig. 5a). We also identified a single possible ancestral match in the Coelacanth genome (Fig. 5c). This hit occurs in the *NRF1* orthologue, a gene that mediates neurite outgrowth²². Interestingly, despite the weak sequence

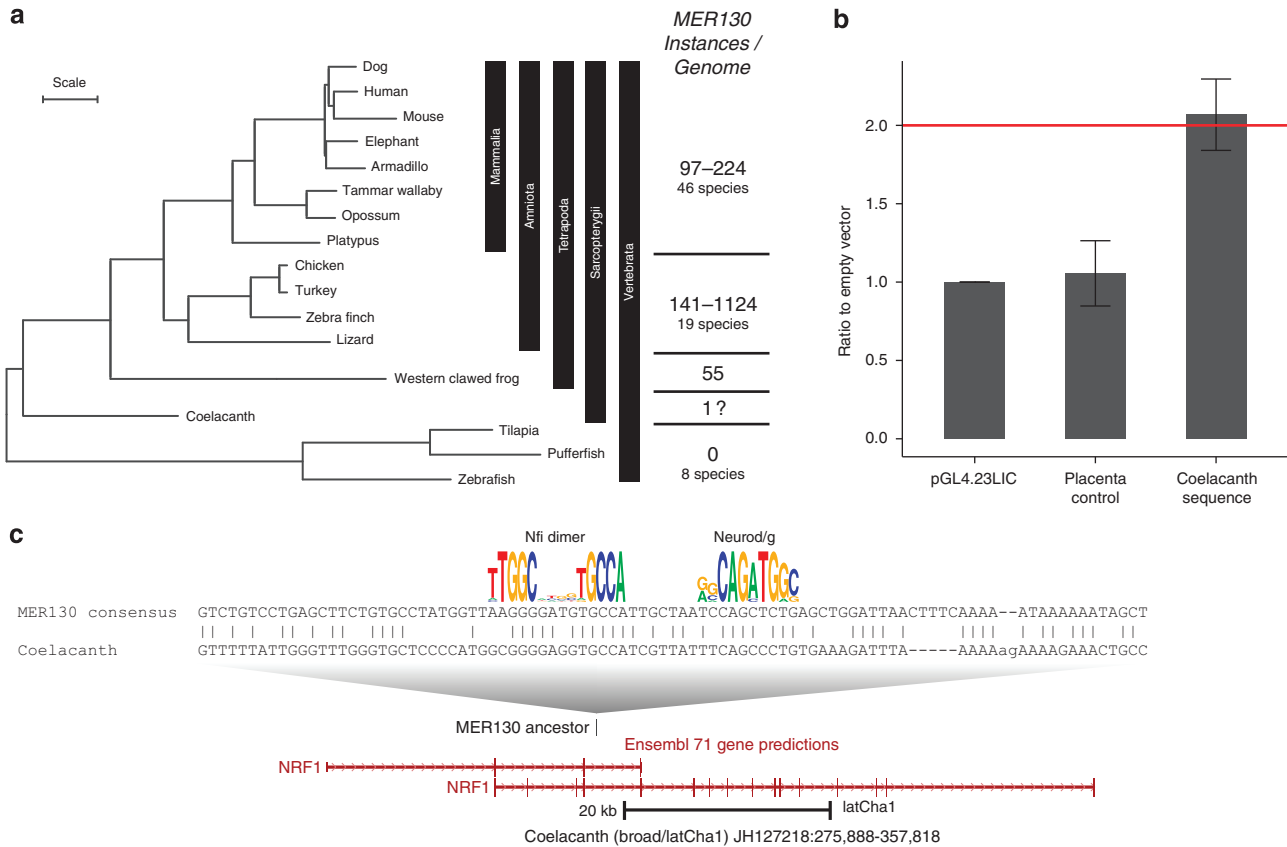


Figure 5 | MER130 originated in the tetrapod or Sarcopterygii ancestor. (a) Phylogenetic tree (adapted from ref. 43) with representative species and the range of MER130 instances in species from each clade. Scale, 0.1 substitutions per site. (b) Average fold activity relative to the empty vector when the coelacanth sequence (c) was transfected into dissociated cortical neurons. Error bars represent s.d. (c) Alignment between Dfam MER130 model and Coelacanth genomic locus. Aligning bases preserve Nfi dimer and Neurod/g motifs (Fig. 2). The putative MER130 instance in coelacanth is located in an intron of *NRF1*, a gene that mediates neurite outgrowth.

