

Online Research @ Cardiff

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <https://orca.cardiff.ac.uk/142799/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Wright, Caroline F., Quaife, Nicholas M., Ramos-Hernández, Laura, Danecek, Petr, Ferla, Matteo P., Samocha, Kaitlin E., Kaplanis, Joanna, Gardner, Eugene J., Eberhardt, Ruth Y., Chao, Katherine R., Karczewski, Konrad J., Morales, Joannella, Gallone, Giuseppe, Balasubramanian, Meena, Banka, Siddharth, Gompertz, Lianne, Kerr, Bronwyn, Kirby, Amelia, Lynch, Sally A., Morton, Jenny E.V., Pinz, Hailey, Sansbury, Francis H., Stewart, Helen, Zuccarelli, Britton D., Cook, Stuart A., Taylor, Jenny C., Juusola, Jane, Retterer, Kyle, Firth, Helen V., Hurles, Matthew E., Lara-Pezzi, Enrique, Barton, Paul J.R. and Whiffin, Nicola 2021. Non-coding region variants upstream of MEF2C cause severe developmental disorder through three distinct loss-of-function mechanisms. *American Journal of Human Genetics* 108 (6) , pp. 1083-1094. 10.1016/j.ajhg.2021.04.025 file

Publishers page: <http://dx.doi.org/10.1016/j.ajhg.2021.04.025>
<<http://dx.doi.org/10.1016/j.ajhg.2021.04.025>>

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies.

See

<http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



1 **MEF2C m/sNon-coding region variants upstream of *MEF2C* cause severe**
2 **developmental disorder through three distinct loss-of-function mechanisms**

3

4 Caroline F. Wright¹, Nicholas M. Quaife^{2,3}, Laura Ramos-Hernández⁴, Petr Danecek⁵, Matteo
5 P. Ferla⁶, Kaitlin E. Samocha⁵, Joanna Kaplanis⁵, Eugene J. Gardner⁵, Ruth Y. Eberhardt⁵,
6 Katherine R. Chao^{7,8}, Konrad J. Karczewski^{7,8}, Joannella Morales⁹, Giuseppe Gallone^{5†},
7 Meena Balasubramanian^{10,11}, Siddharth Banka^{12,13}, Lianne Gompertz¹², Bronwyn Kerr¹³,
8 Amelia Kirby¹⁴, Sally A. Lynch¹⁵, Jenny E.V. Morton¹⁶, Hailey Pinz¹⁷, Francis H. Sansbury¹⁸,
9 Helen Stewart¹⁹, Britton D. Zuccarelli²⁰, Genomics England Research Consortium, Stuart A.
10 Cook², Jenny C. Taylor⁶, Jane Juusola²¹, Kyle Retterer²¹, Helen V. Firth^{5,22}, Matthew E.
11 Hurles⁵, Enrique Lara-Pezzi^{4,23}, Paul J.R. Barton^{2,3} and Nicola Whiffin^{5,8,24*}

12

13 ¹Institute of Biomedical and Clinical Science, University of Exeter Medical School, Royal Devon &
14 Exeter Hospital, Exeter, EX2 5DW, UK

15 ²National Heart & Lung Institute and MRC London Institute of Medical Sciences, Imperial College
16 London, London, W12 0NN, UK

17 ³Cardiovascular Research Centre, Royal Brompton & Harefield Hospitals NHS Trust, London, SW3
18 6NP, UK

19 ⁴Centro Nacional de Investigaciones Cardiovasculares Carlos III (CNIC), 28029 Madrid, Spain

20 ⁵Human Genetics Programme, Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton,
21 CB10 1RQ, UK

22 ⁶National Institute for Health Research Oxford Biomedical Research Centre, Wellcome Centre for
23 Human Genetics, University of Oxford, Oxford, OX3 7BN, UK

24 ⁷Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA 02114, USA

25 ⁸Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA
26 02142, USA

27 ⁹European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Cambridge,
28 CB10 1SD, UK

29 ¹⁰Sheffield Clinical Genetics Service, Sheffield Children's NHS Foundation Trust, Sheffield, S10 2TH,
30 UK

31 ¹¹Academic Unit of Child Health, Department of Oncology & Metabolism, University of Sheffield,
32 Sheffield, S10 2TH, UK

33 ¹²Manchester Centre for Genomic Medicine, St Mary's Hospital, Manchester University Hospitals NHS
34 Foundation Trust, Health Innovation Manchester, Manchester, M13 9WL, UK

35 ¹³Division of Evolution and Genomic Sciences, School of Biological Sciences, University of
36 Manchester, Oxford Road, Manchester, M13 9PL, UK

37 ¹⁴Department of Pediatrics, Wake Forest School of Medicine, Winston-Salem, NC 27101, USA

38 ¹⁵UCD Academic Centre on Rare Diseases, School of Medicine and Medical Sciences, University
39 College Dublin, and Clinical Genetics, Temple Street Children's University Hospital, Dublin, D01
40 XD99, Ireland

41 ¹⁶West Midlands Regional Clinical Genetics Service and Birmingham Health Partners, Birmingham
42 Women's and Children's Hospitals NHS Foundation Trust, Birmingham, B4 6NH, UK

43 ¹⁷Department of Pediatrics, Saint Louis University School of Medicine, Saint Louis, MO 63104, USA

44 ¹⁸All Wales Medical Genomics Service, NHS Wales Cardiff and Vale University Health Board, Institute
45 of Medical Genetics, University Hospital of Wales, Cardiff, CF14 4AY, UK

46 ¹⁹Oxford Centre for Genomic Medicine, Oxford University Hospitals NHS Foundation Trust, Oxford,
47 OX3 7LE, UK

48 ²⁰Department of Neurology, University of Kansas School of Medicine-Salina Campus, Salina, KS
49 67401, USA

50 ²¹GeneDx, Gaithersburg, MD 20877, USA

51 ²²East Anglian Medical Genetics Service, Cambridge University Hospitals NHS Foundation Trust,
52 Cambridge, CB2 0QQ, UK

53 ²³CIBER de enfermedades CardioVasculares (CIBERCV), 28029 Madrid, Spain

54 ²⁴Wellcome Centre for Human Genetics, University of Oxford, Oxford, OX3 7BN, UK

55 †Current address: Max Planck Institute for Molecular Genetics, Ihnestraße 63-73, 14195 Berlin,
56 Germany

57
58 *Correspondence should be addressed to nwhiffin@well.ox.ac.uk
59

60 **Abstract**

61

62 Clinical genetic testing of protein-coding regions identifies a likely causative variant in only
63 around half of developmental disorder (DD) cases. The contribution of regulatory variation in
64 non-coding regions to rare disease, including DD, remains very poorly understood. We
65 screened 9,858 probands from the Deciphering Developmental Disorders (DDD) study for *de*
66 *novo* mutations in the 5'untranslated regions (5'UTRs) of genes within which variants have
67 previously been shown to cause DD through a dominant haploinsufficient mechanism. We
68 identified four single nucleotide variants and two copy number variants upstream of *MEF2C*
69 in a total of 10 individual probands. We developed multiple bespoke and orthogonal
70 experimental approaches to demonstrate that these variants cause DD through three distinct
71 loss-of-function mechanisms, disrupting transcription, translation, and/or protein function.
72 These non-coding region variants represent 23% of likely diagnoses identified in *MEF2C* in
73 the DDD cohort, but these would all be missed in standard clinical genetics approaches.
74 Nonetheless, these variants are readily detectable in exome sequence data, with 30.7% of
75 5'UTR bases across all genes well covered in the DDD dataset. Our analyses show that
76 non-coding variants upstream of genes within which coding variants are known to cause DD
77 are an important cause of severe disease and demonstrate that analysing 5'UTRs can
78 increase diagnostic yield. We also show how non-coding variants can help inform both the
79 disease-causing mechanism underlying protein-coding variants, and dosage tolerance of the
80 gene.

81

82 **Introduction**

83

84 The importance of non-coding regulatory variation in common diseases and traits has long
85 been appreciated, however, the contribution of non-coding variation to rare disease remains
86 poorly understood¹⁻⁴. Consequently, current clinical testing approaches for rare disease
87 focus almost exclusively on regions of the genome that code directly for protein, within which

88 we are able to relatively accurately estimate the effect of any individual variant. Using this
89 approach, however, disease-causing variants are only identified in around 36% of individuals
90 with developmental disorders (DD)⁵ using exome sequencing, with a further 15-20%
91 diagnosed through chromosomal microarrays⁶. In previous work, we assessed the role of *de*
92 *novo* mutations (DNMs) in distal regulatory elements and estimated that 1-3% of
93 undiagnosed DD cases carry pathogenic DNMs in these regions¹.

94

95 Untranslated regions (UTRs) at the 5' and 3' end of genes present a unique opportunity to
96 expand genetic testing outside of protein coding regions given they have important
97 regulatory roles in controlling both the amount and location of mRNA in the cell, and the rate
98 at which it is translated into protein^{7,8}. Crucially, we also know the genes/proteins that these
99 regions regulate. Given that UTRs account for around the same genomic footprint as
100 protein-coding exons, they have substantial potential to harbour novel Mendelian
101 diagnoses^{9,10}. UTRs are, however, not regularly included in exome sequence capture
102 regions, and are excluded in most analysis pipelines. This is primarily due to a lack of
103 guidance on how to determine when UTR variants are likely to be pathogenic.

