

1 **The novel lncRNA *lnc-NR2F1* is pro-neurogenic and mutated in human**
2 **neurodevelopmental disorders**

3

4

5 Cheen Euong Ang^{1,2,†}, Qing Ma^{3,4,†}, Orly L. Wapinski^{3,4,†}, Shenghua Fan⁵, Ryan A.
6 Flynn^{3,4}, Qian Yi Lee^{1,2}, Bradley Coe⁶, Masahiro Onoguchi^{3,4}, Victor H. Olmos¹, Brian T.
7 Do^{3,4}, Lynn Dukes-Rimsky⁵, Jin Xu^{3,4}, Koji Tanabe¹, Liangjiang Wang⁷, Ulrich Elling⁸,
8 Josef Penninger⁸, Yang Zhao^{3,4}, Kun Qu^{3,9}, Evan E. Eichler⁶, Anand Srivastava^{5,7,*},
9 Marius Wernig^{1,*}, Howard Y. Chang^{3,4,*}

10

11 ¹Institute for Stem Cell Biology and Regenerative Medicine and Department of
12 Pathology

13 ²Department of Bioengineering

14 ³Center for Personal Dynamic Regulomes

15 ⁴Department of Dermatology and Department of Genetics

16 Stanford University, Stanford, CA 94305, USA

17 ⁵J.C. Self Research Institute of Human Genetics, Greenwood Genetic Center,
18 Greenwood, S.C. 29646

19 ⁶Howard Hughes Medical Institute and Department of Genome Sciences, University of
20 Washington, Seattle, WA 98195

21 ⁷Department of Genetics and Biochemistry, Clemson University, Clemson, SC 29634

22 ⁸Institute of Molecular Biotechnology of the Austrian Academy of Science (IMBA),
23 Vienna Biocenter (VBC), Dr. Bohr Gasse 3, 1030 Vienna, Austria.

24 ⁹Present address: CAS Key Laboratory of Innate Immunity and Chronic Diseases, School
25 of Life Sciences and Medical Center, University of Science and Technology of China,
26 Hefei 230027, China.

27

28

29 †Co-first authors

30 *Co-Senior authors

31 Correspondence to: howchang@stanford.edu (H.Y.C), wernig@stanford.edu (M.W.)

32

33

34

35

36

37

38

39

40

41

42

43

44

45 **Abstract**

46 Long noncoding RNAs (lncRNAs) have been shown to act as important cell
47 biological regulators including cell fate decisions but are often ignored in human
48 genetics. Combining differential lncRNA expression during neuronal lineage induction
49 with copy number variation morbidity maps of a cohort of children with autism spectrum
50 disorder/intellectual disability versus healthy controls revealed focal genomic mutations
51 affecting several lncRNA candidate loci. Here we find that a t(5:12) chromosomal
52 translocation in a family manifesting neurodevelopmental symptoms disrupts specifically
53 *lnc-NR2F1*. We further show that *lnc-NR2F1* is an evolutionarily conserved lncRNA
54 functionally enhances induced neuronal cell maturation and directly occupies and
55 regulates transcription of neuronal genes including autism-associated genes. Thus,
56 integrating human genetics and functional testing in neuronal lineage induction is a
57 promising approach for discovering candidate lncRNAs involved in neurodevelopmental
58 diseases.

59 **(129 words)**

60 Eukaryotic genomes are extensively transcribed to produce long non-coding
61 RNAs (lncRNAs) in a temporally and spatially regulated manner¹. Until recently,
62 lncRNAs were often dismissed as lacking functional relevance. However, lncRNAs are
63 emerging as critical regulators of diverse biological processes and have been increasingly
64 associated with a wide range of diseases, based primarily on dysregulated expression².
65 lncRNAs represent a new layer of complexity in the molecular architecture of the
66 genome, and strategies to validate disease relevant lncRNAs are much needed. High-
67 throughput analyses have shown that lncRNAs are widely expressed in the brain and may
68 contribute to complex neurodevelopmental processes²⁻⁹. However, few studies have
69 examined the role of lncRNAs in brain development mostly due to technical difficulties.
70 Direct lineage conversion by the transcription factors Brn2, Ascl1 and Myt1l (termed
71 BAM factors in combination) into induced neuronal (iN) cells, recapitulates significant
72 events controlling neurogenesis programs¹⁰⁻¹², and therefore, it is a facile and informative
73 system to study the role of lncRNAs in the establishment of neuronal identity.

74 The noncoding genome has emerged as a major source for human diversity and
75 disease origins. Given that less than 2% of the genome encodes protein-coding genes, the
76 majority of the genomic landscape is largely encompassed by non-coding elements.
77 Efforts to identify genetic variation linked to human disease through genome-wide
78 association studies revealed a significant majority affecting the non-coding landscape.
79 Based on their expression and diversity in the mammalian brain, we postulate neuronal
80 lncRNAs may be recurrently affected by mutations that disrupt normal brain function.
81 Neurodevelopmental disorders manifest as a spectrum of phenotypes particularly early in
82 life¹³. Recent studies suggest that this diversity is the result of different combinations of
83 mutations in multiple genes, often impacting key pathways such as synapse function and
84 chromatin regulation. Nonetheless, despite recent findings that have greatly increased the
85 number of protein coding genes implicated in human intellectual disability and autism, a
86 majority of patients lack well-understood genetic lesions which include a large number of
87 inherited variants occur in noncoding regions that could not be interpreted¹⁴⁻²¹.

88 In this study, we used an integrative approach to identify lncRNA genes
89 important for human disease by incorporating high throughput cell fate reprogramming,
90 human genetics, and lncRNA functional analysis. In addition, we developed a pipeline to

91 enrich for lncRNAs with neuronal function and are associated with disease through focal
92 mutations in patients with autism spectrum disorder and intellectual disability (ASD/ID).
93 Furthermore, we show that one of these lncRNAs, *lnc-NR2F1* participates in neuronal
94 maturation programs *in vitro* by regulating the expression of a network of genes
95 previously linked to human autism.

96

97 **Results**

98 **LncRNA candidate loci are recurrently mutated in patients with** 99 **neurodevelopmental disorders**

100 LncRNAs have been associated with human diseases primarily through alterations
101 in expression levels²²⁻²⁴. However, little is known about mutations affecting the genomic
102 loci that encode lncRNAs. We previously profiled mouse embryonic fibroblasts (MEFs)
103 expressing doxycycline-induced BAM factors after 48 hours, 13 and 22 days of
104 expression (GSE43916). Surprisingly, annotation of the iN cell reprogramming
105 transcriptome revealed that the majority of regulated transcripts were in fact non-coding
106 elements (**Figure 1 – figure supplements 1A**). Specifically, 58% of the changed
107 transcripts corresponded to non-coding genes while only 42% of them corresponded to
108 coding genes. About two thirds of these non-coding transcripts were composed of novel
109 lncRNAs (**Figure 1 – figure supplements 1B**).

110 To study the vast non-coding iN cell reprogramming transcriptome, we developed
111 a rigorous pipeline to select lncRNAs with strong neuronal association (**Figure 1 – figure**
112 **supplements 1C**). We considered expression pattern during MEF-to-iN cell
113 reprogramming and across mouse brain development, protein coding potential, chromatin
114 enrichment, and “guilt-by-association” with neuronal Gene Ontology (GO) terms. We
115 observed 287 non-coding transcripts significantly changed expression during this time
116 course (RPKM >1, fold change >2, p-value<0.05) (**Figure 1 – figure supplements 1D**).
117 Notably, lncRNAs that increased expression during early stages of iN cell
118 reprogramming are more highly expressed in embryonic mouse brain, specifically in
119 ventricular and subventricular zones where neurogenesis occurs; whereas lncRNAs that
120 increased expression during intermediate to late stages of iN cell reprogramming were
121 more highly expressed in adult mouse brain, including in mature cortical layers (**Figure 1**

122 – **figure supplements 1D and E**). Furthermore, robust expression of iN cell lncRNAs in
123 the mouse brain confirmed that these transcripts are indeed bona fide neuronal
124 transcripts.

125 We next assessed lncRNA association with chromatin, reasoning that such RNAs
126 are more likely to exert gene regulatory function as non-coding RNAs. We performed
127 histone H3 immunoprecipitation, followed by deep sequencing of associated RNAs
128 (histone H3 RIP-seq), and discovered some of these 287 iN lncRNAs are chromatin-
129 associated compared to IgG control and input in neural precursor cells (NPCs) or adult
130 mouse brain tissue, suggesting that some of them may have more direct roles in gene
131 expression control at the chromatin (**Figure 1 – figure supplements 1F**). These 287
132 lncRNAs were then selected for further investigation. To further prioritize candidate
133 lncRNAs, we determined regulatory modules based on patterns of co-expression between
134 mRNAs and lncRNAs, inferring co-regulation from co-expression. Among the three
135 predominant modules found, one was strongly associated with neuronal GO terms, such
136 as *neurogenesis*, *axonogenesis*, and *synaptic organization and biogenesis* (**Figure 1 –**
137 **figure supplements 1G**). The remaining two modules consisted of broad non-neuronal
138 biological functions. Based on the criteria included in the pipeline, we nominated 35 iN
139 cell lncRNAs as most promising for possessing functions in the brain and confirmed their
140 expression qRT-PCR (**Figure 1 – figure supplements 1H, S2**). Collectively, these
141 results suggest that MEF-to-iN cell reprogramming can be used to identify lncRNAs
142 expressed in the developing and adult brain. (**Figure 1 – figure supplements 1A-H,**
143 **Figure 1 – figure supplements 2**).

144 We next interrogated these 35 mouse lncRNA loci in patients with autism
145 spectrum disorder and intellectual disability (ASD/ID). Firstly, we found that 28 of the 35
146 mouse lncRNA candidates have human synteny, and 10 loci were already annotated as
147 non-coding RNAs (**Fig. 1A and Figure 1 – figure supplement 3A**). We next overlapped
148 the 28 human lncRNA candidates and the remaining iN cell-lncRNAs coordinates to a
149 CNV morbidity map recently built from 29,085 patients diagnosed with a spectrum of
150 neurodevelopmental disorders and craniofacial congenital malformations, and 19,584
151 controls^{25,26}. This approach was motivated by the fact that the CNV morbidity map has
152 successfully identified novel syndromes characterized by recurrent mutations affecting

153 protein-coding genes of ASD/ID patients, and has offered mechanistic insight into the
154 drivers of the pathogenesis²⁵⁻²⁷.

155 Intersecting genomic coordinates of human lncRNA candidate loci to the CNV
156 morbidity map revealed 7 focal CNVs enriched in disease that overlap with 5 candidate
157 lncRNA loci: E (FLJ42709), H (LOC339529), Z (LOC100630918), D (LINC00094) and
158 O (LO467979) (Sequences in supplementary documents). Among these 7 focal CNVs, 5
159 events corresponded to small deletions in human lncRNA candidates E, H, Z, D and O.
160 Two lncRNA loci corresponding to lncRNAs H and Z were affected by two independent
161 and different focal CNVs. We verified that all five human lncRNAs are expressed during
162 human brain development (**Figure 1 – figure supplement 3A**). We then designed a
163 custom tiling array with dense coverage of the affected loci for comparative genomic
164 hybridization (CGH) to validate the focal CNVs in the genomic DNA from affected
165 individuals. 5 of 7 focal CNVs affecting the lncRNA loci were tested and validated (**Fig.**
166 **1A, 1B, Figure 1 – figure supplement 4A**). We could not test the last two CNVs
167 because patient DNA was no longer available.

168 One of the focally deleted lncRNA was NR_033883 (also known as or
169 LOC339529). This lncRNA locus is disrupted by two focal CNVs in two distinct ASD/ID
170 patients: 990914 and 9900589 (**Figure 1 – figure supplement 4A**). The NR_033883
171 locus neighbors the coding genes *ZFP238* (also known as *ZBTB18*, *ZNF238*, and *RP58*)
172 and *AKT3*. Because the human NR_033883 locus is most proximal but does not overlap
173 *ZFP238*, we propose to refer to this lncRNA as *lnc-ZFP238*. Intriguingly, we previously
174 identified *ZFP238* as a key downstream target of the *Ascl1* network during MEF to iN
175 cell reprogramming¹¹. Additionally, *ZFP238* has an important role in neuronal
176 differentiation during brain development²⁸⁻³⁰, and thus, *lnc-ZFP238* could have a
177 promising neurogenic role given its high expression pattern during the early stages of
178 direct neuronal reprogramming, as well as in postnatal mouse and human brain (**Figure 1**
179 **– figure supplement 2**).

180 Another focally deleted lncRNA was the locus harboring human FLJ42709 (or
181 NR2F1-AS1) (**Fig. 1B**) which is adjacent to the protein-coding gene *NR2F1* (also known
182 as *COUP-TF1*), encoding a transcription factor involved in neurogenesis and patterning
183^{7,31-38}. This lncRNA was previously annotated as “NR2F1-antisense 1” (NR2F1-AS1).

184 However, our RNA-seq analysis showed that at least one isoform of the mouse lncRNA
185 and all detected isoform of the human lncRNA are transcribed divergently from *NR2F1*
186 without antisense overlap⁷. For scientific accuracy, we therefore propose the name *lnc-*
187 *NR2F1*. We first asked whether the coding gene *NR2F1* could also be affected by the
188 focal CNVs in ASD/ID patients. Detailed statistical analysis of the primary data taking
189 into account the relative probe density suggested that inclusion of *NR2F1* is not
190 statistically significant compared to the control group. Moreover, we precisely mapped
191 the independent focal deletion found in patient 9900850 by CGH analysis and found only
192 the *lnc-NR2F1* locus to be disrupted (**Fig. 1B**). These results implicated the genetic
193 disruption of *lnc-NR2F1* as likely contributor to complex neurodevelopmental disorders.

194 Chromosomal aberrations encompassing the *lnc-NR2F1* locus and
195 additional genes on chromosome 5q14 have been previously reported in several patients
196 with neurodevelopmental deficits and congenital abnormalities (**Figure 1 – figure**
197 **supplement 4B and Supplement File 1**)³⁹⁻⁴². However, given that several genes are
198 affected by the deletions, the particular contribution of each gene was difficult to resolve.
199 Three patients (5-year-old girl, 5 and 7-year old boys) with a *de novo* deletion of
200 chromosome 5q14.3-15 were diagnosed with epileptic episodes, intellectual disability,
201 bilateral periventricular heterotopia in the temporal and occipital horns of the lateral
202 ventricles, minor dysmorphic facial features, developmental delay, and impaired to
203 negligible language skills (**Figure 1 – figure supplement 4B and Supplement File 1**).
204 The shared minimal deleted region between the patients spans 5.8Mb and encompasses
205 several annotated genes, amongst them *lnc-NR2F1*⁴¹. Clinical examination of one of the
206 patients that harbors a finer 6.3Mb interstitial deletion, showed macrocephaly (>98th
207 centile) and brain MRI revealed polymicrogyria. No cortical abnormalities were detected
208 on the brain MRI for the other two patients (**Supplement File 1**). More recently, an 8
209 year-old and 3 month-old boy with a *de novo* 582kb deletion was diagnosed with global
210 developmental delay, dysmorphic features, visual motor integration deficit, visual
211 perception disorder, mild conductive hearing loss, and severe fine motor skills
212 abnormalities³⁹ (**Figure 1 – figure supplement 4B and Supplement File 1**). Head
213 circumference was 8th centile. Brain MRI revealed bilateral optic nerve atrophy. The
214 582kb deletion affects the genes *NR2F1*, *lnc-NR2F1*, *FAM172A*, *POU5F2*, and

215 *MIR2277*. A four year-old girl with a balanced de novo paracentric chromosome 5
216 inversion, *inv(5)(q15q33.2)*, and microdeletions near the rearrangement breaking points
217 completely removed *NR2F1* and *lnc-NR2F1*⁴⁰ (**Figure 1 – figure supplement 4B and**
218 **Supplement File 1**). Additional genes on different chromosomes are affected by
219 microdeletions and could potentially contribute to the phenotype. The patient was
220 diagnosed with syndromic deafness, feeding difficulties, dysmorphism, strabism, and
221 developmental delay.

222 Across those patients with structural variation encompassing the *lnc-NR2F1* locus
223 in the literature, the minimal deleted region is approximately 230kb, a small area
224 encompassing the genes *NR2F1* and *lnc-NR2F1* (**Figure 1 – figure supplement 4B and**
225 **Supplement File 1**). The most notable overlapping phenotype consists of global
226 developmental delay, facial dysmorphism, and hearing loss. Hypotonia and
227 ophthalmological abnormalities are also common diagnoses³⁹ (**Figure 1 – figure**
228 **supplement 4B and Supplement File 1**). Phenotypic heterogeneity amongst patients
229 could be the result of dosage sensitive genes, polymorphisms on the unaffected allele,
230 genomic variability, gender, and age, amongst others.

231 Independent of the patients previously reported³⁹⁻⁴², we identified a paternally
232 inherited balanced translocation *t(5;12)(q15;q15)* in a 2 year and 7-month-old male
233 patient (CMS12200) (**Fig. 1C-D and figure 1 – figure supplement 4B**). Patient
234 CMS12200 was diagnosed with developmental delay, speech delay, significant
235 expressive language delay, and congenital infantile left eye esotropia (**Fig. 1C**). Physical
236 examination revealed small head size (head circumference 48.5 cm; 5th-10th centile), and
237 mild fifth finger clinodactyly bilaterally. Other physical features were normal. The
238 patient's father was diagnosed with dyslexia and stutters, and carried the identical *t(5;12)*
239 translocation. The patient's mother had a normal 46, XX karyotype (**Fig. 1C**).
240 Fluorescence *in situ* hybridization and whole genome sequencing defined the
241 chromosomal breakpoints with high precision and revealed that only the *lnc-NR2F1* gene
242 is disrupted in this patient (**Fig. 1D and 1E, Figure 1 – figure supplement 4C**). In
243 humans, three predominant isoforms of *lnc-NR2F1* have been detected in neuronal tissue.
244 The long isoforms (1 and 2) are affected by the chromosomal break, while the gene
245 structure of the short isoform (3) could remain unaffected based on the location of the

246 break (**Fig. 1D**). Further studies by Sanger sequencing of the 5q15 and 12q15 breakpoint-
247 specific junction fragments showed the identical breakpoints in the patient's father and
248 revealed a loss of 9 nucleotides at the 5q15 chromosome and a loss of 12 nucleotides at
249 the 12q15 chromosome in the patient and his father (**Fig. 1D**). The breakpoint at 12q15
250 occurred in a coding gene desert and did not disrupt any coding genes, and is predicted to
251 destabilize the affected transcript due to loss of 3' splice or polyadenylation signals (**Fig.**
252 **1D and Figure 1 – figure supplement 4D-E**). Importantly, whole genome sequencing
253 data indicated the absence of other deleterious mutations known to be associated with
254 autism or intellectual disability, to our knowledge (**Fig. 1D**). Also, the genes adjacent to
255 the break point (*FAM172A*, *ARRDC3*, *KIAA0625*, *USP15*) are not significantly changed
256 (**Figure 1 – figure supplement 4E**). Given that *lnc-NR2F1* is the only disrupted gene in
257 this patient family, it is possible that haploinsufficiency is the primary cause for this
258 syndrome and contribute to the phenotypes manifested across patients mentioned above.
259 Further studies including a larger sample size and independent cases are required to
260 conclusively link *lnc-NR2F1* mutations to multiple clinical symptoms described above.

261

262 **Molecular and functional characterization of *lnc-Nr2f1***

263 Given its potential involvement in neurodevelopmental disease, we next sought to
264 investigate the function of *lnc-Nr2f1*. We focused on mouse *lnc-Nr2f1* as experimental
265 approaches are more tractable in mouse models given the availability of a plethora of
266 genetic tools. Remarkably, and in contrast to many other lncRNAs, *lnc-Nr2f1* was not
267 only syntenically conserved in its genomic context (**Figure 2 – figure supplement 1A,**
268 **Supplemental File 6**)⁴³, but also highly sequence conserved among all human *lnc-*
269 *NR2F1* isoforms and the three exons in mouse *lnc-Nr2f1* (**Fig. 2A, Figure 2 – figure**
270 **supplement 1B**). In addition, we identified short stretches of sequence homology (termed
271 microhomology) near the conserved exons (exon 2 and 3 of human *lnc-NR2F1*) across
272 different species, with recurrent sequence motifs and motif order conserved across
273 different species (**Figure 2 – figure supplement 1C**). All of the above are features hinted
274 at lncRNA functional conservation across different species⁴³.

275 Mouse *lnc-Nr2f1* is induced as early as 48 hours after BAM factors are expressed
276 during MEF-to-iN cell conversion and peaks during mid-to-late stages of reprogramming

277 (Fig. 2B, Figure 1 – figure supplement 2 and Figure 2 – figure supplement 2A). In
278 the developing and adult mouse brain, *lnc-Nr2f1* showed a distinct region-specific pattern
279 of expression (Fig. 2C and Figure 2 – figure supplement 2B). In the developing
280 telencephalon at E13.5, *in situ* hybridization with a probe against *lnc-Nr2f1* revealed
281 strong expression in the caudolateral part of the mouse cortex and ganglionic eminences
282 (GE), similar to *Nr2f1* expression⁴⁴ (Fig. 2C and Figure 2 – figure supplement 2B).

283 To determine *lnc-Nr2f1*'s cellular localization, we performed single molecule
284 RNA-FISH in MEFs ectopically expressing *lnc-Nr2f1*, which revealed a nuclear and
285 cytoplasmic but predominantly nuclear localization (Fig. 2D and Figure 2 – figure
286 supplement 2C). Consistently, cellular fractionation from primary neurons dissected
287 from caudal region of the cortex showed endogenous localization of *lnc-Nr2f1* in both
288 nuclear and cytoplasmic fractions (Fig. 2E, Figure 3 – figure supplement 1H-I). Within
289 the nuclear fraction, *lnc-Nr2f1* is enriched in chromatin as assayed by histone H3 RNA
290 Immunoprecipitation followed by qRT-PCR (histone H3 RIP-qRT-PCR) in brain-derived
291 NPCs, postnatal and adult mouse brain (Fig. 2F).

292 We next wanted to explore potential functional roles of *lnc-Nr2f1* and assessed its
293 role during neuronal induction (Fig. Figure 3 – figure supplement 1). We therefore co-
294 expressed *lnc-Nr2f1* (NR_045195.1 or A830082K12Rik) with *Ascl1* and asked whether
295 it could promote neuronal conversion over *Ascl1* alone as previously observed with other
296 transcription factors (*Brn2*, *Myt1l*)^{45,46}. To that end we infected MEFs with *Ascl1* with or
297 without *lnc-Nr2f1*, and determined the ratio of TauEGFP-positive cells with neuronal
298 processes over the total number of TauEGFP-positive cells at day 7. We chose day 7 to
299 perform the experiment as it is an early time point for reprogramming. Indeed, the
300 addition of *lnc-Nr2f1* showed an approximately 50% (1.5-fold) significant increase in the
301 number of TauEGFP positive cells with neurites relative to *Ascl1* alone (Figure 3 –
302 figure supplement 1A-B). This surprising morphological maturation phenotype were
303 only previously observed only with co-expression of transcription factors (*Brn2* and
304 *Myt1l*) with *Ascl1* demonstrating a role of *lnc-Nr2f1* in neuronal morphological
305 maturation⁴⁷.

306 RNA-seq in sorted 7d MEF-iN cells expressing *Ascl1*, with and without co-
307 expression of *lnc-Nr2f1*, revealed 343 genes significantly changed expression between

308 data sets (RPKM>1, FDR corrected $p < 0.001$) (**Figure 3– figure supplement 1C and E**).
309 The vast majority of these genes were induced in expression upon *lnc-Nr2f1* expression,
310 with 311 genes up- and 32 down-regulated, suggesting *lnc-Nr2f1* may positively enhance
311 transcription. Gene ontology (GO) term enrichment of up regulated genes showed
312 significant enrichment in biological functions related to plasma membrane (*extracellular*
313 *region, cell adhesion, and transmembrane receptor tyrosine kinase activity*) and neuronal
314 function (*neuron projection, calcium binding, neuron differentiation, and axonogenesis*)
315 (**Figure 3 – figure supplement 1D**). These pathways are consistent with phenotype
316 observed for *lnc-Nr2f1* during iN cell reprogramming of promoting precocious
317 maturation programs (**Figure 3 – figure supplement 1D**). Amongst the up-regulated
318 genes are well-characterized neuronal and axon guidance genes such as *NeuroD1, Gap43,*
319 *Tubb4a, Ntf3, Nlgn3, Efnb3, Ntrk3, Bmp4, Sema3d, Slc35d3, Ror1, Ror2, Fgf7.*
320 Additionally, genes previously associated with neurological disorders were similarly up
321 regulated, such as *Mdga2, Clu, EphA3, Chl1, Cntn4, Cdh23, and Pard3b*⁴⁸.