conservation, this coelacanth element drives greater than twofold activity when tested in dissociated cortical neurons (Fig. 5b). No MER130 matches were found in seven species of bony fish or lamprey.

Forebrain expression of genes near MER130 instances. Finally, we used the fact that MER130 is not present in fish to ask whether the appearance of a MER130 instance next to a tetrapod gene makes it more likely to be expressed in the forebrain than its orthologue in fish. Using the GREAT default gene-region association rule, our 22 active MER130 instances can be associated with 37 immediately adjacent genes. Seventeen of these genes, adjacent to 13 MER130 instances, are annotated with ‘TS21–TS23 forebrain expression’ in GXD²¹. Of the same 37 genes, we map 28 to orthologous genes in zebrafish using the Ensembl orthologue mappings. Nineteen of these are annotated for gene expression in any context in ZFIN²³, and of them, only 3 genes are annotated for forebrain expression. While anecdotal, these results suggest that new gene expression domains in the developing neocortex are associated with the emergence of the MER130 family.

Discussion

The 73-fold enrichment of co-opted MER130 elements among E14.5 DCW enhancers is striking. Similar magnitude enrichments have previously only been reported in the context of embryonic stem cells^{5,24}, with maximal folds in the twenties for non-ES cell types or lines^{6,24}. When we examined the MER130 elements, we

observed a clear core whose pattern of preservation recapitulates the known biophysical binding preferences of two important transcription factor families. The common activation code of Neurod/g and Nfi binding sites are the likely reason for their recruitment as a transcriptional network of neocortical enhancers. Together, these features represent a striking example in support of the Britten and Davidson³ hypothesis.

Neurog and Nfi family members are promising candidates for transcriptional regulators of MER130, given that *Neurog2*-mutant mice, as well as mice null for *Nfia*, *Nfib* and *Nfix*, all exhibit abnormal telencephalon morphology, the same phenotype exhibited by mice null for multiple putative MER130 target genes^{25–28}. The developing neocortex has been microdissected at E14.5 to isolate the three layers present at this time point, and gene expression in all three layers has been measured using RNA-seq²⁹. Among *Neurog* family members, *Neurog2* is the most highly expressed in all three layers at E14.5, and among *Nfi* family members, *Nfib* is the most highly expressed in all three layers. Both genes are in the first percentile of all expressed genes at this time point. The presence of four different Nfi binding sites in the MER130 consensus further supports its importance³⁰.

While our experiments suggest that most MER130 instances are capable of driving activity in the cortex, likely by virtue of their shared code, the MER130 instances marked by p300 at E14.5 drive higher relative activity when transfected in dissociated cortical neurons. In contrast with the constructs we tested *in vitro*, endogenous enhancer activity also depends on the chromatin state³¹. The lack of DNase I cleavage in both mouse

and human brains suggests that the MER130 instances not bound by p300 in our measurements are not only weaker drivers in cortical neurons due to their sequence, but also likely not accessible in their endogenous context.

MER130 is more ancient than previously thought. Our search of multiple vertebrate genomes places its origin at least as early as the common ancestor of tetrapods, and possibly as early as our common ancestor with coelacanth. The MER130 family exhibits an LF-SINE-like evolutionary behaviour, where exapted copies shed revealing 5' and 3' marks, leaving behind what is very likely only the body of the co-opted transposon¹³. We did not identify any instances co-opted into coding exons, suggesting MER130 may not have ended in a poly-A tail¹³. The sequencing of additional tetrapod genomes, especially genomically isolated ones, may yet reveal a species with a recently jumping copy, allowing for the classification of this repeat family³².

