1    **A conserved role for SFPQ in repression of pathogenic cryptic last exons**

2    Patricia M. Gordon[1,*], Fursham Hamid[1], Eugene V. Makeyev[1], and Corinne Houart[1]

3

4    [1] Centre for Developmental Neurobiology and MRC Centre for Neurodevelopmental

5    Disorders, IoPPN, Guy's Campus, King's College London, London SE1 1UL, UK

6    * Corresponding author

7

8    **Abstract**

9    The RNA-binding protein SFPQ plays an important role in neuronal development and

10   has been associated with several neurodegenerative disorders, including ALS, FTLD, and

11   Alzheimer's Disease. Here, we report that loss of *sfpq* leads to premature termination of

12   multiple transcripts due to widespread activation of previously unannotated cryptic last

13   exons (CLEs). These CLEs appear preferentially in long introns of genes with neuronal

14   functions and dampen gene expression outputs and/or give rise to short peptides

15   interfering with the normal gene functions. We show that one such peptide encoded by

16   the CLE-containing *epha4b* mRNA isoform is responsible for neurodevelopmental

17   defects in the *sfpq* mutant. The uncovered CLE-repressive activity of SFPQ is conserved

18   in mouse and human, and SFPQ-inhibited CLEs are found across ALS iPSC-derived

19   neurons. These results greatly expand our understanding of SFPQ function and uncover

20   a new gene regulation mechanism with wide relevance to human pathologies.

21

22   **Keywords**: SFPQ, neurodevelopment, zebrafish, alternative polyadenylation, cryptic

23   exons, ALS, neurodegeneration

24

25

26

27  **Introduction**

28  Neurons are highly polarized cells with specialized compartments that must be

29  able to respond to growth cues as well as to form and modify their synapses in an

30  activity-dependent manner. Each compartment of a neuron is able to achieve functional

31  specificity by maintaining a unique proteome (Holt & Schuman, 2013; Hanus &

32  Schuman, 2013; Cagnetta *et al*, 2018). Protein localization in neurons has been shown

33  to be driven largely by RNA transportation and local translation (Zappulo *et al*, 2017),

34  suggesting that neuronally-expressed genes must have special regulatory mechanisms

35  to ensure proper transcription, localization, and translation of each RNA. Indeed, RNAs

36  from neuronal tissue are regulated by a complex array of alternative splicing, intron

37  retention, and alternative cleavage and polyadenylation (Mauger *et al*, 2016;

38  Traunmüller *et al*, 2016; Furlanis *et al*, 2019; Iijima *et al*, 2019; Taliaferro *et al*, 2016;

39  Ciolli Mattioli *et al*, 2019; Guvenek & Tian, 2018; Tushev *et al*, 2018).

40  Splicing Factor Proline/Glutamine Rich (SFPQ) is a ubiquitously expressed RNA

41  binding protein of the DBHS family with diverse roles in alternative splicing,

42  transcriptional regulation, microRNA targeting, paraspeckle formation, and RNA

43  transport into axons (Patton *et al*, 1993; Dye & Patton, 2001; Kim *et al*, 2011; Cosker *et*

44  *al*, 2016; Bottini *et al*, 2017; Mora Gallardo *et al*, 2019; Takeuchi *et al*, 2018; Knott *et al*,

45  2016). Inactivation of the *sfpq* gene causes early embryonic lethality in mouse and

46  zebrafish as well as impaired cerebral cortex development, reduced brain boundary

47  formation, and axon outgrowth defects (Lowery *et al*, 2007; Thomas-Jinu *et al*, 2017;

48  Takeuchi *et al*, 2018; Saud *et al*, 2017). In humans, *sfpq* mutations have been linked to

49  neurodegenerative diseases such as Alzheimer's, ALS, and FTD, and SFPQ interacts with

2

50   the ALS-associated RNA binding proteins TDP-43 and FUS (Ke *et al*, 2012; Wang *et al*,

51   2015; Ishigaki *et al*, 2017; Luisier *et al*, 2018; Tyzack *et al*, 2019; Lu *et al*, 2018).

52        While SFPQ is known to play a role in alternative splicing, only a few RNA targets

53   of SFPQ have been identified. Intriguingly, SFPQ has opposing effects on splicing,

54   depending on the target: it represses inclusion of exon 10 of tau and exon 4 of CD45, but

55   conversely it promotes inclusion of the N30 exon of non-muscle myosin heavy-chain II-B

56   (Ray *et al*, 2011; Ishigaki *et al*, 2017; Heyd & Lynch, 2010; Yarosh *et al*, 2015; Kim *et al*,

57   2011). In addition to its role in splicing, SFPQ has been shown to be part of the 3'-end

58   processing complex, where it enhances cleavage and polyadenylation at suboptimal

59   polyadenylation sites (Hall-Pogar *et al*, 2007; Rosonina *et al*, 2005; Shi *et al*, 2009). The

60   mechanisms by which SFPQ regulates mRNA processing are still unclear, however, and

61   more work is necessary to understand its contribution to normal and pathological cell

62   states.

63        To understand the molecular functions of SFPQ in developing neurons, we

64   performed an RNA-seq analysis of *sfpq* homozygous null mutant zebrafish embryos at

65   24 hpf, the stage of phenotypic onset. Our results reveal a novel role for the protein: loss

66   of SFPQ causes premature termination of transcription as a result of previously

67   unannotated pre-mRNA processing events that we refer to as Cryptic Last Exons (CLEs).

68   Here we describe the formation of CLEs and show that not only do the truncated

69   transcripts act as a form of negative regulation of gene expression levels, but they also

70   directly contribute to the *sfpq* pathology. This function of SFPQ is conserved across

71   vertebrates and may be implicated in human SFPQ-mediated disease states.

72
73   **Results**
74

75   **Identification of the SFPQ-dependent splicing regulation program**

3

76    To examine the effect of SFPQ on gene expression and RNA splicing, we analyzed

77    total RNA extracted from 24 hpf *sfpq*$^{-/-}$ zebrafish embryos and their heterozygous or

78    wildtype siblings by RNA sequencing (RNA-seq). Differential gene expression analysis

79    using Cufflinks RNA-seq workflow (Trapnell *et al*, 2012) uncovered 189 genes that were

80    upregulated and 1044 genes that were downregulated in the mutant samples by a factor

81    of at least 1.3-fold with q≤0.05 (Figure 1a). These results are consistent with our

82    previous microarray study, which showed the vast majority of genes with differential

83    expression in *sfpq*$^{-/-}$ embryos as being downregulated (Thomas-Jinu *et al*, 2017). Gene

84    ontology (GO) analysis of the new dataset, using total transcribed genes as a background

85    gene set, showed enrichment for neuron-specific terms, including neuronal

86    differentiation and axon guidance (Tables S1 and S2). Using Cufflinks' differential

87    isoform switch analysis, we identified 112 genes with significant change in the relative

88    expressions of splice variants in the mutants (q≤0.05; Table S3). GO analysis of these

89    regulated transcripts again showed an over-representation of neuron-specific terms

90    including axonogenesis, axon guidance, and dendrite formation (Table S4). Surprisingly,

91    thorough comparison and annotation of these transcripts revealed that 46% of these

92    genes express a splice variant containing a cryptic alternate last exon, not annotated in

93    the zebrafish assembly (Figure S1a). To verify this, we analyzed the dataset with

94    Whippet (Sterne-Weiler *et al*, 2018), a tool that sensitively detects changes in the usage

95    of alternative exons and additionally allows quantitation of gene expression changes.

96    Whippet also uncovered a high proportion of downregulated genes in *sfpq*$^{-/-}$ embryos

97    (Figure 1b). More importantly, the analysis confirmed that splicing of alternate last

98    exons is the most abundant (18.5%) category of SFPQ-regulated splicing events (Figure

99    1c). Systematic classification of these exons into "known" and "cryptic" events

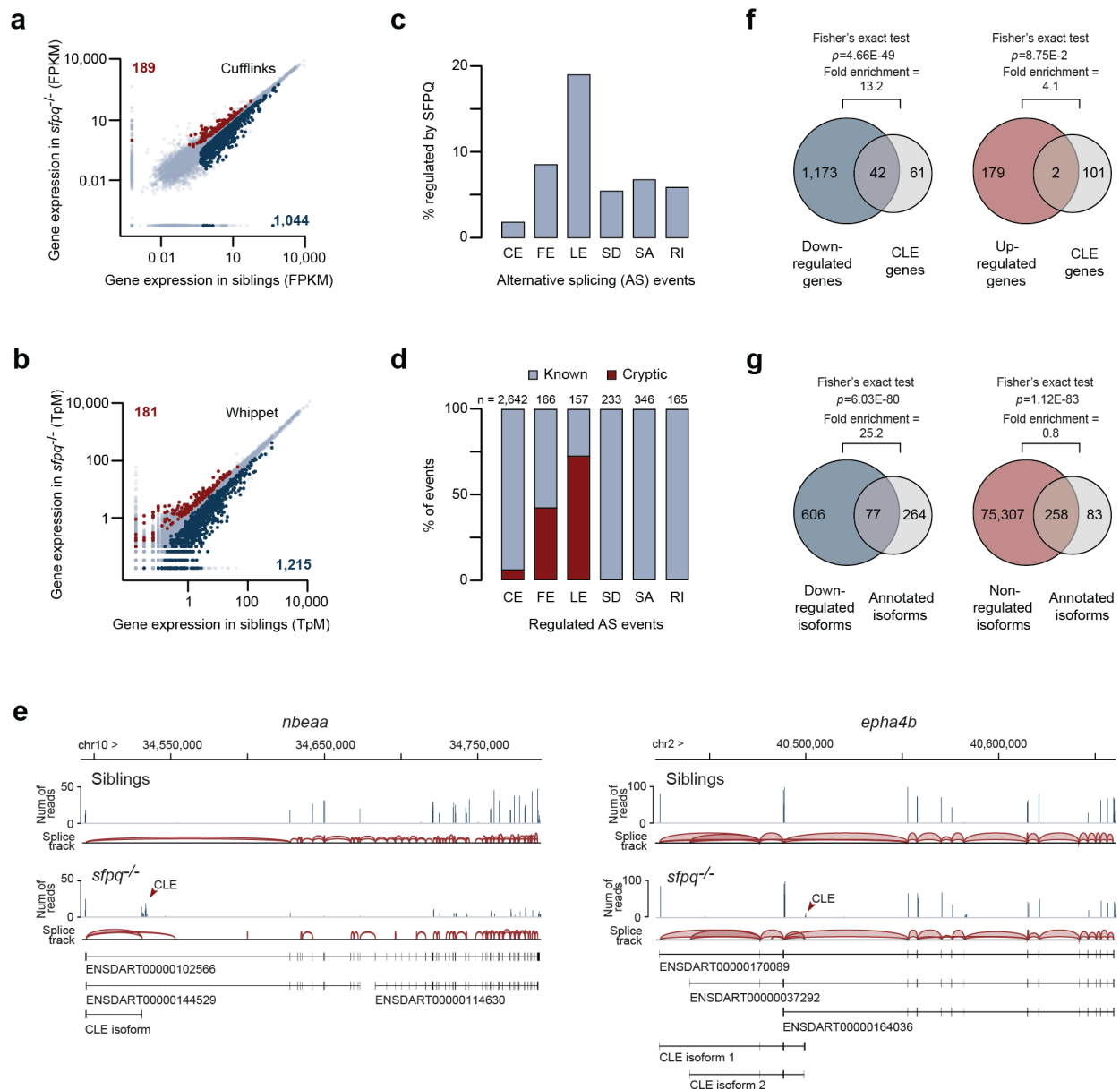100   corroborates that the majority (113 out of 157) of these last exons have not been

4

**Fig. 1**

**Figure 1: SFPQ regulates the formation of cryptic last exons (CLEs)**

**a-b**, Scatter plot showing expression values of genes in sfpq-/- and siblings, analyzed using Cufflinks (a) or Whippet (b) pipelines.