104

105 Recently, we demonstrated that variants creating upstream start codons (uAUGs) in 5'UTRs
106 are under strong negative selection, and are an important cause of Mendelian diseases,
107 including neurofibromatosis and Van der Woude syndrome^{11,12}. Initiation of translation at a
108 newly created uAUG can decrease translation of the downstream coding sequence (CDS).
109 The strength of negative selection acting on uAUG-creating variants varies depending on
110 both the match of the sequence surrounding the uAUG to the Kozak consensus, which is
111 known to regulate the likelihood that translation is initiated^{13,14}, and the nature of the
112 upstream open reading frame (uORF) that is created. Variants that result in ORFs which
113 overlap the CDS have a larger impact on CDS translation and hence are more
114 deleterious^{11,15}.

115

116 Here, we screened 9,858 probands from the Deciphering Developmental Disorders (DDD)⁵
117 study for DNMs in the 5'UTRs of genes within which variants have previously been shown to
118 cause DD through a dominant haploinsufficient mechanism (defined using the clinically-
119 curated Developmental Disorders Genotype to Phenotype (DDG2P) database and
120 henceforth referred to as 'DDG2P haploinsufficient genes'). We uncover likely disease-
121 causing variants that are entirely non-coding and show how these variants cause disease
122 through three distinct loss-of-function mechanisms. We further show how disease-causing
123 missense variants in MEF2C [MIM:600662] are clustered at the N-terminus and likely also
124 cause loss-of-function by disrupting binding of MEF2C protein to DNA. Finally, we analyse
125 the coverage across all UTRs in the DDD exome sequencing dataset to demonstrate how
126 these regions can be readily screened in existing datasets to increase diagnostic yield and
127 glean insight into disease causing mechanisms.

128

129 **Materials and Methods**

130 *Recruitment, sample collection and clinical data*

131

132 The DDD Study has UK Research Ethics Committee approval (10/H0305/83, granted by the
133 Cambridge South REC, and GEN/284/12 granted by the Republic of Ireland REC).

134 Individuals with severe, undiagnosed developmental disorders and their parents were
135 recruited and systematically phenotyped by the 24 Regional Genetics Services within the
136 United Kingdom (UK) National Health Service and the Republic of Ireland. Saliva samples
137 were collected from probands and parents, and DNA extracted as previously described¹⁶;
138 blood-extracted DNA was also collected for probands where available. Clinical data (growth
139 measurements, family history, developmental milestones, etc.) were collected using a
140 standard restricted-term questionnaire within DECIPHER¹⁷. Informed consent was obtained
141 for all participants.

142

143 *Genetic data*

144

145 Array-CGH analysis was performed using 2 x 1M probe custom designed microarrays
146 (Agilent; Amadid No.s 031220/031221) as described previously¹⁶. Exome sequencing was
147 performed using Illumina HiSeq (75-base paired-end sequencing) with SureSelect baits
148 (Agilent Human All-Exon V3 Plus and V5 Plus with custom ELID C0338371) and variants
149 were called and annotated as described previously¹⁶. We used DeNovoGear¹⁸ (version 0.54)
150 to detect likely DNMs from trio exome BAM files and Ensembl Variant Effect Predictor¹⁹ was
151 used to annotate predicted consequences. The data are available under managed access
152 from the European Genome-phenome Archive (Study ID EGAS00001000775), and likely
153 diagnostic variants are available open access in DECIPHER.

154

155 *Defining a gene-set of interest*

156

157 We limited our analysis to 359 DDG2P²⁰ genes with a confirmed or probable role in
158 developmental disorders and with a dominant (including X-linked dominant) loss-of-function
159 disease mechanism (downloaded on 21st July 2020 - see Web Resources section for link;
160 Table S1). We refer to these genes as 'DDG2P haploinsufficient genes'.

161

162 *Identifying uAUG-creating variants in DDD*

163

164 We defined high-confidence DNMs in DD as previously²¹, using the following criteria: minor
165 allele frequency < 0.01 in our cohort and reference databases, depth in the child > 7, depth
166 in both parents > 5, Fisher strand bias p-value > 10⁻³, and a posterior probability of being a
167 DNM from DeNovoGear > 0.00781¹⁸. Additionally, we filtered out DNMs with some evidence
168 of an alternative allele in one of the parents and indels with a low variant allele fraction
169 (<30% of the reads support the alternative) that had a minor allele frequency > 0. We cross-
170 referenced this list of high-confidence DNMs with a list of all possible uAUG-creating SNVs

171 from previous work¹¹. We also assessed any small insertions and deletions that could form
172 uAUGs.

173

174 The strength of the Kozak consensus surrounding each uAUG was assessed as described
175 previously¹¹. Specifically, we assessed the positions at -3 and +3 relative to the A of the

176 AUG, requiring both the -3 base to be either A or G and the +3 to be G for an annotation of

177 'Strong'. If only one of these conditions was true, the strength was deemed to be 'Moderate'

178 and if neither was the case 'Weak'.

179

180 *Defining the 5'UTR of MEF2C*

181

182 We used the MANE Select transcript ENST00000504921.7 for which the 5'UTR was defined
183 using CAGE data from the FANTOM5 project²², RNA-seq supported intron data from the

184 Intropolis resource²³, and exon level expression from the GTEx project²⁴. The Matched

185 Annotation from the NCBI and EMBL-EBI (MANE) is a collaborative project that aims to

186 define a representative transcript (MANE Select) for each protein-coding locus across the

187 genome. The MANE set perfectly aligns to the GRCh38 reference assembly and includes

188 pairs of 100% identical RefSeq and Ensembl/GENCODE transcripts²⁵. The 5'UTR of *MEF2C*

189 was therefore defined as two exons: chr5:88178772-88179001 and chr5:88119606-

190 88119747 on GRCh37, or chr5:88882955-88883184 and chr5:88823789-88823930 on

191 GRCh38.

192

193 *Searching for MEF2C 5'UTR variants in external datasets*

194

195 We queried the regions corresponding to the *MEF2C* 5'UTR for DNMs in (1) a set of 18,789

196 DD trios sequenced by the genetic testing company GeneDx⁵, (2) 13,949 rare disease trios

197 from the main programme v9 release of the UK 100,000 Genomes Project from Genomics
198 England²⁶, and (3) variants in the v3.0 dataset of the Genome Aggregation Database
199 (gnomAD)²⁷.

200

201 *Assessing 5'UTR coverage*

202

203 Regions corresponding to 5'UTRs were extracted from the .gff file from the MANE project
204 v0.91 (see Web Resources; MANE Select transcripts). For each base, we calculated the
205 mean coverage across 1,000 randomly selected samples from DDD. A mean coverage of
206 >10x was used to call a base 'covered'. Analysis was limited to genes with a defined MANE
207 Select transcript. For our DDG2P haploinsufficient genes this was 345/359 genes (96.1%).

208

209 To identify all possible uAUG-creating variants in DDG2P haploinsufficient genes, we
210 extracted the 5'UTR sequence from the MANE rna.fna file and used the UTRannotator²⁸ to
211 find all possible uAUG-creating sites and annotate their consequence.

212

213 *Functional validation of variants creating out-of-frame ORFs (oORFs): by MEF2C 5'UTR-* 214 *luciferase translation assay*

215

216 Expression constructs: WT and variant MEF2C 5'UTRs were cloned directly upstream of
217 Gaussia luciferase (GLuc) in the pEZX-GA02 backbone (Labomics) and sequenced to
218 confirm integrity. Secreted alkaline phosphatase (SEAP) was expressed on the same
219 construct for normalisation of transfection efficiency.

220

221 Cull culture, transfection and analysis: HEK293T cells were purchased from ATCC and
222 cultured in Dulbecco's Modified Eagle Medium (glutamine+, pyruvate+) supplemented with
223 10% foetal bovine serum and 1% penicillin/streptomycin. Cells were transfected with MEF2C
224 5'UTR-luciferase constructs using Lipofectamine 3000, following manufacturer's protocols.

225 After 24h, culture medium was sampled and GLuc and SEAP were simultaneously quantified
226 using the Secrete-Pair Dual Luminescence assay (Genecopoeia). Fifteen technical
227 replicates were performed across three independent experiments.

228

229 qPCR: RNA was purified from cells using phenol-chloroform extraction and the Qiagen
230 RNeasy Miniprep kit. RNA quantity was normalised and cDNA generated using IV VILO
231 reverse transcriptase following manufacturer's protocols. Quantitative PCR was performed
232 using SYBR green master mix on a Quantstudio 7 Real-time PCR system and results
233 normalised to co-amplified GAPDH. The following primers were used: GLUC F: 5'
234 CTGTCTGATCTGCCTGTCCC 3', GLUC R: 5' GGA CTCTTTGTCGCCTTCGT 3', SEAP F:
235 5' ACCTTCATAGCGCACGTCAT 3' and SEAP R: 5' TCTAGAGTAACCCGGGTGCG 3',
236 GAPDH F: 5' GGAGTCAACGGATTTGGTTCG 3', GAPDH R: ATCGCCCCACTTGATTTTGG
237 3'.

238

239 Kozak mutagenesis: The kozak context of the c.-103G>A MEF2C 5'UTR-luciferase construct
240 was modified using the Quikchange II mutagenesis kit, following manufacturers protocols.

241 The following PAGE-purified mutagenesis primers were used: F:

242 5'CTCCTTCTTCAGCATTTTCACAGCTCAGTTCCCAA 3', R: 5'

243 TTGGGAACTGAGCTGTGAAAATGCTGAAGAAGGAG 3'. Constructs were fully sequenced
244 to verify mutation and construct integrity in each case

245

246 *Functional validation of CDS-elongating variants: by MEF2 binding site-luciferase*
247 *transactivation assay*

248

249 Expression and reporter constructs: WT and variant MEF2C 5' UTR+CDS oligos were
250 cloned into the pReceiver-M02 expression construct (Labomics) and sequenced to confirm
251 integrity. For normalisation of transfection efficiency, cells were co-transfected with pRL-
252 Renilla. A desMEF2-luciferase reporter construct was used to quantify the transactivational

253 efficiency of each MEF2C expression construct, and consisted of three copies of a high-
254 affinity MEF2 binding site²⁹, linked to an hsp68 minimal promoter in pGL3 (Promega)³⁰.