Methods

All repeats—p300 sets correlation. We downloaded the RepeatMasker track from the pre-annotated mm9 genome (RepeatMasker open-3.2.8 Repeat Library)³³ and removed low-complexity, satellite and simple repeat classes, while retaining DNA, LINE, SINE, LTR and RNA classes for a total of 1,118 repeat and RNA families. We next downloaded p300 ChIP-seq results for E11.5 forebrain³⁴, midbrain, heart and limb³⁵, and removed proximal peaks (<2.5 kb away from any transcription start site in UCSC knownGene track). We shuffled each p300 set across the genome 10,000 times (excluding the UCSC mm9 gap track). For each p300 set shuffle, we counted the number of p300 set instances overlapped by each repeat family. We also computed the number of actual p300 set instances overlapped by each repeat family. Fold enrichment was computed as observed divided by mean for each repeat family—p300 set combination. To better estimate the *P* value for the enrichment of the MER130 family and determine statistical significance, we performed an additional 1,000,000 shuffles of our neocortex p300 set against this repeat family alone.

Tissue specificity. We downloaded all 24 available H3K27ac tissue and cell-line data from ENCODE and removed all peaks overlapping an exon (UCSC mm9.knownGene) and proximal peaks (<2.5 kb away from any transcription start site). Observed and expected overlaps were computed as above.

MER130 instances identification and intersection with p300. We used nhmmer (from hmmer3.1-snap20121011 using --cut-ga threshold)¹² to identify 107 instances of the MER130 repeat family in the mouse genome (NCBI37/mm9), using Dfam's MER130 DF0000726.2 profile HMM model². Each MER130 instance was semi-automatically classified as high, medium or low confidence based on computing maximum p300 signal intensity over each peak, followed by manual inspection. Read pileups were constructed using MACS1.4, and the maximum peak height was extracted using UCSC bigWigAverageOverBed. Maximum p300 read pileup from the p300 ChIP-seq was used as the primary selection attribute, and instances classified as high confidence had a punctuated peak over the MER130 family member, as well as low input control signal. This resulted in 23 well-conserved MER130 instances with maximum p300 read pileup height ≥ 20 (high), 28 instances marked by intermediate levels of p300 ($7 \leq p300 \leq 24$; medium) and 56 instances marked by background levels of p300 ($0 \leq p300 \leq 6$; low).

Multiple alignment. Sequences aligning to the DFAM DF0000726.2 profile HMM model (corresponding to the 23 MER130 instances strongly enriched for p300 signal) were extracted from the mouse genome (NCBI37/mm9). A multiple alignment of these sequences was then constructed using MUSCLE (v3.8.31; ref. 36).

Ethics. All animal work was carried out in compliance with the Stanford University Institutional Animal Care and Use Committee under approved protocols #18487 and #21758. Timed pregnant female mice were adults >6 weeks of age and maintained on a Swiss Webster background (Charles River). Embryos at the time of harvest were 14.5 days post coitum.

Cloning and transfections. Cloning and transfections were performed as in ref. 15. Briefly, inserts were amplified from mouse genomic DNA (Clontech Laboratories Inc.) using Phusion High Fidelity DNA Polymerase (NEB Inc.). For transfections using the coelacanth match to the MER130 consensus, a 196 bp insert containing the MER130 core match was amplified from coelacanth (*Latimeria menadoensis*) genomic DNA. A Ligation Independent Cloning (LIC) site was ligated to the firefly luciferase vector, PGL4.23 (Promega Corp.), at the 5' KpnI and 3' HindIII sites, and all inserts were cloned into the PGL4.23 LIC vector using an LIC method³⁷. All positive clones were identified by colony PCR and sequenced.