**c**, Alternative last exon splicing is highly regulated by SFPQ. CE: cassette exon, FE: first exon, LE: last exon, SD: splice donor, SA: splice acceptor, RI: retained intron.

**d**, Majority of SFPQ-regulated last exon events are cryptic.

**e**, Sashimi plots showing example CLE formation in nbeaa and epha4b. Top tracks: plot of read coverage from siblings (upper) and sfpq-/- (lower). Bottom tracks: isoforms discovered for each gene.

**f**, Genes expressing CLE-containing isoforms tend to be down-regulated in sfpq-/-.

**g**, Normal long isoforms (annotated isoforms) from CLE-expressing genes tend to be down-regulated in sfpq-/-.

101   previously annotated (Figure 1D, Table S5). These last exons were expressed from 106

102   genes, 25 of which were also detected by the Cufflinks pipeline (Figure S1b). We refer to

103   this pervasive splicing defect, in which the transcript undergoes premature termination

104   after the inclusion of a cryptic exon, as Cryptic Last Exons (CLEs) (Figure 1e).

105

106   **The use of CLEs inversely correlates with expression of full-length transcripts**

107      Of the 106 CLE-expressing genes, 97% exhibited increased splicing of CLEs in

108   *sfpq*$^{-/-}$ (Table S5). Notably, more than half of these genes were downregulated in mutants

109   (~13 fold enrichment over the number of genes expected by chance; Fisher's exact test

110   p=4.66 x 10$^{-49}$) indicating concurrent alterations in expression level and splicing for

111   these genes (Figure 1f). In line with this finding, the full-length (non-CLE) isoforms from

112   these genes showed an even stronger enrichment for the downregulation effect

113   (exceeding the expectation ~25-fold; Fisher's exact test p=6.03 x 10$^{-80}$) (Figure 1g). To

114   verify these results, we performed RT-qPCR on five selected CLE-containing genes:

115   *nbeaa*, *gdf11*, *epha4b*, *trip4*, and *b4galt2*. In all cases, cryptic exons showed a substantial

116   increase in expression level in *sfpq*$^{-/-}$ mutants compared to siblings (Figure S1c-g).

117   Additionally, we detected a strong downregulation of the full-length isoforms in four of

118   the five genes, suggesting that the loss of *sfpq* causes upregulation of the CLE isoforms at

119   the expense of their normal counterparts (Figure S1d-g). These results argue that SFPQ

120   is required to repress CLE splicing in order to maintain stable gene expression.

121

122   **CLEs tend to occur in long introns and show evidence of interspecies conservation**

123      In order to understand under what conditions CLEs form, we examined CLE-

124   containing introns and compared them to all other introns from the same genes. We

125   first asked where CLE-containing introns are found within their genes and found no bias

126    (Figure 2a). However, when we ranked the introns by length, we found that CLEs are

127    frequently located in the longest intron of the gene (Figure 2b). CLE-containing introns

128    are also significantly longer than the average intron size in the entire zebrafish

129    transcriptome (Figure 2c). Consistent with these results, CLE-containing genes are

130    significantly longer than average zebrafish genes (Figure S2). Within the intron, location

131    of the CLE is biased toward the 5' end, with most appearing approximately 22.4% (95%

132    confidence interval of 18.1% to 26.7%) of the way into the intron (Figure 2d). The

133    distance between the CLE and the upstream exon is generally <10 kb (Figure 2e). We

134    next asked whether the sequences within and neighboring these CLEs are conserved. To

135    this end, we calculated the mean conservation scores of 1 kb sequences (sliding window,

136    1 bp steps) along these CLE-containing introns, using the PhastCons analysis method

137    (Siepel *et al*, 2005). Our analyses showed that sequences containing CLEs tend to have

138    higher conservation scores as compared to sequences within the same intron that do not

139    contain CLEs (Figure S2b). In fact, 18% of these sequences displayed a mean PhastCons

140    score of at least 0.5 (as opposed to 12% of non-CLE sequences; Fisher's exact test: 5.71 x

141    $10^{-51}$) (Figure S2c). Next, we calculated the mean base conservation scores of each CLE

142    together with 250 bp flanking sequences (Figure 2f and S2d). Although only 35% of the

143    CLEs had a PhastCons score of at least 0.5, the sequences near its 3' acceptor site

144    showed the highest conservation (Figure 2f). Together, these data indicate that CLEs are

145    often found close to the 5' ends of very long introns and that at least some of these exons

146    are evolutionarily conserved.

147

148    **CLE-terminated transcripts are cleaved and polyadenylated at the 3' end**

149        Our data thus far suggests that these cryptic transcripts are stably expressed and

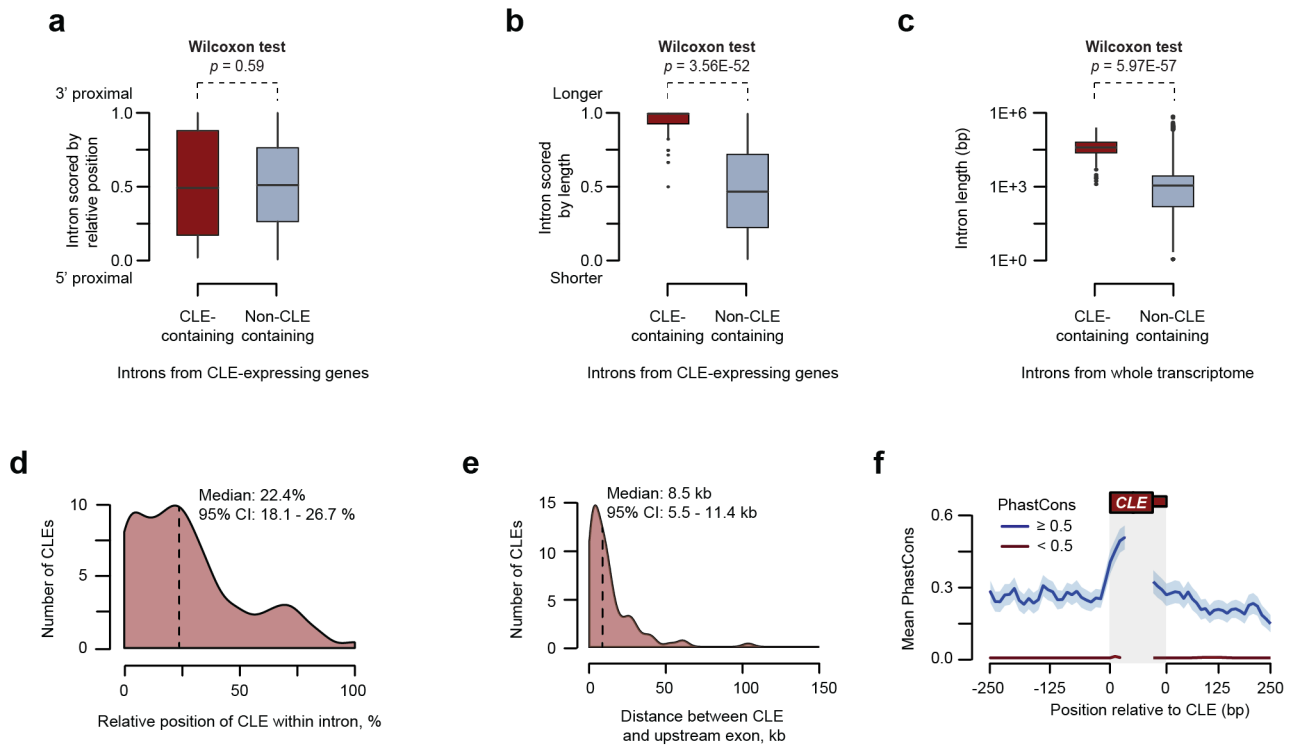150    detectable using RNA-seq and RT-qPCR techniques. Sequence analyses revealed that

**Fig. 2 (a-f)**

**Figure 2 (a-f): Molecular properties of CLEs**

**a-b**, Introns from CLE-expressing genes were scored by its relative position (a) and by its relative length (b), and the distribution of these scores were plotted. Note that introns containing CLE tend to be long and sparsely distributed.

**c**, CLE-containing introns are longer than average introns. Length of CLE-containing introns is compared to all other introns from the zebrafish transcriptome.

**d**, CLEs tend to be found closer to the 5' end of its intron.

**e**, CLEs are found within 10 kb of the upstream constituitive exon.

**f**, Line-plot showing the conservation score of sequences surrounding conserved (blue) and non-conserved (red) CLEs. 280 bp of surrounding intron/CLE junction sequence (250 bp intron and 30bp exon) were binned into 10 bp windows and the mean PhastCons score for each bins were shown (± SEM).
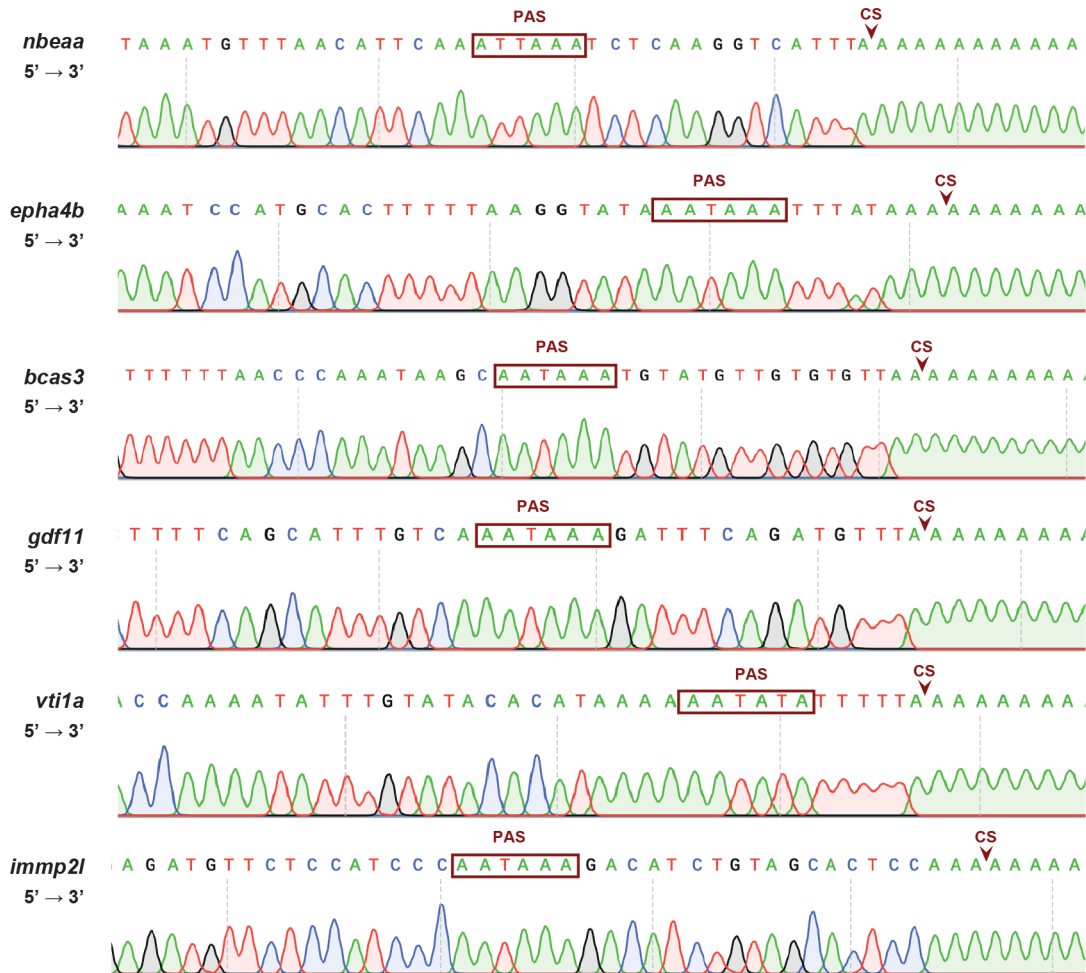
**g**



**Fig. 2 (g)**

**Figure 2 (g): Molecular properties of CLEs**

**g**, Sanger sequencing of 3'RACE PCR products of CLE isoforms. PAS hexamers are shown within red boxes and the predicted cleavage site are marked by arrowheads.