322

323 ***Lnc-Nr2f1* is required for proper neuronal gene expression**

324 To investigate the contribution of *lnc-Nr2f1* to overall gene regulation, we sought to
325 achieve *lnc-Nr2f1* gain and loss-of-function in one experimental system. We reasoned
326 that mouse ES cells were the best way to accomplish loss-of-function. Since the
327 functional domains of *lnc-Nr2f1* RNA are unknown, non-coding sequences cannot be
328 turned into missense information by frame-shift mutations. Also, a large deletion
329 encompassing the entire 20kb *lnc-Nr2f1* locus may also inactivate interweaved intronic
330 regulatory elements and chromatin structure. Instead, we chose to insert a polyA
331 transcriptional termination signal to eliminate *lnc-Nr2f1* transcripts. We obtained mouse
332 ES cells that were previously genetically characterized in a genome-saturating haploid ES
333 cell mutagenesis screen⁴⁹. One of those ES cell clones had the mutagenesis cassette
334 containing an inverted (therefore inactive) splicing acceptor and polyA site inserted after
335 the first exon of *lnc-Nr2f1* (“Control” thereafter, **Fig. 3A**). The mutagenesis cassette was
336 designed to be conditionally reversible as it is flanked by combinations of loxP sites.

337 To achieve gain-of-function in the same cell system, we first overexpressed *lnc-*
338 *Nr2f1* in mES cells together with the proneural Ngn2 because that was shown to

339 efficiently and rapidly induce neurons from ES cells⁵⁰. Again, we quantified the neurite
340 length and neuron count, to test for its effect on increasing maturation kinetics.
341 Consistent with the fibroblast reprogramming results, we saw a significant increase in
342 neurite length upon *lnc-Nr2f1* overexpression though the number of neurons remained the
343 same (**Fig. 3B-D**). In contrast, overexpression of the coding gene *Nr2f1* did not induce
344 these phenotypes, and instead caused a drastic reduction in the number of neurons (**Fig.**
345 **3B-D**). These divergent results suggest that *lnc-Nr2f1* functions independently of *Nr2f1*.
346 RNA-seq analysis showed *lnc-Nr2f1* overexpression led to the induction of genes with
347 functions in axon guidance (*Sema5d*, *Epha1*) and neuronal projection development
348 (*Tubb6*, *Stmn2*, *Dtnbp1*) (**Fig. 3E**), confirming the cell biology phenotype (FDR
349 corrected $p < 0.10$, fold change > 1.5).

350 Next we turned to inactivate *lnc-Nr2f1* function in mouse ES cells. To generate an
351 isogenic knockout line, we treated the conditionally mutant ES cell line with Cre
352 recombinase, which resulted in an inversion of the polyA cassette, which in turn
353 terminates *lnc-Nr2f1* transcription (“*lnc-Nr2f1* KO” thereafter) (**Fig. 3F**). Since *lnc-Nr2f1*
354 is not expressed in ES cells, we differentiated *lnc-Nr2f1* KO and control mES cells into
355 induced neuronal cells by Ngn2 overexpression as above to assess the transcriptional
356 consequences of loss of *lnc-Nr2f1* in a neuronal context. RNA-seq showed that 348 genes
357 were differentially expressed between the control and the *lnc-Nr2f1* KO neurons, which
358 can be subsequently rescued with *lnc-Nr2f1* overexpression (FDR corrected $p < 0.10$) (**Fig.**
359 **3G and Figure 3 – figure supplement 1F**). Consistent with target genes in our *lnc-Nr2f1*
360 overexpression study, we found *lnc-Nr2f1* KO led to down regulation of neuronal
361 pathfinding and axon guidance genes such as *Sema6d* and proneural bHLH transcription
362 factor *Neurod2* as well as deregulation of genes associated with autism spectrum disorder
363 such as *Bdnf*, *Dcx* and *Nlgn3*⁴⁸ (**Fig. 3G**). The transcriptional abnormalities were
364 reversed by enforced expression of *lnc-Nr2f1* from a heterologous construct via lentiviral
365 transduction. The rescue data indicate that the downregulation of neuronal genes and the
366 upregulation of ectopic genes are caused by the loss of *lnc-Nr2f1* expression in the knock
367 out cells and unlikely by disruption of the nearby DNA regulatory elements due to the
368 insertion of the targeting cassette. Gene Ontology analysis of the downregulated genes in
369 *lnc-Nr2f1* KO neurons revealed enrichment for terms related to neural functions

370 (*regionalization, central nervous system development and neural precursor cell*
371 *proliferation*), whereas the upregulated genes are enriched in development of non-
372 neuronal tissues such as *circulatory system and skin development* (**Fig. 3H and 3I**).
373 Finally, the rescued genes overlap significantly with curated autism risk genes by Basu et.
374 al.⁴⁸ ($p=0.0012$, Chi-square) (**Figure 3 – figure supplement 1G**).

375 To further distinguish the function of *lnc-Nr2f1* vs *Nr2f1*, we generated *Nr2f1*
376 heterozygous and homozygous KO mouse ES cell lines using CRISPR/Cas9 (**Figure 3 –**
377 **figure supplement 2A**). When we performed qRT-PCR on the day 4 iN cells generated
378 from the Ctrl, *Nr2f1* heterozygous and homozygous null mES cells, we found that the
379 level of both *Nr2f1* and *lnc-Nr2f1* RNA transcripts did not change (**Figure 3 – figure**
380 **supplement 2B**). However, protein quantitation confirmed that *Nr2f1* protein level was
381 reduced or eliminated in the heterozygous lines (clone 21 and 44) or homozygous null
382 clones (Clone 2, 11 and 18), respectively (**Figure 3 – figure supplement 2C**). The *Nr2f1*
383 KO did not affect *lnc-Nr2f1* expression (**Figure 3 – figure supplement 2B**) and had no
384 impact on the neurite length or number in mES-iN cells (**Figure 3 – figure supplement**
385 **2D-E**). In summary, both gain and loss of function studies demonstrated that *lnc-Nr2f1*
386 plays a role in the transcriptional regulation of a gene network involved in neuronal
387 maturation pathways that ultimately resulted in faster acquisition of a mature neuronal
388 identity in both MEFs and mES cells and is functionally distinct from its neighboring
389 coding gene, *Nr2f1*.

390

391 **Mouse *lnc-Nr2f1* binds to distinct genomic loci regulating neuronal genes**

392 As described above, histone pull-down experiments suggested an association of *lnc-Nr2f1*
393 with chromatin. We therefore sought to map the precise *lnc-Nr2f1* genome wide
394 occupancy and performed Chromatin Isolation by RNA precipitation followed by
395 sequencing (ChIRP-seq) on day 4 mES-iN (**Fig. 4A, Figure 4 – figure supplement 1A-**
396 **B**). To minimize background sequencing we used even and odd probes targeting *lnc-*
397 *Nr2f1* in replicate experiments and only considered the overlapping peaks from both
398 experiments (**Fig. 4B**). Both even and odd probes pulled down *lnc-Nr2f1* efficiently.
399 There are 14975 peaks called by MACS, and the peak signals are consistent between
400 replicates (n=4) with little signal in RNase-treated control samples (**Fig. 4B**). We

401 obtained 1092 high confidence peaks with further filtering for the most significant and
402 reproducible binding events (see methods for filtering criteria). As an example, *lnc-Nr2f1*
403 binds to the intronic region of *Nrp2*, a gene with known roles in neuronal pathfinding
404 (**Fig. 4C**). To understand which gene ontology terms are enriched in the genes adjacent to
405 the mES-iN ChIRP peaks, we performed GREAT analysis and associated the 1092 *lnc-*
406 *Nr2f1* binding sites to 1534 genes. GO term analysis revealed that these genes that
407 enriched for neuronal terms such as *central nervous system development*, *synapse*
408 *organization and chemical synaptic transmission* (**Fig. 4D**). Using the publicly available
409 ChIP-seq peak sets for CTCF, enhancer, H3K27ac, H3K4me3 and PolIII obtained from
410 mouse adult cortex and E14.5 brain, we found significant enrichments of those peaks co-
411 localizing with the enhancers, H3K27ac and H3K4me3 marks relatively to the
412 background (**Fig. 4E and Figure 4 – figure supplement 1D**)⁵¹. DNA motif analysis of
413 *lnc-Nr2f1* binding sites revealed several basic helix loop helix (bHLH) factor motifs that
414 are significantly enriched (*NeuroD1*, *Atoh1*, *Olig2* and *Ptf1a*) (**Fig. Figure 4 – figure**
415 **supplement 1C**). To test whether *lnc-Nr2f1* and the bHLH factors are binding to same
416 genomic regions, we compared ChIP-seq data of bHLH neurogenic factors (Ngn2 and
417 *Ascl1*) with *lnc-Nr2f1* ChIRP-seq data (**Figure 4 – figure supplement 1F, G**). We found
418 that direct overlap between *lncNr2f1* ChIRP-seq peaks and ChIP-seq peaks are low but
419 statistically significant. 3.2% and 13.2% *lncNr2f1* ChIRP-seq peaks overlapped with
420 Ngn2 and *Ascl1* ChIP peaks, respectively (**Figure 4 – figure supplement 1F**). For peak-
421 associated genes, *lncNr2f1* target genes overlap very significantly with Ngn2 (37.9%) or
422 *Ascl1* (67.3%) target genes (**Figure 4 – figure supplement 1G**). These results suggest
423 that *lncNr2f1* and bHLH transcription factors such as *Ascl1* and Ngn2 may coordinately
424 regulate the same set of neuronal genes, and the majority of instances occur with *lnc-*
425 *Nr2f1* and the bHLH factors binding nearby but non-overlapping sites. Finally, to
426 understand whether the *lnc-Nr2f1* regulates the genes listed in figure 3E, we overlapped
427 the 1534 genes adjacent to the mouse *lnc-Nr2f1* ChIRP peaks with the genes up or
428 downregulated upon *lnc-Nr2f1* overexpression, we found 177 common genes between the
429 two lists (p<0.0001, Chi Square), suggesting that those genes might be direct *lnc-Nr2f1*
430 targets since they are occupied by *lnc-Nr2f1* RNA and are significantly altered when
431 manipulating *lnc-Nr2f1* (**Figure 4 – figure supplement 1E**)

432

433 **Human *lnc-NR2F1* shows isoform-specific chromatin binding**

434 The balanced chromosomal translocation t(5;12) detected in patient CMS12200
435 disrupts the long *lnc-NR2F1* while the short isoform appears unaffected. We therefore
436 hypothesized that the *lnc-NR2F1* might have isoform-specific functions and the long
437 isoforms are contributing to the phenotype observed in patient CMS12200. Given that
438 *lnc-NR2F1* is highly expressed in human brain tissue and it has high sequence
439 conservation between mouse and human (**Fig. 5A**), we next sought to determine whether
440 human *lnc-NR2F1* had a similar role in neuronal reprogramming as the mouse transcript
441 and whether the different isoforms may have distinct functions. Therefore, we
442 individually expressed each of the three human *lnc-NR2F1* isoforms in MEFs, and
443 measured their ability to enhance Ascl1- mediated neuronal reprogramming, as judged by
444 morphological complexity of TauEGFP cells (**Fig. 5B**). The long human *lnc-NR2F1*
445 isoform 2 significantly increased the proportion of TauEGFP cells with projections, albeit
446 with a slight smaller magnitude than the mouse lncRNA. Intriguingly, long isoform 1
447 inhibited neuronal maturation while the short isoform 3 had no significant effect (**Fig.**
448 **5B**). Thus, different isoforms of *lnc-NR2F1* may possess differential regulatory activity.
449 Long *lnc-NR2F1* isoforms disrupted by chromosomal translocation in patient CMS12200
450 can impact neuronal maturation, while the short *lnc-NR2F1* isoform remaining intact in
451 patient CMS12200 did not have a detectable effect on neuronal complexity.

452 Due to *lnc-NR2F1*'s strong association with chromatin and isoform-specific
453 function, we hypothesized that different domains of *lnc-NR2F1* may have differential
454 chromatin localization. To test this idea, we mapped the genome-wide localization of
455 different RNA domains in *lnc-NR2F1* using domain-specific Chromatin Isolation by
456 RNA precipitation followed by sequencing (domain ChIRP-seq)^{51,52}. We performed *lnc-*
457 *NR2F1* ChIRP-seq in human neural progenitor cell (hNPC) differentiated 12d from
458 human embryonic stem cells using dual SMAD inhibition protocol⁵³ (**Figure 5 – figure**
459 **supplement 1A**). We performed ChIRP-seq separately with two orthogonal probe sets
460 (termed odd and even sets) against two different domains of *lnc-NR2F1* (long isoform-
461 specific exon 11 and the short isoform-specific exon 7) and only accepted concordant
462 results between the odd and even probe sets. There are approximately 10-fold more

463 genomic occupancy for the long vs. short isoform of *lnc-NR2F1*: 4404 ChIRP-seq peaks
464 for exon 11 (n=4) vs. 415 ChIRP-seq peaks for exon 7 (n=4), respectively. It is unlikely
465 that low expression or inefficient pulldown of the short isoform are the cause of the
466 difference given that we detected comparable level of long and short isoform in hNPC
467 and obtained similar recovery of the long and short RNA isoform (**Figure 5 – figure**
468 **supplement 1B and C**). Consistent with our hypothesis of domain specific chromatin
469 localization, genomic occupancy sites of different *lnc-NR2F1* exons showed limited
470 overlap of peaks (**Figure 5C**), suggesting that the long- and short- specific exon might be
471 each binding to different genomic loci and regulating different subsets of downstream
472 genes. For example, only the long isoform-specific exon has a distinct binding site
473 surrounding *POLR1A* (**Figure 5D**).

474 We then assessed the transcriptional response which the three isoforms of human
475 *lnc-NR2F1* varied quantitatively: The shortest isoform, human *lnc-NR2F1* isoform 3, had
476 the lowest number of differentially regulated genes (5 downregulated and 4 upregulated)
477 compared to isoform 1 (1147 downregulated and 414 upregulated) and isoform 2 (141
478 downregulated and 45 upregulated) (**Figure 5 – figure supplement 1F**). This is
479 consistent with our hypothesis that the long isoform is the functional one in neurogenesis
480 and in patients with neurodevelopmental disorders including ASD.

481 After further filtering the peaks for high confidence peaks, we obtained 913 high
482 confidence peaks for long-isoform domain ChIRP and no peaks for short-isoform domain
483 ChIRP (see methods for filtering criteria). To further characterize *lnc-NR2F1* occupancy
484 patterns, high confidence ChIRP peaks were classified according to distance to putative
485 *cis*-genes (**Figure 5E and Figure 5 – figure supplement 1D**). Relative to the human
486 genome, the long isoform-specific peaks are more enriched in the exonic, intronic,
487 enhancer and promoter regions and depleted in the intergenic regions (**Figure 5E**).
488 Furthermore, ChIRP peaks are characterized by chromatin state model, which defines
489 human genome with 25 chromatin states using 12 biochemical features (histone
490 modifications, DNA accessibility, DNA methylation, RNA-seq and other epigenetic
491 signals)⁵⁴. *lnc-NR2F1* ChIRP peaks are enriched in promoter, enhancer, and transcribed
492 regions, compared to the whole genome which majority is in the quiescent state (**Figure**
493 **5 – figure supplement 1D**).

494 The 1361 high confidence peak associated genes by the long isoform-specific
495 domain ChIRP are all enriched for neuronal specific biological terms such *central*
496 *nervous system development*, *cell-cell adhesion* and *regulation of nervous system*
497 *development* (**Figure 5F**). There is also a significant overlap between the genes adjacent
498 to peaks from the long isoform-specific domain ChIRP and the genes adjacent to the
499 mouse *lnc-Nr2f1* ChIRP, indicating a possible conserved functions of *lnc-Nr2f1* between
500 mouse and human (207 genes overlapped, Chi-square test, $p < 0.0001$) (**Figure 5G**). We
501 also observed a significant overlap between the genes adjacent to the peaks from the long
502 isoform-specific domain ChIRP and the autism risk genes (114 genes, Chi-square test,
503 $p < 0.0001$) (**Figure 5H**). Notably, peaks of long isoform-specific exon 11 ChIRP are
504 enriched for multiple basic-helix-loop-helix motifs (*Ascl1*, *NeuroD1*, *Olig2* and *Atoh1*)
505 which all share the CANNTG motif⁵⁵ (**Figure 5 – figure supplement 1E**). Mouse ES-iN
506 cell ChIRP peaks are also enriched for similar motifs suggesting possibly conserved
507 mechanisms of *lnc-NR2F1* in human and mouse. Given the pervasive roles of bHLH
508 proteins in neuronal development and in induced neuronal reprogramming, the binding
509 preference of *lnc-NR2F1* suggests a biochemical basis for the functional cooperativity
510 with proneural bHLH factors (*Ngn2* and *Ascl1*). In summary, we conclude that the *lnc-*
511 *NR2F1* isoforms have different genomic occupancy and transcriptional effects. The long
512 isoform showed most biological activity and chromatin binding and is also solely affected
513 in patient CMS12200.

514

515 **Discussion**

516

517 Given the stringency required to rewire a cellular state from an unrelated lineage,
518 factors expressed during direct reprogramming likely have an active role in the
519 establishment of the new cell identity. Direct neuronal lineage induction represents a
520 synchronized and streamlined conversion of cell fate, and should be a powerful system to
521 enrich for lineage-specific regulatory factors. In this study, direct conversion of
522 fibroblasts to induced neuronal cells enabled the identification of lncRNAs with unique
523 properties in establishing neuronal fate via neurogenesis or maturation programs. The
524 pipeline described in this study may be extrapolated to identify potential lncRNA
525 regulators of other specific cellular states.

526 Because the brain is the organ with the greatest number of unique cell-type
527 specific lncRNAs⁵⁶, our approach may be useful to identify lncRNAs with roles in neural
528 lineage specification. Indeed, we identified *lnc-Nr2f1* as functional player in neuronal
529 maturation and pathfinding. Remarkably its sequence is remarkably conserved between
530 the first few exons of mouse and human *lnc-NR2F1* which is an atypical pattern for
531 lncRNAs^{57,58}. Consequently, we found that there is high synteny, sequence and
532 microdomain conservation between mouse and human *lnc-NR2F1*. These observations
533 suggest that some lncRNAs may have been functionally conserved throughout evolution.

534 In this study, we focus on the functional characterization of *lnc-NR2F1* locus
535 because it is recurrently mutated in human patients with ASD/ID. We identified a patient,
536 whose genome harbors a balanced translocation disrupting the *lnc-NR2F1* locus without
537 any other detectable pathogenic genetic variations and shows abnormal
538 neurodevelopmental symptoms; therefore, implicating this lncRNA as a critical regulator
539 of brain development and function. The father of the proband carries the same
540 translocation and suffers from dyslexia and stuttering, suggesting that the phenotype may
541 be transmitted in a Mendelian manner. However, the much milder phenotype of the father
542 implies that additional genes, environmental factors, or compensatory neuronal circuitry
543 acquired during adulthood may influence the severity of the outcome.

544 Given the close genomic proximity of *lnc-NR2F1* and *NR2F1* increased attention
545 must be devoted to consider the possibility of a contribution of the coding gene. Since
546 unknown regulatory elements for the coding gene could be affected the human genetics
547 data are not decisive. However, several functional experiments point to a contribution of
548 *lnc-NR2F1* rather than the coding gene. First, gain and loss of function studies as well as
549 chromatin localization clearly show that *lnc-NR2F1* acts in trans to affect gene
550 expression. Second, we can rescue the phenotype of the *lnc-Nr2f1* KO by overexpression
551 of *lnc-Nr2f1* mRNA. Third, *lnc-Nr2f1* overexpression in *lnc-Nr2f1* KO cells does not
552 affect *Nr2f1* expression. Fourth, ChIRP-sequencing of *lnc-Nr2f1* in induced neurons
553 derived from mES cells does not show binding of *lnc-Nr2f1* in the *Nr2f1* promoter region
554 (**Fig. 4F**). The only definitive answer may be obtained from human postmortem tissue
555 analysis of affected patients.

556 Neurodevelopmental and neuropsychiatric disorders are complex diseases
557 manifesting in a spectrum of phenotypes. We integrated lncRNA expression pattern, *in*
558 *vitro* functional screen, and human genetic data to pinpoint potentially causal lncRNAs.
559 We concentrated on genomic lesions affecting lncRNAs, which have been largely
560 understudied regulatory factors in these diseases, and connected them to specific
561 phenotypes. We found several of the lncRNA candidates were disrupted by focal
562 chromosomal aberrations in patients diagnosed with ASD/ID, establishing a link between
563 human disease and lncRNA function. The advent of next generation sequencing has
564 greatly improved the ability to pinpoint causal disease mutations in protein coding genes
565 including the discovery of novel autism genes. Most of the other functional regions of the
566 genome, however, have largely been ignored as part of exome sequencing approaches.
567 While full genome sequencing of patients is beginning, functional interpretation remains
568 a daunting challenge. We present a strategy to begin to characterize the functionally
569 important non-coding regions as it relates to disease. Our work highlights lncRNA
570 mutations as an understudied and important potential next frontier in human genetics
571 related to neurodevelopmental disease.

572

573 **Acknowledgements**

574 We thank Cindy Skinner, Mrs. Sydney Ladd, Dr. Barbara R. DuPont, Dr. Katie R.
575 Clarkson for patient recruitment and evaluation and members of our labs for discussion
576 and advice. We thank the Stanford Functional Genomic Facility especially Vanita for her
577 assistant in the project. This project is supported by NIH RC4-NS073015 (H.Y.C.,
578 M.W.), P50-HG007735 (H.Y.C.), California Institute for Regenerative Medicine (M.W.,
579 H.Y.C.), NIH R01 HD39331 (A.K.S.) and Self Regional Healthcare Foundation Funds
580 (A.K.S.). C.E.A. was supported by California Institute of Regenerative Medicine
581 Training Grant and Siebel Foundation. Q.M was supported by Stanford Dean's
582 Fellowship. O.L.W. was supported by a NSF fellowship. M.W. is a NYSCF–Robertson
583 Stem Cell Investigator. Haplobank is generously funded by Nestlé Institute of Health
584 Science NIHS as well as the Austrian National Bank (OeNB) and Era of Hope/National
585 Coalition against Breast Cancer/DoD. U.E. is supported by the Austrian Academy of
586 Sciences, the Austrian National Bank (OeNB), and is a Wittgenstein Prize fellow. J.M.P.
587 is supported by an Advanced ERC grant and an Era of Hope/DoD grant. E.E. is an
588 Investigator of the Howard Hughes Medical Institute.