E14.5 embryos were harvested from Swiss Webster mice (Charles River). The dorsal telencephalon was removed and stored briefly on ice in a Hank's balanced salt solution dissection solution. Cells were dissociated using 0.25% trypsin and $10 \mu\text{g} \mu\text{l}^{-1}$ DNase. Nucleofections were performed using the experimental luciferase construct and a pRL-CMV Renilla control (Promega) and a P3 Primary Cell 96-well Nucleofector Kit (Lonza). Nucleofected cells were plated onto poly-D-lysine coated 96-well plates (NUNC) containing PNB (Lonza). Luciferase assays were performed 48 h after nucleofection using the Dual-Luciferase Reporter 1000 Assay System (Promega) according to the manufacturer's protocol and read using a Promega Glomax luminometer. An enhancer sequence that drove activity in placental cells, but not in cortical neurons in luciferase reporter assays³⁸, was used as a negative control. Primers are listed in Supplementary Table 3.

Site-directed mutagenesis. Site-directed mutagenesis was performed as in ref. 39. Briefly, the PGL4.23 LIC vector containing the normal insert served as the template for short cycle PCR using specially designed primers and Phusion High Fidelity DNA Polymerase. DpnI was used to digest the template. All clones were identified by colony PCR and sequenced. Primers are listed in Supplementary Table 3.

Chromatin state. Day 85 human foetal brain tissue DNase-seq intensities were downloaded from the Roadmap Epigenomics Project¹⁸. Mouse sequences were lifted to human using the UCSC liftOver tool (-minMatch = 0.2), and the maximum peak height was measured using UCSC bigWigAverageOverBed.

Genomic region enrichment analysis. Enrichment analysis was performed using GREAT (<http://great.stanford.edu/>) version 2.0.2 (ref. 20) with default parameters.

MER130 origins and activity. Sequences were queried using nhmmer¹² (from hmmer3.1-snap20121011 using --cut-ga threshold) with the Dfam² MER130 DF0000726.2 profile HMM model.

We first searched the RepBase catalogue of repeats⁹ to identify homologous repeat families, as well as the eukaryotic tRNA database⁴⁰ and Rfam⁴¹ to probe for RNAs from which a SINE element may have been derived¹³. We then downloaded vertebrate genome sequences from UCSC (<http://hgdownload.cse.ucsc.edu/downloads.html>) and vertebrate scaffolds from Genbank Release 196 (ftp://ftp.ncbi.nlm.nih.gov/genbank/genomes/Eukaryotes/vertebrates_other/).

A total of 75 available vertebrate genome drafts were searched. We observed repeat densities of 0.0–0.5 copies per Mb (Fig. 5a, Supplementary Table 2). The green sea turtle, with the highest repeat density of a modest 0.5 copies per Mb, had 1,124 MER130 copies present in its genome (Fig. 5a, Supplementary Table 2). However, these instances were quite diverged from each other, suggesting that MER130 has long been inactive in this lineage as well^{13,42}.