151     63% of these transcripts contain an open reading frame predicted to express truncated

152     peptides with missing C-terminal domains (Table 6). To test if CLE transcripts are

153     polyadenylated, we performed 3' RACE on six CLE-containing transcripts: *bcas3*, *epha4b*,

154     *gdf11*, *immp2l*, *nbeaa*, and *vti1a*. We found that all six showed elements of strong

155     polyadenylation sites (Shi & Manley, 2015): four of the six exons had canonical AAUAAA

156     hexamers just upstream of the cleavage site, while the other two had common one-base

157     substitutions of AUUAAA and AAUAUA. In addition, five of the six contained

158     downstream GUGU sequences, while two also had an upstream UGUA. Although none of

159     the exons had a canonical CA sequence directly 5' of the cleavage site, overall the cryptic

160     exons displayed strong polyadenylation sequences.

161

162     **SFPQ directly binds to sequences adjacent to CLEs**

163     The accumulation of CLE-terminated transcripts in *sfpq* mutants raises the

164     question of whether SFPQ represses CLEs in a direct manner. SFPQ binds promiscuously

165     to a wide range of RNA sequences (Yarosh *et al*, 2015; Knott *et al*, 2016), making binding

166     prediction difficult. Using a binding motif produced by a recent *in vitro* study (Ray *et al*,

167     2013), however, we found a significant enrichment in predicted SFPQ binding sites

168     upstream of cryptic exon sequences compared to control last exons (Figure 3a). To

169     validate this, we purified SFPQ-RNA complexes in 24 hpf embryos using standard CLIP

170     protocol and quantified the relative amount of bound CLE RNA fragments using RT-

171     qPCR. Our results confirmed that SFPQ binds either within the CLE or in adjacent 5' or 3'

172     intronic regions of at least three CLE transcripts (Figure 3b-d). These results support

173     the idea that SFPQ directly binds to region surrounding CLEs to regulate their inclusion.

174

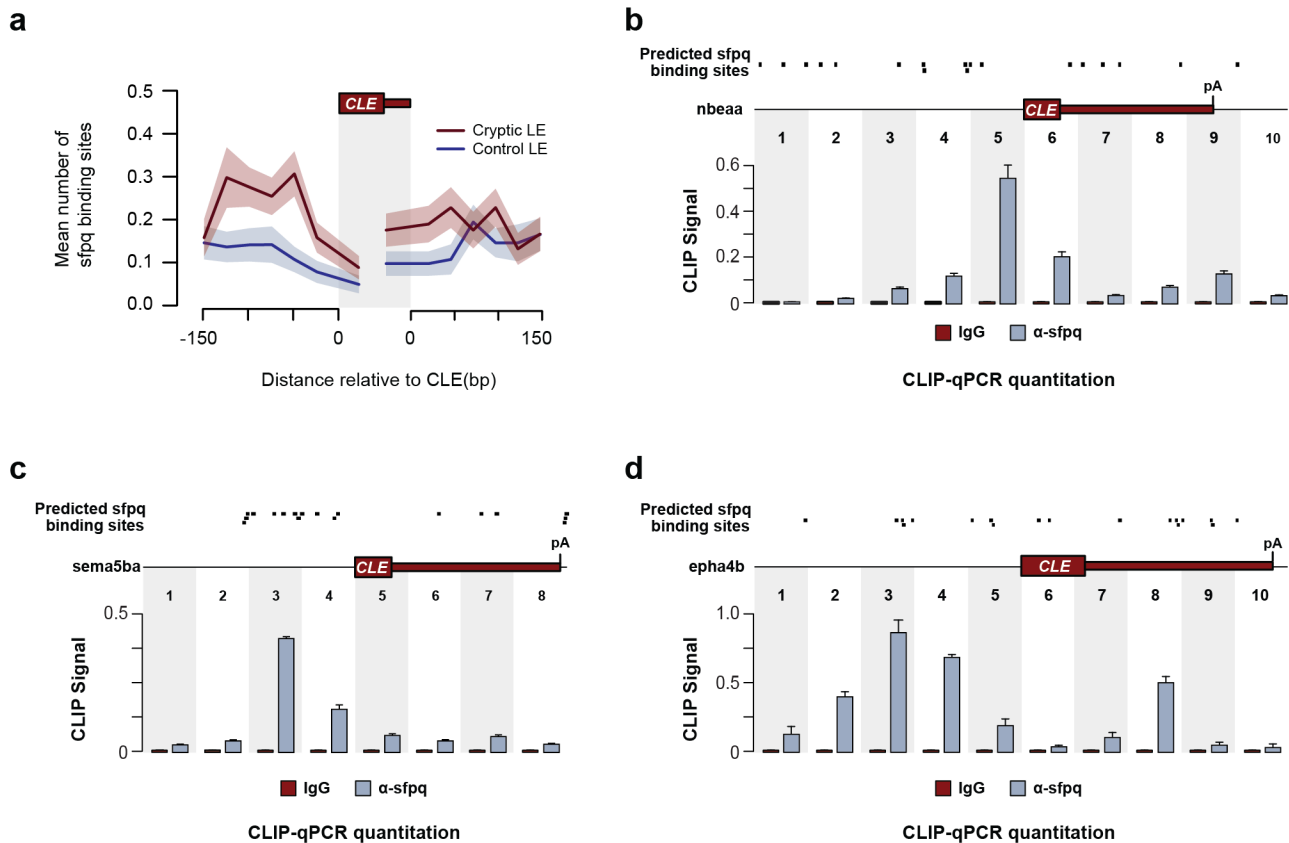175     **CLEs can dampen the expression of full-length transcripts**

7

**Fig. 3**

**Figure 3: SFPQ directly binds to RNA adjacent to CLE sequences**

**a**, Line plot showing the distribution of predicted SFPQ-binding sites surrounding CLEs (red) and constitutive last exons of each CLE-containing gene (blue). 200 bp of surrounding intron/CLE junction sequence (150 bp intron and 50bp exon) were binned into 50 bp windows and the mean number of predicted motifs were shown (± SEM).

**b-d**, Top: Location of SFPQ binding motifs predicted using MEME suite. Bottom: RT-qPCR quantitation showing the relative enrichment of SFPQ-interacting regions surrounding CLEs. Abundance of SFPQ- or IgG(control)-crosslinked RNAs were normalized to input and the mean value from three replicates were shown (± SD).

176    The reciprocal relationship between CLEs and the abundance of full-length

177    transcripts (Figure 1f) suggests that these exons may act as negative regulators of gene

178    expression.  If production of CLE transcripts is a mechanism for down-regulating the

179    normal full-length transcripts, then eliminating the cryptic exon in *sfpq*$^{-/-}$ mutants

180    should rescue their expression.  To test this possibility, we used the gene *b4galt2* as case

181    study, as it shows a very strong loss of expression of its three normal isoforms in the

182    mutant (Figure S1f). We used CRISPR/Cas9 to delete the *b4galt2* CLE, injecting Cas9

183    along with two guide RNAs that targeted directly upstream of the cryptic exon and at the

184    3' end of the exon (Figure 4a).  Injected founder embryos (crispants) will show

185    mosaicism, so a complete loss of the cryptic exon would not be expected in every cell of

186    the embryo.  Despite mosaicism, PCR analysis of the "crispants" showed a strong

187    deletion band for six out of eight tested embryos (Figure 4a).  Encouraged by the high

188    efficiency of the gRNAs, we performed RT-qPCR on pooled injected *sfpq*$^{-/-}$ embryos to

189    measure the expression levels of the normal *b4galt2* transcripts and saw a significant

190    rescue of the longer transcripts compared to the uninjected *sfpq*$^{-/-}$ control (Figure 4a).

191    This result did not hold true with two other CLEs we deleted (example *gdf11* CLE, Figure

192    S4a). We concluded that CLEs can regulate expression levels of at least some of the

193    genes containing these exons.

194

195    **Truncated protein derived from CLE-containing *epha4b* transcripts accounts for**

196    **the boundary defects in *sfpq*$^{-/-}$ brain**

197    In addition to affecting the expression levels of normal isoforms, CLE transcripts

198    could impact the *sfpq* phenotype through aberrant functions of the truncated RNAs or

199    the short peptides they produce.  We focused on the candidate gene *epha4b*, which

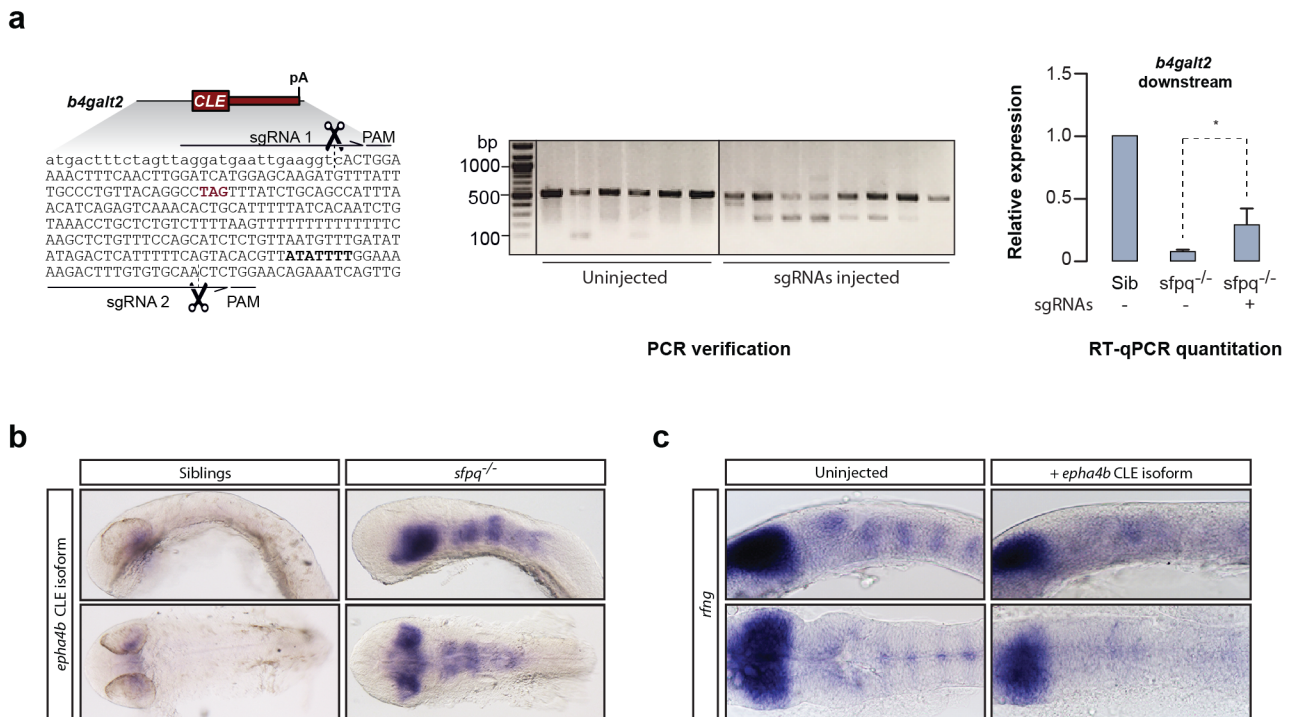200    expresses a CLE-containing short mRNA in *sfpq* null embryos, while showing no change