255

256 Cell culture and transfection: HL1 cardiomyocytes were cultured in Claycomb medium,
257 supplemented with 2 mM L-glutamine, 10% FBS and 100 g/ml Penicillin/Streptomycin.

258 Culture surfaces were pre-treated with gelatin/fibronectin. Cells were co-transfected with 1)

259 desMEF2-luciferase reporter construct, 2) pRL-Renilla transfection control, and 3)

260 expression construct of either: i) empty pcDNA3.1 (negative control), ii) WT MEF2C 5'

261 UTR+CDS, iii) MEF2C -26C>T, or iv) MEF2C -8C>T. Transfection was with Lipofectamine

262 2000, following manufacturers protocols. 48h after transfection, firefly and Renilla

263 Luciferases were quantified by the Promega Dual-Luciferase Reporter Assay System.

264 Eighteen technical replicates were performed across three independent experiments.

265

266 Western blot: HL1 cells were lysed in RIPA buffer in the presence of protease and

267 phosphatase inhibitors (04693159001 and 04906845001, Roche Diagnostics). Lysates were

268 separated on SDS-PAGE gels and transferred to PVDF membranes, which were blocked

269 with 3% skimmed milk in TBS. The primary antibody was anti-MEF2C (ab211493, Abcam),

270 and the secondary antibody was anti-mouse P0447 from Dako. The membrane was

271 developed using ECL reagent (AC2204, Azure Biosystems) and intensity of the bands

272 quantified using ImageJ software.

273

274 Statistical analysis for all assays: Data were analysed for statistical significance using 1-way

275 ANOVA followed by Tukey's post-test, using GraphPad Prism 8.0.

276

277 *CNV calling*

278

279 Four CNV detection algorithms (XHMM³¹, CONVEX¹⁶, CLAMMS³² and CANOES³³) were
280 used to ascertain CNVs from exome data, followed by a random forest machine learning
281 approach to integrate and filter the results (manuscript in preparation).

282

283 Layered H3K4me3 data (to visualise active promoter regions) was downloaded from the
284 UCSC table browser for GN12878 as a representative cell line and plotted alongside the
285 identified CNVs in Figure S1.

286

287 *Modelling missense disruption to DNA-binding*

288

289 We collated a set of missense variants identified in *MEF2C* in DD cases comprising all *de*
290 *novo* variants from trios in DDD and GeneDx published previously⁵, and variants from
291 ClinVar either flagged as being identified as *de novo*, or with functional evidence (Table S3).

292

293 As a comparator, we used missense variants from gnomAD v2.1.1²⁷. Given that there are
294 only three variants in the N-terminal region of *MEF2C* in gnomAD, but the sequence of the
295 N-terminal region is near identical across the four MEF2 proteins (Figure S4), we used
296 missense variants from all four genes (*MEF2A-D*; Table S4).

297

298 Based on structures of the N-terminal MADS-box of *MEF2A* homodimer(1egw, 3kov and
299 6byy, residues 1-92) bound to its DNA consensus sequence³⁴, we categorised residues into
300 one of four categories: (1) in N-terminal random coil and in contact with the DNA (2) in N-
301 terminal alpha-helix pointing towards the DNA; (3) in N-terminal alpha-helix pointing away
302 from the DNA; or (4) distal to the DNA contact surface (Table S5). We used a two-sided
303 Fisher's exact test to assess for an enrichment of variants in contact or pointing towards the
304 DNA helix in DD cases (Table S6).

305

306 The Swissmodel threaded model of MEF2C based upon PDB:6BYY (89% identity)^{35,36} was
307 energy minimised using Pyrosetta³⁷ with 15 FastRelax cycles³⁸ against the electron density
308 of PDB:6BYY and 5 unconstrained. The DNA was extended on both ends due to the
309 proximity of R15. Mutations were introduced and the 10 Å neighbourhood was energy
310 minimised. Gibbs free energy was calculated using the Rosetta ref2015 scorefunction³⁹.
311 Gibbs free energy of binding was calculated by pulling away the DNA and repacking
312 sidechains and, in the case of residues in the N-terminal loop, thoroughly energy minimising
313 the backbone of the loop as this is highly flexible when unbound. N-terminal extensions were
314 made using the RemodelMover⁴⁰ with residues 2-5 also remodelled as determined by
315 preliminary test. Closest distance of each residue to the DNA was calculated with the Python
316 PyMOL module. Code used for this analysis can be found at the link in Web Resources. This
317 interactive page was made in Michelangelo⁴¹.

318

319 All missense variants are annotated with respect to the Ensembl canonical transcript
320 ENST00000340208.5.

321

322 *Calculating regional missense constraint and de novo enrichment*

323

324 We determined regional missense constraint by (1) extracting observed variant counts from
325 the 125,748 samples in gnomAD v2.1.1, (2) calculating the expected variant count per
326 transcript, and (3) applying a likelihood ratio test to search for significant breaks that split a
327 transcript into two or more sections of variable missense constraint.

328

329 Observed missense variants were extracted from the gnomAD exomes Hail Table (version
330 2.1.1) as described previously²⁷, using the following criteria:

- 331 ● Annotated as a missense change in a canonical transcript of a protein-coding gene in
332 Gencode v19 by Variant Effect Predictor (VEP, version 85)
- 333 ● Median coverage greater than zero in the gnomAD exomes data

- 334 • Passed variant filters
- 335 • Adjusted allele count of at least one and an allele frequency less than 0.1% in the
- 336 gnomAD exomes

337

338 To calculate the expected variant count, we extended methods described previously²⁷ to
339 compute the proportion of expected missense variation per base. Briefly, we annotated each
340 possible substitution with local sequence context, methylation level (for CpGs), and
341 associated mutation rate from the table computed in Karczewski *et al.*²⁷ We aggregated
342 these mutation rates across the transcript and calibrated models based on CpG status and
343 median coverage. To determine the expected variants for a given section of the transcript,
344 we calculated the fraction of the overall the mutation rate represented by the section and
345 multiplied it by the aggregated expected variant count for the full transcript.

346

347 We defined missense constraint by extending the methods from Samocha *et al.*⁴² We
348 employed a likelihood ratio test to compare the null model (transcript has no regional
349 variability in missense constraint) with the alternative model (transcript has evidence of
350 regional variability in missense constraint). We required a χ^2 value above a threshold of 10.8
351 to determine significance for each breakpoint, and in the case of multiple breakpoints,
352 retained the breakpoint with the maximum χ^2 . This approach defined a single breakpoint in
353 the *MEF2C* canonical transcript at chr5:88057138 (GRCh37).

354

355 To evaluate the enrichment of DNMs in the transcript when removing the N-terminal section,
356 we determined the probability of a missense mutation in that region and then compared the
357 observed number of DNMs (n=3) with the expected count in 28,641 individuals using a
358 Poisson test. Specifically, we took the probability of a missense mutation (μ_{mis}) as
359 provided in the gnomAD v2.1 constraint files for *MEF2C* and adjusted it for the fraction of
360 mutability represented in the latter section of the gene (~79.5%).

361

362 **Results**

363 *Identifying de novo 5'UTR variants in DD cases*

364

365 To investigate the contribution of uAUG-creating variants to severe DD cases, we analysed
366 29,523 high-confidence DNMs identified in exome sequencing data from 9,858 parent-
367 offspring trios in the DDD study⁵. Although the majority of DNMs identified are coding, as
368 expected with exome sequencing data, many non-coding variants are also detectable,
369 particularly near exon boundaries. Given that uAUG-creating variants that decrease CDS
370 translation would only be expected to be deleterious in genes that are dosage sensitive, we
371 restricted our analysis to the 5'UTRs of 359 haploinsufficient genes from the curated DDG2P
372 database²⁰ (Table S1).

373

374 We identified five unique uAUG-creating *de novo* single nucleotide variants (SNVs) in five
375 unrelated probands upstream of two different genes. All of these variants are absent from
376 the Genome Aggregation Database (gnomAD) population reference dataset (both v2.1.1 and
377 v3.0)²⁷. Notably, four of the five variants were found in the 5'UTR of *MEF2C* in probands with
378 phenotypes consistent with *MEF2C* haploinsufficiency (Table 1; [MIM: 613443])⁴³. Two of
379 these DNMs create uAUGs out-of-frame with the *MEF2C* CDS, which are expected to
380 reduce downstream protein translation, whilst the other two create uAUGs in-frame with the
381 CDS, which are expected to elongate the protein (Figure 1). The fifth variant was located in a
382 strong Kozak consensus upstream of *STXBP1* (ENST00000373302.8:c.-26C>G), creating
383 an uAUG out-of-frame with the *STXBP1* CDS; the phenotype of the proband with this variant
384 is consistent with *STXBP1* haploinsufficiency⁴⁴, including global developmental delay,
385 microcephaly, and delayed speech and language development.

386

387 Given the identification of multiple uAUG-creating *de novo* SNVs in *MEF2C* in the DDD
388 study, we subsequently queried high-confidence DNMs identified in 18,789 trios with DD that
389 were exome sequenced by GeneDx⁵ for additional *MEF2C* DNMs. We uncovered three

390 additional *de novo* occurrences of two of the uAUG-creating variants observed in the DDD
391 study. In addition, we identified a further *de novo* occurrence of one of these variants in a DD
392 proband in the UK 100,000 Genomes Project²⁶(Table 1).