References

- Levo, M. & Segal, E. In pursuit of design principles of regulatory sequences. *Nat. Rev. Genet.* **15**, 453–468 (2014).
- Wheeler, T. J. *et al.* Dfam: a database of repetitive DNA based on profile hidden Markov models. *Nucleic Acids Res.* **41**, D70–D82 (2012).
- Britten, R. J. & Davidson, E. H. Repetitive and non-repetitive DNA sequences and a speculation on the origins of evolutionary novelty. *Q. Rev. Biol.* **46**, 111–138 (1971).
- Lowe, C. B., Bejerano, G. & Haussler, D. Thousands of human mobile element fragments undergo strong purifying selection near developmental genes. *Proc. Natl Acad. Sci. USA* **104**, 8005–8010 (2007).
- Kunaro, G. *et al.* Transposable elements have rewired the core regulatory network of human embryonic stem cells. *Nat. Genet.* **42**, 631–634 (2010).
- Chuong, E. B., Rumi, M. A. K., Soares, M. J. & Baker, J. C. Endogenous retroviruses function as species-specific enhancer elements in the placenta. *Nat. Genet.* **45**, 325–329 (2013).
- May, D. *et al.* Large-scale discovery of enhancers from human heart tissue. *Nat. Genet.* **44**, 89–93 (2012).
- Wenger, A. M. *et al.* The enhancer landscape during early neocortical development reveals patterns of dense regulation and co-option. *PLoS Genet.* **9**, e2003728 (2013).
- Jurka, J. *et al.* Repbase Update, a database of eukaryotic repetitive elements. *Cytogenet. Genome Res.* **110**, 462–467 (2005).
- Creyghton, M. P. *et al.* Histone H3K27ac separates active from poised enhancers and predicts developmental state. *Proc. Natl Acad. Sci. USA* **107**, 21931–21936 (2010).
- ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
- Wheeler, T. J. & Eddy, S. R. nhmmer: DNA homology search with profile HMMs. *Bioinformatics* **29**, 2487–2489 (2013).
- Bejerano, G. *et al.* A distal enhancer and an ultraconserved exon are derived from a novel retroposon. *Nature* **441**, 87–90 (2006).
- Spitz, F. & Furlong, E. E. M. Transcription factors: from enhancer binding to developmental control. *Nat. Rev. Genet.* **13**, 613–626 (2012).