8

**Fig. 4 (a-c)**

**Figure 4 (a-c): CLE formation is functionally relevant**

**a**, Deletion of the b4galt2 CLE using CRISPR/Cas9 rescues expression of downstream exons. Left: cut sites of the b4galt2 sgRNAs. CLE is indicated by capital letters. Center: PCR verification of Cas9 cleavage after injection of sgRNAs. Right: RT-qPCR quantitation of the relative expression of the downstream b4galt2 exons in sfpq-/- embryos compared to siblings.

**b**, in-situ hybridization of the epha4b CLE at 24 hpf, displaying strong expression in the midbrain and hindbrain of sfpq-/- embryos.

**c**, in-situ hybridization of rfng shows rhombomere boundary defects at 22ss after injection of the epha4b cryptic transcript into WT embryos

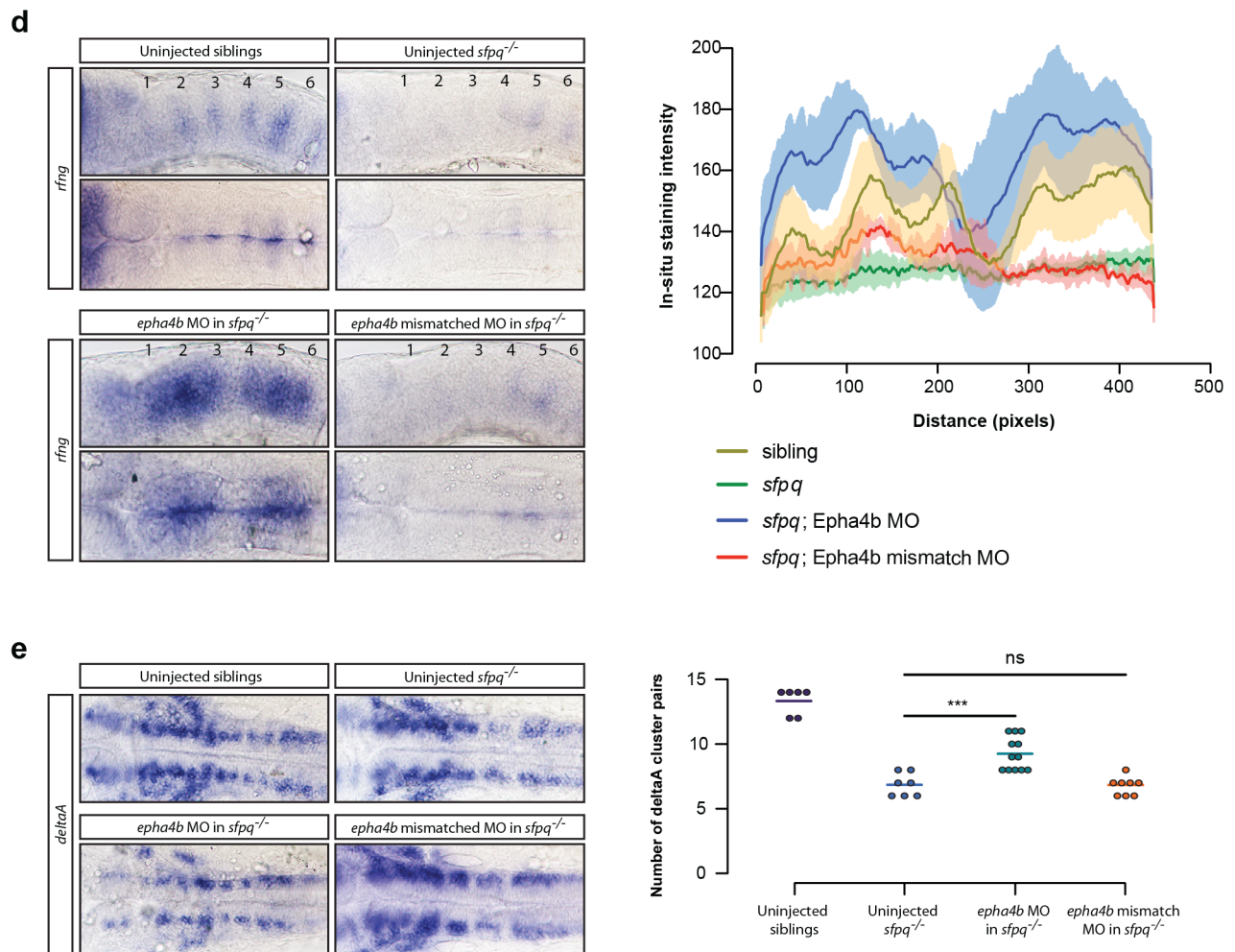**b-d**, Upper: lateral view. Lower: dorsal view.

Fig. 4 (d-e)

**Figure 4 (d-e): CLE formation is functionally relevant**

**d**, Left: in-situ hybridization of rfng shows rhombomere boundary defects of sfpq-/- embryos are rescued by injection of the epha4b cryptic splice junction morpholino but not a mismatch morpholino. Rhombomere boundaries are numbered. Right: quantification of staining in rhombomeres in three lateral view samples for each condition

**e**, Left: in-situ hybridization of DeltaA shows a loss of discrete neuronal clusters in sfpq-/- which is rescued by injection of the epha4b cryptic splice junction morpholino but not a mismatch morpholino. Right: Quantification of number of DeltaA clusters in each condition.

**b-d**, Upper: lateral view. Lower: dorsal view.

201    in expression of the normal transcripts (Figure S1c).  This gene is one of two zebrafish

202    paralogues of the human ALS-associated gene Epha4 (Van Hoecke *et al*, 2012; Wu *et al*,

203    2017), coding for a protein-tyrosine kinase of the Ephrin receptor family known to

204    regulate hindbrain boundary formation (Cooke *et al*, 2005; Kemp *et al*, 2009). Truncated

205    forms of EPH receptors have been shown to act as dominant negatives by competing

206    with full-length versions of the protein for ligand binding (Smith *et al*, 2004).  The

207    predicted peptide produced by the *epha4b* CLE-containing short transcript would

208    contain the ligand binding domain but not the transmembrane and intracellular

209    domains and thus would be predicted to be a dominant negative (Table S6).

210         To assess possible effects of the shortened *epha4b*, we first performed an *in-situ*

211    hybridization using a probe for the cryptic exon.  We found that in *sfpq*$^{-/-}$ embryos, but

212    not in siblings, the *epha4b* CLE was expressed strongly in the midbrain and hindbrain

213    (Figure 4b), where the gene is normally transcribed at that developmental stage.  We

214    then tested whether, in wildtype fish, injection of the CLE transcript would induce

215    defects in the midbrain or hindbrain.  Using the early hindbrain boundary marker *rfng*,

216    we found that injection of the short *epha4b* transcript did not affect formation of the

217    midbrain but did cause a loss of hindbrain rhombomere boundaries similar to that seen

218    in the *sfpq*$^{-/-}$ mutant (Figure 4c).

219         We then asked whether repressing the *epha4b* CLE in *sfpq*$^{-/-}$ embryos could

220    rescue the *sfpq* hindbrain defect.  We used a splice junction morpholino (MO) that

221    targeted the 3' splice acceptor site of the CLE to prevent the cryptic exon from being

222    used in *sfpq* mutants.  Although MOs frequently have off-target effects, those effects are

223    generally the opposite of what we would expect to see from a rescue (i.e. increased cell

224    death and off-target phenotypes, never rescue of phenotypes).  However, as MOs have

225    been shown to have some phenotypic effects on the hindbrain (Gerety & Wilkinson,

9

226    2011), we used mismatch controls to ensure that our results were specific to the *epha4b*

227    CLE splice-MO.  We tested the MO efficiency using RT-PCR with primers both within the

228    cryptic exon and across the exon junction (Figure S4b).  We then examined the effects of

229    the MO on hindbrain development using both the boundary-specific *rfng* marker (Figure

230    4d) and the pan-neuronal marker DeltaA (Figure 4e).  We saw that the CLE splice

231    junction MO, but not the mismatch control, rescued formation of rhombomere

232    boundaries in *sfpq*$^{-/-}$ mutants.  Taken together, these results indicate that the hindbrain

233    boundary defect in *sfpq*$^{-/-}$ embryos can be explained by the dominant-negative effects of

234    the *epha4b* CLE transcript.

235

236    **Repression of CLEs by SFPQ is conserved across vertebrates and relevant to**

237    **human neuropathologies.**

238          As our analysis of the *sfpq* loss-of-function phenotype was performed solely in

239    zebrafish, we wondered whether SFPQ repressed CLEs in other organisms.  Accordingly,

240    we turned to publicly available RNA-seq datasets from *sfpq* loss-of-function

241    experiments.  A conditional mouse knock-out model (Takeuchi *et al*, 2018) inactivated

242    *Sfpq* in the cerebral cortex.  Examining mouse orthologs of zebrafish CLE-containing

243    genes, we were able to identify CLE formation in mouse *Sfpq*-null brains for *Epha4b*,

244    *Cpped1*, *Fam172a*, and *Exoc4* (Figure 5a and S5).  Overall, we identified 144 instances of

245    upregulation of CLEs in the cortical *Sfpq* knockout (Table S7).  Examination of the CLE-

246    containing introns showed results similar to those for the zebrafish CLEs: the CLE-

247    containing introns have a bias towards appearing earlier in the gene, they are often

248    embedded within the largest intron in a gene, and their host introns are significantly