393

394 In a separate analysis, we analysed copy number variants (CNVs) identified in the DDD
395 study using exome sequencing data and identified five *de novo* CNVs overlapping *MEF2C*
396 (Figure S1). Two of these CNVs (each found in a single additional proband) overlap the
397 5'UTR of *MEF2C* without impacting any of the coding exons (Table 1). These two non-
398 coding CNVs delete the first exon of the *MEF2C* 5'UTR and >40kb of immediately upstream
399 sequence (294kb and 97kb, respectively), removing the entire promoter (as defined by the
400 Ensembl regulatory build⁴⁵ and H3K4me3 peaks from ENCODE⁴⁶) and likely abolishing
401 transcription of this allele (Figure S1). There are no large deletions (>600bps) in this
402 upstream region in the gnomAD structural variant dataset (v2.1)⁴⁷. Both coding *MEF2C*
403 disruptions and non-coding deletions further upstream of *MEF2C* that are predicted to
404 disrupt enhancer function have been identified in DD probands previously^{48,49}.

405

406 *De novo 5'UTR variants cause phenotypes consistent with MEF2C haploinsufficiency*

407

408 We collated all available clinical data for the ten probands with *MEF2C* 5'UTR *de novo*
409 variants and in each case the observed phenotype is consistent with previously reported
410 *MEF2C* haploinsufficiency^{50,51} (Table S2). Specifically, of the nine individuals for which
411 detailed phenotypic information was available, the following features were noted: global
412 developmental delay (9/9) with delayed or absent speech (9/9), seizures (8/9), hypotonia
413 (5/9) and stereotypies (2/9). These probands had no other likely disease-causing variants in
414 the coding sequence of *MEF2C*, or in any other DDG2P genes following exome sequencing.

415

416 *uAUG-creating SNVs cause loss-of-function by reducing translation or disrupting protein*
417 *function*

418

419 The four uAUG-creating SNVs identified in *MEF2C* result in two different downstream
420 effects. We used two distinct experimental approaches to evaluate the impact of i) out-of-
421 frame uAUG-creating variants on downstream translation and ii) CDS-elongating variants on
422 *MEF2C*-dependent transactivation.

423

424 Two of the variants (c.-66A>T and c.-103G>A), each found in a single proband, create
425 uAUGs that are out-of-frame with the coding sequence (CDS), creating an overlapping ORF
426 (oORF) that terminates 128 bases after the canonical start site (Figure 1b). Using a
427 translation assay, with wild-type or mutant 5'UTR sequence cloned upstream of a luciferase
428 reporter gene, we show that both variants result in a significant decrease in translational
429 efficiency (Figure 2a; Figure S2a). The amount by which translation is reduced appears to be
430 dependent on the uAUG match to the Kozak consensus sequence, consistent with previous
431 observations¹¹. The c.-103G>A variant, which creates an uAUG with a weak Kozak
432 consensus, results in only a moderate decrease in luciferase expression, and the proband
433 with this variant displays a milder phenotype on clinical review. To validate that this
434 difference in effect is indeed due to the differing Kozak strengths, in the c.-103G>A
435 translation assay, we mutated a single base to alter the oORF start context to a moderate
436 Kozak consensus match (see methods). This modification resulted in significantly decreased
437 translational efficiency compared to the unmodified c.-103G>A variant, to a level equivalent
438 to the c.-66A>T variant (Figure S3). The individual carrying the c.-103G>A does not have
439 any other 5'UTR variants that could similarly modify the variant's effect. These data suggest
440 that *MEF2C* is sensitive to even partial loss-of-function.

441

442 The other two variants (c.-8C>T and c.-26C>T) are both observed recurrently *de novo*, each
443 in three unrelated probands (Table 1). Both variants create uAUGs that are in-frame with the
444 CDS, resulting in N-terminal extensions of three and nine amino acids respectively (Figure
445 1c). *MEF2C* is a transcription factor, and critical to its function is the DNA-binding domain

446 located at the extreme N-terminal region⁵². Although no structure is available for the MEF2C
447 protein, numerous crystal and NMR structures of the N-terminal DNA-binding domain of
448 human MEF2A are available, which is 96% identical in sequence to MEF2C. These
449 structures show clearly that the extreme N-terminus of the protein is in direct contact with
450 DNA^{34,53}, and that the first few residues bind directly into the minor groove (Figure 3). We
451 assayed MEF2C-dependent transactivation using MEF2C expression constructs with wild-
452 type and mutant 5'UTR sequences. These data demonstrate significantly reduced activation
453 of target gene transcription from the variants (Figure 2b; Figure S2b and c), compared to
454 wild-type MEF2C. Once again, the strength of the effect is dependent on the uAUG context,
455 with the c.-8C>T variant that creates a strong Kozak consensus having a larger effect,
456 almost abolishing transactivation activity.

457

458 We looked in the gnomAD dataset²⁷ for uAUG-creating variants that might have similar
459 impacts. Across the exome (v2.1.1) and genome (v3.0) sequencing datasets, there are only
460 two uAUG-creating variants in the *MEF2C* 5'UTR. Crucially, neither of these fall into the
461 proximal 5'UTR exon and neither create ORFs overlapping the CDS. In both instances, the
462 uAUGs are created into weak Kozak-consensus contexts, and they have in-frame stop
463 codons after 6bps (allele count = 6) and 57bps (allele count = 1) respectively (Table 1;
464 Figure 1d). These variants would therefore not be expected to have substantial, if any, effect
465 on *MEF2C* translation.

466

467 *Pathogenic de novo missense variants likely cause loss-of-function of MEF2C through*
468 *disrupting DNA-binding*

469

470 Whilst the major recognised mechanism through which pathogenic variants in *MEF2C* lead
471 to severe developmental phenotypes is loss-of-function, *de novo* missense variants are also
472 significantly enriched in DD trios ($P=1.3 \times 10^{-14}$)⁵ and multiple pathogenic missense variants
473 are reported in ClinVar⁵⁴. These variants are almost exclusively found at the extreme N-

474 terminus of the protein (Table S3), in the DNA-binding region, which is also highly
475 constrained for missense variants in gnomAD (obs/exp=0.069; calculated on 125,748 exome
476 sequenced samples in v2.1.1; Figure 3a). We hypothesised that these pathogenic missense
477 variants are also causing loss-of-function by disrupting DNA-binding of MEF2C as has been
478 demonstrated for random disruptions to the N-terminal region⁵² and two proband variants⁴⁹
479 previously. Using the structure of the N-terminal MEF2A homodimer bound to DNA, we
480 modelled the location of pathogenic missense variants in MEF2C, as well as missense
481 variants in gnomAD v2.1.1 across all members of the myocyte enhancer factor 2 protein
482 family (MEF2A-D; 84% N-terminal domain sequence identity; Table S4; Figure S4), and saw
483 a significant enrichment of pathogenic variants interacting directly with DNA via both the N-
484 terminal loop and DNA-binding helix (Fisher's $P=2.6 \times 10^{-5}$, Figure 3b; Tables S5 & S6). We
485 further calculated the change in Gibbs free energy ($\Delta\Delta G$) of both the protein-DNA interaction
486 and the complex stability for each missense change. Variants found in DD cases have
487 significantly increased $\Delta\Delta G$ scores compared to gnomAD variants (Wilcoxon $P=2.7 \times 10^{-4}$;
488 Figure 3c) and are significantly closer to the bound DNA (Wilcoxon $P=1.5 \times 10^{-5}$; Figure 3d;
489 Table S7). Together, these data suggest that disease-causing missense variants in *MEF2C*
490 act through a loss-of-function mechanism, as has been experimentally demonstrated for two
491 proband variants previously⁴⁹. Indeed, excluding the N-terminal DNA-binding domain, the
492 remainder of *MEF2C* shows much weaker constraint against missense variants in gnomAD
493 (obs/exp=0.41), and only nominal enrichment for *de novo* missense variants in DD cases
494 ($P=0.041$).

495

496 *Disease-causing 5'UTR variants can be detected in exome sequencing data*

497

498 Given our ability to identify 5'UTR variants in *MEF2C*, we investigated the extent to which
499 these regions are captured across all genes in the exome sequencing dataset from the DDD
500 study. We find that 30.7% of all gene 5'UTR bases and 20.4% of 5'UTR bases of our
501 DDG2P haploinsufficient genes (average of 73 bps per gene; n=345 with MANEv0.91

502 transcripts) are covered at a mean coverage threshold of >10x. The average length of
503 5'UTRs in DDG2P haploinsufficient genes is 356 bps (Figure 4a), with 42.0% containing
504 multiple exons (Figure 4b). As expected, 5'UTR coverage decays as distance from the CDS
505 increases (Figure 4c), with distal exons very poorly covered (6.7% of bases >10x). In
506 comparison, a much lower proportion of 3'UTR bases (6.0%) are covered at >10x, which is
507 unsurprising given that 3'UTRs are much longer than 5'UTRs, at an average of 2,652 bps for
508 our DDG2P haploinsufficient genes.

509

510 To determine the proportion of all possible uAUG-creating variants that are sufficiently
511 covered in the DDD exome sequence data, we computationally identified 3,962 possible
512 uAUG-creating variants in DDG2P haploinsufficient genes that would create out-of-frame
513 overlapping ORFs (n=2,782) or CDS-elongations (n=1,180). Of these, 42.4% are sequenced
514 at >10x coverage across the DDD study dataset (40.2% of out-of-frame and 47.6% of CDS-
515 elongating). However, we would not expect CDS-elongating variants to cause a loss-of-
516 function for the majority of genes. Rather, we expect this to be limited to genes with
517 important functional domains at the extreme N-terminus that would be adversely affected by
518 the addition of extra N-terminal amino acids, either through disrupting binding or altering
519 protein structure. Based on Pfam domain predictions, only three of the proteins encoded by
520 our 359 DDG2P haploinsufficient genes, including *MEF2C*, have DNA-binding domains that
521 start within 10 bps of the N-terminus (Figure 4d); the other two (*ZNF750* and *SIM1*) encode
522 an N-terminal zinc-finger and basic helix-loop-helix, respectively, and although no structures
523 are available, these bind DNA via specific motifs that are unlikely to include the extreme N-
524 terminal residues.