589

590 **Data availability statement**

591 A summary table containing all the lnc-Nr2f1 mutation reported in the literature the
592 ChIRP-sequencing probes (**supplemental file 1**), the ChIRP-sequencing probes
593 (**supplemental file 2**) Datasets used in this paper and their corresponding source

594 (supplemental file 3), the qRT-PCR primer sequences (supplemental file 4), RNA FISH
 595 (supplemental file 5) and the sequence conservation (supplemental file 6) can be found
 596 in the supplementary documents. The datasets generated during and/or analyzed during
 597 the current study are available in the NIH GEO repository (GSE115079)
 598
 599
 600
 601
 602

603 Key resource table

604

Reagent (species) resource	type or	Designation	Source	Identifier	Additional information
Antibody		Rabbit polyclonal anti-H3	Abcam	ab1791 (RRID:AB_302613)	
Antibody		Goat polyclonal anti-Sox1	R&D	AF3369 (RRID:AB_2239879)	IHC (1:50)
Antibody		Rabbit polyclonal anti- β -III-tubulin	Covance	Discontinued	IHC (1:1000)
Antibody		Mouse monoclonal anti-Nestin	R&D	MAB1259 (RRID:AB_2251304)	IHC (1:1000)
Antibody		Rabbit monoclonal anti-HSP90	Cell Signalling	4877 (RRID:AB_2233307)	WB (1:2500)
Antibody		Rabbit monoclonal anti-Nr2f1	Cell Signalling	6364 (RRID:AB_11220432)	WB (1:1000)
Chemical compound, drug		SB431542	Tocris	1614	
Chemical compound, drug		LDN198189	MiliporeSigma	5.09882.0001	
Chemical compound, drug		CHIR99021	StemGen	04-0004	
Chemical compound, drug		PD0325901	Axon	1408	
Chemical compound, drug		Leukemia Inhibitory Factor	Generated in the lab		
Cell line (H. Sapiens)		Human: 293T	ATCC	CRL-3216 (RRID:CVCL_0063)	
Cell line (H.		Human: H9 hESC line	UWisco	H9	

Sapiens)		nsin	(RRID:CVCL_9773)	
Cell line (H. Sapiens)	Human: SK-N-SH	ATCC	HTB-11 (RRID:CVCL_0531)	
Cell line (Mus musculus)	Tau: EGFP Mouse embryonic fibroblasts	Generated in the lab		
Cell line (Mus musculus)	Mouse: Haploid ES cells	Obtained from Penninger lab (Elling et al. 2011)		
Genetic reagent (<i>M. musculus</i>)	B6.129S4(Cg)- <i>Mapt</i> ^{tm1(EGFP)Klt/J}	Jackson	29219 (RRID:IMSR_JAX:004779)	
Recombinant DNA reagent	TetO- <i>lnc-Nr2f1</i> (Mouse)	This paper		
Recombinant DNA reagent	TetO- <i>lnc-NR2F1-I</i> PGK blast (Human)	This paper		
Recombinant DNA reagent	TetO- <i>lnc-NR2F1-II</i> PGK blast (Mouse)	This paper		
Recombinant DNA reagent	TetO- <i>lnc-NR2f1-III</i> PGK blast (Mouse)	This paper		
Recombinant DNA reagent	TetO-NR2F1 (Mouse)	This paper		

605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621

622

623

624

625

626

627

628 **Materials and Methods**

629 This method section was organized into four categories: animal and human protocols, cell
630 culture, computational and sequencing methods and biochemistry. Within each category,
631 method descriptions were arranged in the order they appear in figures. The ChIRP
632 probes, public datasets, qRT-PCR primers and RNA FISH probes are available in the
633 supplemental File 2, 3, 4 and 5 respectively.

634

635 **Animal and human protocols**

636

637 Animal

638 All mouse work was performed according to IACUC approved protocols at Stanford
639 University. Samples in the paper were obtained without determining their sex. All
640 animals were housed in an animal facility with a 12hr light/dark cycle.

641

642 Human subjects

643 The study protocol, consent form, consent to publish and privacy practices were reviewed
644 and approved by the Institutional Review Board of the Self Regional Healthcare,
645 Greenwood, SC (Reference number Pro00074882).

646

647 **Cell culture and tissue dissection**

648

649 Cell culture

650 All cell lines (SK-N-SH, 293T) were purchased from ATCC and were verified by the
651 manufacturer by STR profiling. They were also screened for mycoplasma and cultured
652 using recommended conditions. Mouse embryonic fibroblasts (MEF) were derived from
653 E13.5 Tau::EGFP embryos and cultured in MEF media [500ml of DMEM (Gibco), 50ml
654 of Cosmic Calf Serum (Thermo Scientific), 5ml of Non-essential amino acid, 5ml of
655 Sodium Pyruvate, 5ml of Penicillin/Streptomycin (Gibco), 4ul of β -mercaptoethanol
656 (Sigma)].

657

658 Mouse haploid embryonic stem cells were cultured in mouse embryonic stem cell media
659 [341.5ml DMEM (Gibco), 50ml Knockout Serum Replacement (Gibco), 12.5ml of
660 Cosmic Calf Serum (ThermoScientific), 4.2ml of Penicillin/Streptomycin, 4.2ml of Non-
661 essential amino acid, 4.2ml of Sodium Pyruvate (Gibco), 4ul of β -mercaptoethanol
662 (Sigma) with leukemia inhibitory factor, 3 μ M of CHIR99021 and 1 μ M of PD3259010
663 (Both Tocris, Final concentration)].

664

665 Human embryonic stem cells (HUES9) were cultured in mTESR media (Stem Cell
666 Technologies). The experiments were performed in accordance with California State
667 Regulations, CIRM Regulations and Stanford's Policy on Human Embryonic Stem Cell
668 Research.

669

670 Mouse postnatal/adult brain dissection

671 Briefly, forebrains were dissected from TauGFP heterozygous E13.5 embryos in cold
672 HBSS, triturated in DMEM/F12 media, filtered through a 70um filter and cultured in
673 monolayer. Neural stem cells (NSC) were propagated in DMEM/F12 with N2 and B27

674 supplements (Invitrogen) with 20ng/ml of FGF2 and 10ng/ml of EGF. Postnatal brains
675 (Postnatal day 0) and adult brains (three weeks old) were obtained from C57BL6 mice.
676 To obtain postnatal brains, pups were anaesthetized in an ice bath before the whole brain
677 was removed. To obtain adult brains, mice were euthanized using cervical dislocation
678 before dissecting the whole brain out. For both adult and postnatal brains, they were
679 manually dissociated to fine pieces before being digested in 0.25% trypsin for 30
680 minutes. They were triturated from time to time until a clear suspension was obtained.
681 The cells were spun down at 1000rpm for 5 minutes before proceeding to glutaraldehyde
682 fixation.

683

684 Reprogramming of mouse fibroblasts to induced neuronal cells (iN cells)

685 We followed protocols previously described (Wapinski et al., 2013). Briefly, mouse
686 embryonic fibroblasts harvested from E13.5 Tau::EGFP embryos were plated at a density
687 of 25000 cells/cm². The next day, lentiviruses carrying TetO-FUW-ASCL1 and FUW-
688 rtTA were added. Doxycycline (Final concentration: 2µg/ml, Sigma) in MEF media was
689 added to the wells. Media was changed to neuronal media [N2 + B27 + DMEM/F12
690 (Invitrogen) + 1.6ml Insulin (6.25mg/ml, Sigma)] + doxycycline two days after
691 doxycycline induction. Subsequently, media was changed every three days. To obtain a
692 pure population of day 7 TauEGFP positive neurons, the cells were digested using 0.25%
693 trypsin (Invitrogen) and subjected to FACS. Forward and side scatters were used to
694 exclude doublets and dead cells. The gating for GFP was set with a negative control
695 (MEF).

696

697 Maturation screen for *lnc-NR2F1*

698 To examine whether the lncRNA candidates can facilitate mouse embryonic fibroblasts
699 (MEFs) to induced neuronal cells reprogramming, mouse and the three human isoforms
700 for *lnc-NR2F1* were synthesized (sequences in the supplementary documents) and cloned
701 into the TetO-FUW or TetO-PGK-blast^R backbone respectively (available from
702 Addgene). To examine whether those lncRNA candidates can help facilitate the
703 maturation process, the number of MAP2 positive neuronal cells with projections three
704 times the diameter of the cell body was counted at day 7 and normalized to the total
705 number of MAP2 positive cells. For neurite length measurement, Simple Neurite Tracer
706 (ImageJ) was used manually to track neurite.

707

708 Reprogramming of mouse embryonic stem cells to induced neurons

709 We followed the protocol previously described⁵⁰. Mouse embryonic stem cells were
710 plated single cell and infected the next day with TetO-NGN2-T2A-PURO^R and FUW-
711 rtTA. Doxycycline was added to the wells the next day. To select for only Ngn2
712 transducing cells, puromycin (Final concentration: 2µg/ml, Sigma) was added in addition
713 to doxycycline the next day and kept for 3 days.

714

715 Generating *lnc-Nr2f1* KO ES-iN cells

716 We obtained mouse ES cells that were previously generated in a genome-saturating
717 haploid ES cell mutagenesis screen⁵⁹. We identified one ES cell clone had the
718 mutagenesis cassette containing a splicing acceptor and polyA site inserted after the first
719 exon of *lnc-Nr2f1*. The orientation of the polyA site is in reverse from the transcription

720 direction of *lnc-Nr2f1* so it's non-disruptive. The insertion is confirmed by PCR and
721 sequencing. For generating *lnc-Nr2f1* KO clones, we did nucleofection of cre
722 recombinase to invert polyA cassette since the polyA cassette is flanked by loxP sites.
723 After nucleofection, we plated cells at low density and picked single colonies for testing
724 the polyA inversion. The control and KO clones were then expanded for a few passages,
725 allowing majority of them to become diploid cells. The homozygous diploid cells were
726 then plated at 300K cells/6 well in mES media at day 0. They were then infected with
727 TetO-Ngn2-T2a-puro, FUW-rtTA and TetO-GFP the next day. At day 2, the media was
728 changed to neuronal media [N2 + B27 + DMEM/F12 (Invitrogen) + 1.6ml Insulin
729 (6.25mg/ml, Sigma)] and doxycycline (Sigma, Final conc: 2µg/ml) was added. At day 3,
730 puromycin (Sigma, 2ug/ml) was added to the neuronal + dox media. RNA was harvested
731 four days after dox induction.

732

733 Generating *Nr2f1* KO ES-iN cells

734 For CRISPR/Cas9 genome editing of *Nr2f1*, gRNAs targeting second exon of *Nr2f1* are
735 cloned to a plasmid (pSpCas9(BB)-2A-Puro, pX459, Addgene #62988) expressing both
736 the Cas9 protein and the gRNA. gRNA sequences were designed using the online tool
737 (<http://crispr.mit.edu/>) provided by the Zhang lab (gRNA sequence used:
738 CATGTCCGCGGACCGCGTCG). 24-48 hours after ES cell nucleofection, puromycin
739 was added to select for 2-3 days. The remaining cells were plated at plate 100, 300, 1000,
740 3000 cells per plate for picking single colonies. Genomic DNA of each single colony was
741 extracted using QuickExtract™ DNA Extraction Solution (Epicentre, QE09050). This
742 extract was then used in a PCR of the genomic region that had been targeted for knock
743 out (Fwd primer: AGAGACACCTGGTCCGTGAT. Reverse primer:
744 GAGCCGGTGAAGGTAGATGA). PCR products were then Sanger sequenced to
745 identify clones that would result in frameshifts and truncated *Nr2f1*. Sequence alignment
746 and genomic PCR primer design was carried out using SnapGene software and cutting
747 efficiency is calculated using web tool TIDE (<https://tide-calculator.nki.nl/>).

748

749 Computational and sequencing methods

750 LncRNA discovery pipelines (related to Figure 1 – figure supplement 1)

751 TopHat was used for de novo alignment of paired-end reads for each of the samples. An
752 assembled transcriptome was built from merged time points using Cuffmerge function.
753 The de novo iN transcriptome was compared to RefSeq genes and annotated protein
754 coding genes were removed, while non-coding genes annotated as “NR” were kept.
755 Expression level of genes was calculated in unit of fragments per kilobase of exon model
756 per million mapped fragments (FPKM). Genes with low FPKM (average log₂ FPKM
757 across all samples less than 1) were removed. Genes with p-value<0.05 and at least two-
758 fold expression change during iN reprogramming were defined as significant.

759

760 Histone H3 RIP-seq (related to Figure 1 -figure supplement 1F)

761 RNA isolated from H3 RIP was amplified and converted to cDNA using Nugen
762 Ovation RNA-Seq System V2. The product was sheared using Covaris to 100-300bp.
763 Libraries were prepared using SPRIworks system for Illumina sequencing. The following
764 antibodies were used for RIP: Rabbit anti-H3 (Abcam ab1791) and rabbit IgG (Abcam
765 ab37415). For the H3 RIP-seq analysis we used a similar pipeline to bulk RNAseq

766 assays. We first remove duplicate reads, clip adaptor sequences, discard short reads.
767 Reads were then aligned to mm9 using Tophat2. Using samtools we convert the files
768 from sam file to bam format. Filtered reads are normalized to sequence depth, and
769 subsequently we calculate RPKM. The RIP-seq experiment was conducted with H3 and
770 IgG antibodies. We sequenced both libraries and also 1% input material. To determine
771 whether a lncRNA is enriched we compute the number of reads from H3 relative to input,
772 and IgG relative to input using an in-house Perl script ([rnaexp_rpkm.pl](#)). We then
773 calculated the fold-enrichment between H3 and IgG RIPs. Since the background was very
774 low, anything greater than 2-fold and $p < 0.05$ between H3 and IgG was considered
775 enriched. Experiments in NPCs, MEFs, and adult brain were conducted in biological
776 replicates. Only lncRNAs reproducibly enriched in the H3 RIP from biological replicates
777 were considered chromatin associated and display as binary in the Figure 1 – figure
778 supplement 1F.

779

780 Co-expression analysis for lncRNAs (related to Figure 1 – figure supplement 1G)

781 We first obtained mouse RNAseq data from ENCODE, and calculated the RPKMs for all
782 transcripts including coding and non-coding. We then for each non-coding RNA,
783 calculated the Pearson correlation of the non-coding RNA with every coding transcript. If
784 the correlation is greater than 0.5, this non-coding RNA was defined as positively
785 correlated with the coding gene, and if the correlation is less than 0.5, it was defined as
786 negatively correlated. We then obtained a matrix of coding genes versus non-coding
787 genes, with positive (+1) and negative (-1) correlations as values in the matrix. We then
788 use the GeneSets function in Genomica software from (<http://genomica.weizmann.ac.il/>),
789 and generated a enrichment ($-\log(p\text{-values})$) matrix for lncRNAs associated with Gene
790 Ontology terms based on similar expression pattern with mRNAs. The default settings set
791 in the software were used.

792

793 Overlap with CNV morbidity map (synteny of coordinates, significance calculation)
794 (related to Figure 2A)

795 To find syntenic conservation from mouse to human, UCSC Genome Browser tool
796 Lifter was used. To determine the potential role of lncRNAs in neurodevelopmental
797 disorders, we analyzed array CGH profiles from 29,085 children with intellectual
798 disability and developmental delay that were submitted to Signature Genomics
799 Laboratories, LLC, for clinical microarray-based CGH. The CNV map intersecting
800 lncRNAs was compared with that of 19,584 healthy controls²⁵(dbVar nstd100). Focal
801 enrichment was calculated using fishers exact test statistics and odds ratios comparing
802 cases and control CNV counts at each locus. Validation of focal CNVs affecting
803 lncRNAs of interest was performed on a custom 8-plex Agilent CGH array using
804 standard methods²⁵.

805

806 Ingenuity Variant Analysis (related to Figure 2A)

807 Using Ingenuity Variant Analysis software, we filtered 4,038,671 sequencing variants
808 and obtained a list of 45 variants possibly related with the patient's phenotype. This list
809 included three structural variants (deletions) and one gene fusion. We verified by Sanger
810 sequencing each of the variants associated with disease and found them to be either false
811 positives or non-causative.

812 The 16.8 Mb large deletion on Chromosome 11 was found to be false positive
813 based on the observation of heterozygosity in the deleted region in whole genome
814 sequencing data. The same false-positive was also shown in the whole genome
815 sequencing data of other two translocation patients. The gene fusion between *ARHGEF3*
816 on Chromosome 3 and *TRIO* on Chromosome 5 was determined as false positive by
817 Sanger sequencing. One fragment of *ARHGEF3* (132bp) intron sequence was inserted
818 into an intron of *TRIO*. The insertion led to the false detection of gene fusion. RT-PCR
819 proved that mRNA splicing of *TRIO* was not affected by this insertion and qRT-PCR
820 proved that the expression level of *TRIO* was not affected. For the rest of variants, we
821 reviewed the reported functions of genes having these variants and copy number variation
822 information in these regions in Database of Genome Variants. We found seven variants
823 occurred in the genes having closely related functions with patient's phenotype and also
824 not extensively covered by CNVs in Database of Genome Variants. We performed
825 Sanger sequencing for these seven variants in patient's family (the father and son having
826 the same translocation and the father also having dyslexia and stutter). Four variants were
827 found to be false positive. Both the patient and healthy mother possessed two variants.
828 All three family members possessed one variant. Overall, we have not found a promising
829 disease causing variants other than the translocation found in patient CMS12200.

830

831 Histone H3 RIP-qRT-PCR (related to Figure 2F)

832 Approximately $20\text{-}50 \times 10^6$ cells were used for each experiment. Cells were crosslinked
833 with 1% formaldehyde. Cell pellet was resuspended in equivalent volume of Nuclear
834 Lysis Buffer (i.e. 100mg- 1mL buffer) (50 mM Tris-Cl pH 7.0, 10 mM EDTA, 1% SDS,
835 100x PMSF, 50x protease inhibitors, and 200x Superase inhibitor). Chromatin was
836 sheared using Covaris sonicator until DNA was fragmented to 200-1000 bp range and
837 diluted 2-fold using Dilution Buffer (0.01% SDS, 1.1% Triton X 100, 1.2 mM EDTA,
838 16.7 mM Tris-Cl pH 7.0, 167 mM NaCl, 100x PMSF, 50x protease inhibitors, and 200x
839 Superase inhibitor). Samples were incubated with 5 μg of H3 or IgG antibody rotating
840 overnight at 4°C. Protein A dynabeads (50uL) were washed in Dilution buffer and added
841 to the chromatin for 2 hours rotating at room temperature. Immunoprecipitate fraction
842 was washed four times with Wash buffer (100 mM Tris-Cl pH 7.0, 500mM LiCl, 1%
843 NP40, 1% sodium deoxycholate, and PMSF). Subsequently, the immunoprecipitate
844 fraction was eluted from beads by vortexing for 30 minutes at room temperature using
845 elution buffer (50mM sodium bicarbonate, and 1% SDS). Immediately after 5% of 3M
846 Sodium Acetate was added to neutralize pH. Proteinase K treatment proceeded for 45
847 minutes at 45C, followed by RNA extraction using Trizol. Isolated RNA was subjected to
848 DNase treatment using TurboDNase and purified by phenol-chloroform extraction and
849 ethanol precipitation. For qRT-PCR analysis we used Roche's Lightcycler and
850 Stratagene's RT kit.

851

852 RNA-seq library preparation (related to Figure 3E, 3G, Figure 3 – figure supplement 1C)

853 We followed protocols previously described (Wapinski et al., 2013). Briefly, for the
854 RNA-sequencing experiment in figure 3, 4 and S7, libraries were produced from poly-A
855 enriched mRNA using TruSeq kit (Illumina). They were subsequently sequenced using
856 the NextSeq or HiSeq platform producing paired ends reads.

857

858 RNA-seq analysis for loss and gain of function analysis (related to Figure 3E, 3G, Figure
859 3 – figure supplement 1C)

860 Reads obtained were first mapped using Tophat. Expression for each gene was calculated
861 using Cuffdiff (Figure 3E, Figure 3G) or DEGSeq (Figure 3 – figure supplement 1C)
862 using default settings. For DEGSeq briefly, only properly paired mapped reads were used
863 ⁶⁰. DEGSeq selected longest transcript for each gene, when multiple isoforms were
864 found. Raw counts for each sample were merged into a table and transformed to
865 logarithmic scale. Batch effect among samples was removed using ComBat method in
866 *sva* package in R ⁶¹. Subsequently, expression values were transformed raw counts and
867 differentially expressed genes were identified by DESeq2 package by comparing
868 different conditions using default parameters ⁶¹. Gene ontology analyses were performed
869 using PANTHER/DAVID.

870

871 ChIRP-seq and data analysis (related to Figure 4)

872 To determine the genome-wide localization of *lnc-Nr2f11* we followed protocols
873 previously described (Chu et al., 2011)⁵¹. ChIRP was performed using biotinylated
874 probes designed against mouse *lnc-Nr2f1* using the ChIRP probes designer (Biosearch
875 Technologies). Independent even and odd probe pools were used to ensure lncRNA-
876 specific retrieval (Refer to separate document for odd and even sequences targeting
877 human and mouse *lnc-Nr2F1*, **Supplemental File 2**). Mouse ES-iN samples are
878 crosslinked in 3% formaldehyde. RNase pre-treated samples are served as negative
879 controls for probe-DNA hybridization. ChIRP libraries are constructed using the
880 NEBNext DNA library preparation kit (New England Biolabs). Sequencing libraries were
881 barcoded using TruSeq adapters and sequenced on HiSeq or NextSeq instruments
882 (Illumina). Reads were processed using the ChIRP-seq pipeline⁵⁹. Even-odd ChIRP-seq
883 tracks are merged as previously described ⁵⁹. Peaks were called from the merged tracks
884 over RNase control tracks using MACS14. Overlapping peaks from all replicates were
885 final peaks. High confidence peaks were then filtered by their significance [$-\log_{10}$ (p-
886 value) ≥ 100] and correlation between even/odd probes >0 , average coverage (>2 for
887 mES-iN, >1 for hNPC). For hNPC ChIRP of long and short isoforms, the analysis
888 pipeline and filtering criteria are the same. Sequence motifs were discovered using
889 Homer in 200-bp windows. Peak associated gene sets were obtained through GREAT ⁵²
890 (<http://great.stanford.edu/>). Peaks are assigned to genes according to whether peaks are in
891 gene's regulatory domain. Gene regulatory domain is defined as: Each gene is assigned a
892 basal regulatory domain of a distance 5kb upstream and 1kb downstream of the TSS
893 (regardless of other nearby genes). The gene regulatory domain is extended in both
894 directions to the nearest gene's basal domain but no more than the 1000kb extension in
895 one direction. Gene Ontology of gene sets were performed using Metascape
896 (<http://metascape.org/>). For overlapping mouse ChIRP-seq peak with chromatin features
897 (CTCF, enhancer, H3K27ac, H3K4me3 and PolII) in mouse cortex and E14.5 brain,
898 chromatin annotation files are obtained from public available Chip-seq data from Bing
899 Ren lab ⁵¹. For overlapping human ChIRP-seq peak with chromatin features,
900 ChromHMM model of 25 chromatin states and 12 histone modification marks in neuron
901 cells was used⁶². In addition, peaks are annotated according to distance to genes in figure
902 6 (promoter: -2kb to +1kb of TSS, enhancer: -2kb to -10kb of TSS, exon: exon of a gene,

903 intron: intro of a gene, gene tail: 0 to 2kb downstream of the end of a gene, intergenic:
904 none of the above).

905

906 Compare ChIP-seq and IncNr2f1 ChIRP-seq

907 The fastq files of ChIP-seq data were first aligned to mm9 genome using bowtie2⁶³. Then
908 the reads with alignment score lower than 10 were removed. The aligned sam files were
909 converted to bam files and sorted by Samtools. Picard
910 (<http://broadinstitute.github.io/picard/>) were used to remove duplicates with
911 MarkDuplicates module. After that, Samtools was used to index the bam files. MACS2
912 was used to call peaks with “-f BAM -g mm -B -p 0.005” options⁶⁴. Each ChIP-seq peak
913 was annotated by its closest gene using R package “ChIPSeeker”⁶⁵, and the number of
914 overlap peaks and genes between mouse ChIRP-seq peaks and ChIP-seq peaks were
915 reported. Random peaks were sampled with the same size of mouse ChIRP-seq peaks and
916 annotated by ChIPSeeker package. The number of overlap peaks and genes between the
917 random peaks and ChIP-seq peaks were recorded. The random sampling procedure was
918 conducted 1000 times to construct null distributions at both peak level and gene level,
919 and the empirical p values were then computed respectively.

920

921

922 Biochemistry

923 Single molecule RNA FISH protocol and probes

924 Probes were designed using Stellaris probe designer tool and synthesized by Stellaris.
925 Adherent cells were grown in 12mm coverglass, fixed in 1% formaldehyde for 10
926 minutes at room temperature, washed twice with phosphate buffer saline (PBS), and
927 permeabilized using 70% ethanol at 4C overnight. Fixed cells were subjected to RNase
928 treatment for 30 minutes at 37C with 0.1mg/mL RNase A. After washing (2x SSC, 10%
929 Formamide), hybridization with 250nM probes in hybridization buffer (10% dextran
930 sulfate, 10% formamide, 2x SSC) at 37C overnight in a coverglass protected from light.
931 The next day, washed (2x SSC, 10% Formamide) at 37C for 30 minutes. DAPI staining
932 was added to a clean coverslip and coverglass mounted. Slides were images using
933 confocal microscopy.