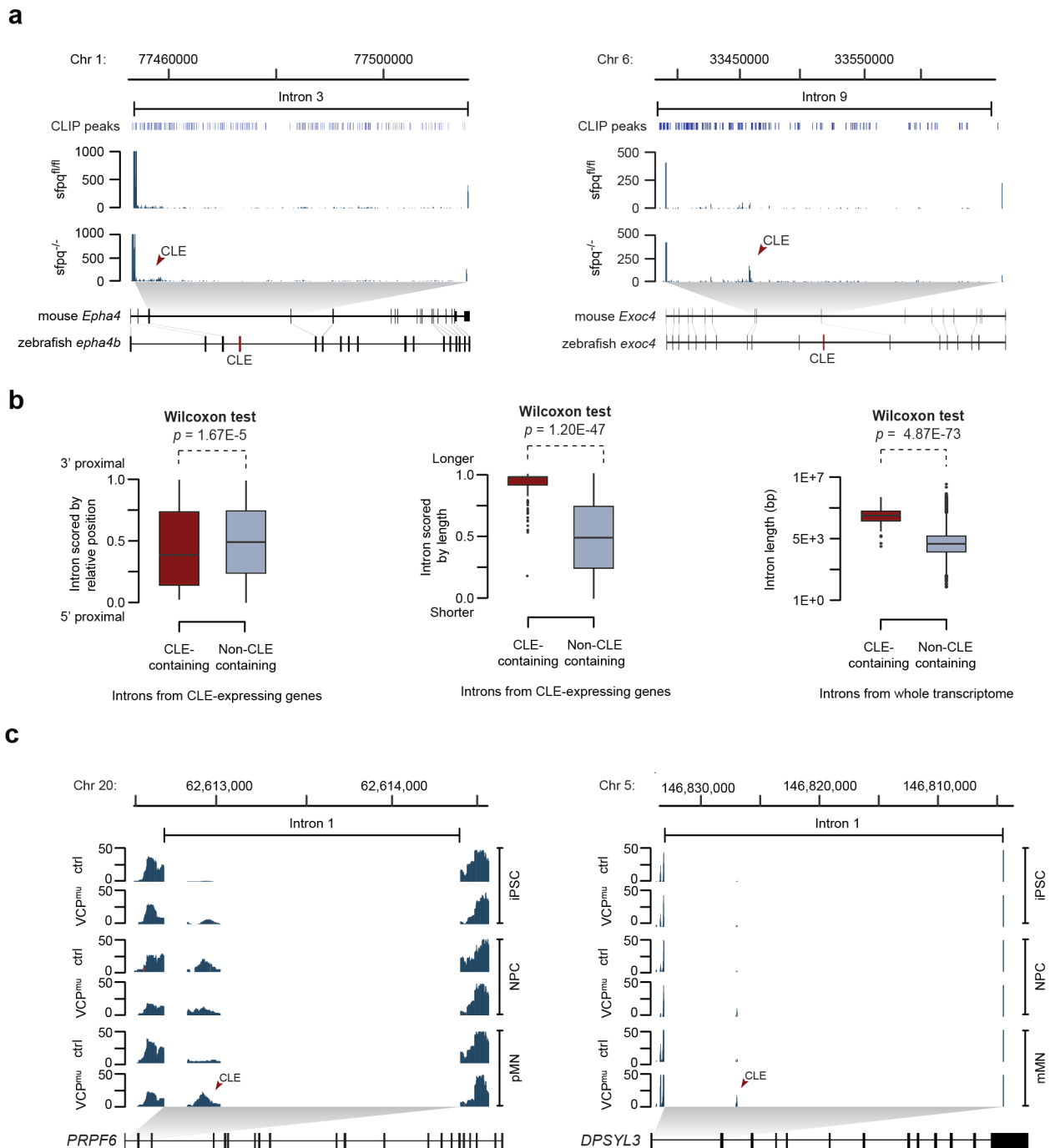249    larger than the average mouse intron (Figure 5b).

**Fig. 5**

**Figure 5: The CLE-repressing function of SFPQ is conserved in mouse and human**

**a**, Meta-analysis of RNA-seq and CLIP-seq dataset from conditional Sfpq knockout mice for cryptic last exons. Top: distribution of Sfpq CLIP peaks within the CLE-containing intron. Middle: tracks showing read coverage plots and "sashimi" plots from Sfpqfl/fl and Sfpq-/- mice. Bottom: exon architecture of orthologous CLE-expressing genes. Homologous regions between orthologues are shown as connecting lines.

**b**, Introns from mouse CLE-expressing genes were scored by its relative position (left) and by its relative length (mid), and the distribution of these scores were plotted. Note that introns containing CLE tend to be long and sparsely distributed. Right: CLE-containing introns are longer than average introns. Length of CLE-containing introns is compared to all other introns from the mouse transcriptome.

**c**, Representative RNA-seq coverage plots from ALS-derived iPSC dataset of CLEs up-regulated in VCPmu samples.

250    As SFPQ has been recently linked with ALS in human, we also examined RNA-seq

251    results from iPSCs derived from ALS patients, which show loss of nuclear SFPQ

252    expression (Luisier *et al*, 2018). In total, we found 76 CLE events up-regulated in ALS-

253    mutant backgrounds across the neuronal differentiation stages (Table S8). This is

254    probably an underestimation since the sequencing depth in this dataset was somewhat

255    lower than that in the mouse knockout study. Interestingly, CLEs spliced from PRPF6

256    and DPYSL3 genes showed consistent up-regulation in three time-points (Figure 5c).

257    The latter gene is involved in positive regulation of axon guidance and genetic variants

258    of this gene have been previously implicated in ALS patients (Blasco *et al*, 2013). These

259    results indicate that CLE repression is a conserved function of SFPQ, and that CLE-

260    dependent short transcripts may have a substantial impact on SFPQ-mediated disease

261    states.

262

263

**Discussion**

264

265 Our study uncovers a critical role of SFPQ in repression of cryptic last exons

266 (CLEs). We show that truncated transcripts appearing as a result of increased use of

267 CLEs are functionally relevant both as regulators of gene expression output and as a

268 source of interfering protein isoforms. Moreover, the CLE-repressing function of SFPQ is

269 conserved in mouse and human, indicating an important developmental role, with

270 implications for human pathology.

271

**Mechanism of CLE formation**

272

273 The presence of strong polyadenylation sites in CLE sequences suggests that the

274 paucity of CLE-containing isoforms under normal physiological conditions is due to

275 active suppression of CLE cleavage/polyadenylation or/and splicing. Our CLIP-seq

276 experiments provide evidence for SFPQ binding within or directly adjacent to CLEs.

277 Moreover, the bias of CLEs towards the 5' end of long introns is consistent with previous

278 analyses of SFPQ localization on RNA (Takeuchi *et al*, 2018). These data argue that SFPQ

279 may play a direct role in repressing cryptic exon formation. However, further work will

280 be required to distinguish between suppression of splicing versus blocking of the

281 polyadenylation site.

282 The relationship between SFPQ and CLEs extends our understanding of the

283 regulation possibilities afforded by long introns. Indeed, long introns have been

284 previously shown to control gene expression through interplay between premature

285 cleavage/polyadenylation and the U1 snRNP-dependent antitermination mechanism

286 known as "telescripting" (Langemeier *et al*, 2013; Oh *et al*, 2017; Venters *et al*, 2019;

287 Kainov & Makeyev, 2020). SFPQ-mediated CLE repression also operates in long introns

288 (Figure 2b) but, unlike telescripting, CLE involves definition of a last exon possibly via

12

289　　interactions between the U2 snRNP and U2AF with the cleavage/polyadenylation

290　　machinery (Martinson, 2011). Moreover, inactivation of U1 often promotes

291　　cleavage/polyadenylation relatively close to the 5' end of the gene, whereas CLEs do not

292　　show such a gene location bias.

293　　　　Long introns have been also shown to be subject to recursive splicing (RS), a

294　　multistep process promoting accuracy and efficiency of intron excision (Sibley *et al*,

295　　2015; Blazquez *et al*, 2018). Like CLEs, RS-sites appear primarily in long introns in

296　　genes with neuronal function. RS-sites initially produce an RS-exon that is spliced to the

297　　upstream exon prior to being excised at the subsequent round of splicing reactions.

298　　However, RS-exons do not contain polyadenylation sequences, so inclusion of the exon

299　　would not lead to truncation of the transcript. In addition, recursive splicing creates a

300　　stereotypical saw-tooth pattern of RNA-seq reads, which is not seen in the *sfpq*$^{-/-}$ RNA-

301　　seq data set. Therefore, SFPQ and CLEs provide a distinct regulation modality compared

302　　to telescripting and recursive splicing.

303

304　　**Pathology of cryptic transcripts**

305　　　　A notable feature of the SFPQ-repressed CLEs is the detrimental effect that they

306　　have on the function of their host genes. We previously showed that loss of *sfpq* leads to

307　　an array of morphological and neurodevelopmental abnormalities in zebrafish embryos,

308　　including loss of brain boundaries and altered motor axon morphology (Thomas-Jinu *et*

309　　*al*, 2017). However, the mechanism by which those abnormalities formed was

310　　unresolved. Here, we found that CLEs contribute to at least one aspect of the *sfpq*

311　　phenotype: the dominant negative *epha4b* truncated transcript induces hindbrain

312　　boundary defects. Moreover, a subset of the identified CLE-dependent short transcripts

313　　identified in *sfpq*$^{-/-}$ is predicted to affect axon growth and connectivity.

13

314       While CLE formation is clearly detectable under pathological conditions of loss of

315 *sfpq*, our data do not preclude the possibility of CLEs being expressed under non-

316 pathological conditions. Although the CLEs are not annotated in the current zebrafish,

317 mouse, and human genomes, it is possible that they may be regulated in a spatio-

318 temporal manner such that they only appear in specific tissues and/or at specific

319 developmental time points. Indeed, this possibility is supported by the relatively low

320 expression of SFPQ in non-neuronal tissue (Thomas-Jinu *et al*, 2017; Lowery *et al*, 2007),

321 and by low-level detection of the *epha4b* CLE transcript in siblings by PCR (Figure S4B).

322 Early termination of long pre-mRNAs has been shown to be a developmentally

323 controlled regulatory mechanism: the RNA-binding protein Sex-lethal promotes the

324 formation of truncated transcripts during short nuclear cycles in *Drosophila* (Sandler *et*

325 *al*, 2018), and downregulation of the cleavage and polyadenylation factor PCF11 during

326 differentiation of mouse C2C12 myoblast cells suppresses intronic polyadenylation to

327 promote long gene expression (Wang *et al*, 2019). Further examination of CLE

328 expression in wildtype animals across development may identify possible role of these

329 truncated transcripts in normal tissues.

330

331 **Cryptic exons in neurodegenerative disease**

332       Neurodegenerative diseases such as Alzheimer's, ALS, and FTD are frequently

333 characterized by altered localization and function of splicing factors (Tyzack *et al*, 2019;

334 Ling *et al*, 2013; Neumann *et al*, 2006; Nag *et al*, 2018). The ALS-associated proteins

335 transactivation response element DNA-binding protein 43 (TDP-43) and fused in

336 sarcoma (FUS) regulate alternative splicing and alternative polyadenylation (Ishigaki *et*

337 *al*, 2012; Masuda *et al*, 2016; Deshaies *et al*, 2018; Melamed *et al*, 2019; Klim *et al*, 2019;

338 Ling *et al*, 2015). TDP-43 has been shown to act as a repressor of cryptic exons, a

339  minority of which contain polyadenylation sites and thus would form CLEs (Ling *et al*,

340  2015).  Stathmin-2 is one of the latter and rescue of its normal full-length expression in

341  TDP-43-knockdown cell culture improves axonal growth in this model (Melamed *et al*,

342  2019; Klim *et al*, 2019), indicating that CLEs are pathogenic across various splicing

343  protein-dependent pathologies. These findings place CLEs at the center of priority for

344  understanding molecular mechanisms of neurodegenerative diseases and developing

345  new ways to diagnose and treat these increasingly prevalent disorders.