525

526 **Discussion**

527

528 Here, we have identified six unique non-coding, pathogenic DNMs in *MEF2C* in ten
529 individuals with severe developmental disorders (six in the DDD study, three in a cohort from

530 GeneDx, and one in the UK 100,000 Genomes Project). These variants act via three distinct
531 loss-of-function mechanisms at different stages of expression regulation: (1) two large
532 deletions remove the promoter and part of the 5'UTR and are predicted to abolish normal
533 transcription of *MEF2C*; (2) two SNVs create out-of-frame uAUGs and reduce normal
534 translation of the *MEF2C* coding sequence; and (3) two SNVs create in-frame uAUGs that
535 elongate the *MEF2C* coding sequence, disrupting binding of the MEF2C protein to DNA and
536 reducing subsequent transactivation of gene-expression. We also identified a single uAUG-
537 creating variant in *STXBP1* in a proband whose phenotype was consistent with *STXBP1*
538 haploinsufficiency. This variant is predicted to create an out-of-frame oORF into a strong
539 Kozak consensus, thus decreasing normal *STXBP1* translation (as ribosomes first
540 encounter, and begin to translate from this new uAUG), leading to reduced levels of *STXBP1*
541 protein.

542

543 These observations demonstrate the importance of screening 5'UTRs of genes known to
544 harbour disease-causing coding variants in individuals that remain genetically undiagnosed.
545 We have previously identified 20 probands with diagnostic DNMs (15 SNVs and 5 CNVs)
546 impacting *MEF2C* protein-coding regions in the 9,858 family trios analysed in the DDD
547 study. The six additional non-coding DNMs described here (4 SNVs and 2 CNVs) therefore
548 comprise 23% of diagnoses impacting *MEF2C* in this cohort.

549

550 Our data show that 5'UTR variants can be identified in existing datasets that were primarily
551 designed to capture coding sequences, with 30.7% of 5'UTR bases having sufficient (>10x)
552 coverage in exome sequencing data from the DDD study. However, exome sequencing data
553 is likely to only identify UTR variants that are proximal to the first and last exons of genes,
554 and whole genome or expanded panel sequencing will be required to assay distal or poorly
555 covered UTRs. Furthermore, given their large size, 3'UTRs are particularly poorly covered in
556 exome sequencing datasets. There are examples of disease-causing variants within 3'UTRs,

557 including those impacting polyA signals and microRNA binding^{9,55-57}, which will not be
558 detected using these methodologies but that could increase diagnostic yield.

559

560 Although we screened DNMs in the 5'UTRs of a set of 359 haploinsufficient DDG2P genes,
561 four of the five identified *de novo* uAUG-creating variants were found in *MEF2C*. This
562 enrichment in a single gene is likely due to a combination of factors (Figure S5). Firstly,
563 *MEF2C* has a proximal 5'UTR exon that is very well covered in the DDD exome sequencing
564 data. Secondly, this 5'UTR exon contains a large number of sites where a variant could
565 create an uAUG, with only two DDG2P haploinsufficient genes having more well-covered
566 possible uAUG-creating sites. Thirdly, unlike the other genes with well-covered possible
567 uAUG-creating sites, *MEF2C* haploinsufficiency is a recurrent cause of DD within the DDD
568 study (Figure S5). Finally, due to the direct interaction of the extreme N-terminus of *MEF2C*
569 with DNA, CDS-elongating variants are also likely to be pathogenic, which is unlikely to be
570 the case in the vast majority of other haploinsufficient genes. As a result, *MEF2C* may be
571 unusual in its potential for pathogenic mutations in the 5'UTR and similarly large increases in
572 diagnostic yield are unlikely across most DDG2P haploinsufficient genes. Nethertheless the
573 enrichment of uAUG-creating variants in *MEF2C* is striking: only 14 of 426 possible variants
574 create uAUGs (at 142 5'UTR bases that are well-covered in the DDD study exome
575 sequencing data), yet all four DNMs observed in the DDD study in the *MEF2C* 5'UTR are
576 uAUG-creating (binomial $P=1.2\times 10^{-6}$).

577

578 In our functional data, we see a difference in the size of variant effects dependent on the
579 strength of the Kozak consensus surrounding the newly created uAUG. The Kozak
580 sequence is known to influence the likelihood of a ribosome initiating translation at any given
581 AUG as it scans along the 5'UTR from the 5' cap¹³. Our four uAUG-creating variants each
582 generate a new uORF that overlaps the coding sequence. Ribosomes that initiate translation
583 at these uAUGs will not be available to translate from the wild-type coding start site (which

584 itself has a strong Kozak consensus), resulting in reduced translation of the CDS. The
585 stronger the Kozak consensus around the uAUG, the greater this effect will be.

586

587 As we extend our analyses to detect non-coding variants, we caution that interpretation of
588 UTR variants still remains a critical challenge. Every 5'UTR has a unique combination of
589 regulatory elements tightly regulating RNA stability and protein expression^{58,59}, and the
590 impact of any variant will vary with the gene-specific context. Functional validation of
591 identified variants will therefore be crucial to prove (or reject) causality. Some variants may
592 have only a partial regulatory effect, but these variants can nonetheless be harnessed to
593 assess the extent to which perturbation of protein levels or function is tolerated, potentially
594 leading to reduced expressivity and/or lower penetrance. In the case of *MEF2C*, our results
595 suggest that even partial reductions in protein expression lead to severe disease.

596

597 Finally, we note how the mechanism of action of non-coding variants can inform the
598 mechanisms underlying protein-coding variants. Identification and characterisation of the
599 effect of the CDS-elongating *MEF2C* variants led us to analyse the domain structure of
600 *MEF2C* protein and confirm that all the currently identified missense variants likely also act
601 via disrupting DNA-binding, leading to a loss-of-function.

602

603 In conclusion, our results further highlight the important contribution of non-coding regulatory
604 variants to rare disease and underscore the huge promise of large whole-genome
605 sequencing datasets to both find new diagnoses and further our understanding of regulatory
606 disease mechanisms.

607

608 **Supplementary Data**

609

610 Supplementary data include five figures and seven tables. Also included is the Genomics
611 England Research Consortium author list.

612

613 **Declaration of Interests**

614

615 K.J.K. is a consultant for Vor Biopharma. J.J. and K.R. are employees of GeneDx, Inc. K.R.

616 holds shares in Opko Health, Inc. B.D.Z. is a member of the speakers bureau for Biogen,

617 Neurelis, and Supernus. S.A.C. is co-founder and shareholder of Enleofen Bio Pte Ltd.

618 M.E.H. is co-founder, shareholder, consultant, and non-executive director of Congenica Ltd.

619 All other authors declare no competing interests.

620

621 **Acknowledgements**

622

623 NW is currently supported by a Sir Henry Dale Fellowship jointly funded by the Wellcome

624 Trust and the Royal Society (Grant Number 220134/Z/20/Z). Initial work was completed

625 whilst NW was supported by a Rosetrees and Stoneygate Imperial College Research

626 Fellowship. NQ is supported by the Imperial College Academic Health Science Centre. This

627 work is additionally supported by The Rosetrees Trust (Grant Number H5R01320), the

628 Wellcome Trust [WT200990/Z/16/Z, WT200990/A/16/Z], Fondation Leducq [16 CVD 03], the

629 National Institute for Health Research (NIHR) Imperial College Biomedical Research Centre,

630 the Cardiovascular Research Centre, Royal Brompton & Harefield NHS Trust, and the NIHR

631 Oxford Biomedical Research Centre Programme. The views expressed are those of the

632 author(s) and not necessarily those of the NHS, the NIHR or the Department of Health. The

633 DDD study presents independent research commissioned by the Health Innovation

634 Challenge Fund (Grant Number HICF-1009-003) a parallel funding partnership between the

635 Wellcome Trust and the Department of Health, and the Wellcome Trust Sanger Institute

636 (Grant Number WT098051). See Nature 2015;519:223-8 or www.ddduk.org/access.html for

637 full acknowledgement. This research was made possible through access to the data and

638 findings generated by the 100,000 Genomes Project. The 100,000 Genomes Project is

639 managed by Genomics England Limited (a wholly owned company of the Department of

640 Health and Social Care). The 100,000 Genomes Project is funded by the National Institute
641 for Health Research and NHS England. The Wellcome Trust, Cancer Research UK and the
642 Medical Research Council have also funded research infrastructure. The 100,000 Genomes
643 Project uses data provided by patients and collected by the National Health Service as part
644 of their care and support. This research was funded, in whole or in part, by Wellcome. A CC
645 BY or equivalent licence is applied to the Author Accepted Manuscript, in accordance with
646 Wellcome open access conditions.