934

935 In situ hybridization

936 E13.5 mouse embryos were fixed at 4°C with 4% (weight/volume) paraformaldehyde in
937 PBS overnight. Samples were cryoprotected overnight with 30% (weight/volume)
938 sucrose in PBS, embedded in OCT (Tissue-Tek), and frozen on dry ice. Frozen embryos
939 were sectioned on a cryostat at 16µm. Sections were processed for in situ hybridization.
940 Frozen sections were treated sequentially with 0.3% (volume/volume) Triton-X in PBS
941 and RIPA buffer (150mM NaCl, 50mM Tris-HCl (pH 8.1), 1mM EDTA, 1% NP-40,
942 0.5% Sodium Deoxycholate, 0.1% SDS). Sections were postfixed in 4%
943 paraformaldehyde at room temperature for 15 minutes and washed with PBS.
944 Subsequently, the sections were treated with 0.25% acetic anhydride in 0.1M
945 triethanolamine for 15 minutes and washed with PBS. Sections were incubated in
946 hybridization buffer (50% formamide deionized, 5× SSC, 5× Denhardt's, 500µg/mL
947 Salmon Sperm DNA, 250µg/mL yeast tRNA) containing DIG-labeled probes at 65°C
948 overnight. Hybridized sections were washed two times in washing solution (2× SSC,

949 50% formamide, 0.1% Tween 20) at 65°C for 60 minutes. After washing, sections were
950 incubated for 1 hour in 1% (weight/volume) blocking reagent (0.1M Maleic Acid, 0.15M
951 NaCl (pH 7.5), 0.1% Tween 20, Roche). Subsequently, incubated with an alkaline
952 phosphatase (AP)-coupled antibody (Roche) at 4°C overnight. After rinsing, the signals
953 were visualized with nitro-blue tetrazolium chloride (NBT)/5-Bromo-4-Chloro-3'-In-
954 dolyolphosphatase p-Toluidine salt (BCIP) (Sigma). The DIG-labeled antisense RNA
955 probe for detecting mouse NR2F1 corresponds to the CDS region and for *lnc-Nr2f1*
956 corresponds to 800bp region. The DIG-labeled sense RNA probe for both, NR2F1 and
957 *lnc-Nr2f1*, corresponds to the same region as the antisense probe in the reverse direction.
958 Probes were generated by in vitro transcription with T7 RNA polymerase (Roche) using
959 the DNA templates containing a promoter sequence of T7 RNA polymerase promoter
960 (TAA TAC GAC TCA CTA TAG GG) followed by a complimentary sequence of target
961 RNA. DNA templates were amplified by PCR with the following primers: For *lnc-Nr2f1*
962 probe: (F) GTG GCC ATG GAA TGG TGT AGC AGA, and (R) GTC TGA GTG TTT
963 GTT TGA CTG AAT GT; NR2F1 probe: (F) CGG TTC AGC GAG GAA GAA TGC
964 CT, and (R) CTA GGA ACA CTG GAT GGA CAT GTA AG.

965 Cellular fractionation

966 Cell fractionation of primary neocortical cells (prepared from E12.5 mouse cerebral
967 cortex) into cytoplasmic and nuclear RNA fractions was performed with a
968 nuclear/cytoplasm fractionation kit (PARIS kit, Ambion) following the instructions of the
969 manufacturer. The chromatin/cytoplasm fractionation was performed following the
970 published protocol⁶⁶. The amount of RNA in each fraction was determined by qRT-PCR
971 in a Roche LightCycler with Brilliant III Ultra-Fast SYBR Green QRT-PCR Master Mix
972 (Agilent). For primer sequences refer to separate document.

973 Immunofluorescence

974 Cells were fixed with 4% paraformaldehyde for 15 minutes and subsequently lysed and
975 blocked with blocking buffer [PBS + 0.1% Triton X (Sigma Aldrich) + 5% Cosmic calf
976 serum (Thermo Scientific)] for 30 minutes. Primary antibodies diluted with blocking
977 buffer were added to the wells and left for an hour. The following antibodies were used
978 for immunostaining: mouse anti-MAP2 (Sigma, 1:500), rabbit anti-Tuj1 (Covance,
979 1:1000), goat anti-Sox1 (R&D, 1:100) and mouse anti-hNestin (R&D, 1:1000). The wells
980 were subsequently washed three times with the blocking buffer. Secondary antibodies
981 conjugated with Alexa dyes (1:1000, Invitrogen) diluted with the blocking buffer were
982 added to the wells and left for an hour. The wells were again washed three times with the
983 blocking buffer. 4',6-Diamidino-2-phenylindole (DAPI) (Life Technologies, 1: 10,000)
984 diluted in PBS was added for 1 minute for nuclear staining.

985 Western blotting

986 Cells were lysed with 1 volume of RIPA buffer with cOmplete protease inhibitors
987 (Sigma-Aldrich) and equivolume of 2x Laemmli buffer was added. The samples were
988 then boiled for 5 minutes at 95°C and subsequently separated in 4-12% Bis-Tris gel with
989 MES buffer (Invitrogen) and transferred onto PVDF membrane for 2 hours at 4°C. Blots
990 were then blocked in blocking buffer (PBS + 0.1% Tween-20 (Sigma-Aldrich) + 5% fat-
991 free milk) for 30 minutes and subsequently incubated overnight with primary antibodies
992 at 4°C. The primary antibodies used are rabbit anti-HSP90 (Cell Signalling) and rabbit
993

994 anti-NR2F1 (Cell Signalling). The blots were washed three times in PBS + 0.1% Tween-
 995 20 for 10 minutes each. Next, the blots were incubated with secondary antibodies
 996 conjugated to horseradish peroxidase (Jackson immunoresearch) were diluted in blocking
 997 buffer for 1 hour. The blots were washed three times in PBS + 0.1% Tween-20 for 10
 998 minutes each and once with PBS before adding chemiluminescence substrates (Perkin
 999 Elmer) for signal detection on films.

1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007

1008 **Figure legends**

1009

1010 **Fig. 1. lncRNA loci are recurrently mutated in patients with neurodevelopmental** 1011 **disorders**

1012

(A) Schematic representation of CNV morbidity map analysis for candidate lncRNAs
 1013 and all other in lncRNAs loci. The 35 mouse lncRNA candidates (28 human loci)
 1014 is from Figure 1 – figure supplement 1H.

1015

(B) Top: Representative tracks for lncRNA E locus, also known as *lnc-NR2F1*.
 1016 Depicted in blue are deletions and in red duplications. Arrow points to patient
 1017 with focal deletion affecting the *lnc-NR2F1* locus only. Bottom: Custom CGH
 1018 arrays used to validate chromosomal aberration in patient 9900850 harboring
 1019 focal deletion represented in green signal.

1020

(C) Genetic pedigree analysis for family with paternally inherited balanced
 1021 chromosomal translocation (5;12)(q15;q15), including a summary of clinical
 1022 features for patient CMS12200 and father. The mother has a normal karyotype.
 1023 Listed in the box are the symptoms of the patients.

1024

(D) Top: Circa plot representing the pathogenic chromosomal event for patient
 1025 CMS12200 involving chromosomes 5 and 12. Bottom: Representative
 1026 chromosome ideogram and track of the balanced chromosomal break affecting
 1027 patient CMS12200. Below the ideoplots is the schematic representation of
 1028 predominant human isoforms for *lnc-NR2F1* and the site of the break site
 1029 disrupting the long isoforms.

1030

(E) The locations of the probes are in figure 1 – figure supplement 4C. Left:
 1031 Metaphase spread from patient CMS12200 with the t(5;12) translocation showing
 1032 FISH signals obtained with the clone RP11-608G16 (green) spanning
 1033 Chromosome 5 breakpoint, and a Chromosome 5 telomere-specific probe (red).
 1034 Middle: Metaphase spread from patient CMS12200 with the t(5;12) translocation
 1035 showing FISH signals obtained with the clone RP11-597C7 (green) proximal to
 1036 Chromosome 12 breakpoint, and a Chromosome 12 centromere-specific probe
 1037 (red). Right: Metaphase spread from patient CMS12200 with the t(5;12)
 1038 translocation showing FISH signals obtained with the clone RP11-641O3 (green)
 1039 distal to Chromosome 12 breakpoint, and a Chromosome 12 centromere-specific
 1040 probe (red).

1041

1042

Fig. 2. Molecular characterization of mouse *lnc-Nr2f1*

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

(A) Schematic showing the different isoforms reported by Refseq of the human *lnc-NR2F1* and mouse *lnc-Nr2f1*. Exons highlighted in red are conserved among human and mouse. The table at the bottom right corner shows the sequence similarity as reported by T-COFFEE. The sequence alignment is available as data S7.

(B) *lnc-Nr2f1* expression measured by qRT-PCR across stages of mouse brain development and early stages of iN cell reprogramming. Results show early detection in E13.5 brain, peak expression at postnatal stages, and continued expression through adulthood.

(C) In situ hybridization for *Nr2f1* and *lnc-Nr2f1* shows similar expression pattern in E13.5 mouse brain. Highlighted by arrows are neocortex (NCX) and ganglionic eminences (GE) with high expression levels.

(D) Cellular localization of *lnc-Nr2f1* by single molecule RNA-FISH in MEFs ectopically expressing the lncRNA 48 hours after dox induction reveals nuclear and cytoplasmic localization, with slight nuclear enrichment. Green arrow points at *lnc-Nr2f1* in the nucleus and red arrows point at the uninfected nuclei.

(E) Cellular fractionation of primary neurons derived from E13.5 caudal cortex dissection shows nuclear and cytoplasmic localization of *lnc-Nr2f1*.

(F) Chromatin enrichment of *lnc-Nr2f1* by histone H3 RIP-qRT-PCR in brain derived neuronal precursor cells (NPCs), postnatal and adult mouse brain.

Fig. 3. Mouse *lnc-Nr2f1* KO reveals *lnc-Nr2f1* regulates neuronal genes

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1081

1082

1083

1084

1085

(A) Schematic showing the experimental strategy for *lnc-Nr2f1* overexpression. In control mouse ES cells, an inverted construct with a splice acceptor (marked in yellow) and a polyadenylation signal (marked in red) are added after the first exon of the *lnc-Nr2f1*. The mouse ES were infected with rtTA and Ngn2-T2A-puro and mES derived induced neurons (mES-iN) were assayed after 3 or 4 days after dox induction.

(B) Graph showing that overexpression of *lnc-Nr2f1* increased the neurite length in day 3 Ngn2 mouse ES derived iN cells relative to the Ctrl. The same effect was not seen with *Nr2f1* overexpression. For each replicate, the individual neurite length for all neurons in each of the five 20x field was manually traced in Fiji. The sequence used for mouse *lnc-Nr2f1* overexpression is available in the supplementary document (n=3, Student t-test, Two-tailed, * indicates $p=0.048<0.05$). Error bars show s.e.m.

(C) Graph showing that overexpression of *Nr2f1* decreased the neurite number in day 3 Ngn2 mouse ES-iN cells relative to the Ctrl. The same effect was not seen with *lnc-Nr2f1* overexpression. (n=3, 10 field per replicate, Student t-test, Two-tailed, ** indicates $p=0.0022<0.01$). Error bars show s.e.m.

(D) β -III-tubulin staining of the day 3 Ngn2 mouse ES derived iN cells for Ctrl, *lnc-Nr2f1* overexpression and *Nr2f1* overexpression. Scale bar = 50 μ m. Red arrow pointed at immature induced neuronal cells with short projection. Green arrow

- 1086 pointed at mature induced neuronal cells with longer projection. Note that the *lnc-Nr2f1*
 1087 *lnc-Nr2f1* overexpression condition have more mature induced neuronal cells.
- 1088 (E) Hierarchical clustering heatmap of day 4 Ngn2 ES-iN cells between control and
 1089 *lnc-Nr2f1* overexpression (OE). There are 1912 genes differentially expressed
 1090 (n=2, FDR corrected p<0.10, Fold change>1.5 fold). Listed to the right are genes
 1091 which are upregulated upon *lnc-Nr2f1* overexpression.
- 1092 (F) Schematic showing the knocking out strategy for *lnc-Nr2f1*. The *lnc-Nr2f1*
 1093 knockout mouse ES cells are generated after Cre recombinase introduction to the
 1094 Ctrl line in Fig.4A. The mouse ES were infected with rtTA and Ngn2-T2A-puro
 1095 and mES derived induced neurons were assayed after 3 or 4 days after dox
 1096 induction.
- 1097 (G) Hierarchical clustering heatmap of day 4 Ngn2 ES-iN cells between wild type,
 1098 *lnc-Nr2f1* knockout (KO) and *lnc-Nr2f1* knockout with *lnc-Nr2f1* overexpression
 1099 (OE). There are 348 genes differentially expressed and can be subsequently
 1100 rescued with *lnc-Nr2f1* overexpression (n=2, FDR corrected p<0.10). Listed to
 1101 the right are genes which are upregulated upon *lnc-Nr2f1* KO.
- 1102 (H) Gene ontology of the upregulated genes in *lnc-Nr2f1* knockout day 4 Ngn2 mouse
 1103 ES- iN cells as compared to the Ctrl.
- 1104 (I) Gene ontology of the downregulated genes in *lnc-Nr2f1* knockout day 4 Ngn2
 1105 mouse ES- iN cells as compared to the Ctrl.

1106
 1107 **Fig. 4. *lnc-Nr2f1* binds to distinct genomic loci regulating neuronal genes**

- 1108 (A) Schematic showing the location of ChIRP probe for mouse *lnc-Nr2f1* (highlighted
 1109 in red). Yellow lines represent the conserved exons between mouse and human
 1110 *lnc-Nr2f1*.
- 1111 (B) Heatmaps representing genome-wide occupancy profile for mouse *lnc-Nr2f1* in
 1112 day 4 Ngn2 mouse ES- iN cells and the RNase control obtained by ChIRP. There
 1113 are 14975 significant peaks called with respect to the RNase treated control. E and
 1114 O represents even and odd probes respectively.
- 1115 (C) UCSC browser track showing the binding site within the intronic region of *Nrp2*.
 1116 The “R” represents the RNase treated control.
- 1117 (D) Gene ontology terms associated with genes adjacent to the high confident mES-iN
 1118 ChIRP-seq peaks. Terms highlighted in red are terms related to nervous system
 1119 development.
- 1120 (E) Percentage of mES-iN ChIRP-seq peaks which overlap with CTCF, enhancer,
 1121 H3K27ac, H3K4Me3 and PolII defined in mouse cortex relative to the control.
 1122 (***) represents p<0.0001, ** represents p<0.01, Chi-square test)
- 1123 (F) Proposal mechanism of *lnc-Nr2f1* action. *lnc-Nr2f1* binds to the genomic region
 1124 enriched with bHLH motif and regulates the downstream neuronal genes.

1125
 1126 **Fig. 5. Human *lnc-NR2F1* shows isoform-specific chromatin binding**

- 1127 (A) Schematic showing the location of the ChIRP probes target the short isoform-
 1128 specific exon (exon 7) and long isoform-specific exon (exon 11). The red line
 1129 denotes the break point for the patient.

- 1130 (B) Overexpression of human *lnc-NR2F1* isoforms in combination with *Ascl1* relative
 1131 to MEFs expressing *Ascl1* alone. The graph quantifies the proportion of TauGFP
 1132 positive cells with projections normalized to number of TauGFP cells. TauGFP
 1133 cells with projections longer than three times the diameter of the cell body were
 1134 counted and normalized to the total number of TauGFP positive cells. The
 1135 sequences for human *lnc-Nr2f1* isoforms is available in the supplementary
 1136 documents (n=3, Student t-test, two tailed, scale bar= * represents $p < 0.05$, **
 1137 represents $p < 0.01$). Error bars show s.e.m.
- 1138 (C) Heatmap representing genome-wide occupancy profile for domain ChIRP
 1139 performed using probes specific to the long and short isoform-specific exon of
 1140 *lnc-NR2F1* in human ES derived neural progenitor cells (NPC). There are 4404
 1141 and 415 significant peaks called relative to the RNase control for the long and
 1142 short isoform respectively. E, O and M represents even, odd and merge track
 1143 respectively.
- 1144 (D) UCSC browser track showing the site within the promoter region of LOC90784
 1145 bound by the long isoform-specific exon (exon 11) but not the short isoform-
 1146 specific exon (exon 7)
- 1147 (E) Bar graph showing the distribution of the 913 high confident long isoform-
 1148 specific peaks. The long isoform-specific peaks are enriched in the introns, exons,
 1149 promoters and enhancers but depleted in the intergenic regions.
- 1150 (F) Gene ontology terms associated with genes adjacent to the human ES derived
 1151 NPC ChIRP-seq high confident peaks. Terms highlighted in red are terms related
 1152 to nervous system development.
- 1153 (G) Venn diagram representing the peak associated gene overlap between the domain
 1154 ChIRP of the long isoform-specific exon (exon 11) from human ES derived NPC
 1155 and mouse mES-iN ChIRP. ($p < 0.0001$, $\chi^2 = 239.921$, DF=1, Chi square test)
- 1156 (H) Venn diagram representing overlap between genes involved in the autism risk and
 1157 genes identified by the domain ChIRP of the long isoform-specific exon (exon 11)
 1158 from human ES derived NPC. ($p < 0.0001$, $\chi^2 = 71.670$, DF=1, Chi square test).

1162 **Figure 1 – figure supplement 1. Molecular profiling of direct fibroblast to iN cell**
 1163 **reprogramming nominates functional lncRNAs involved in neurogenesis**

- 1164 (A) Schematic representation of experimental time points generated for this study.
- 1165 (B) Classification of the iN cell transcriptome consisting of 51,470 transcripts based
 1166 on coding genes and non-coding RNAs.
- 1167 (C) Diagram depicting pipeline derived in this study to enrich for candidate lncRNAs
 1168 with strong neuronal association.
- 1169 (D) Hierarchical clustering heatmap of lncRNA expression during iN cell
 1170 reprogramming by RNA-seq across indicated time points (n=2 biological
 1171 replicates). Shown are 287 lncRNAs that changed expression at least two-fold at
 1172 any time point ($p < 0.05$). Fold change is represented in logarithmic scale
 1173 normalized to the mean expression value of a gene across all samples. The green
 1174 box highlights the genes which are upregulated in the MEF + BAM 48 hours. The
 1175 same set of genes are upregulated in embryonic mouse brain (see Figure 1 –

- 1176 figure supplement 1E). The pink box highlights the genes which are upregulated
 1177 in the MEF + BAM 22 hours. The same set of genes are upregulated in adult
 1178 mouse brain (see Figure 1 – figure supplement 1E).
- 1179 (E) Hierarchical clustering heatmap of iN lncRNAs from Figure 1 – figure
 1180 supplement 1D across mouse brain tissues from publicly available data.
 1181 Expression levels are represented in logarithmic scale normalized to the mean
 1182 expression value of a gene across all samples.
- 1183 (F) Chromatin association of iN lncRNAs determined by histone H3-RIP-seq in
 1184 Neuronal Progenitor Cells (NPCs) and total adult brain (n=2 biological replicates
 1185 for NPCs and n=3 for mouse brain) presented in binary format. Shown are
 1186 lncRNAs that have significant enrichment over background (>2-fold, p<0.05) and
 1187 consistent enriched in the chromatin among biological replicates.
- 1188 (G) Co-expression analysis using Genomica between iN lncRNAs and mRNAs
 1189 associated with Gene ontology (GO) terms. Highlighted in the yellow box are
 1190 lncRNA candidates which are associated with neuronal GO terms.
- 1191 (H) RNA-seq heatmap of 35 filtered candidate lncRNAs across MEF-to-iN cell
 1192 reprogramming. In brackets are the Refseq ID for the annotated lncRNAs. See the
 1193 supplementary documents for the sequences and coordinates of the 35 candidates.
 1194

1195 **Figure 1 – figure supplement 2. QRT-PCR validation of candidate lncRNAs**
 1196 **expression**

- 1197 (A) Expression detection of candidate lncRNAs by qRT-PCR across early stages of
 1198 iN cell reprogramming and mouse brain development.
 1199

1200 **Figure 1 – figure supplement 3. Conserved lncRNAs have distinct pattern of**
 1201 **expression across different stages of the human brain**

- 1202 (A) Heatmap of annotated syntenic lncRNAs across human brain development
 1203 detected by RNA-seq. Expression is represented in logarithmic scale normalized
 1204 to the mean RPKM value of a gene across all samples. See Supplementary
 1205 documents for the lncRNA sequences.
 1206

1207 **Figure 1 – figure supplement 4. Other reports of CNVs affecting *lnc-NR2F1* and an**
 1208 **example of focal deletion affecting *lnc-ZFP238* and characterization fo patient**
 1209 **CMS12200,**

- 1210 (A) Top: Representative tracks for lncRNA H locus, also known as *lnc-ZFP238*.
 1211 Depicted in blue are deletions and in red duplications. Arrow points to two
 1212 significant focal deletions in two distinct patients. Bottom: Custom CGH arrays
 1213 used to validate chromosomal aberration in patients 9909124 and 9900584
 1214 harboring focal deletions represented in green signal.
- 1215 (B) Representative tracks of previously reported deletions affecting chromosome
 1216 region 5q15 in patients with neurodevelopmental and neuropsychiatric disorders.
 1217 The black rectangle represents affected genomic region corresponding to the
 1218 patient described in a previous publication to the left. Top panel: light blue box
 1219 represents minimal common region amongst three patients in Cardoso et al, 2009
 1220 report. There are multiple genes affected in the locus. Middle panel: purple box
 1221 depicts common deleted region reported amongst two patients from Brown, et al

1222 2009, and Al-Kateb, et al 2013 studies. The region encompasses only two genes:
 1223 *NR2F1* and *lnc-NR2F1*. Bottom panel: Patient CMS 12200 described in this study
 1224 harboring a balanced chromosomal translocation disrupting *lnc-NR2F1* only as
 1225 depicted by the pink box. In red, *lnc-NR2F1* is highlighted.

1226 (C) Schematic representation of breakpoint region between chromosome 5 and 12 for
 1227 patient CMS12200. Illustrated is also the probe design to confirm event. Probe
 1228 RP11-608G16 spans the breakpoint in chromosome 5. Probe RP11-597C7 is
 1229 proximal to the breakpoint on chromosome 12. Probe RP11-64103 is distal to the
 1230 breakpoint on chromosome 12.

1231 (D) Genomic PCR to confirm t(5;12) translocation using primers spanning control
 1232 region in chromosome 5 (left), translocation between chromosome 5 and 12
 1233 (middle left), unaffected region in ARGH gene (middle right), and non-
 1234 pathological gene duplication for TRIO exon in ARGH intronic loci (right).
 1235 Samples from control (GM12878), mother (does not harbor translocation), father
 1236 (t(5;12)), and CMS12200 patient (t(5;12)).

1237 (E) Expression of genes proximal to breakpoint is unaffected as measured by RT-
 1238 qPCR in CMS12200 patient lymphocytes and three distinct control samples.
 1239

1240 **Figure 2 – figure supplement 1. Synteny, sequence and microdomain conservation of**
 1241 ***lnc-NR2F1*.**

1242 (A) UCSC browser track for human (top) and mouse (bottom) showing synteny
 1243 conservation around the *lnc-Nr2f1 – Nr2f1* locus. The same gene between the two
 1244 species is color-coded

1245 (B) UCSC browser track showing the sequence conservation of the three exons
 1246 (highlighted in red) across different species.

1247 (C) Sequences around the conserved exons show short sequence homology from
 1248 different species. MEME (<http://meme-suite.org/tools/meme>) is used to discover
 1249 motif homology. See Supplementary documents for the sequences used.
 1250

1251 **Figure 2 – figure supplement 2. Characterization of *lnc-Nr2f1* localization.**

1252 (A) RPKM counts for *lnc-Nr2f1* during MEF-to-iN-reprogramming and across
 1253 different stages and tissues of mouse brain development.

1254 (B) Control in situ hybridization for NR2F1 and *lnc-Nr2f1* using sense probe.

1255 (C) Control single molecule RNA-FISH in MEFs infected with rtTA alone in order to
 1256 determine background signal.
 1257

1258 **Figure 3 – figure supplement 1. Characterization of the roles of *lnc-Nr2f1* during iN**
 1259 **reprogramming.**

1260 (A) Percentage of TauGFP positive cells with projections normalized to number of
 1261 TauGFP cells. TauGFP cells with projections longer than three times the diameter
 1262 of the cell body were counted and normalized to the total number of TauGFP
 1263 positive cells. The sequence for mouse *lnc-Nr2f1* is available in the
 1264 supplementary documents (n=4, Student t-test, Two tailed, Error bars show
 1265 s.e.m).