346

347

353

354  **Materials and Methods**

355  *Zebrafish husbandry*

356      Zebrafish (*Danio rerio*) were reared in accordance with the Animals (Scientific

357  Procedures) Act 1986.  Fish were maintained on a 14 hr light/10 hr dark cycle at 28°C.

358  Embryos were cultured in fish water containing 0.01% methylene blue to prevent fungal

359  growth.  Wildtype fish were AB strain from the Zebrafish International Resource Center

360  (ZIRC), while *sfpq* null mutants were *sfpq$^{kg41}$* (Thomas-Jinu *et al*, 2017).

361

362  *RNA-seq*

15

363    RNA was extracted from 24 hpf *sfpq⁻ᐟ⁻* embryos and their heterozygous or

364    homozygous wildtype siblings using the RNeasy Mini Kit (Qiagen). RNA was sequenced

365    using the Illumina HiSeq 2500 with 50bp paired-end reads.

366

367    *3' RACE*

368    RNA was extracted from 24 hpf *sfpq⁻ᐟ⁻* embryos using the RNease Mini Kit

369    (Qiagen). Reverse transcription was performed using the 3' RACE System for Rapid

370    Amplification of cDNA Ends (ThermoFisher). cDNA was amplified in two subsequent

371    PCR reactions, using the adapter primer as a reverse primer and the following primers

372    as forward primers:

373    *nbeaa:* AGAGAGGGACCGTGTAGAC, AAGGCACATCGAGCCCATATTG

374    *epha4b:* ATGGCAACCCTTTGGATTTATCCG, CTACTGTCAGGCTGTTTCGG

375    *bcas3:* CGCTGCATGTCAGCTTCAC, CTTCAGGAAACTGACAAACGCGAG

376    *gdf11:* GGGAGCATTATAGGCATCGGTAC , TGGCTTCAGAGCGAGTCATAG

377    *vti1a:* CGTCAATAAGAAGCAGACACAAGCAAC, GATTTTGTGGTCACATTTTCGTG

378    *immp2l:* CGACGCACAGCACGTACATAAG, GCGACATTGTGTCAGTTTTAACATC

379

380    *CLIP-qPCR*

381    Dechorionated 24 hpf wildtype fish were irradiated (twice at 0.8 J/cm² , 254nm)

382    and deyolked using high calcium Ringer's solution (116 mM NaCl, 2.9 mM KCl, 10 mM

383    CaCl2, 5 mM HEPES, pH 7.2) with 0.3 mM PMSF and 1 mM EDTA. After several washes,

384    embryos were lysed using PXL buffer (0.1% SDS, 0.5% deoxycholate, 0.5% NP-40) and

385    homogenized using a plastic pestle. Lysates were treated with 10 μL diluted RNAseI

386    (1:500 dilution; Thermo Fisher) and 2 μL Turbo DNase (Thermo Fisher) at 37°C for 3

387    minutes on a shaking incubator. Protein-RNA complexes were purified by centrifugation

16

388   and 5% of the lysate was retained as input. The remaining lysate were split into two and

389   its volume were topped up to 100 µL using PXL buffer. 100 µL of protein A Dynabeads

390   (Thermo Fisher) primed with either anti-SFPQ antibody (ab38148) or anti-IgG antibody

391   (MA5-14453) were added to each lysate and incubated at 4°C for an hour on a rotator.

392   Bound SFPQ-RNA complexes were purified and washed thrice in high salt wash buffer

393   (50 mM Tris-HCL, pH 7.4, 1M NaCl, 1mM EDTA, 1% Igepal, 0.1% SDS, 0.5% sodium

394   deoxycholate). Subsequently, bound complexes were washed twice in PNK wash buffer

395   (20 mM Tris-HCL, pH 7.4, 10 mM $MgCl_2$, 0.2% Tween-20) and followed by proteinase K

396   digestion (Thermo Fisher). Bound RNAs were purified using phenol-chloroform

397   extraction followed by reverse-transcription to generate cDNAs. Relative amounts of

398   SFPQ-bound RNAs was quantified by qPCR using primers: (Table S9).

399

400

401   *CRISPR/Cas9*

402       gRNAs were formed from chemically synthesized Alt-R®-modified crRNAs from

403   Integrated DNA Technologies (IDT).  Each crRNA was suspended in duplex buffer to

404   100µM concentration, then a crRNA:tracrRNA duplex was formed by combining 3µl

405   crRNA, 3µl 100µM tracrRNA, and 19 µl duplex buffer at 95°C for five minutes, then

406   cooled to room temperature and stored at -20°C.  To make gRNA:Cas9 RNP complexes, a

407   mix was formed as follows: 1.5 µl each gRNA, 0.75 µl 2M KCl, 1.25 µl EnGen Spy Cas9

408   NLS (NEB).  The mix was incubated at 37°C for five minutes, then brought to room

409   temperature.  One nanoliter of the gRNA:Cas9 complex was injected into embryos at the

410   1-cell stage.  The following gRNAs were used:

411   *b4galt2*: AAGGATGAATTGAAGGTCAC, AAAGACTTTGTGTGCAACTC

412   *gdf11*: GTAGAGAGTAGGTTCAGAGT, GACCAAATGTTGTTAGAAAG

413

414 *RNA and morpholino injections*

415       The *epha4b* cryptic transcript was amplified from cDNA and inserted into the

416 multi-cloning site of plasmid pCS2+ (Addgene). The *in-vitro* transcription reaction was

417 performed on linearized plasmid using the mMessage mMachine SP6 Transcription Kit

418 (ThermoFisher), and the RNA was purified using a Mini Quick Spin Column (Roche).

419 100 pg RNA was injected into the embryo at the one-cell stage.

420       For morpholino knockdown of the *epha4b* cryptic exon, embryos were injected

421 into the yolk at the one-cell stage with 0.1 pmol of Epha4b splice junction morpholino or

422 mismatch.

423 Epha4b splice junction morpholino: ACAGCTGAGAAAAAAACACGGATAT

424 Epha4b splice junction mismatch morpholino: ACAcCTcAGAAAtAAAgACcGATAT

425

426 *In-situ hybridization*

427       Linearized plasmids containing the antisense sequence for *rfng* (Cheng *et al*,

428 2004), *deltaA* (Allende & Weinberg, 1994), or the *epha4b* cryptic exon were transcribed

429 into RNA probes using DIG labeling mix (Roche) according to the manufacturer's

430 instructions. Probes were purified using Mini Quick Spin Columns (Roche). *In-situ*

431 hybridization reaction was performed as described elsewhere (Thomas-Jinu & Houart,

432 2013).

433

434 *qPCR*

435       RNA was extracted from 24-28 hpf *sfpq-/-* embryos and heterozygous or WT

436 siblings using the RNease Mini Kit (Qiagen). 1 ug of extracted RNA was used in a reverse

437 transcriptase reaction using the Superscript III First Strand cDNA Synthesis Kit

18

438    (Invitrogen).  250 ng of cDNA was used in qPCR reactions with the LightCycler 480 SYBR

439    Green I Master Mix (Roche).  Each sample was compared against a B-actin control

440    reaction.

441

442    *Bioinformatics*

443         For analyses of 24 hpf *sfpq⁻/⁻* RNA-seq data using Cufflinks package (Trapnell *et*

444    *al*, 2012), reads were mapped to zebrafish GRCz9 assembly and differential expression

445    analysis were carried out using default settings.

446

447         For analyses of 24 hpf *sfpq⁻/⁻* RNA-seq data using Whippet pipeline (Sterne-

448    Weiler *et al*, 2018), a GRCz10 Ensembl-based index was generated using Whippet's

449    index      building      function      from      the      Ensembl-based      fasta

450    (ftp://ftp.ensembl.org/pub/release-

451    91/fasta/danio_rerio/dna/Danio_rerio.GRCz10.dna.toplevel.fa.gz) and gene annotation

452    files                                                    (ftp://ftp.ensembl.org/pub/release-

453    91/gtf/danio_rerio/Danio_rerio.GRCz10.91.gtf.gz). Quantification of aligned RNA-seq

454    reads were done as follows:

$$whippet-quant.jl\ fwd_{file}.fastq.gz\ rev_{file}.fastq.gz\ --biascorrect-x\ index_{graph}.jls$$

$$-o < output\_directory > --sam < SAM\_output\_directory >$$

455

456         The above quantification function outputs several tables containing read counts

457    at the gene and isoform level. Differential gene and isoform expression analyses were

458    identified  using  the  edgeR  package  with  the  estimateGLMRobustDisp  function

459    (Robinson *et al*, 2010). Differential splicing events were identified using Whippet's delta

460    analysis function with default parameters. An event with a "Probability" score exceeding

19

461    80% is classified as significantly regulated. Cryptic splicing events were annotated using

462    custom R-scripts.

463

464        For analyses of conditional Sfpq knock-out mouse model (Takeuchi *et al*, 2018)

465    dataset, the above Whippet pipeline was carried out using Ensembl's GRCm38 fasta

466    (ftp://ftp.ensembl.org/pub/release-

467    99/fasta/mus_musculus/dna/Mus_musculus.GRCm38.dna.toplevel.fa.gz        )        and

468    annotation                                        (ftp://ftp.ensembl.org/pub/release-

469    99/gtf/mus_musculus/Mus_musculus.GRCm38.99.gtf.gz)    files.      For    analyses    of

470    conditional ALS-derived iPSC differentiation dataset (Luisier *et al*, 2018), the above

471    Whippet    pipeline    was    carried    out    using    Ensembl's    GRCh37    fasta

472    (ftp://ftp.ensembl.org/pub/grch37/current/fasta/homo_sapiens/dna/Homo_sapiens.G

473    RCh37.dna.toplevel.fa.gz)    and    annotation    (ftp://ftp.ensembl.org/pub/release-

474    75/gtf/homo_sapiens/Homo_sapiens.GRCh37.75.gtf.gz) files.

475

476        To construct CLE-containing transcripts, read alignments from Whippet were

477    sorted , indexed and assembled using the StringTie program (Kovaka *et al*, 2019).

478    Ensembl's GRCz10 transcriptome was used as reference and assembly was done for

479    each biological replicate as follows:

$$stringtie < file1.bam > -p < num\_threads > -o < file1.gtf > -G < reference >$$

480

481        Assembled transcripts from each sample were subsequently combined using

482    StringTie's merge function using GRCz10 annotations as reference. CLE-containing

483    isoforms were identified by intersecting exon coordinates from the merged transcript

484    assembly with CLE coordinates from Whippet delta analysis output. Intersection

20

485    operation was done in R using Bioconductor's GenomicRanges package (Lawrence *et al*,

486    2013).  Analyses on the coding potential of CLE isoforms and its functional loss of

487    protein domains were carried out using custom R-scripts.

488

489        For the analyses of introns from which the CLEs were spliced from, intronic

490    features were extracted from the custom-assembled transcript in R using

491    Bioconductor's GenomicFeatures package (Lawrence *et al*, 2013). A list of the largest,

492    non-overlapping introns was generated and annotated for an overlap with a CLE

493    segment using GenomicRanges' reduce and subsetByOverlaps functions respectively.

494    The relative position of CLEs within its intron was determined using psetdiff operation

495    followed by extracting the width of the upstream intronic segment.