15. Wenger, A. M. *et al.* PRISM offers a comprehensive genomic approach to transcription factor function prediction. *Genome Res.* **23**, 889–904 (2013).
16. Guturu, H., Doxey, A. C., Wenger, A. M. & Bejerano, G. Structure-aided prediction of mammalian transcription factor complexes in conserved non-coding elements. *Philos. Trans. R. Soc. Lond., B, Biol. Sci.* **368**, 20130029–20130029 (2013).
17. Portales-Casamar, E. *et al.* JASPAR 2010: the greatly expanded open-access database of transcription factor binding profiles. *Nucleic Acids Res.* **38**, D105–D110 (2010).
18. Bernstein, B. E. *et al.* The NIH roadmap epigenomics mapping consortium. *Nat. Biotechnol.* **28**, 1045–1048 (2010).
19. Clancy, B., Finlay, B. L., Darlington, R. B. & Anand, K. J. Extrapolating brain development from experimental species to humans. *Neurotoxicology* **28**, 931–937 (2007).
20. McLean, C. Y. *et al.* GREAT improves functional interpretation of cis-regulatory regions. *Nat. Biotechnol.* **28**, 495–501 (2010).
21. Eppig, J. T. *et al.* The Mouse Genome Database (MGD): facilitating mouse as a model for human biology and disease. *Nucleic Acids Res.* **43**, D726–D736 (2014).
22. Tong, C.-W. *et al.* Novel genes that mediate nuclear respiratory factor 1-regulated neurite outgrowth in neuroblastoma IMR-32 cells. *Gene* **515**, 62–70 (2013).
23. Sprague, J. *et al.* The Zebrafish Information Network: the zebrafish model organism database. *Nucleic Acids Res.* **34**, D581–D585 (2006).
24. Jacques, P. É., Jeyakani, J. & Bourque, G. The majority of primate-specific regulatory sequences are derived from transposable elements. *PLoS Genet.* **9**, e1003504 (2013).
25. Fode, C. *et al.* A role for neural determination genes in specifying the dorsoventral identity of telencephalic neurons. *Genes Dev.* **14**, 67–80 (2000).
26. Neves, das, L. *et al.* Disruption of the murine nuclear factor I-A gene (*Nfia*) results in perinatal lethality, hydrocephalus, and agenesis of the corpus callosum. *Proc. Natl Acad. Sci. USA* **96**, 11946–11951 (1999).
27. Steele-Perkins, G. *et al.* The transcription factor gene *Nfib* is essential for both lung maturation and brain development. *Mol. Cell Biol.* **25**, 685–698 (2005).
28. Campbell, C. E. *et al.* The transcription factor *Nfix* is essential for normal brain development. *BMC Dev. Biol.* **8**, 52 (2008).
29. Ayoub, A. E. *et al.* Transcriptional programs in transient embryonic zones of the cerebral cortex defined by high-resolution mRNA sequencing. *Proc. Natl Acad. Sci. USA* **108**, 14950–14955 (2011).
30. Gotea, V. *et al.* Homotypic clusters of transcription factor binding sites are a key component of human promoters and enhancers. *Genome Res.* **20**, 565–577 (2010).
31. Arnold, C. D. *et al.* Genome-wide quantitative enhancer activity maps identified by STARR-seq. *Science* **339**, 1074–1077 (2013).
32. Lowe, C. B., Bejerano, G., Salama, S. R. & Haussler, D. Endangered species hold clues to human evolution. *J. Hered.* **101**, 437–447 (2010).
33. *RepeatMasker Open-3.0*. <http://www.repeatmasker.org> (1996).
34. Visel, A. *et al.* A high-resolution enhancer atlas of the developing telencephalon. *Cell* **152**, 895–908 (2013).
35. Blow, M. J. *et al.* ChIP-Seq identification of weakly conserved heart enhancers. *Nat. Genet.* **42**, 806–810 (2010).
36. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).
37. Du, R., Li, S. & Zhang, X. A modified plasmid vector pCMV-3Tag-LIC for rapid, reliable, ligation-independent cloning of polymerase chain reaction products. *Anal. Biochem.* **408**, 357–359 (2011).
38. Tuteja, G. *et al.* Automated discovery of tissue-targeting enhancers and transcription factors from binding motif and gene function data. *PLoS Comput. Biol.* **10**, e1003449 (2014).
39. Camp, J. G., Jazwa, A. L., Trent, C. M. & Rawls, J. F. Intronic cis-regulatory modules mediate tissue-specific and microbial control of *angptl4/fiaf* transcription. *PLoS Genet.* **8**, e1002585 (2012).
40. Chan, P. P. & Lowe, T. M. GtRNAdb: a database of transfer RNA genes detected in genomic sequence. *Nucleic Acids Res.* **37**, D93–D97 (2009).
41. Burge, S. W. *et al.* Rfam 11.0: 10 years of RNA families. *Nucleic Acids Res.* **41**, D226–D232 (2013).
42. Kazazian, H. H. Mobile elements: drivers of genome evolution. *Science* **303**, 1626–1632 (2004).
43. Amemiya, C. T. *et al.* The African coelacanth genome provides insights into tetrapod evolution. *Nature* **496**, 311–316 (2013).

Acknowledgements

We thank Chris Amemiya for coelacanth genomic DNA; Travis Wheeler for guidance in using nhmmer; Craig Lowe for insights into repeat homology; Dr Susan McConnell and members of the Bejerano lab for helpful comments. This work was supported by a National Science Foundation Fellowship DGE-1147470 (J.H.N.), a Bio-X Stanford Interdisciplinary Graduate Fellowship (J.H.N.) and an NIH U01MH105949 Grant (G.B.). G.B. is a Packard Fellow and Microsoft Research Fellow.

Author contributions

J.H.N. and G.B. designed the study. J.H.N., T.C. and W.H. performed the experiments. J.H.N. and G.B. analyzed the data. J.H.N., W.H. and G.B. wrote the manuscript.

Additional information

Supplementary Information accompanies this paper at <http://www.nature.com/naturecommunications>.

Competing financial interests: The authors declare no competing financial interests.

Reprints and permission information is available online at <http://npg.nature.com/reprintsandpermissions>.

How to cite this article: Notwell, J. H. *et al.* A family of transposable elements co-opted into developmental enhancers in the mouse neocortex. *Nat. Commun.* 6:6644 doi: 10.1038/ncomms7644 (2015).