647

648 **Data and Code Availability**

649

650 The DDD study data are available under managed access from the European Genome-
651 phenome Archive (Study ID EGAS00001000775), and likely diagnostic variants are available
652 open access in DECIPHER. Code used for modelling case and population variants on the
653 MEF2C protein structure can be found here: https://github.com/matteoferla/MEF2C_analysis

654

655

656

657 **Web resources**

- 658 Online Mendelian Inheritance in Man (<http://www.omim.org>)
- 659 Gene-2-phenotype (<https://www.ebi.ac.uk/gene2phenotype/downloads>)
- 660 Matched Annotation between NCBI and EBI project information (MANE;
661 <https://www.ncbi.nlm.nih.gov/refseq/MANE/>)
- 662 MANE data download (ftp://ftp.ncbi.nlm.nih.gov/refseq/MANE/MANE_human/release_0.91/)
- 663 Genomics England 100,000 Genomes Project de novo call set
664 (<https://cnfl.extge.co.uk/display/GERE/De+novo+variant+research+dataset>)
- 665 UCSC table browser (<https://genome.ucsc.edu/cgi-bin/hgTables>)
- 666 Code for MEF2C protein modelling (https://github.com/matteoferla/MEF2C_analysis)
- 667 Interactive protein structure browser (<https://michelangelo.sgc.ox.ac.uk/r/mef2c>)

668

669 **Declaration of Interests**

670

- 671 K.J.K. is a consultant for Vor Biopharma. J.J. and K.R. are employees of GeneDx, Inc. K.R.
672 holds shares in Opko Health, Inc. B.D.Z. is a member of the speaker's bureau for Biogen,
673 Neurelis, and Supernus. S.A.C. is co-founder and shareholder of Enleofen Bio Pte Ltd.
674 M.E.H. is co-founder, shareholder, consultant, and non-executive director of Congenica Ltd.
675 All other authors declare no competing interests.

676

677 **References**

- 678 1. Short, P.J., McRae, J.F., Gallone, G., Sifrim, A., Won, H., Geschwind, D.H., Wright, C.F.,
679 Firth, H.V., FitzPatrick, D.R., Barrett, J.C., et al. (2018). De novo mutations in regulatory
680 elements in neurodevelopmental disorders. *Nature* 555, 611–616.
- 681 2. Spielmann, M., and Mundlos, S. (2016). Looking beyond the genes: the role of non-coding
682 variants in human disease. *Hum. Mol. Genet.* 25, R157–R165.
- 683 3. An, J.-Y., Lin, K., Zhu, L., Werling, D.M., Dong, S., Brand, H., Wang, H.Z., Zhao, X.,
684 Schwartz, G.B., Collins, R.L., et al. (2018). Genome-wide de novo risk score implicates
685 promoter variation in autism spectrum disorder. *Science* 362,.
- 686 4. Chong, J.X., Buckingham, K.J., Jhangiani, S.N., Boehm, C., Sobreira, N., Smith, J.D.,

- 687 Harrell, T.M., McMillin, M.J., Wiszniewski, W., Gambin, T., et al. (2015). The Genetic Basis
688 of Mendelian Phenotypes: Discoveries, Challenges, and Opportunities. *Am. J. Hum. Genet.*
689 *97*, 199–215.
- 690 5. Kaplanis, J., Samocha, K.E., Wiel, L., Zhang, Z., Arvai, K.J., Eberhardt, R.Y., Gallone, G.,
691 Lelieveld, S.H., Martin, H.C., McRae, J.F., et al. (2020). Evidence for 28 genetic disorders
692 discovered by combining healthcare and research data. *Nature* *586*, 757–762.
- 693 6. Srivastava, S., Love-Nichols, J.A., Dies, K.A., Ledbetter, D.H., Martin, C.L., Chung, W.K.,
694 Firth, H.V., Frazier, T., Hansen, R.L., Prock, L., et al. (2019). Meta-analysis and
695 multidisciplinary consensus statement: exome sequencing is a first-tier clinical diagnostic
696 test for individuals with neurodevelopmental disorders. *Genet. Med.* *21*, 2413–2421.
- 697 7. Mignone, F., Gissi, C., Liuni, S., and Pesole, G. (2002). Untranslated regions of mRNAs.
698 *Genome Biol* *3*, reviews0004.1.
- 699 8. Mortimer, S.A., Kidwell, M.A., and Doudna, J.A. (2014). Insights into RNA structure and
700 function from genome-wide studies. *Nature Reviews Genetics* *15*, 469–479.
- 701 9. Wanke, K.A., Devanna, P., and Vernes, S.C. (2018). Understanding Neurodevelopmental
702 Disorders: The Promise of Regulatory Variation in the 3'UTRome. *Biol. Psychiatry* *83*, 548–
703 557.
- 704 10. Chatterjee, S., and Pal, J.K. (2009). Role of 5'- and 3'-untranslated regions of mRNAs in
705 human diseases. *Biology of the Cell* *101*, 251–262.
- 706 11. Whiffin, N., Genome Aggregation Database Production Team, Karczewski, K.J., Zhang,
707 X., Chothani, S., Smith, M.J., Gareth Evans, D., Roberts, A.M., Quaife, N.M., Schafer, S., et
708 al. (2020). Characterising the loss-of-function impact of 5' untranslated region variants in
709 15,708 individuals. *Nature Communications* *11*,.
- 710 12. de Lima, R.L.L.F., Hoper, S.A., Ghassibe, M., Cooper, M.E., Rorick, N.K., Kondo, S.,
711 Katz, L., Marazita, M.L., Compton, J., Bale, S., et al. (2009). Prevalence and nonrandom
712 distribution of exonic mutations in interferon regulatory factor 6 in 307 families with Van der
713 Woude syndrome and 37 families with popliteal pterygium syndrome. *Genet. Med.* *11*, 241–
714 247.
- 715 13. Kozak, M. (1987). An analysis of 5'-noncoding sequences from 699 vertebrate
716 messenger RNAs. *Nucleic Acids Res.* *15*, 8125–8148.
- 717 14. Noderer, W.L., Flockhart, R.J., Bhaduri, A., Diaz de Arce, A.J., Zhang, J., Khavari, P.A.,
718 and Wang, C.L. (2014). Quantitative analysis of mammalian translation initiation sites by
719 FACS-seq. *Mol. Syst. Biol.* *10*, 748.
- 720 15. Sample, P.J., Wang, B., Reid, D.W., Presnyak, V., McFadyen, I.J., Morris, D.R., and
721 Seelig, G. (2019). Human 5' UTR design and variant effect prediction from a massively
722 parallel translation assay. *Nat. Biotechnol.* *37*, 803–809.
- 723 16. Deciphering Developmental Disorders Study (2015). Large-scale discovery of novel
724 genetic causes of developmental disorders. *Nature* *519*, 223–228.
- 725 17. Firth, H.V., Richards, S.M., Bevan, A.P., Clayton, S., Corpas, M., Rajan, D., Van Vooren,
726 S., Moreau, Y., Pettett, R.M., and Carter, N.P. (2009). DECIPHER: Database of
727 Chromosomal Imbalance and Phenotype in Humans Using Ensembl Resources. *Am. J.*
728 *Hum. Genet.* *84*, 524–533.

- 729 18. Ramu, A., Noordam, M.J., Schwartz, R.S., Wuster, A., Hurles, M.E., Cartwright, R.A.,
730 and Conrad, D.F. (2013). DeNovoGear: de novo indel and point mutation discovery and
731 phasing. *Nat. Methods* 10, 985–987.
- 732 19. McLaren, W., Gil, L., Hunt, S.E., Riat, H.S., Ritchie, G.R.S., Thormann, A., Flicek, P.,
733 and Cunningham, F. (2016). The Ensembl Variant Effect Predictor. *Genome Biol.* 17, 122.
- 734 20. Thormann, A., Halachev, M., McLaren, W., Moore, D.J., Svinti, V., Campbell, A., Kerr,
735 S.M., Tischkowitz, M., Hunt, S.E., Dunlop, M.G., et al. (2019). Flexible and scalable
736 diagnostic filtering of genomic variants using G2P with Ensembl VEP. *Nat. Commun.* 10,
737 2373.
- 738 21. Deciphering Developmental Disorders Study (2017). Prevalence and architecture of de
739 novo mutations in developmental disorders. *Nature* 542, 433–438.
- 740 22. (dgt), T.F.C.A.T.R.P.A.C., and The FANTOM Consortium and the RIKEN PMI and CLST
741 (DGT) (2014). A promoter-level mammalian expression atlas. *Nature* 507, 462–470.
- 742 23. Nellore, A., Jaffe, A.E., Fortin, J.-P., Alquicira-Hernández, J., Collado-Torres, L., Wang,
743 S., Phillips, R.A., III, Karbhari, N., Hansen, K.D., Langmead, B., et al. (2016). Human splicing
744 diversity and the extent of unannotated splice junctions across human RNA-seq samples on
745 the Sequence Read Archive. *Genome Biol.* 17, 266.
- 746 24. Consortium, T.G., and The GTEx Consortium (2020). The GTEx Consortium atlas of
747 genetic regulatory effects across human tissues. *Science* 369, 1318–1330.
- 748 25. Yates, A.D., Achuthan, P., Akanni, W., Allen, J., Allen, J., Alvarez-Jarreta, J., Amode,
749 M.R., Armean, I.M., Azov, A.G., Bennett, R., et al. (2020). Ensembl 2020. *Nucleic Acids Res.*
750 48, D682–D688.
- 751 26. Caulfield, M., Davies, J., Dennys, M., Elbahy, L., Fowler, T., Hill, S., Hubbard, T., Jostins,
752 L., Maltby, N., Mahon-Pearson, J., et al. (2019). The National Genomics Research and
753 Healthcare Knowledgebase (figshare).
- 754 27. Karczewski, K.J., Francioli, L.C., Tiao, G., Cummings, B.B., Alföldi, J., Wang, Q., Collins,
755 R.L., Laricchia, K.M., Ganna, A., Birnbaum, D.P., et al. (2020). The mutational constraint
756 spectrum quantified from variation in 141,456 humans. *Nature* 581, 434–443.
- 757 28. Zhang, X., Wakeling, M., Ware, J., and Whiffin, N. (2020). Annotating high-impact
758 5'untranslated region variants with the UTRannotator. *Bioinformatics*.
- 759 29. Naya, F.J., Wu, C., Richardson, J.A., Overbeek, P., and Olson, E.N. (1999).
760 Transcriptional activity of MEF2 during mouse embryogenesis monitored with a MEF2-
761 dependent transgene. *Development* 126, 2045–2052.
- 762 30. Wu, H., Rothermel, B., Kanatous, S., Rosenberg, P., Naya, F.J., Shelton, J.M.,
763 Hutcheson, K.A., DiMaio, J.M., Olson, E.N., Bassel-Duby, R., et al. (2001). Activation of
764 MEF2 by muscle activity is mediated through a calcineurin-dependent pathway. *EMBO J.* 20,
765 6414–6423.
- 766 31. Fromer, M., Moran, J.L., Chambert, K., Banks, E., Bergen, S.E., Ruderfer, D.M.,
767 Handsaker, R.E., McCarroll, S.A., O'Donovan, M.C., Owen, M.J., et al. (2012). Discovery
768 and statistical genotyping of copy-number variation from whole-exome sequencing depth.
769 *Am. J. Hum. Genet.* 91, 597–607.
- 770 32. Packer, J.S., Maxwell, E.K., O'Dushlaine, C., Lopez, A.E., Dewey, F.E., Chernomorsky,