496

497        For the analyses of CLE conservation, 8-way PhastCons data were downloaded

498    from                                                                      UCSC

499    (http://hgdownload.soe.ucsc.edu/goldenPath/danRer7/phastCons8way/fish.phastCons

500    8way.bw). Coordinates of CLE containing intronswere converted to GRCv9 using UCSC's

501    LiftOver function and binned into 1 kb sequence using a sliding window technique (1 bp

502    steps). Average PhastCons score of each bin was calculated using bedtools' "map"

503    function and bins containing CLE were annotated through intersection. Conservation

504    scores of each CLE and 250 nt of its flanking introns were calculated using the same tool.

505    To refine the conservation regions surrounding the intron-CLE borders, average

506    PhastCons score were calculated for 10 nt windows including 30 nt of each exonic ends.

507

508        For the analyses of SFPQ binding motifs within sequences surrounding CLEs, its

509    Position-Specific      Scoring      Matrix      was      downloaded      from      RBPmap

21

510 (http://rbpmap.technion.ac.il/download.html) and manually converted into a MEME

511 motif format (http://meme-suite.org/doc/meme-format.html). Occurrence of SPFQ

512 binding sites was analyzed using MEME's FIMO program (http://meme-

513 suite.org/doc/fimo.html) using the following parameters:

$$fimo --thresh\ 0.005 --o < output\_directory > < SFPQ\_PSSM > < CLE\_fasta >$$

514 The average number of SFPQ binding motifs were calculated for 25 nt windows of

515 flanking intronic sequence including 25 nt of each exonic ends.

516

517

518 **References**

519 Allende ML & Weinberg ES (1994) The expression pattern of two zebrafish achaete-
520     scute homolog (ash) genes is altered in the embryonic brain of the cyclops mutant.
521     *Dev. Biol.* **166:** 509–530
522 Blasco H, Bernard-Marissal N, Vourc'h P, Guettard YO, Sunyach C, Augereau O,
523     Khederchah J, Mouzat K, Antar C, Gordon PH, Veyrat-Durebex C, Besson G, Andersen
524     PM, Salachas F, Meininger V, Camu W, Pettmann B, Andres CR & Corcia P (2013) A
525     Rare Motor Neuron Deleterious Missense Mutation in the *DPYSL3* ( *CRMP4* ) Gene is
526     Associated with ALS. *Hum. Mutat.* **34:** 953–960
527 Blazquez L, Emmett W, Faraway R, Pineda JMB, Bajew S, Gohr A, Haberman N, Sibley CR,
528     Bradley RK, Irimia M & Ule J (2018) Exon Junction Complex Shapes the
529     Transcriptome by Repressing Recursive Splicing. *Mol. Cell* **72:** 496-509.e9
530 Bottini S, Hamouda-Tekaya N, Mategot R, Zaragosi LE, Audebert S, Pisano S, Grandjean V,
531     Mauduit C, Benahmed M, Barbry P, Repetto E & Trabucchi M (2017) Post-
532     transcriptional gene silencing mediated by microRNAs is controlled by
533     nucleoplasmic Sfpq. *Nat. Commun.* **8:** 1189
534 Cagnetta R, Frese CK, Shigeoka T, Krijgsveld J & Holt CE (2018) Rapid Cue-Specific
535     Remodeling of the Nascent Axonal Proteome. *Neuron* **99:** 29-46.e4
536 Cheng Y-C, Amoyel M, Qiu X, Jiang Y-J, Xu Q & Wilkinson DG (2004) Notch Activation
537     Regulates the Segregation and Differentiation of Rhombomere Boundary Cells in
538     the Zebrafish Hindbrain. *Dev. Cell* **6:** 539–550
539 Ciolli Mattioli C, Rom A, Franke V, Imami K, Arrey G, Terne M, Woehler A, Akalin A,
540     Ulitsky I & Chekulaeva M (2019) Alternative 3' UTRs direct localization of
541     functionally diverse protein isoforms in neuronal compartments. *Nucleic Acids Res.*
542     **47:** 2560–2573
543 Cooke JE, Kemp HA & Moens CB (2005) EphA4 Is Required for Cell Adhesion and
544     Rhombomere-Boundary Formation in the Zebrafish. *Curr. Biol.* **15:** 536–542
545 Cosker KE, Fenstermacher SJ, Pazyra-Murphy MF, Elliott HL & Segal RA (2016) The RNA-
546     binding protein SFPQ orchestrates an RNA regulon to promote axon viability. *Nat.*
547     *Neurosci.* **19:** 690–696

548    Deshaies J-E, Shkreta L, Moszczynski AJ, Sidibé H, Semmler S, Fouillen A, Bennett ER,
549         Bekenstein U, Destroismaisons L, Toutant J, Delmotte Q, Volkening K, Stabile S,
550         Aulas A, Khalfallah Y, Soreq H, Nanci A, Strong MJ, Chabot B & Vande Velde C (2018)
551         TDP-43 regulates the alternative splicing of hnRNP A1 to yield an aggregation-
552         prone variant in amyotrophic lateral sclerosis. *Brain* **141:** 1320–1333
553    Dye BT & Patton JG (2001) An RNA recognition motif (RRM) is required for the
554         localization of PTB-associated splicing factor (PSF) to subnuclear speckles. *Exp. Cell*
555         *Res.* **263:** 131–144
556    Furlanis E, Traunmüller L, Fucile G & Scheiffele P (2019) Landscape of ribosome-
557         engaged transcript isoforms reveals extensive neuronal-cell-class-specific
558         alternative splicing programs. *Nat. Neurosci.* **22:** 1709–1717
559    Gerety SS & Wilkinson DG (2011) Morpholino artifacts provide pitfalls and reveal a
560         novel role for pro-apoptotic genes in hindbrain boundary development. *Dev. Biol.*
561         **350:** 279–289
562    Guvenek A & Tian B (2018) Analysis of alternative cleavage and polyadenylation in
563         mature and differentiating neurons using RNA-seq data. *Quant. Biol.* **6:** 253–266
564    Hall-Pogar T, Liang S, Hague LK & Lutz CS (2007) Specific trans-acting proteins interact
565         with auxiliary RNA polyadenylation elements in the COX-2 3′-UTR. *RNA* **13:** 1103–
566         1115
567    Hanus C & Schuman EM (2013) Proteostasis in complex dendrites. *Nat. Rev. Neurosci.*
568         **14:** 638–648
569    Heyd F & Lynch KW (2010) Phosphorylation-dependent regulation of PSF by GSK3
570         controls CD45 alternative splicing. *Mol. Cell* **40:** 126–137
571    Van Hoecke A, Schoonaert L, Lemmens R, Timmers M, Staats KA, Laird AS, Peeters E,
572         Philips T, Goris A, Dubois B, Andersen PM, Al-Chalabi A, Thijs V, Turnley AM, van
573         Vught PW, Veldink JH, Hardiman O, Van Den Bosch L, Gonzalez-Perez P, Van Damme
574         P, et al (2012) EPHA4 is a disease modifier of amyotrophic lateral sclerosis in
575         animal models and in humans. *Nat. Med.* **18:** 1418–1422
576    Holt CE & Schuman EM (2013) The central dogma decentralized: New perspectives on
577         RNA function and local translation in neurons. *Neuron* **80:** 648–657
578    Iijima Y, Tanaka M, Suzuki S, Hauser D, Tanaka M, Okada C, Ito M, Ayukawa N, Sato Y,
579         Ohtsuka M, Scheiffele P & Iijima T (2019) SAM68-specific splicing is required for
580         proper selection of alternative 3'UTR isoforms in the nervous system. *ISCIENCE*
581    Ishigaki S, Fujioka Y, Okada Y, Riku Y, Udagawa T, Honda D, Yokoi S, Endo K, Ikenaka K,
582         Takagi S, Iguchi Y, Sahara N, Takashima A, Okano H, Yoshida M, Warita H, Aoki M,
583         Watanabe H, Okado H, Katsuno M, et al (2017) Altered Tau Isoform Ratio Caused by
584         Loss of FUS and SFPQ Function Leads to FTLD-like Phenotypes. *Cell Rep.* **18:** 1118–
585         1131
586    Ishigaki S, Masuda A, Fujioka Y, Iguchi Y, Katsuno M, Shibata A, Urano F, Sobue G & Ohno
587         K (2012) Position-dependent FUS-RNA interactions regulate alternative splicing
588         events and transcriptions. *Sci. Rep.* **2:**
589    Kainov YA & Makeyev E V. (2020) A transcriptome-wide antitermination mechanism
590         sustaining identity of embryonic stem cells. *Nat. Commun.* **11:** 1–18
591    Ke Y, Dramiga J, Schütz U, Kril JJ, Ittner LM, Schröder H & Götz J (2012) Tau-mediated
592         nuclear depletion and cytoplasmic accumulation of SFPQ in Alzheimer's and Pick's
593         disease. *PLoS One* **7:**
594    Kemp HA, Cooke JE & Moens CB (2009) EphA4 and EfnB2a maintain rhombomere
595         coherence by independently regulating intercalation of progenitor cells in the
596         zebrafish neural keel. *Dev. Biol.* **327:** 313–326

597 Kim KK, Kim YC, Adelstein RS & Kawamoto S (2011) Fox-3 and PSF interact to activate
598      neural cell-specific alternative splicing. *Nucleic Acids Res.* **39:** 3064–3078

599 Klim JR, Williams LA, Limone F, Guerra San Juan I, Davis-Dusenbery BN, Mordes DA,
600      Burberry A, Steinbaugh MJ, Gamage KK, Kirchner R, Moccia R, Cassel SH, Chen K,
601      Wainger BJ, Woolf CJ & Eggan K (2019) ALS-implicated protein TDP-43 sustains
602      levels of STMN2, a mediator of motor neuron growth and repair. *Nat. Neurosci.* **22:**
603      167–179

604 Knott GJ, Bond CS & Fox AH (2016) The DBHS proteins SFPQ, NONO and PSPC1: A
605      multipurpose molecular scaffold. *Nucleic Acids Res.* **44:** 3989–4004

606 Kovaka S, Zimin A V., Pertea GM, Razaghi R, Salzberg SL & Pertea M (2019)
607      Transcriptome assembly from long-read RNA-seq alignments with StringTie2.
608      *Genome Biol.* **20:** 278

609 Langemeier J, Radtke M & Bohne J (2013) U1 snRNP-mediated poly(A) site suppression:
610      Beneficial and deleterious for mRNA fate. *RNA Biol.* **10:** 180–184

611 Lawrence M, Huber W, Pagès H, Aboyoun P, Carlson M, Gentleman R, Morgan MT &
612      Carey VJ (2013) Software for Computing and Annotating Genomic Ranges. *PLoS*
613      *Comput. Biol.* **9:** e1003118

614 Ling JP, Pletnikova O, Troncoso JC & Wong PC (2015) TDP-43 repression of
615      nonconserved cryptic exons is compromised in ALS-FTD. *Science (80-. ).* **349:** 650–
616      655

617 Ling S-C, Polymenidou M & Cleveland DW (2013) Converging Mechanisms in ALS and
618      FTD: Disrupted RNA and Protein Homeostasis. *Neuron* **79:** 416–438

619 Lowery LA, Rubin J & Sive H (2007) whitesnake/sfpq is required for cell survival and
620      neuronal development in the zebrafish. *Dev. Dyn.* **236:** 1347–1357

621 Lu J, Shu R & Zhu Y (2018) Dysregulation and Dislocation of SFPQ Disturbed DNA
622      Organization in Alzheimer's Disease and Frontotemporal Dementia. *J. Alzheimer's*
623      *Dis.* **61:** 1311–1321