1266 (B) Immunofluorescence staining depicting how a TauGFP positive cell with
 1267 projections (example highlighted in light green) is differentiated from a TauGFP

- 1268 positive partially reprogrammed iN cell (example highlighted in dark green) at 7
 1269 days. Scale bar = 50 μ m.
- 1270 (C) Hierarchical clustering heatmap of differentially expressed genes detected by
 1271 RNA-seq in MEFs expressing *Ascl1* alone compared to *Ascl1* and *lnc-Nr2f1* after
 1272 7 days (n=2 biological replicates, FDR corrected p<0.001). Shown are 343 genes.
 1273 311 genes are upregulated and 32 genes are downregulated. Fold change is
 1274 represented in logarithmic scale normalized to the mean expression value of a
 1275 gene across all samples. Representative gene names are included. Those with (*)
 1276 have been linked to neurological disorders curated by Basu et al.
- 1277 (D) Gene ontology of the upregulated genes upon *lnc-Nr2f1* overexpression in *Ascl1*
 1278 MEF-iN 7 days. Highlighted in red are neuronal GO terms.
- 1279 (E) qRT-PCR validation downstream neuronal genes of the RNA-sequencing results
 1280 in Figure 1 – figure supplement 3C. Ectopic expression of *lnc-Nr2f1* led to
 1281 upregulation of several downstream neuronal genes. (n=3, * indicates p<0.05).
 1282 Error bars show s.e.m
- 1283 (F) qRT-PCR validation of several target genes in Fig. 3G that go down when *lnc-*
 1284 *Nr2f1* is knocked out which can be subsequently rescued with *lnc-Nr2f1*. (n=4, *
 1285 indicates p<0.05). Error bars show s.e.m.
- 1286 (G) Venn diagram representing the overlap between autism related genes and rescued
 1287 genes from figure 4g. (p=0.0015, $\chi^2=10.097$, DF=1, Chi square test)
- 1288 (H) Non-denaturing TAE agarose gel showing the different species of RNA in
 1289 different fractions. Note the presence of 45S in the chromatin fraction and tRNA
 1290 and 5S in the cytoplasmic fraction.
- 1291 (I) RT-PCR performed on day4 mES iN using different primers targeting GAPDH
 1292 (positive control), Xist (nuclear control), two independent exon-exon/exon-intron
 1293 primers. WT represent *lnc-Nr2f1*. WT and KO represent *lnc-Nr2f1* wild-type and
 1294 knock-out mES-iN respectively.

1295
 1296 **Figure 3 – figure supplement 2. Characterization of the epistasis relationship**
 1297 **between mouse *Nr2f1* and *lnc-Nr2f1*.**
 1298

- 1299 (A) CRISPR knock out strategy to generate *Nr2f1* knockout (Homo) and heterozygous
 1300 null lines (Het) from the control mES cells (Ctrl).
- 1301 (B) qRT-PCR results for *lnc-Nr2f1* and *Nr2f1* in the Ctrl, *Nr2f1* heterozygous null and
 1302 *Nr2f1* knock out day 4 Ngn2 mES-iN. (n=4 for Ctrl and Homo, n=6 for Het; n.s.
 1303 denotes not significant by two tailed t test)
- 1304 (C) Western blot showing the level of NR2F1 for individual clones of Ctrl, Het and
 1305 Homo for day 4 Ngn2 mES-iN.
- 1306 (D) Neurite length measurement of the Ngn2 day 3 mES iN cells generated from the
 1307 *Nr2f1* Ctrl, Het or Homo lines. (n=4 for Ctrl and Homo, n=6 for Het) (n.s.
 1308 indicates p<0.05).
- 1309 (E) Number of neurons per 20x the Ngn2 day 3 mES iN cells generated from the
 1310 *Nr2f1* Ctrl, Het or Homo lines. (n=4 for Ctrl and Homo, n=6 for Het) (n.s.
 1311 indicates p<0.05).
 1312

1313 **Figure 4 – figure supplement 1. Identification of *lnc-NR2F1* role in transcriptional**
 1314 **regulation.**

- 1315 (A) Immunofluorescence staining for the day 4 mouse embryonic stem cell derived
 1316 induced neurons (mES-iN) used in the *lnc-Nr2f1* ChIRP. (Green= β -III-tubulin,
 1317 Blue=DAPI) (Scale bar=50 μ m)
 1318 (B) *Lnc-Nr2f1* RNA pull down efficiency for both even and odd probes
 1319 (C) The enriched motifs with their corresponding p-value and the percentage of peaks
 1320 with the given motifs.
 1321 (D) Percentage of mES-iN ChIRP-seq peaks which overlap with CTCF, enhancer,
 1322 H3K27ac, H3K4Me3 and PolIII defined in mouse E14.5 brain relative to the
 1323 control. (***) represents $p < 0.0001$, * represents $p < 0.05$, Chi-square test)
 1324 (E) Venn diagram showing the overlap between 1912 genes from Figure 4E and the
 1325 1534 genes adjacent to the 1092 high confident mES-iN ChIRP peaks. ($p < 0.0001$,
 1326 $\chi^2 = 21.983$, DF=1, Chi square test)
 1327 (F) Venn diagram and statistical analysis of the overlapped peaks between *lnc-Nr2f1*
 1328 ChIRP-seq peaks and *Ngn2* (top) or *Ascl1* (bottom) ChIP-seq peaks. Right panels
 1329 are the null distribution of overlapped peaks constructed by 1000 random
 1330 sampling. The number of overlap peaks in the observed data is marked in solid
 1331 red line.
 1332 (G) Venn diagram and statistical analysis of the overlapped peak-associated-genes
 1333 between *lnc-Nr2f1* ChIRP-seq and *Ngn2* (top) or *Ascl1* (bottom) ChIP-seq. Right
 1334 panels are the null distribution of overlapped genes constructed by 1000 random
 1335 sampling. The number of overlap genes in the observed data is marked in solid
 1336 red line.
 1337

1338 **Figure 5 – figure supplement 1. Identification of *lnc-NR2F1* role in transcriptional**
 1339 **regulation**

- 1340 (A) Immunofluorescence staining for day 12 human neural progenitor cells (hNPC)
 1341 differentiated from HUES9 cells using dual SMAD protocol. (Green=human
 1342 NESTIN, Red=human SOX1, Blue=DAPI) (Scale bar=50 μ m)
 1343 (B) qRT-PCT using primers specific to the common, long or short exon in day 12
 1344 hNPC. (n=3).
 1345 (C) *Lnc-NR2F1* RNA pull down efficiency for the long and short isoform-specific
 1346 exons of for both even and odd probes. The pull down is specific to *lnc-NR2F1*
 1347 since there is very little *NR2F1* or *GAPDH*.
 1348 (D) ChromHMM model ran on the peaks from the high confident long isoform-
 1349 specific exon (exon 11) ChIRP peaks performed in hNPC. The peaks are enriched
 1350 in the promoters and TSSs, active enhancers and quiescent chromatin regions.
 1351 (E) The enriched transcription factor motifs for the long isoform-specific domain
 1352 ChIRP experiments in hNPC.
 1353 (F) Heatmap showing the gene expression changes in a human neuroblastoma cell
 1354 line (SK-N-SH) upon overexpression of human *lnc-NR2F1* isoform 1, 2 and 3
 1355 respectively normalized to the control (>2-fold, $p < 0.05$, FDR<0.05). The box
 1356 below the heatmap shows the number of genes significantly upregulated or
 1357 downregulated.
 1358

1359
1360 **Supplemental File 1. Diagnostic comparison between studies of patients with**
1361 **affected *lnc-NR2F1* locus. Related to Fig. 2**

1362 (A) Summary of diagnosis for previously reported patients, including patient
1363 CMS12200 described in this study. Highlighted in grey are the shared diagnostic
1364 features across patients. Adapted figure³⁹.

1365

1366 **Supplemental File 2 CHIRP sequencing probes used in the study**

1367 **Supplemental File 3 Public datasets used in the study**

1368 **Supplemental File 4 qRT-PCR primers used in the study**

1369 **Supplemental File 5 RNA FISH primers used in the study**

1370 **Supplemental File 6 Sequence conservation used in the study**

1371 **Supplemental File 7 A list of human lncRNAs reported in the study**

1372 **Supplemental File 8 A list of mouse lncRNAs reported used in the study**

1373

1374

1375

1376

1377

1378

1379

1380

1381 **References**

- 1382 1. Flynn, R.A. & Chang, H.Y. Long noncoding RNAs in cell-fate programming and
1383 reprogramming. *Cell Stem Cell* **14**, 752-61 (2014).
- 1384 2. Wapinski, O. & Chang, H.Y. Long noncoding RNAs and human disease. *Trends*
1385 *Cell Biol* **21**, 354-61 (2011).
- 1386 3. Fertuzinhos, S. *et al.* Laminar and temporal expression dynamics of coding and
1387 noncoding RNAs in the mouse neocortex. *Cell Rep* **6**, 938-50 (2014).
- 1388 4. Valadkhan, S. & Nilsen, T.W. Reprogramming of the non-coding transcriptome
1389 during brain development. *J Biol* **9**, 5 (2010).
- 1390 5. Lv, J. *et al.* Long non-coding RNA identification over mouse brain development
1391 by integrative modeling of chromatin and genomic features. *Nucleic Acids Res*
1392 **41**, 10044-61 (2013).
- 1393 6. Aprea, J. *et al.* Transcriptome sequencing during mouse brain development
1394 identifies long non-coding RNAs functionally involved in neurogenic commitment.
1395 *EMBO J* **32**, 3145-60 (2013).
- 1396 7. Ramos, A.D. *et al.* Integration of genome-wide approaches identifies lncRNAs of
1397 adult neural stem cells and their progeny in vivo. *Cell Stem Cell* **12**, 616-28
1398 (2013).

- 1399 8. Ramos, A.D. *et al.* The long noncoding RNA pnky regulates neuronal
1400 differentiation of embryonic and postnatal neural stem cells. *Cell Stem Cell* **16**,
1401 439-47 (2015).
- 1402 9. Ng, S.Y., Bogu, G.K., Soh, B.S. & Stanton, L.W. The long noncoding RNA RMST
1403 interacts with SOX2 to regulate neurogenesis. *Mol Cell* **51**, 349-59 (2013).
- 1404 10. Vierbuchen, T. *et al.* Direct conversion of fibroblasts to functional neurons by
1405 defined factors. *Nature* **463**, 1035-41 (2010).
- 1406 11. Wapinski, O.L. *et al.* Hierarchical mechanisms for direct reprogramming of
1407 fibroblasts to neurons. *Cell* **155**, 621-35 (2013).
- 1408 12. Ang, C.E. & Wernig, M. Induced neuronal reprogramming. *J Comp Neurol* **522**,
1409 2877-86 (2014).
- 1410 13. Voineagu, I. *et al.* Transcriptomic analysis of autistic brain reveals convergent
1411 molecular pathology. *Nature* **474**, 380-4 (2011).
- 1412 14. Iossifov, I. *et al.* The contribution of de novo coding mutations to autism spectrum
1413 disorder. *Nature* **515**, 216-21 (2014).
- 1414 15. Ronemus, M., Iossifov, I., Levy, D. & Wigler, M. The role of de novo mutations in
1415 the genetics of autism spectrum disorders. *Nat Rev Genet* **15**, 133-41 (2014).
- 1416 16. Gilman, S.R. *et al.* Rare de novo variants associated with autism implicate a
1417 large functional network of genes involved in formation and function of synapses.
1418 *Neuron* **70**, 898-907 (2011).
- 1419 17. Iossifov, I. *et al.* De novo gene disruptions in children on the autistic spectrum.
1420 *Neuron* **74**, 285-99 (2012).
- 1421 18. De Rubeis, S. *et al.* Synaptic, transcriptional and chromatin genes disrupted in
1422 autism. *Nature* **515**, 209-15 (2014).
- 1423 19. O'Roak, B.J. *et al.* Multiplex targeted sequencing identifies recurrently mutated
1424 genes in autism spectrum disorders. *Science* **338**, 1619-22 (2012).
- 1425 20. O'Roak, B.J. *et al.* Sporadic autism exomes reveal a highly interconnected
1426 protein network of de novo mutations. *Nature* **485**, 246-50 (2012).
- 1427 21. Hormozdiari, F., Penn, O., Borenstein, E. & Eichler, E.E. The discovery of
1428 integrated gene networks for autism and related disorders. *Genome Res* **25**, 142-
1429 54 (2015).
- 1430 22. Meng, L. *et al.* Towards a therapy for Angelman syndrome by targeting a long
1431 non-coding RNA. *Nature* **518**, 409-12 (2015).
- 1432 23. Cheetham, S.W., Gruhl, F., Mattick, J.S. & Dinger, M.E. Long noncoding RNAs
1433 and the genetics of cancer. *Br J Cancer* **108**, 2419-25 (2013).

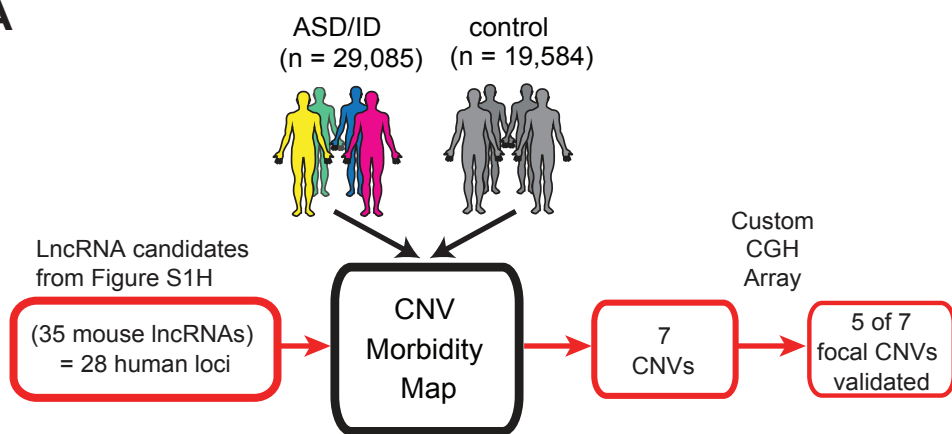
- 1434 24. Gupta, R.A. *et al.* Long non-coding RNA HOTAIR reprograms chromatin state to
1435 promote cancer metastasis. *Nature* **464**, 1071-6 (2010).
- 1436 25. Coe, B.P. *et al.* Refining analyses of copy number variation identifies specific
1437 genes associated with developmental delay. *Nat Genet* **46**, 1063-71 (2014).
- 1438 26. Cooper, G.M. *et al.* A copy number variation morbidity map of developmental
1439 delay. *Nat Genet* **43**, 838-46 (2011).
- 1440 27. Turner, T.N. *et al.* Genomic Patterns of De Novo Mutation in Simplex Autism.
1441 *Cell* **171**, 710-722 e12 (2017).
- 1442 28. Xiang, C. *et al.* RP58/ZNF238 directly modulates proneurogenic gene levels and
1443 is required for neuronal differentiation and brain expansion. *Cell Death Differ* **19**,
1444 692-702 (2012).
- 1445 29. Ohtaka-Maruyama, C. *et al.* RP58 regulates the multipolar-bipolar transition of
1446 newborn neurons in the developing cerebral cortex. *Cell Rep* **3**, 458-71 (2013).
- 1447 30. Baubet, V. *et al.* Rp58 is essential for the growth and patterning of the
1448 cerebellum and for glutamatergic and GABAergic neuron development.
1449 *Development* **139**, 1903-9 (2012).
- 1450 31. Armentano, M., Filosa, A., Andolfi, G. & Studer, M. COUP-TFI is required for the
1451 formation of commissural projections in the forebrain by regulating axonal
1452 growth. *Development* **133**, 4151-62 (2006).
- 1453 32. Borello, U. *et al.* Sp8 and COUP-TF1 reciprocally regulate patterning and Fgf
1454 signaling in cortical progenitors. *Cereb Cortex* **24**, 1409-21 (2014).
- 1455 33. Faedo, A. *et al.* COUP-TFI coordinates cortical patterning, neurogenesis, and
1456 laminar fate and modulates MAPK/ERK, AKT, and beta-catenin signaling. *Cereb*
1457 *Cortex* **18**, 2117-31 (2008).
- 1458 34. Harrison-Uy, S.J., Siegenthaler, J.A., Faedo, A., Rubenstein, J.L. & Pleasure,
1459 S.J. CoupTFI interacts with retinoic acid signaling during cortical development.
1460 *PLoS One* **8**, e58219 (2013).
- 1461 35. Lin, F.J., Qin, J., Tang, K., Tsai, S.Y. & Tsai, M.J. Coup d'Etat: an orphan takes
1462 control. *Endocr Rev* **32**, 404-21 (2011).
- 1463 36. Job, C. & Tan, S.S. Constructing the mammalian neocortex: the role of intrinsic
1464 factors. *Dev Biol* **257**, 221-32 (2003).
- 1465 37. Tsai, S.Y. & Tsai, M.J. Chick ovalbumin upstream promoter-transcription factors
1466 (COUP-TFs): coming of age. *Endocr Rev* **18**, 229-40 (1997).
- 1467 38. O'Leary, D.D., Chou, S.J. & Sahara, S. Area patterning of the mammalian cortex.
1468 *Neuron* **56**, 252-69 (2007).

- 1469 39. Al-Kateb, H. *et al.* NR2F1 haploinsufficiency is associated with optic atrophy,
1470 dysmorphism and global developmental delay. *Am J Med Genet A* **161A**, 377-81
1471 (2013).
- 1472 40. Brown, K.K. *et al.* NR2F1 deletion in a patient with a de novo paracentric
1473 inversion, inv(5)(q15q33.2), and syndromic deafness. *Am J Med Genet A* **149A**,
1474 931-8 (2009).
- 1475 41. Cardoso, C. *et al.* Periventricular heterotopia, mental retardation, and epilepsy
1476 associated with 5q14.3-q15 deletion. *Neurology* **72**, 784-92 (2009).
- 1477 42. Malan, V. *et al.* Molecular characterisation of a prenatally diagnosed 5q15q21.3
1478 deletion and review of the literature. *Prenat Diagn* **26**, 231-8 (2006).
- 1479 43. Quinn, J.J. *et al.* Rapid evolutionary turnover underlies conserved lncRNA-
1480 genome interactions. *Genes Dev* **30**, 191-207 (2016).
- 1481 44. Jonk, L.J. *et al.* Cloning and expression during development of three murine
1482 members of the COUP family of nuclear orphan receptors. *Mech Dev* **47**, 81-97
1483 (1994).
- 1484 45. Chanda, S. *et al.* Generation of induced neuronal cells by the single
1485 reprogramming factor ASCL1. *Stem Cell Reports* **3**, 282-96 (2014).
- 1486 46. Treutlein, B. *et al.* Dissecting direct reprogramming from fibroblast to neuron
1487 using single-cell RNA-seq. *Nature* **534**, 391-5 (2016).
- 1488 47. Mall, M. *et al.* Myt1l safeguards neuronal identity by actively repressing many
1489 non-neuronal fates. *Nature* **544**, 245-249 (2017).
- 1490 48. Basu, S.N., Kollu, R. & Banerjee-Basu, S. AutDB: a gene reference resource for
1491 autism research. *Nucleic Acids Res* **37**, D832-6 (2009).
- 1492 49. Elling, U. *et al.* A reversible haploid mouse embryonic stem cell biobank resource
1493 for functional genomics. *Nature* **550**, 114-118 (2017).
- 1494 50. Zhang, Y. *et al.* Rapid single-step induction of functional neurons from human
1495 pluripotent stem cells. *Neuron* **78**, 785-98 (2013).
- 1496 51. Shen, Y. *et al.* A map of the cis-regulatory sequences in the mouse genome.
1497 *Nature* **488**, 116-20 (2012).
- 1498 52. McLean, C.Y. *et al.* GREAT improves functional interpretation of cis-regulatory
1499 regions. *Nat Biotechnol* **28**, 495-501 (2010).
- 1500 53. Chambers, S.M. *et al.* Highly efficient neural conversion of human ES and iPSC
1501 cells by dual inhibition of SMAD signaling. *Nat Biotechnol* **27**, 275-80 (2009).
- 1502 54. Matsui, T. *et al.* Neural stem cells directly differentiated from partially
1503 reprogrammed fibroblasts rapidly acquire gliogenic competency. *Stem Cells* **30**,
1504 1109-19 (2012).

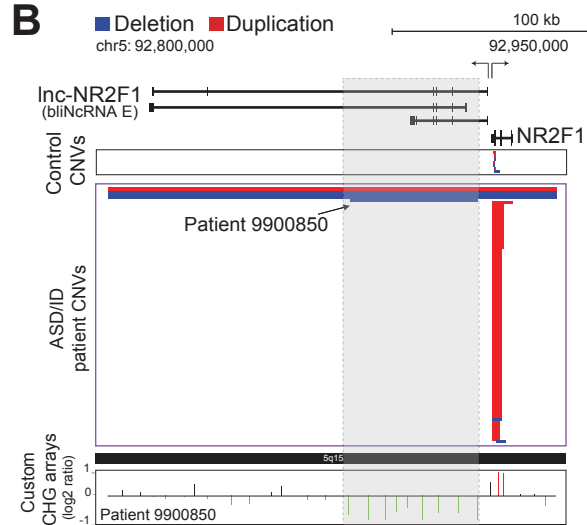
- 1505 55. Kim, J. *et al.* Direct reprogramming of mouse fibroblasts to neural progenitors.
1506 *Proc Natl Acad Sci U S A* **108**, 7838-43 (2011).
- 1507 56. Qureshi, I.A., Mattick, J.S. & Mehler, M.F. Long non-coding RNAs in nervous
1508 system function and disease. *Brain Res* **1338**, 20-35 (2010).
- 1509 57. Ulitsky, I. & Bartel, D.P. lincRNAs: genomics, evolution, and mechanisms. *Cell*
1510 **154**, 26-46 (2013).
- 1511 58. Ulitsky, I., Shkumatava, A., Jan, C.H., Sive, H. & Bartel, D.P. Conserved function
1512 of lincRNAs in vertebrate embryonic development despite rapid sequence
1513 evolution. *Cell* **147**, 1537-50 (2011).
- 1514 59. Chu, C., Qu, K., Zhong, F.L., Artandi, S.E. & Chang, H.Y. Genomic maps of long
1515 noncoding RNA occupancy reveal principles of RNA-chromatin interactions. *Mol*
1516 *Cell* **44**, 667-78 (2011).
- 1517 60. Wang, L., Feng, Z., Wang, X., Wang, X. & Zhang, X. DEGseq: an R package for
1518 identifying differentially expressed genes from RNA-seq data. *Bioinformatics* **26**,
1519 136-8 (2010).
- 1520 61. Leek, J.T., Johnson, W.E., Parker, H.S., Jaffe, A.E. & Storey, J.D. The sva
1521 package for removing batch effects and other unwanted variation in high-
1522 throughput experiments. *Bioinformatics* **28**, 882-3 (2012).
- 1523 62. Ernst, J. & Kellis, M. Discovery and characterization of chromatin states for
1524 systematic annotation of the human genome. *Nat Biotechnol* **28**, 817-25 (2010).
- 1525 63. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nat*
1526 *Methods* **9**, 357-9 (2012).
- 1527 64. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol* **9**,
1528 R137 (2008).
- 1529 65. Yu, G., Wang, L.G. & He, Q.Y. ChIPseeker: an R/Bioconductor package for ChIP
1530 peak annotation, comparison and visualization. *Bioinformatics* **31**, 2382-3 (2015).
- 1531 66. Conrad, T. & Orom, U.A. Cellular Fractionation and Isolation of Chromatin-
1532 Associated RNA. *Methods Mol Biol* **1468**, 1-9 (2017).
1533
1534

Figure 1: LncRNA candidate loci are recurrently mutated in patients with neurodevelopmental disorders