- 771 R., Baras, A., Overton, J.D., Habegger, L., and Reid, J.G. (2016). CLAMMS: a scalable
772 algorithm for calling common and rare copy number variants from exome sequencing data.
773 *Bioinformatics* 32, 133–135.
- 774 33. Backenroth, D., Homsy, J., Murillo, L.R., Glessner, J., Lin, E., Brueckner, M., Lifton, R.,
775 Goldmuntz, E., Chung, W.K., and Shen, Y. (2014). CANOES: detecting rare copy number
776 variants from whole exome sequencing data. *Nucleic Acids Res.* 42, e97.
- 777 34. Santelli, E., and Richmond, T.J. (2000). Crystal structure of MEF2A core bound to DNA
778 at 1.5 Å resolution. *J. Mol. Biol.* 297, 437–449.
- 779 35. Bienert, S., Waterhouse, A., de Beer, T.A.P., Tauriello, G., Studer, G., Bordoli, L., and
780 Schwede, T. (2017). The SWISS-MODEL Repository—new features and functionality. *Nucleic
781 Acids Res.* 45, D313–D319.
- 782 36. Lei, X., Kou, Y., Fu, Y., Rajashekar, N., Shi, H., Wu, F., Xu, J., Luo, Y., and Chen, L.
783 (2018). The Cancer Mutation D83V Induces an α -Helix to β -Strand Conformation Switch in
784 MEF2B. *J. Mol. Biol.* 430, 1157–1172.
- 785 37. Chaudhury, S., Lyskov, S., and Gray, J.J. (2010). PyRosetta: a script-based interface for
786 implementing molecular modeling algorithms using Rosetta. *Bioinformatics* 26, 689–691.
- 787 38. Conway, P., Tyka, M.D., DiMaio, F., Konerding, D.E., and Baker, D. (2014). Relaxation
788 of backbone bond geometry improves protein energy landscape modeling. *Protein Sci.* 23,
789 47–55.
- 790 39. Alford, R.F., Leaver-Fay, A., Jeliaskov, J.R., O’Meara, M.J., DiMaio, F.P., Park, H.,
791 Shapovalov, M.V., Renfrew, P.D., Mulligan, V.K., Kappel, K., et al. (2017). The Rosetta All-
792 Atom Energy Function for Macromolecular Modeling and Design. *J. Chem. Theory Comput.*
793 13, 3031–3048.
- 794 40. Huang, P.-S., Ban, Y.-E.A., Richter, F., Andre, I., Vernon, R., Schief, W.R., and Baker,
795 D. (2011). RosettaRemodel: a generalized framework for flexible backbone protein design.
796 *PLoS One* 6, e24109.
- 797 41. Ferla, M.P., Pagnamenta, A.T., Damerell, D., Taylor, J.C., and Marsden, B.D. (2020).
798 MichelaNglo: sculpting protein views on web pages without coding. *Bioinformatics* 36, 3268–
799 3270.
- 800 42. Samocha, K.E., Kosmicki, J.A., Karczewski, K.J., O’Donnell-Luria, A.H., Pierce-Hoffman,
801 E., MacArthur, D.G., Neale, B.M., and Daly, M.J. (2017). Regional missense constraint
802 improves variant deleteriousness prediction (bioRxiv).
- 803 43. Le Meur, N., Holder-Espinasse, M., Jaillard, S., Goldenberg, A., Joriot, S., Amati-
804 Bonneau, P., Guichet, A., Barth, M., Charollais, A., Journal, H., et al. (2010). MEF2C
805 haploinsufficiency caused by either microdeletion of the 5q14.3 region or mutation is
806 responsible for severe mental retardation with stereotypic movements, epilepsy and/or
807 cerebral malformations. *J. Med. Genet.* 47, 22–29.
- 808 44. Suri, M., Evers, J.M.G., Laskowski, R.A., O’Brien, S., Baker, K., Clayton-Smith, J., Dabir,
809 T., Josifova, D., Joss, S., Kerr, B., et al. (2017). Protein structure and phenotypic analysis of
810 pathogenic and population missense variants in. *Mol Genet Genomic Med* 5, 495–507.
- 811 45. Zerbino, D.R., Wilder, S.P., Johnson, N., Juettemann, T., and Flicek, P.R. (2015). The
812 ensembl regulatory build. *Genome Biol.* 16, 56.

- 813 46. ENCODE Project Consortium, Moore, J.E., Purcaro, M.J., Pratt, H.E., Epstein, C.B.,
814 Shores, N., Adrian, J., Kawli, T., Davis, C.A., Dobin, A., et al. (2020). Expanded
815 encyclopaedias of DNA elements in the human and mouse genomes. *Nature* 583, 699–710.
- 816 47. Collins, R.L., Brand, H., Karczewski, K.J., Zhao, X., Alföldi, J., Francioli, L.C., Khera,
817 A.V., Lowther, C., Gauthier, L.D., Wang, H., et al. (2020). A structural variation reference for
818 medical and population genetics. *Nature* 581, 444–451.
- 819 48. D'haene, E., Bar-Yaacov, R., Bariah, I., Vantomme, L., Van Loo, S., Cobos, F.A.,
820 Verboom, K., Eshel, R., Alatawna, R., Menten, B., et al. (2019). A neuronal enhancer
821 network upstream of MEF2C is compromised in patients with Rett-like characteristics. *Hum.*
822 *Mol. Genet.* 28, 818–827.
- 823 49. Zweier, M., Gregor, A., Zweier, C., Engels, H., Sticht, H., Wohlleber, E., Bijlsma, E.K.,
824 Holder, S.E., Zenker, M., Rossier, E., et al. (2010). Mutations in MEF2C from the 5q14.3q15
825 microdeletion syndrome region are a frequent cause of severe mental retardation and
826 diminish MECP2 and CDKL5 expression. *Hum. Mutat.* 31, 722–733.
- 827 50. Biennu, T., Diebold, B., Chelly, J., and Isidor, B. (2013). Refining the phenotype
828 associated with MEF2C point mutations. *Neurogenetics* 14, 71–75.
- 829 51. Vrečar, I., Innes, J., Jones, E.A., Kingston, H., Reardon, W., Kerr, B., Clayton-Smith, J.,
830 and Douzgou, S. (2017). Further Clinical Delineation of the MEF2C Haploinsufficiency
831 Syndrome: Report on New Cases and Literature Review of Severe Neurodevelopmental
832 Disorders Presenting with Seizures, Absent Speech, and Involuntary Movements. *J. Pediatr.*
833 *Genet.* 6, 129–141.
- 834 52. Molkenin, J.D., Black, B.L., Martin, J.F., and Olson, E.N. (1996). Mutational analysis of
835 the DNA binding, dimerization, and transcriptional activation domains of MEF2C. *Mol. Cell.*
836 *Biol.* 16, 2627–2636.
- 837 53. Potthoff, M.J., and Olson, E.N. (2007). MEF2: a central regulator of diverse
838 developmental programs. *Development* 134, 4131–4140.
- 839 54. Landrum, M.J., Lee, J.M., Benson, M., Brown, G., Chao, C., Chitipiralla, S., Gu, B., Hart,
840 J., Hoffman, D., Hoover, J., et al. (2016). ClinVar: public archive of interpretations of clinically
841 relevant variants. *Nucleic Acids Res.* 44, D862–D868.
- 842 55. Bennett, C.L., Brunkow, M.E., Ramsdell, F., O'Briant, K.C., Zhu, Q., Fuleihan, R.L.,
843 Shigeoka, A.O., Ochs, H.D., and Chance, P.F. (2001). A rare polyadenylation signal
844 mutation of the FOXP3 gene (AAUAAA-->AAUGAA) leads to the IPEX syndrome.
845 *Immunogenetics* 53, 435–439.
- 846 56. Devanna, P., Chen, X.S., Ho, J., Gajewski, D., Smith, S.D., Gialluisi, A., Francks, C.,
847 Fisher, S.E., Newbury, D.F., and Vernes, S.C. (2018). Next-gen sequencing identifies non-
848 coding variation disrupting miRNA-binding sites in neurological disorders. *Molecular*
849 *Psychiatry* 23, 1375–1384.
- 850 57. Reamon-Buettner, S.M., Cho, S.-H., and Borlak, J. (2007). Mutations in the 3'-
851 untranslated region of GATA4 as molecular hotspots for congenital heart disease (CHD).
852 *BMC Med. Genet.* 8, 38.
- 853 58. Araujo, P.R., Yoon, K., Ko, D., Smith, A.D., Qiao, M., Suresh, U., Burns, S.C., and
854 Penalva, L.O.F. (2012). Before It Gets Started: Regulating Translation at the 5' UTR.
855 *Comparative and Functional Genomics* 2012, 1–8.