624 Luisier R, Tyzack GE, Hall CE, Mitchell JS, Devine H, Taha DM, Malik B, Meyer I,
625      Greensmith L, Newcombe J, Ule J, Luscombe NM & Patani R (2018) Intron retention
626      and nuclear loss of SFPQ are molecular hallmarks of ALS. *Nat. Commun.* **9:** 2010

627 Martinson HG (2011) An active role for splicing in 3′-end formation. *Wiley Interdiscip.*
628      *Rev. RNA* **2:** 459–470

629 Masuda A, Takeda J & Ohno K (2016) FUS-mediated regulation of alternative RNA
630      processing in neurons: insights from global transcriptome analysis. *Wiley*
631      *Interdiscip. Rev. RNA* **7:** 330–340

632 Mauger O, Lemoine F & Scheiffele P (2016) Targeted Intron Retention and Excision for
633      Rapid Gene Regulation in Response to Neuronal Activity. *Neuron* **92:** 1266–1278

634 Melamed Z, López-Erauskin J, Baughn MW, Zhang O, Drenner K, Sun Y, Freyermuth F,
635      McMahon MA, Beccari MS, Artates JW, Ohkubo T, Rodriguez M, Lin N, Wu D, Bennett
636      CF, Rigo F, Da Cruz S, Ravits J, Lagier-Tourenne C & Cleveland DW (2019)
637      Premature polyadenylation-mediated loss of stathmin-2 is a hallmark of TDP-43-
638      dependent neurodegeneration. *Nat. Neurosci.* **22:** 180–190

639 Mora Gallardo C, Sánchez de Diego A, Gutiérrez Hernández J, Talavera-Gutiérrez A,
640      Fischer T, Martínez-A C & van Wely KHM (2019) Dido3-dependent SFPQ
641      recruitment maintains efficiency in mammalian alternative splicing. *Nucleic Acids*
642      *Res.***:** 1–14

643 Nag S, Yu L, Boyle PA, Leurgans SE, Bennett DA & Schneider JA (2018) TDP-43 pathology
644      in anterior temporal pole cortex in aging and Alzheimer's disease. *Acta Neuropathol.*
645      *Commun.* **6:** 33

646  Neumann M, Sampathu DM, Kwong LK, Truax AC, Micsenyi MC, Chou TT, Bruce J, Schuck
647      T, Grossman M, Clark CM, McCluskey LF, Miller BL, Masliah E, Mackenzie IR,
648      Feldman H, Feiden W, Kretzschmar HA, Trojanowski JQ & Lee VM-Y (2006)
649      Ubiquitinated TDP-43 in Frontotemporal Lobar Degeneration and Amyotrophic
650      Lateral Sclerosis. *Science (80-. ).* **314:** 130–133
651  Oh JM, Di C, Venters CC, Guo J, Arai C, So BR, Pinto AM, Zhang Z, Wan L, Younis I &
652      Dreyfuss G (2017) U1 snRNP telescripting regulates a size-function-stratified
653      human genome. *Nat. Struct. Mol. Biol.* **24:** 993–999
654  Patton JG, Porro EB, Galceran J, Tempst P & Nadal-Ginard B (1993) Cloning and
655      characterization of PSF, a novel pre-mRNA splicing factor. *Genes Dev.* **7:** 393–406
656  Ray D, Kazan H, Cook KB, Weirauch MT, Najafabadi HS, Li X, Gueroussov S, Albu M,
657      Zheng H, Yang A, Na H, Irimia M, Matzat LH, Dale RK, Smith SA, Yarosh CA, Kelly SM,
658      Nabet B, Mecenas D, Li W, et al (2013) A compendium of RNA-binding motifs for
659      decoding gene regulation. *Nature* **499:** 172–177
660  Ray P, Kar A, Fushimi K, Havlioglu N, Chen X & Wu JY (2011) PSF suppresses tau exon 10
661      inclusion by interacting with a stem-loop structure downstream of exon 10. In
662      *Journal of Molecular Neuroscience* pp 453–466. Humana Press Inc
663  Robinson MD, McCarthy DJ & Smyth GK (2010) edgeR: a Bioconductor package for
664      differential expression analysis of digital gene expression data. *Bioinformatics* **26:**
665      139–140
666  Rosonina E, Ip JYY, Calarco JA, Bakowski MA, Emili A, McCracken S, Tucker P, Ingles CJ &
667      Blencowe BJ (2005) Role for PSF in Mediating Transcriptional Activator-Dependent
668      Stimulation of Pre-mRNA Processing In Vivo. *Mol. Cell. Biol.* **25:** 6734–6746
669  Sandler JE, Irizarry J, Stepanik V, Dunipace L, Amrhein H & Stathopoulos A (2018) A
670      Developmental Program Truncates Long Transcripts to Temporally Regulate Cell
671      Signaling. *Dev. Cell* **47:** 773-784.e6
672  Saud K, Cánovas J, Lopez CI, Berndt FA, López E, Maass JC, Barriga A & Kukuljan M
673      (2017) SFPQ associates to LSD1 and regulates the migration of newborn pyramidal
674      neurons in the developing cerebral cortex. *Int. J. Dev. Neurosci.* **57:** 1–11
675  Shi Y, Di Giammartino DC, Taylor D, Sarkeshik A, Rice WJ, Yates JR, Frank J & Manley JL
676      (2009) Molecular Architecture of the Human Pre-mRNA 3′ Processing Complex.
677      *Mol. Cell* **33:** 365–376
678  Shi Y & Manley JL (2015) The end of the message: Multiple protein–RNA interactions
679      define the mRNA polyadenylation site. *Genes Dev.* **29:** 889–897
680  Sibley CR, Emmett W, Blazquez L, Faro A, Haberman N, Briese M, Trabzuni D, Ryten M,
681      Weale ME, Hardy J, Modic M, Curk T, Wilson SW, Plagnol V & Ule J (2015) Recursive
682      splicing in long vertebrate genes. *Nature* **521:** 371–375
683  Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, Clawson H, Spieth
684      J, Hillier LDW, Richards S, Weinstock GM, Wilson RK, Gibbs RA, Kent WJ, Miller W &
685      Haussler D (2005) Evolutionarily conserved elements in vertebrate, insect, worm,
686      and yeast genomes. *Genome Res.* **15:** 1034–1050
687  Smith A, Robinson V, Patel K & Wilkinson DG (2004) The EphA4 and EphB1 receptor
688      tyrosine kinases and ephrin-B2 ligand regulate targeted migration of branchial
689      neural crest cells. *Curr. Biol.* **7:** 561–570
690  Sterne-Weiler T, Weatheritt RJ, Best AJ, Ha KCH & Blencowe BJ (2018) Efficient and
691      Accurate Quantitative Profiling of Alternative Splicing Patterns of Any Complexity
692      on a Laptop. *Mol. Cell* **72:** 187-200.e6
693  Takeuchi A, Iida K, Tsubota T, Hosokawa M, Denawa M, Brown JB, Ninomiya K, Ito M,
694      Kimura H, Abe T, Kiyonari H, Ohno K & Hagiwara M (2018) Loss of Sfpq Causes

695   Long-Gene Transcriptopathy in the Brain. *Cell Rep.* **23:** 1326–1341

696 Taliaferro JM, Vidaki M, Oliveira R, Olson S, Zhan L, Saxena T, Wang ET, Graveley BR,
697   Gertler FB, Swanson MS & Burge CB (2016) Distal Alternative Last Exons Localize
698   mRNAs to Neural Projections. *Mol. Cell* **61:** 821–833

699 Thomas-Jinu S, Gordon PM, Fielding T, Taylor R, Smith BN, Snowden V, Blanc E, Vance C,
700   Topp S, Wong CH, Bielen H, Williams KL, McCann EP, Nicholson GA, Pan-Vazquez A,
701   Fox AH, Bond CS, Talbot WS, Blair IP, Shaw CE, et al (2017) Non-nuclear Pool of
702   Splicing Factor SFPQ Regulates Axonal Transcripts Required for Normal Motor
703   Development. *Neuron* **94:** 322-336.e5

704 Thomas-Jinu S & Houart C (2013) Dynamic expression of neurexophilin1 during
705   zebrafish embryonic development. *Gene Expr. Patterns* **13:** 395–401

706 Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, Pimentel H, Salzberg SL, Rinn JL
707   & Pachter L (2012) Differential gene and transcript expression analysis of RNA-seq
708   experiments with TopHat and Cufflinks. *Nat. Protoc.* **7:** 562–578

709 Traunmüller L, Gomez AM, Nguyen TM & Scheiffele P (2016) Control of neuronal
710   synapse specification by a highly dedicated alternative splicing program. *Science*
711   *(80-. ).* **352:** 982–986

712 Tushev G, Glock C, Heumüller M, Biever A, Jovanovic M & Schuman EM (2018)
713   Alternative 3′ UTRs Modify the Localization, Regulatory Potential, Stability, and
714   Plasticity of mRNAs in Neuronal Compartments. *Neuron* **98:** 495-511.e6

715 Tyzack GE, Luisier R, Taha DM, Neeves J, Modic M, Mitchell JS, Meyer I, Greensmith L,
716   Newcombe J, Ule J, Luscombe NM & Patani R (2019) Widespread FUS
717   mislocalization is a molecular hallmark of amyotrophic lateral sclerosis. *Brain*

718 Venters CC, Oh J-M, Di C, So BR & Dreyfuss G (2019) U1 snRNP Telescripting:
719   Suppression of Premature Transcription Termination in Introns as a New Layer of
720   Gene Regulation. *cshperspectives.cshlp.org* **11:** a032235

721 Wang G, Yang H, Yan S, Wang C-E, Liu X, Zhao B, Ouyang Z, Yin P, Liu Z, Zhao Y, Liu T, Fan
722   N, Guo L, Li S, Li X-J & Lai L (2015) Cytoplasmic mislocalization of RNA splicing
723   factors and aberrant neuronal gene splicing in TDP-43 transgenic pig brain. *Mol.*
724   *Neurodegener.* **10:** 42

725 Wang R, Zheng D, Wei L, Ding Q & Tian B (2019) Regulation of Intronic Polyadenylation
726   by PCF11 Impacts mRNA Expression of Long Genes. *Cell Rep.* **26:** 2766-2778.e6

727 Wu B, De SK, Kulinich A, Salem AF, Koeppen J, Wang R, Barile E, Wang S, Zhang D, Ethell I
728   & Pellecchia M (2017) Potent and Selective EphA4 Agonists for the Treatment of
729   ALS. *Cell Chem. Biol.* **24:** 293–305

730 Yarosh CA, Tapescu I, Thompson MG, Qiu J, Mallory MJ, Fu XD & Lynch KW (2015)
731   TRAP150 interacts with the RNA-binding domain of PSF and antagonizes splicing of
732   numerous PSF-target genes in T cells. *Nucleic Acids Res.* **43:** 9006–9016

733 Zappulo A, Van Den Bruck D, Ciolli Mattioli C, Franke V, Imami K, McShane E, Moreno-
734   Estelles M, Calviello L, Filipchyk A, Peguero-Sanchez E, Müller T, Woehler A,
735   Birchmeier C, Merino E, Rajewsky N, Ohler U, Mazzoni EO, Selbach M, Akalin A &
736   Chekulaeva M (2017) RNA localization is a key determinant of neurite-enriched
737   proteome. *Nat. Commun.* **8:**

738