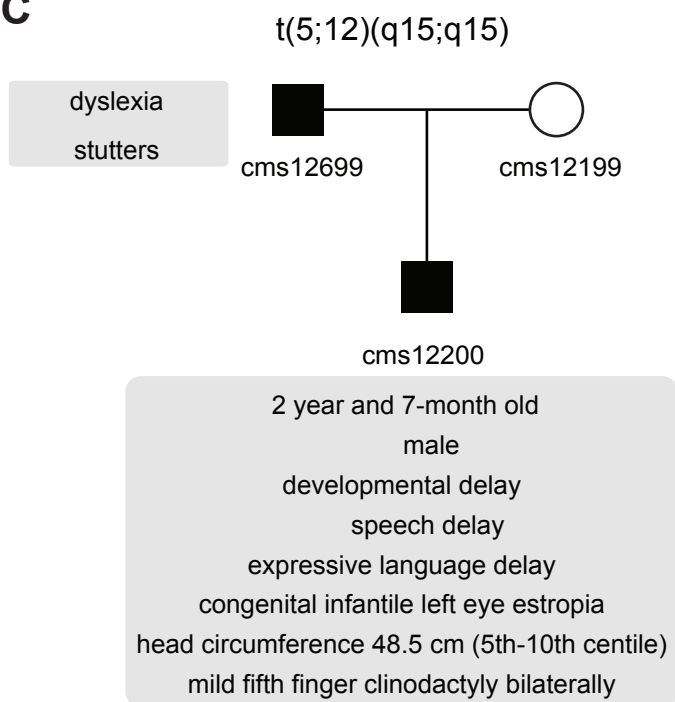
A



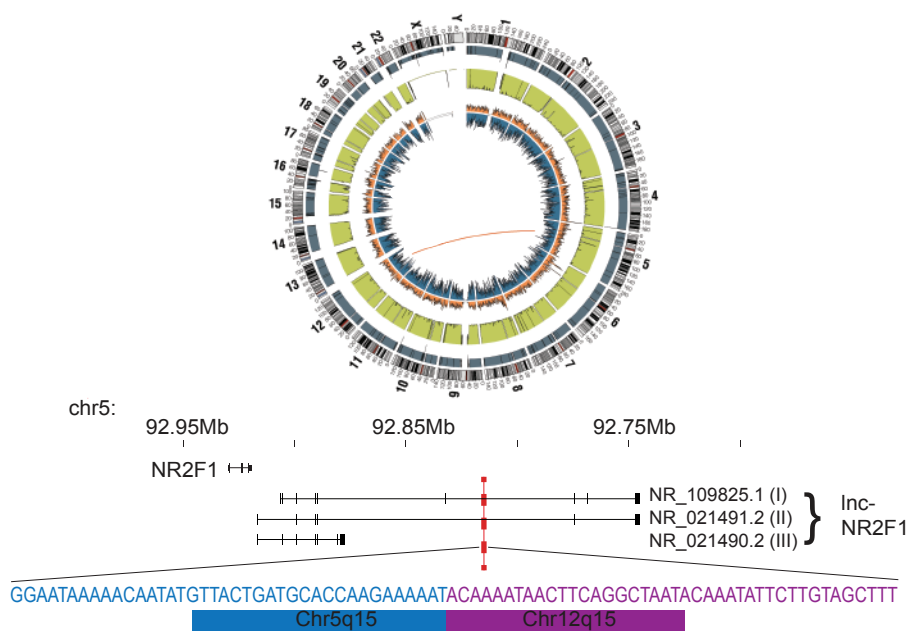
B



C



D



E

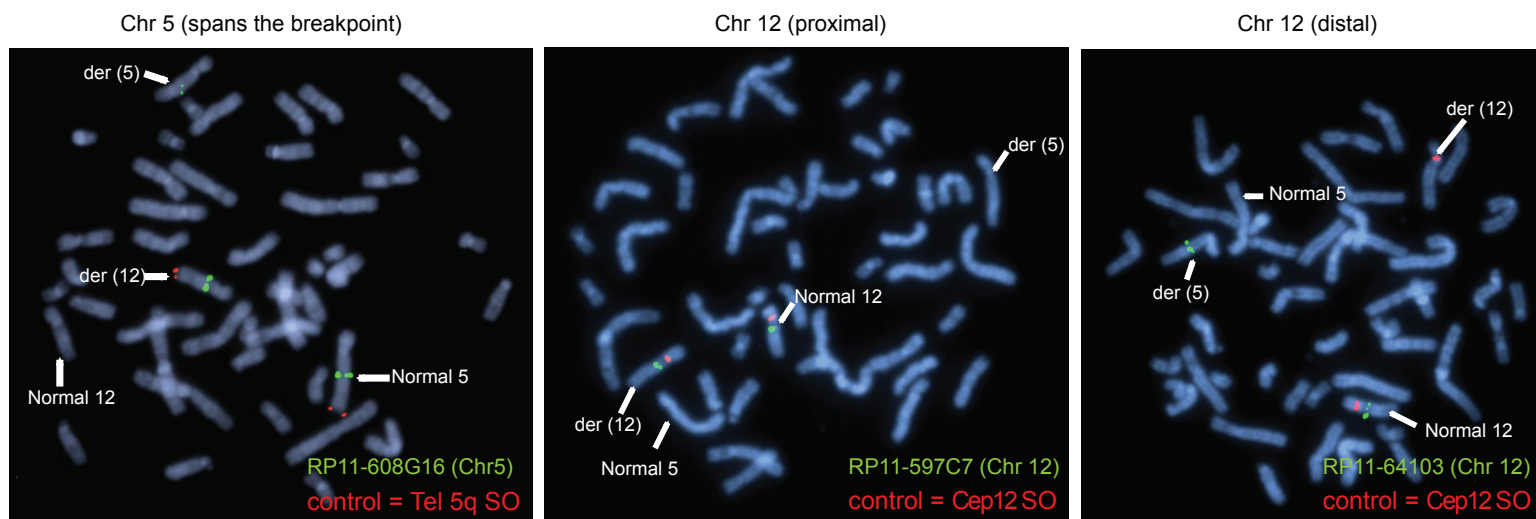
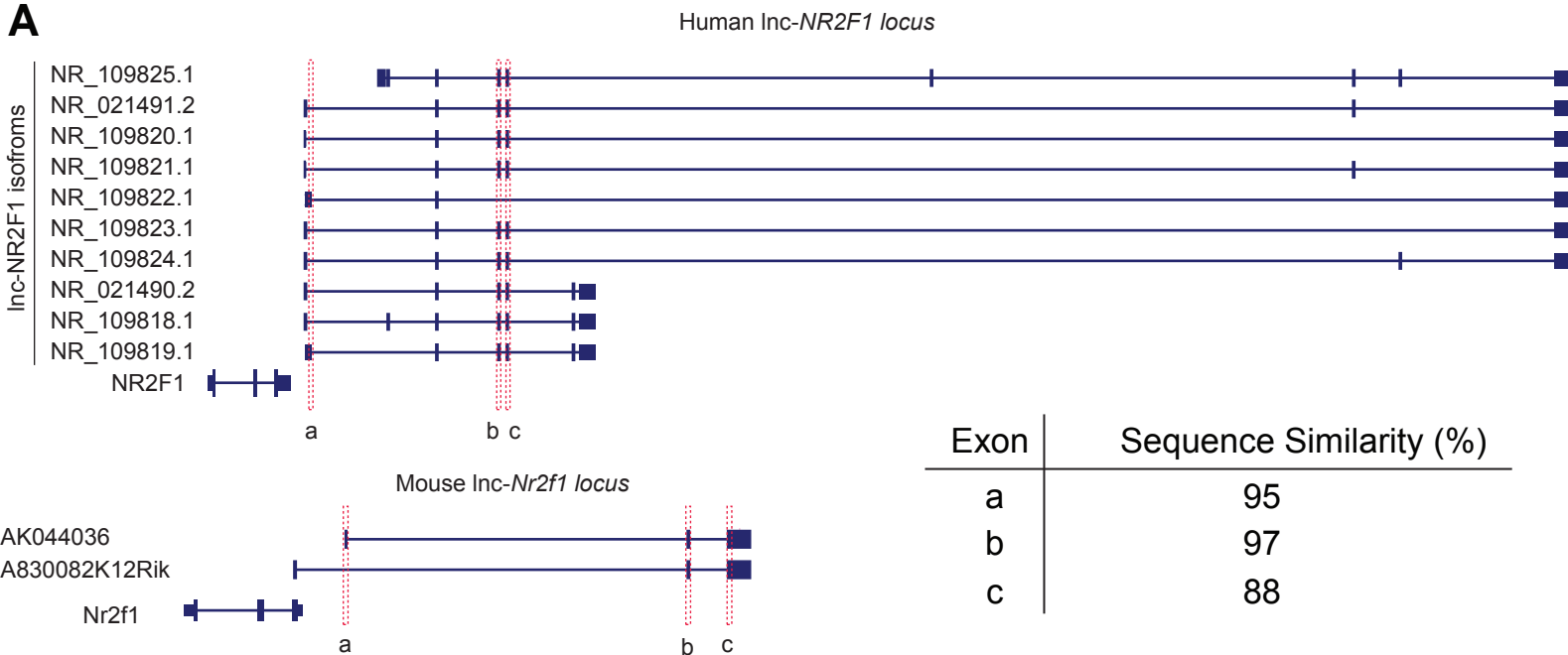
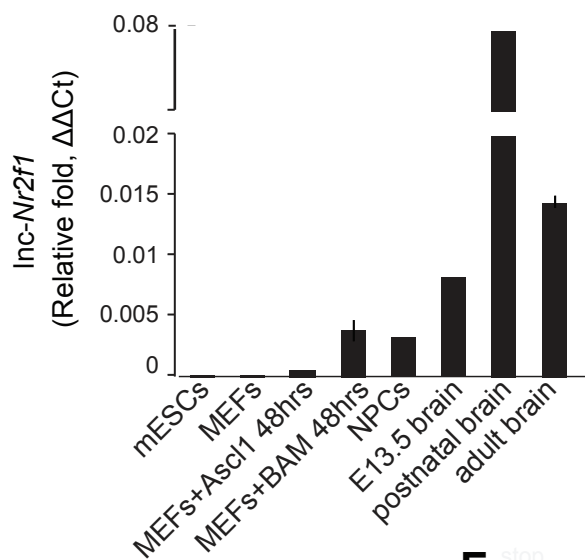


Figure 2: Molecular characterization of mouse *Inc-Nr2f1*

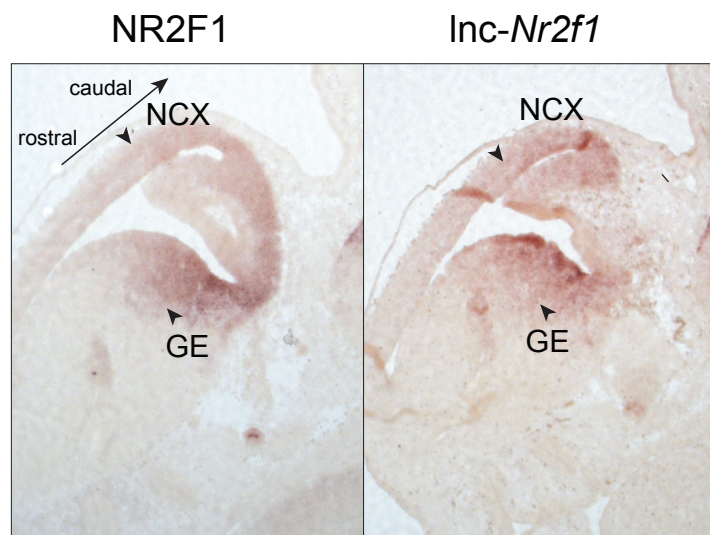
A



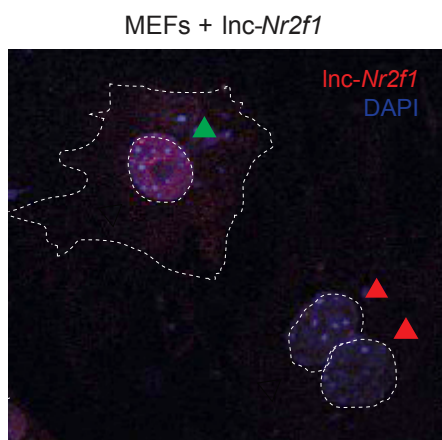
B Mouse brain development expression



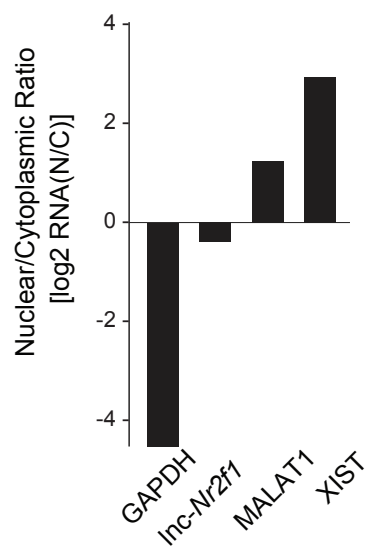
C



D



E ^{stop} *Inc-Nr2f1* are present in nucleus and cytoplasm in E13.5 caudal brain



F

Inc-NR2F1 is enriched in the chromatin

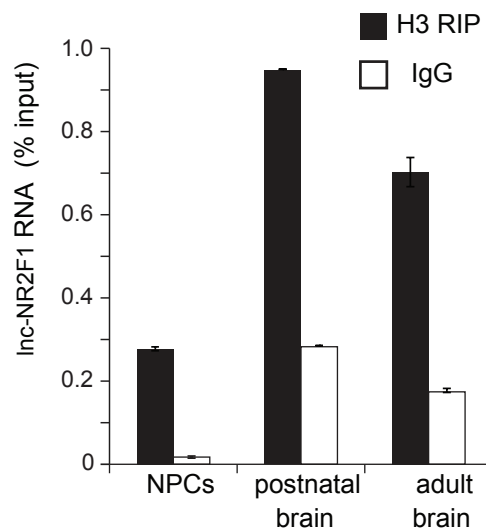


Figure 3: Mouse *Inc-Nr2f1* KO reveals *Inc-Nr2f1* regulates neuronal genes

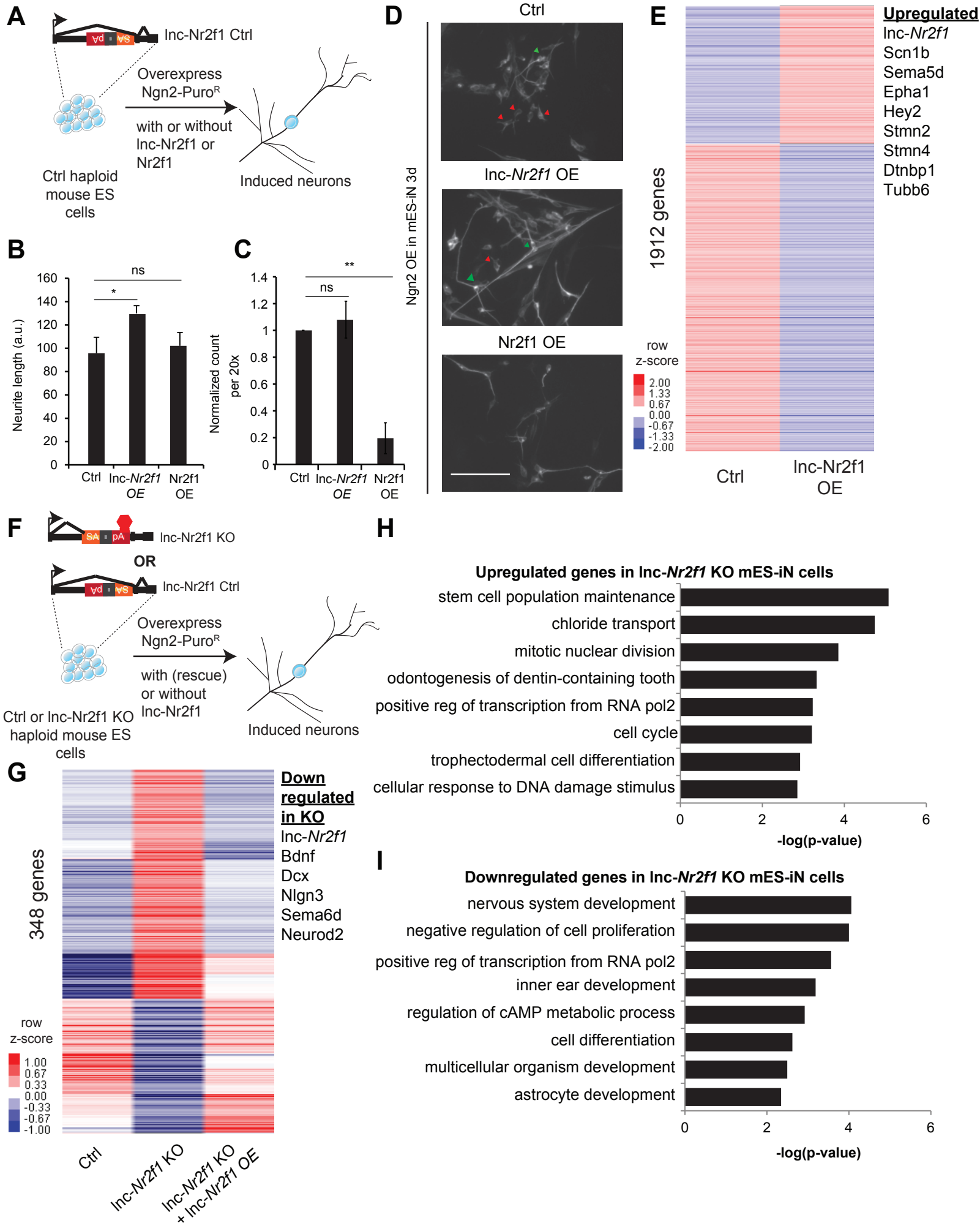


Figure 4: lnc-Nr2f1 binds to distinct genomic loci regulating neuronal genes

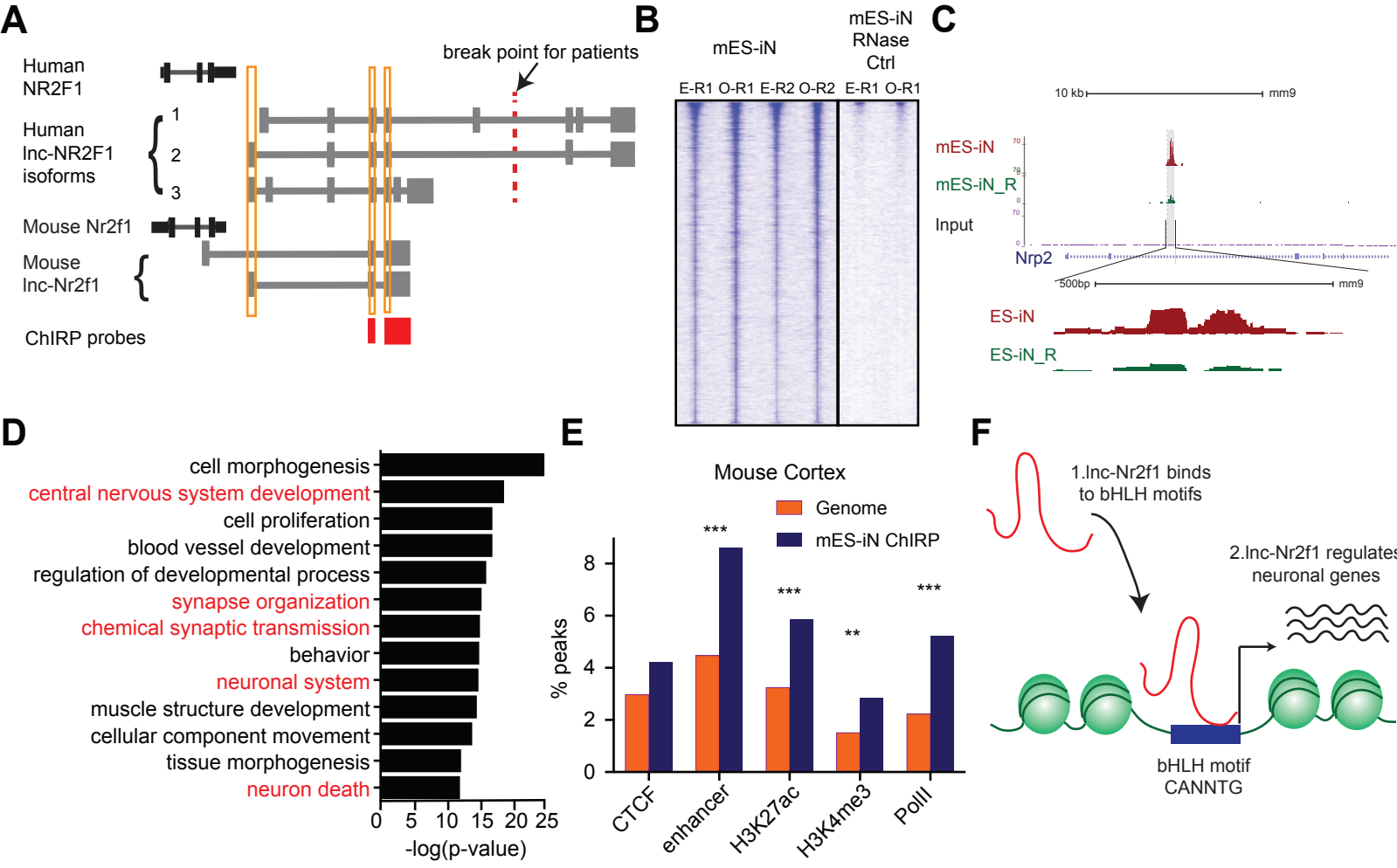


Figure 5: Human Inc-NR2F1 shows isoform specific chromatin binding

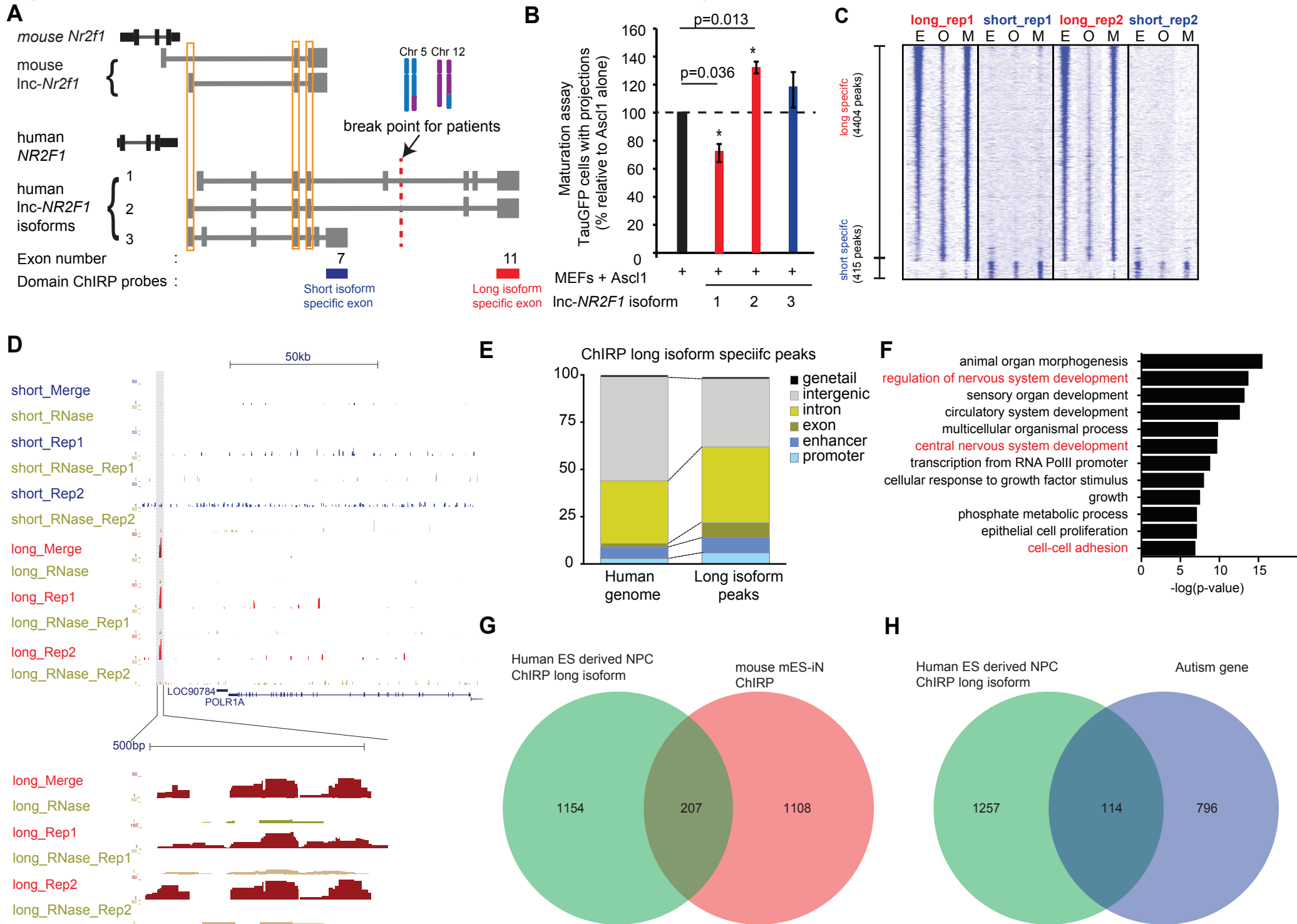


Figure 1 - figure supplement 1

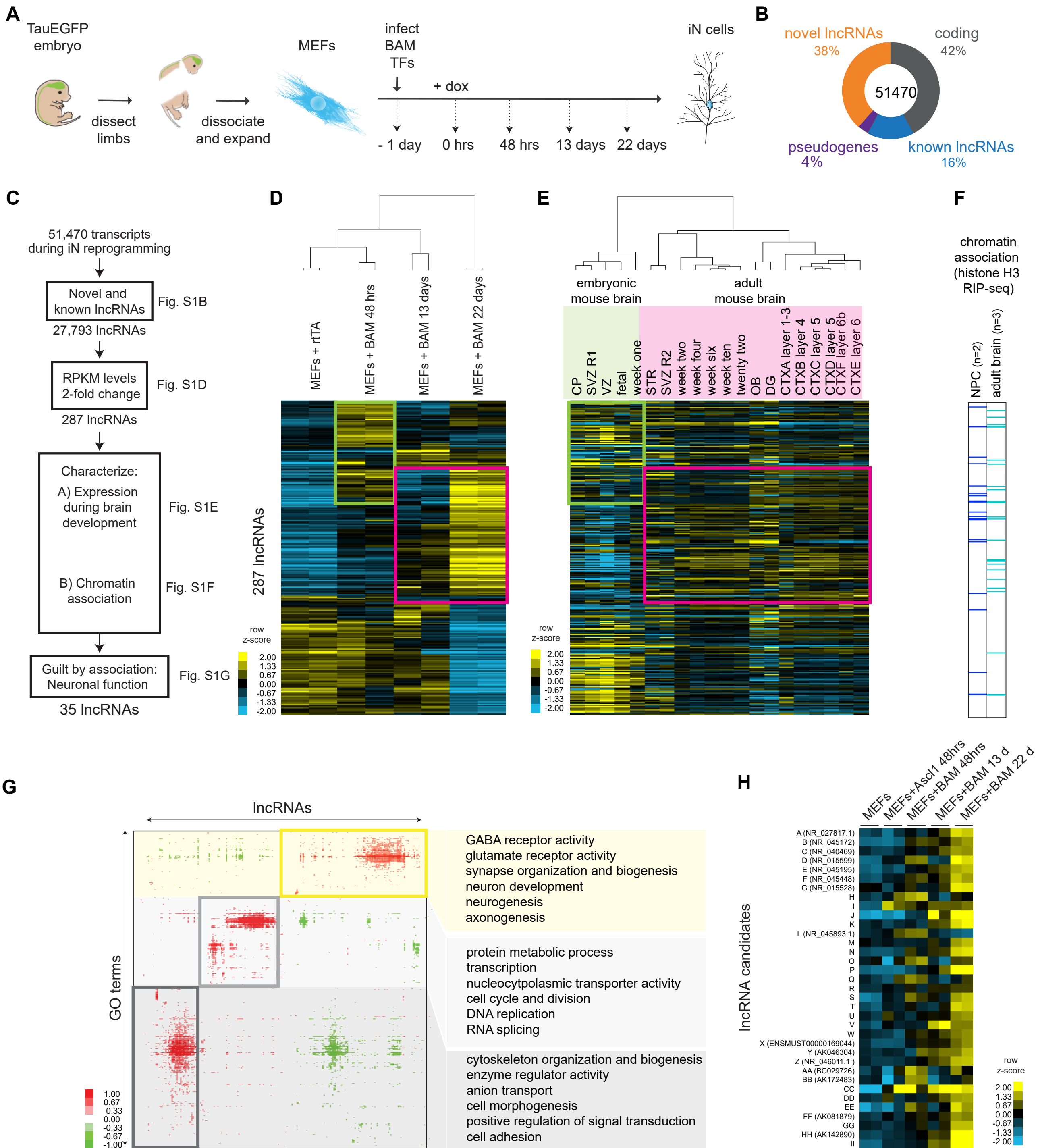


Figure 1 - figure supplement 2

A

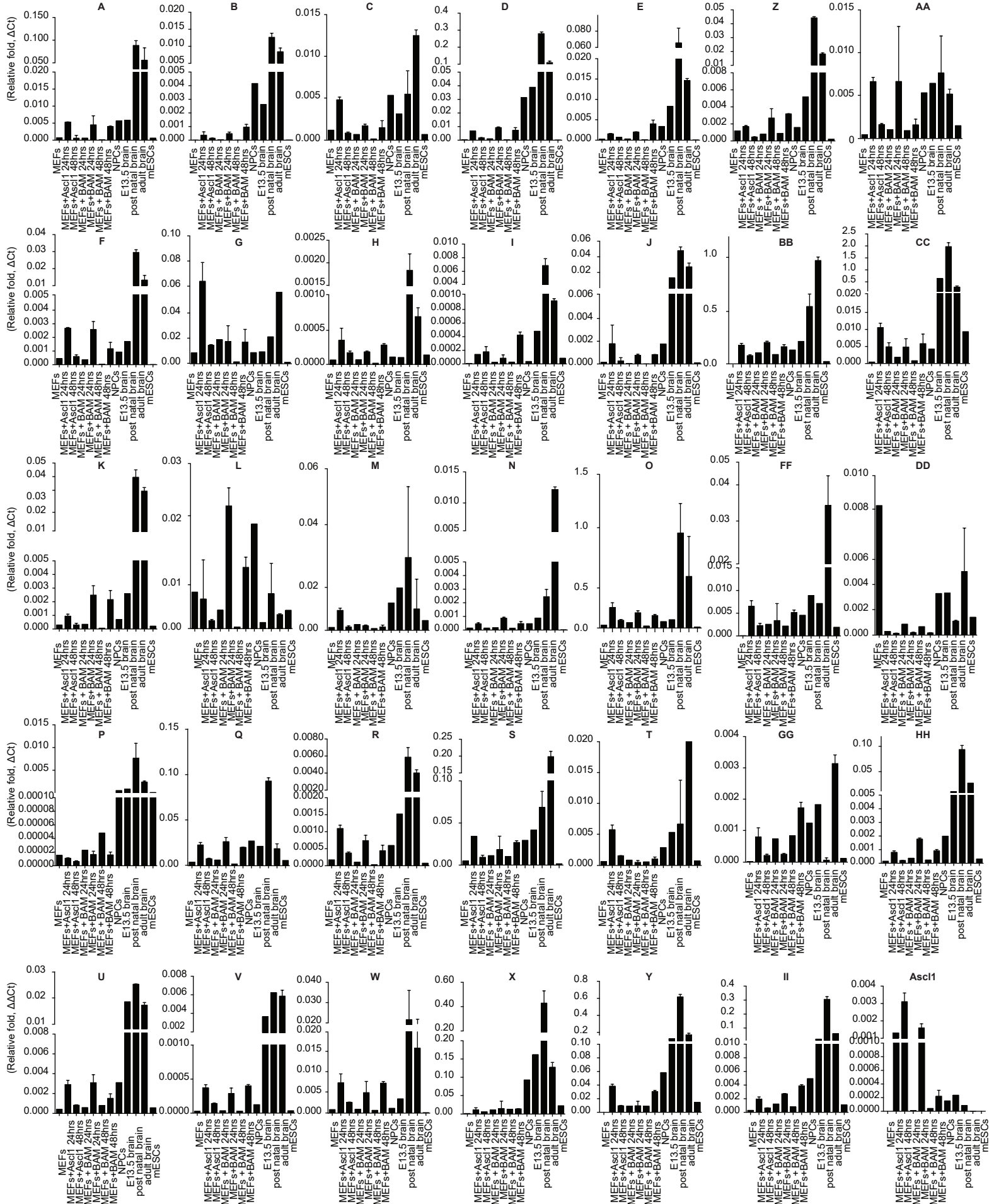


Figure 1 - figure supplement 3

A

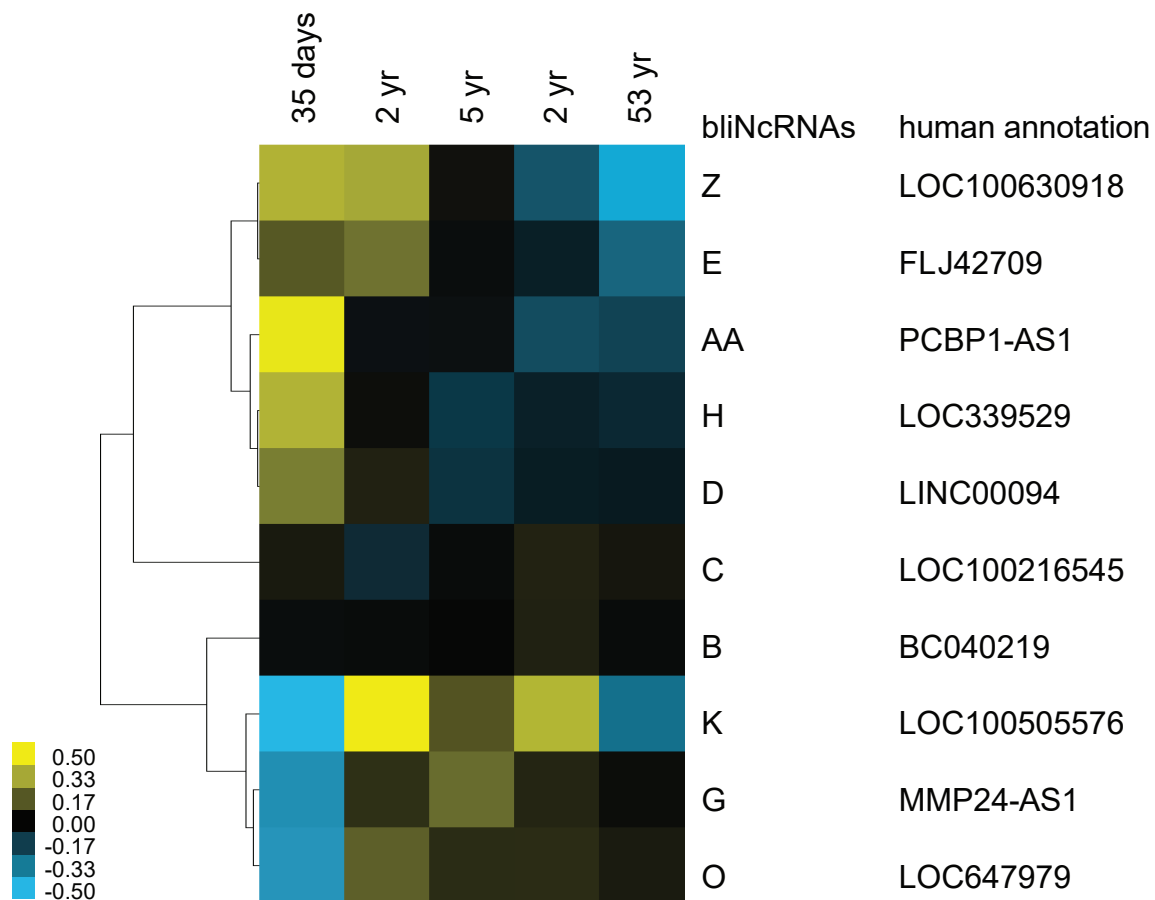
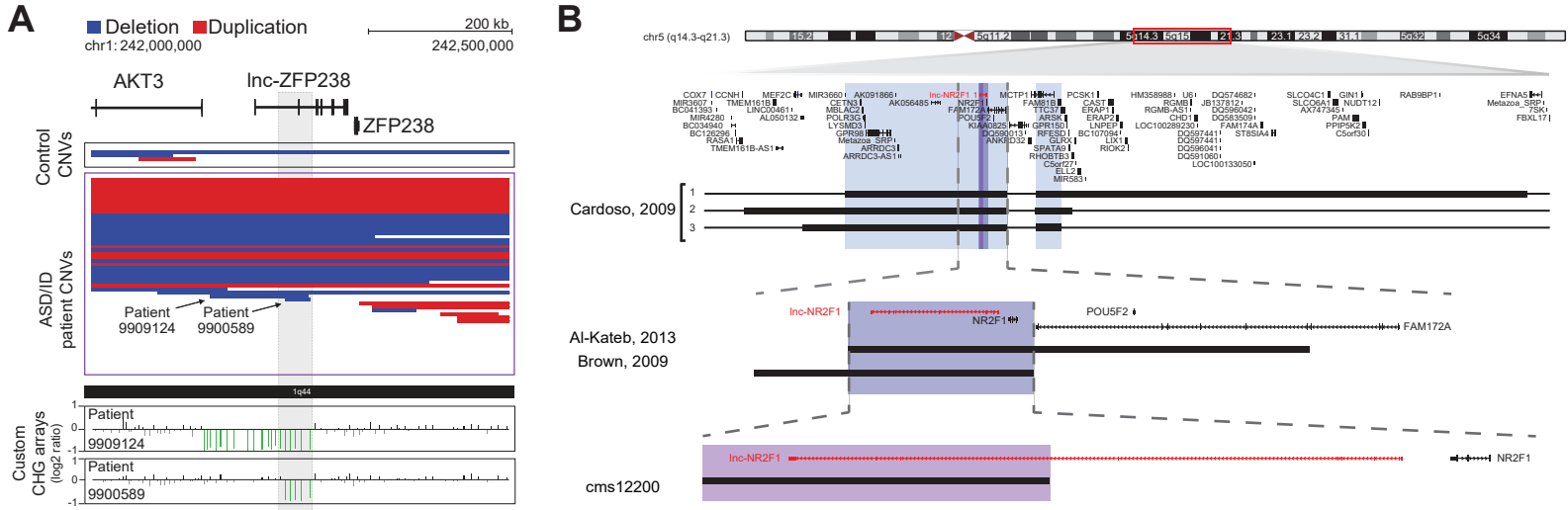
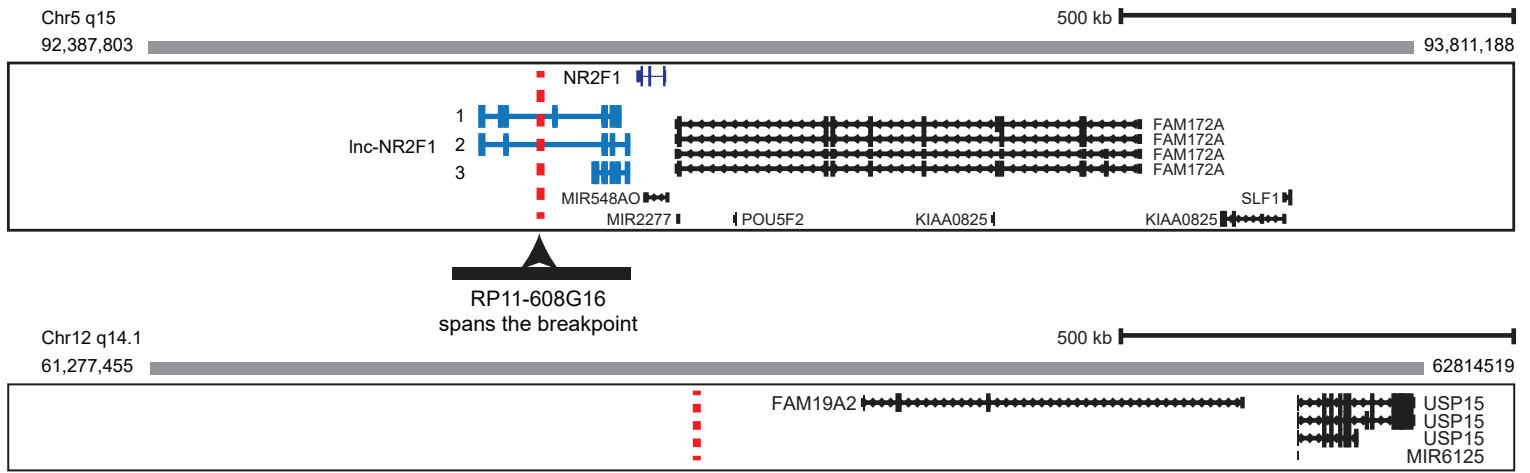


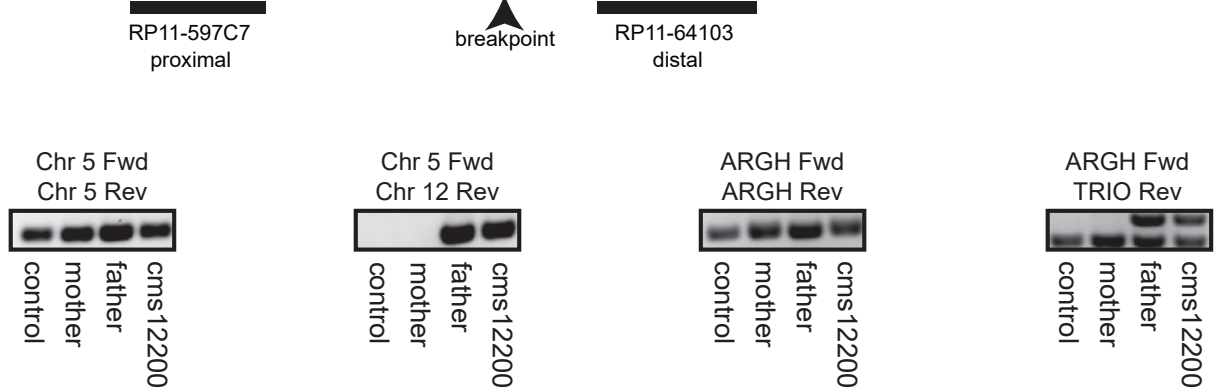
Figure 1 - figure supplement 4



C cms12200: 46, XY, t(5;12)(q15;q15)



D



E

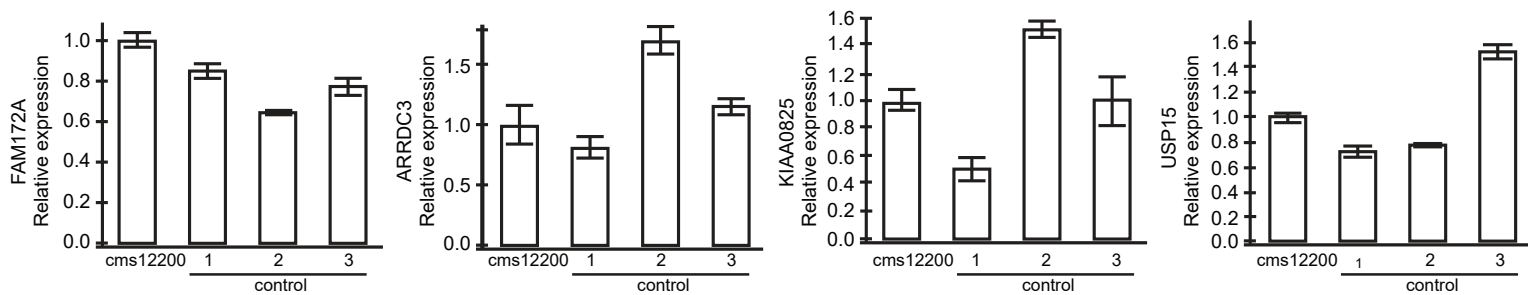
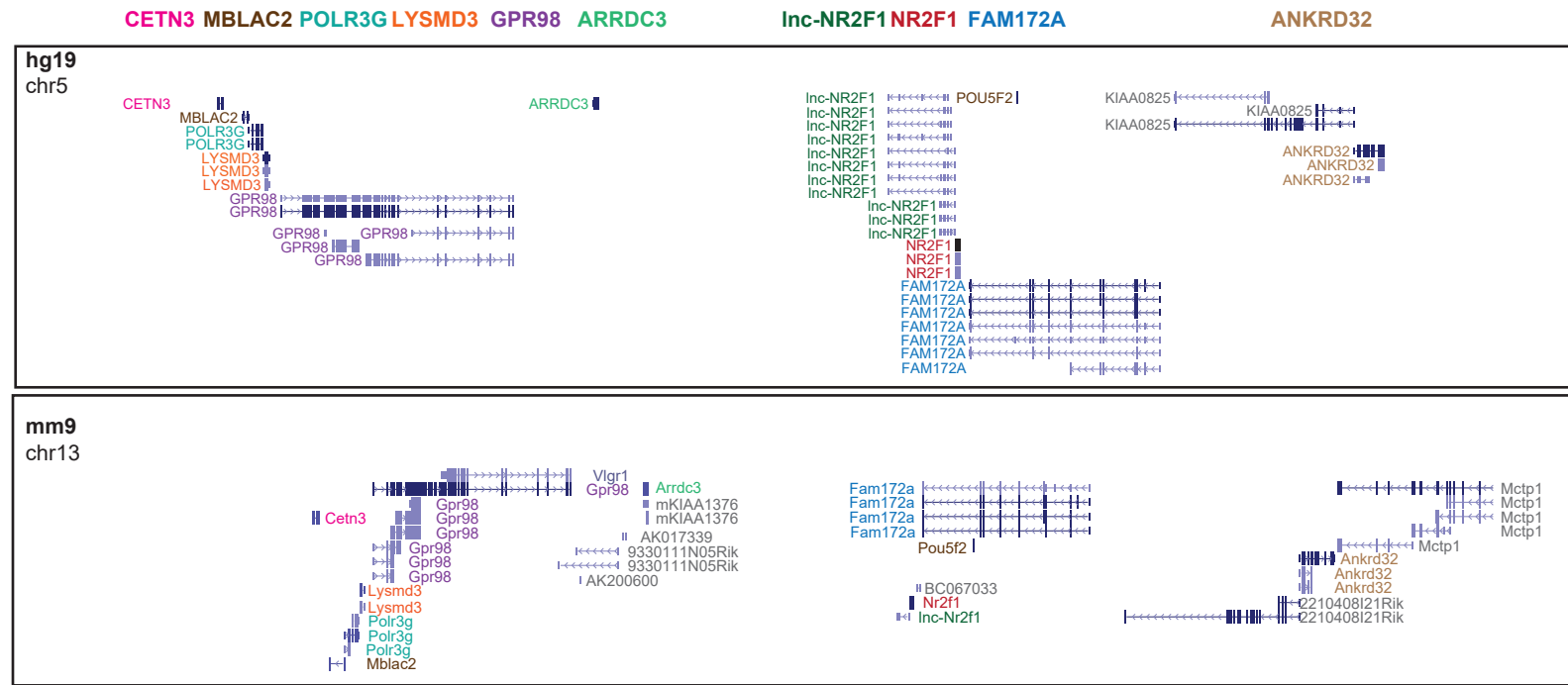
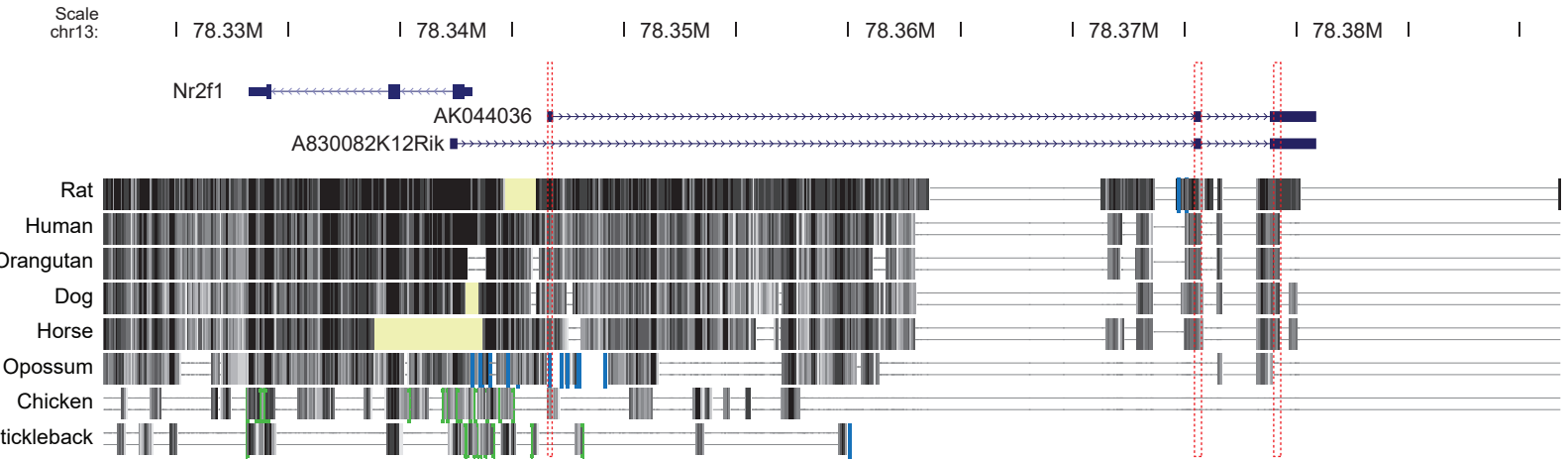


Figure 2 - figure supplement 1

A Synteny conservation



B Sequence conservation



C Microdomain conservation

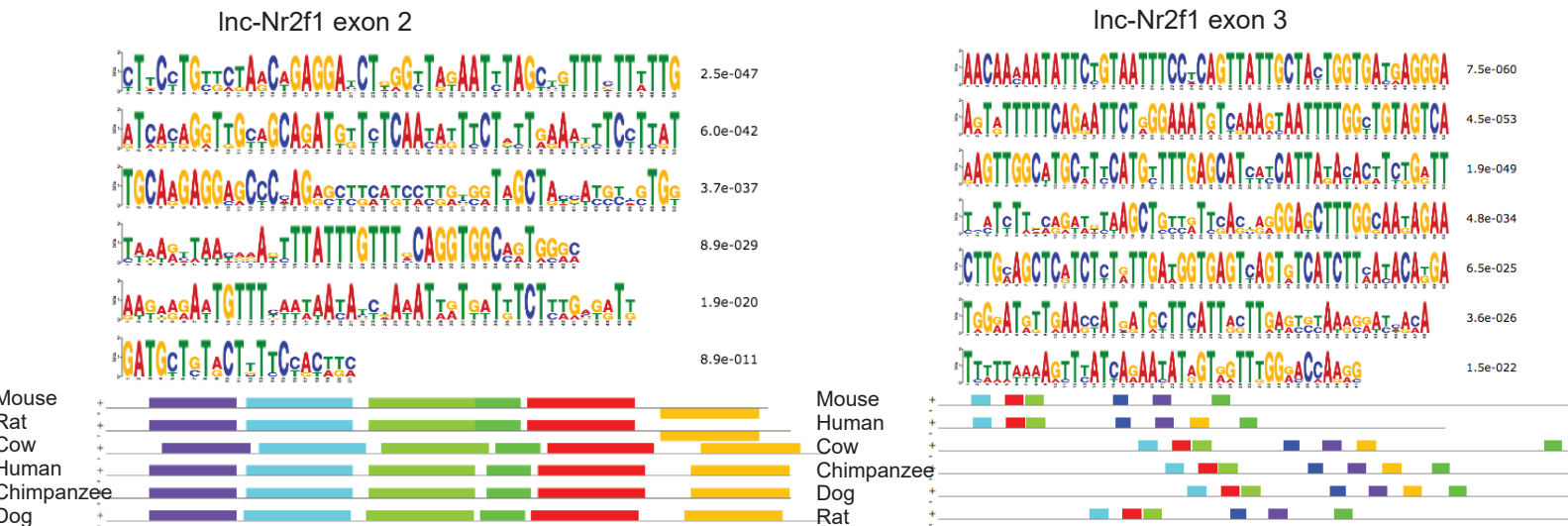
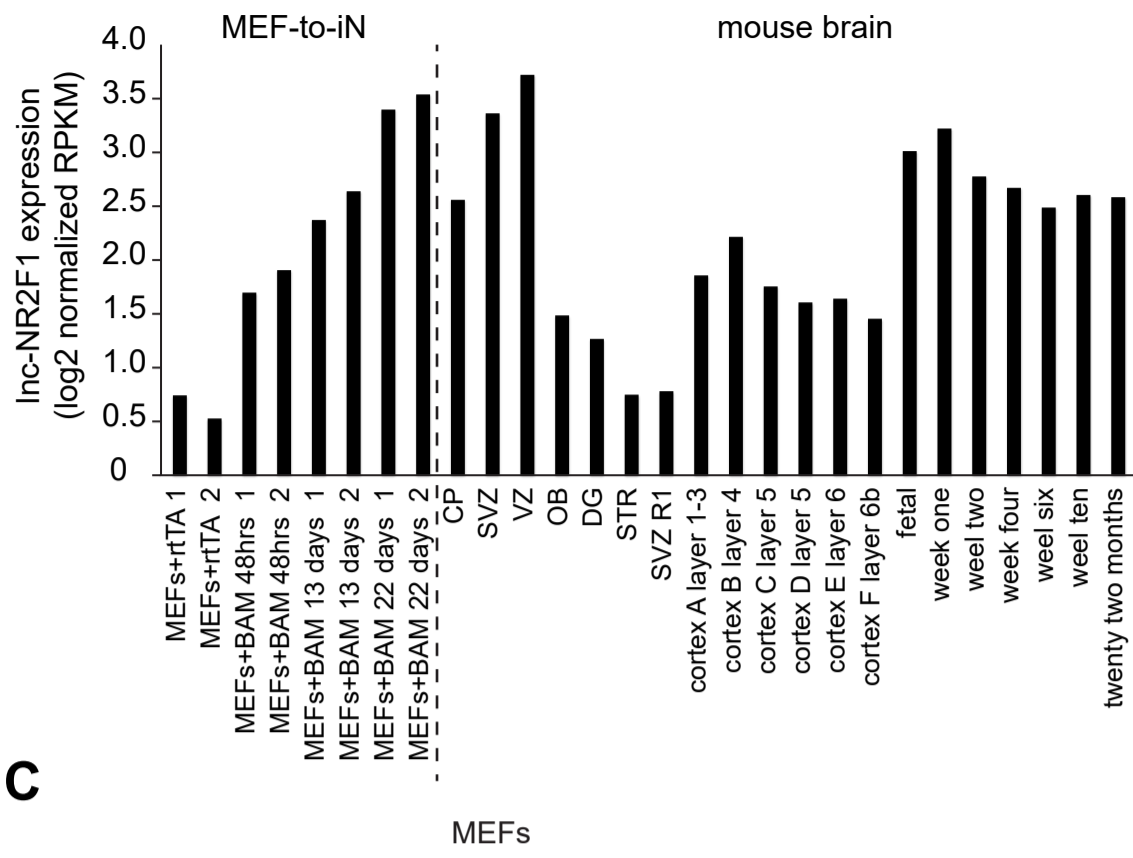
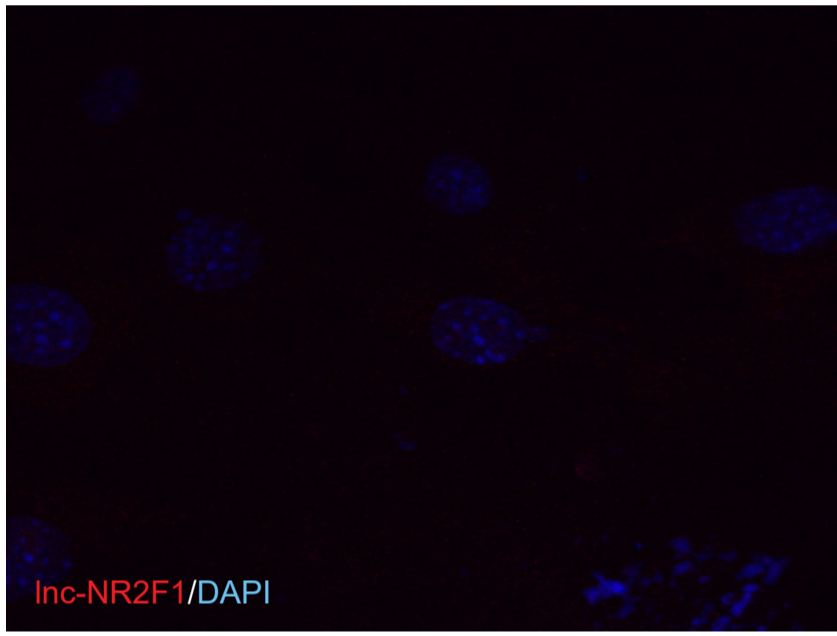


Figure 2 - figure supplement 2

A



C



B

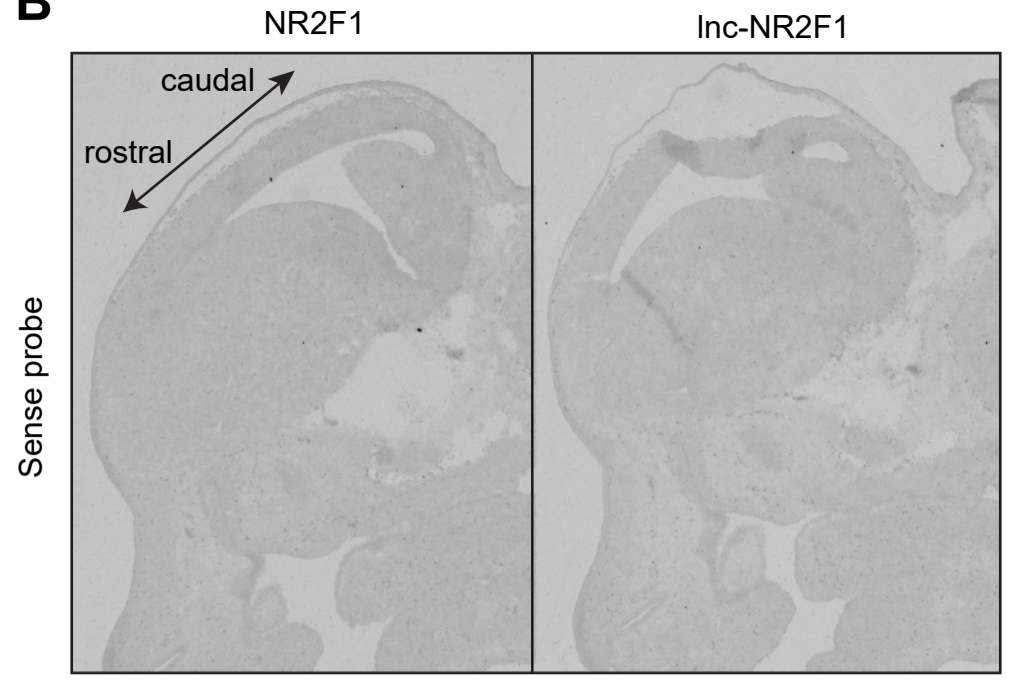


Figure 3 – figure supplement 1

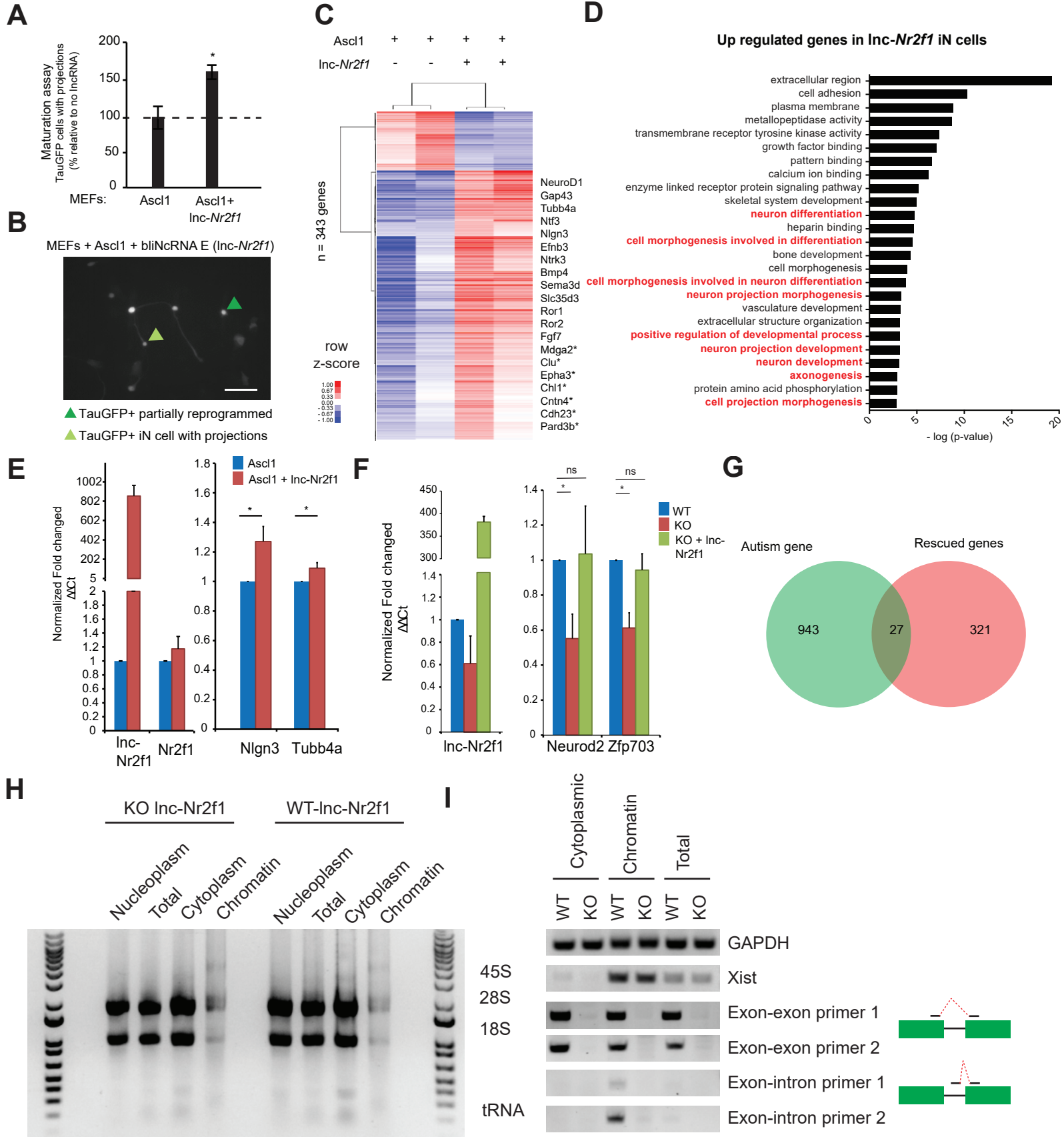


Figure 3 – figure supplement 2

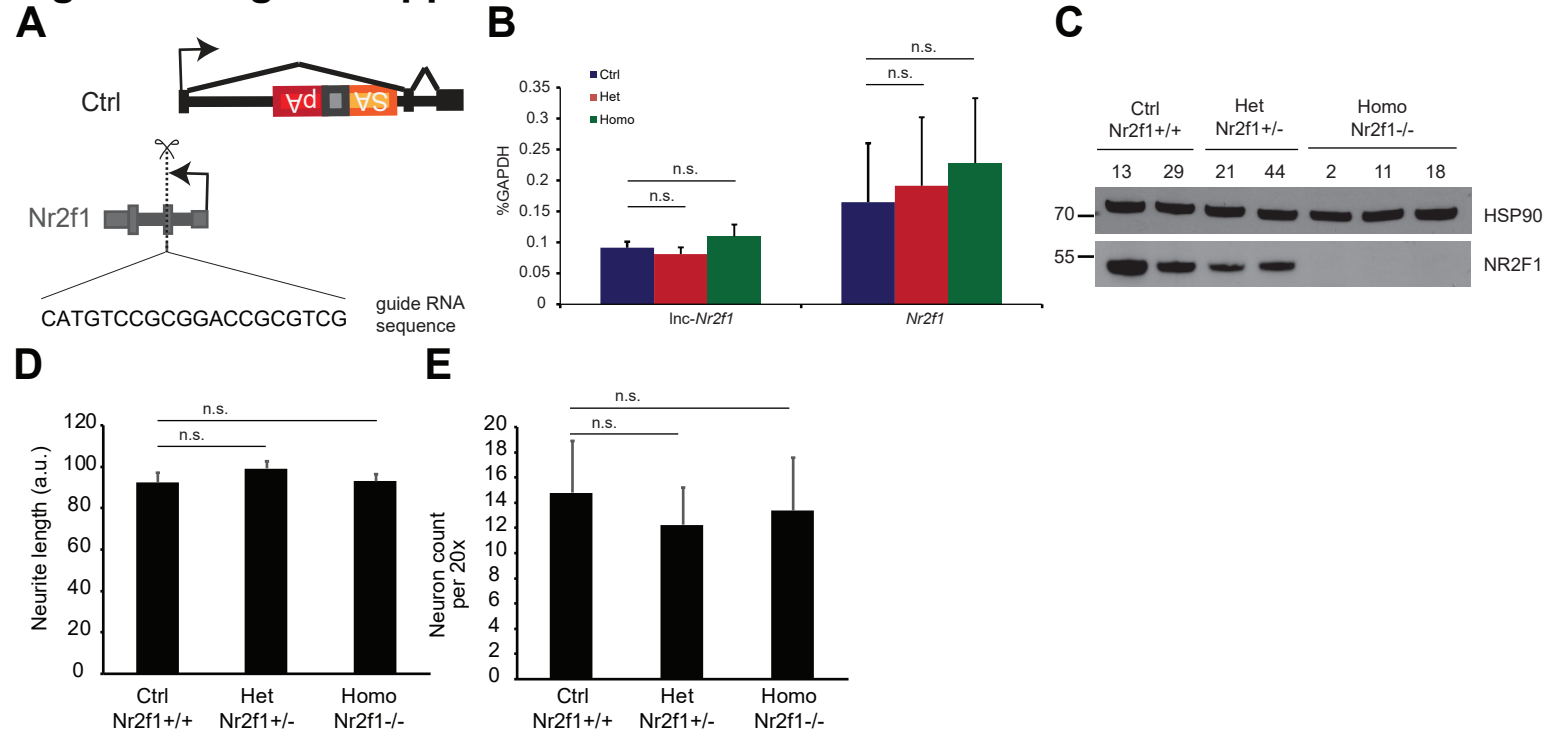
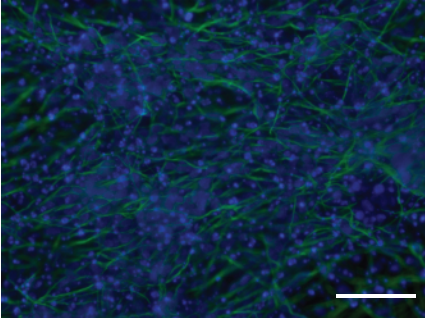
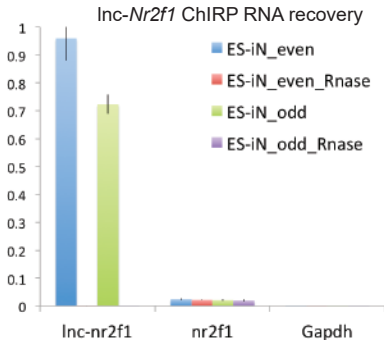


Figure 4 – figure supplement 1

A



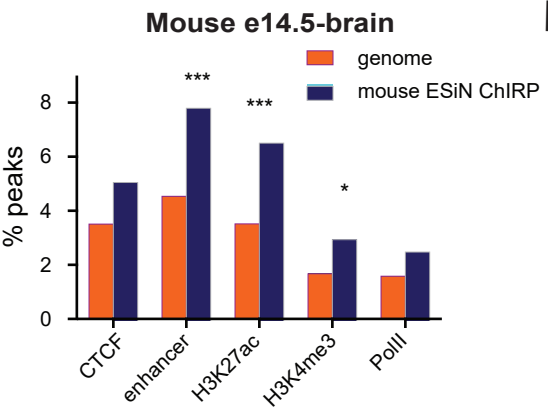
B



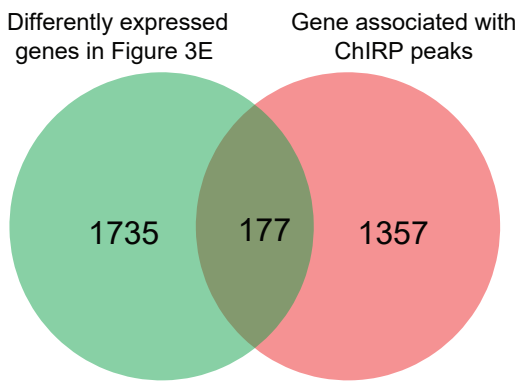
C

Enriched top 10 motifs	mES-iN ChIRP	p-value	% peaks
NeuroD1		1e-8	11.9%
ZBTB18		1e-7	8.1%
Tcf21		1e-7	12.1%
Atoh1		1e-5	13.5%
Ascl1		1e-5	16.4%
Olig2		1e-5	26.1%
Tgif1/2		1e-4	41%
SCL		1e-4	48.2%
Stat6		1e-4	27.9%
Ptf1a		1e-4	8.9%

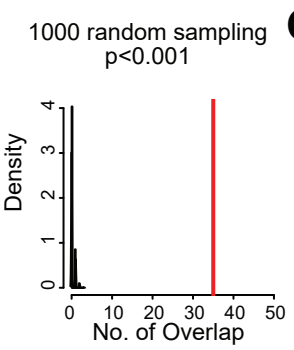
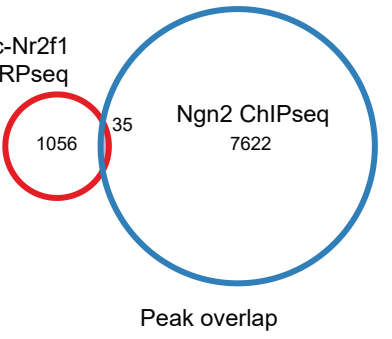
D



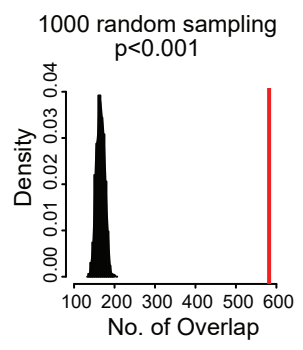
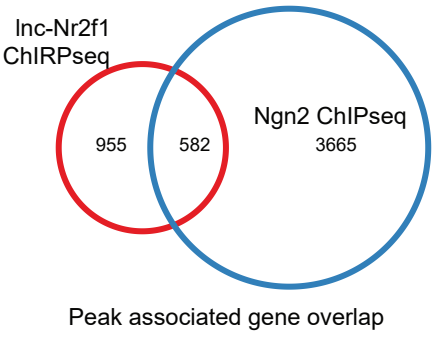
E



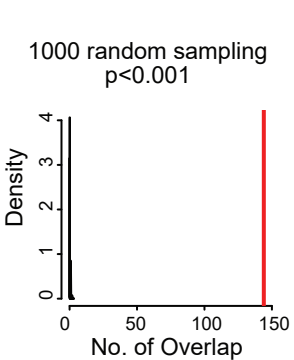
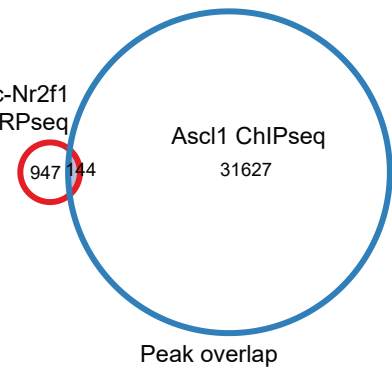
F



G



F



G

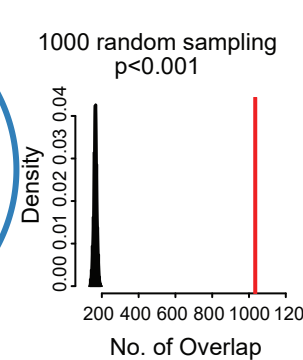
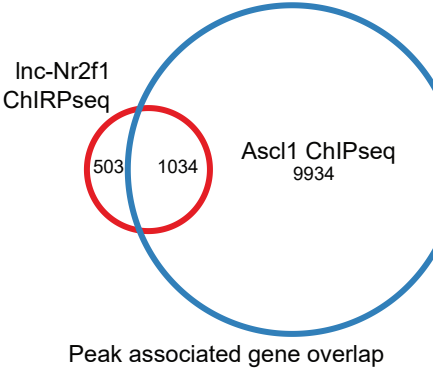


Figure 5 - figure supplement 1

Human Inc-NR2F1 ChIRP for hNPC