856 59. Leppék, K., Das, R., and Barna, M. (2018). Functional 5' UTR mRNA structures in
857 eukaryotic translation regulation and how to find them. *Nature Reviews Molecular Cell*
858 *Biology* 19, 158–174.

859

860

861 **Figure legends**

862

863 **Figure 1:** Schematic of the wild-type *MEF2C* gene (a) and the position and effect of uAUG-
864 creating variants identified as *de novo* in developmental disorder cases (b and c) and in
865 gnomAD population controls (d). The two 5'UTR exons are shown as light grey boxes,
866 separated by an intron shown as a thinner broken grey line. Upstream open reading frames
867 (uORFs) already present in the sequence are shown in green. Variant positions are
868 represented by arrows. New ORFs created by the variants are shown as blue boxes. (b)
869 Two case variants create ORFs that overlap the coding sequence (CDS) out-of-frame
870 (oORF-creating). If translation initiates at the uAUG, the ribosome will not translate the CDS.
871 (c) Two recurrent case variants create uAUGs in-frame with the CDS. If translation initiates
872 at this uAUG, an elongated protein will be translated. (d) Two variants identified in gnomAD
873 create uORFs far upstream of the CDS which would not be predicted to disrupt translation of
874 the normal protein.

875

876 **Figure 2:** uAUG-creating variants decrease translation of *MEF2C* (a) or transactivation of
877 target genes (b). (a) *MEF2C* 5' UTR out-of-frame overlapping ORF (oORF)-creating variants
878 c.-103G>A and c.-66A>T (Figure 1b) reduce downstream luciferase expression relative to
879 wild-type (WT) 5' UTR in a translation reporter assay. Reduction is stronger for c.-66A>T
880 (moderate uAUG Kozak context) than for c.-103G>A (weak Kozak context). (b)
881 Overexpression of *MEF2C* with the WT 5' UTR/CDS induces expression of luciferase from a
882 *MEF2C*-dependent enhancer-luciferase reporter construct, relative to an empty pcDNA3.1
883 construct negative control. The *MEF2C* N-terminus-extending variants c.-26C>T (9 amino
884 acids) and c.-8C>T (3 amino acids; Figure 1c) both reduce transactivation. For (a) and (b)
885 bars are coloured by Kozak consensus: yellow=weak; orange=moderate; red=strong.
886 Luciferase expression was normalised for transfection efficiency.

887

888 **Figure 3:** (a) The N-terminal region of *MEF2C* is highly constrained for missense variants in
889 gnomAD (obs/exp=0.069), with much lower constraint across the rest of the protein

890 (obs/exp=0.41). This region of high constraint correlates with the location of the majority of
891 *de novo* missense variants identified in DD cases (red circles), while gnomAD variants are
892 mostly outside of this N-terminal region (grey circles). (b) The N-terminal portion of the
893 MEF2C dimer [1-92], modelled using structures of the human MEF2A dimer which is 96%
894 identical in sequence to MEF2C, bound directly to its consensus DNA sequence. Side
895 chains of amino acids with pathogenic *de novo* missense variants from DDD, GeneDx and
896 ClinVar are shown in yellow, with gnomAD *MEF2C* missense variants in grey. Most
897 pathogenic missense variants either protrude directly into the DNA or are located in the
898 DNA-binding helix. In particular, the terminal amine (Gly2, top inset) along with Arg3 (bottom
899 inset) act as reader-heads for nucleobase specificity, which is likely disrupted in the N-
900 terminal extension variants (middle inset). All pathogenic and gnomAD variants can be
901 viewed in our interactive protein structure browser (see link in Web Resources). (c-d)
902 Missense variants from DD cases (DDD, GeneDx and ClinVar) are significantly more
903 disruptive to the interaction with DNA as measured by $\Delta\Delta G$ values (c) and closer to the
904 bound DNA molecule (d) than *MEF2A-D* variants in gnomAD (see online methods).

905

906 **Figure 4:** 5'UTRs of DDG2P haploinsufficient genes (red) are longer (a), and a higher
907 proportion have multiple exons (b) compared to 5'UTRs of all genes (light grey), and other
908 DDG2P genes (dark grey). Mean lengths for each gene set in (a) are shown as dotted lines.
909 (c) The coverage of 5'UTRs decays rapidly with distance from the CDS (x-axis truncated at
910 1000 bps). Note that these figures were calculated using exome sequence data from the
911 DDD study and may vary between different exome capture designs. (d) The position of DNA-
912 binding domains (including homeodomains, zinc-fingers, and specific DNA-binding domains)
913 in DDG2P haploinsufficient genes with respect to the N-terminus of the protein; MEF2C is
914 one of three proteins with a DNA-binding domain that starts within 10 bps of the N-terminus.

915

variant (GRCh37)	cDNA description (ENST00000504921.7)	variant effect	deletion size	kozak strength	proband ID(s)	proband count	gnomAD v3 AC
<i>uUAG-creating de novo variants discovered in probands with DD:</i>							
chr5:88119671 T>A	c.-66A>T	out-of-frame oORF created	-	moderate	1	1	-
chr5:88119708 C>T	c.-103G>A	out-of-frame oORF created	-	weak	2	1	-
chr5:88119613 G>A	c.-8C>T	CDS-elongating	-	strong	3,4,5	3	-
chr5:88119631 G>A	c.-26C>T	CDS-elongating	-	moderate	6,7,8	3	-
<i>uAUG-creating variant present in gnomAD:</i>							
chr5:88883052 G>A	c.-240C>T	uORF created	-	weak	-	0	1
chr5:88883059 G>A	c.-247C>T	uORF created	-	weak	-	0	6
-							
chr5:88133089-88427361 del	-	promoter and partial 5'UTR deletion	294kb	-	9	1	-
chr5:88123099-88220350 del	-	promoter and partial 5'UTR deletion	97kb	-	10	1	-

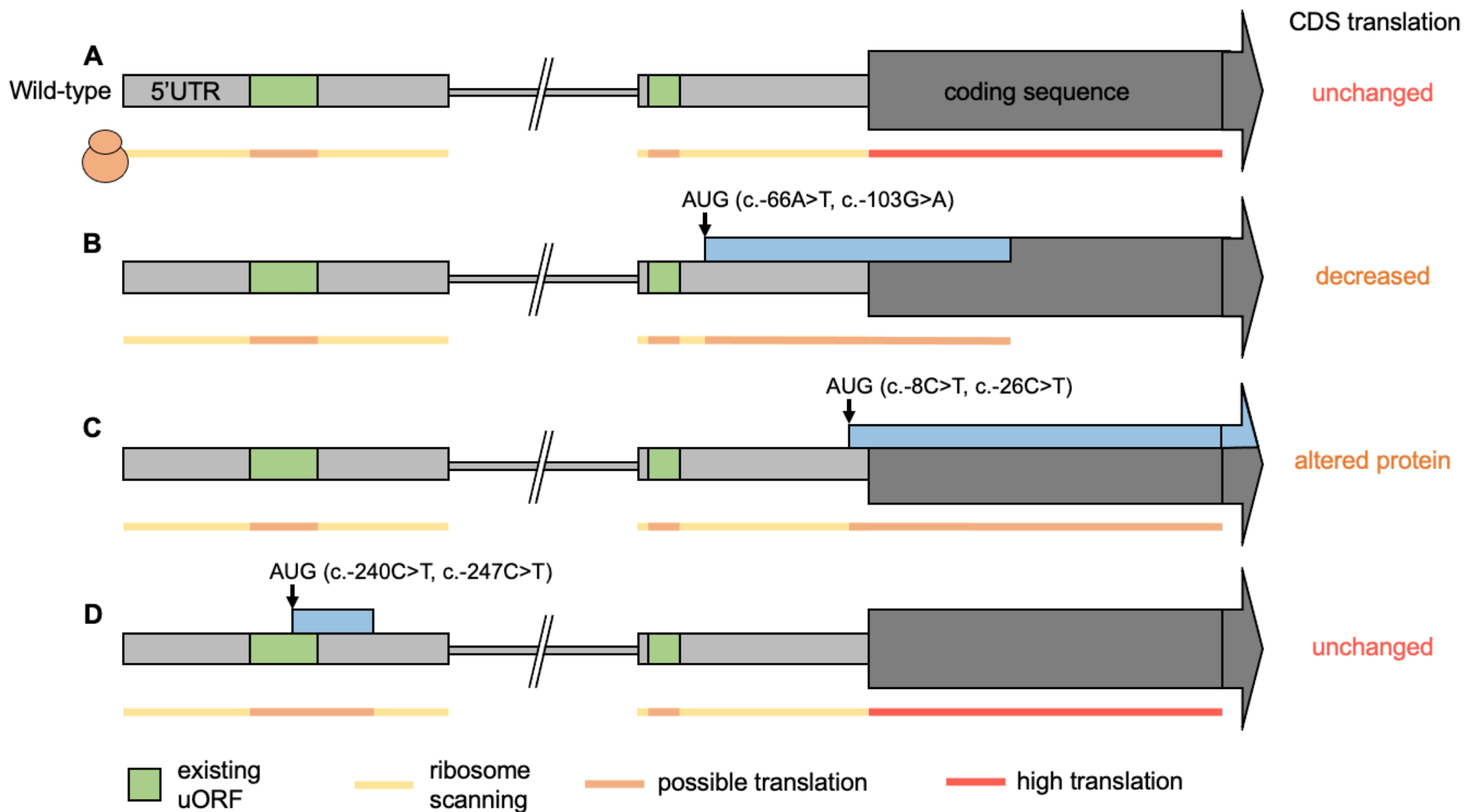
918

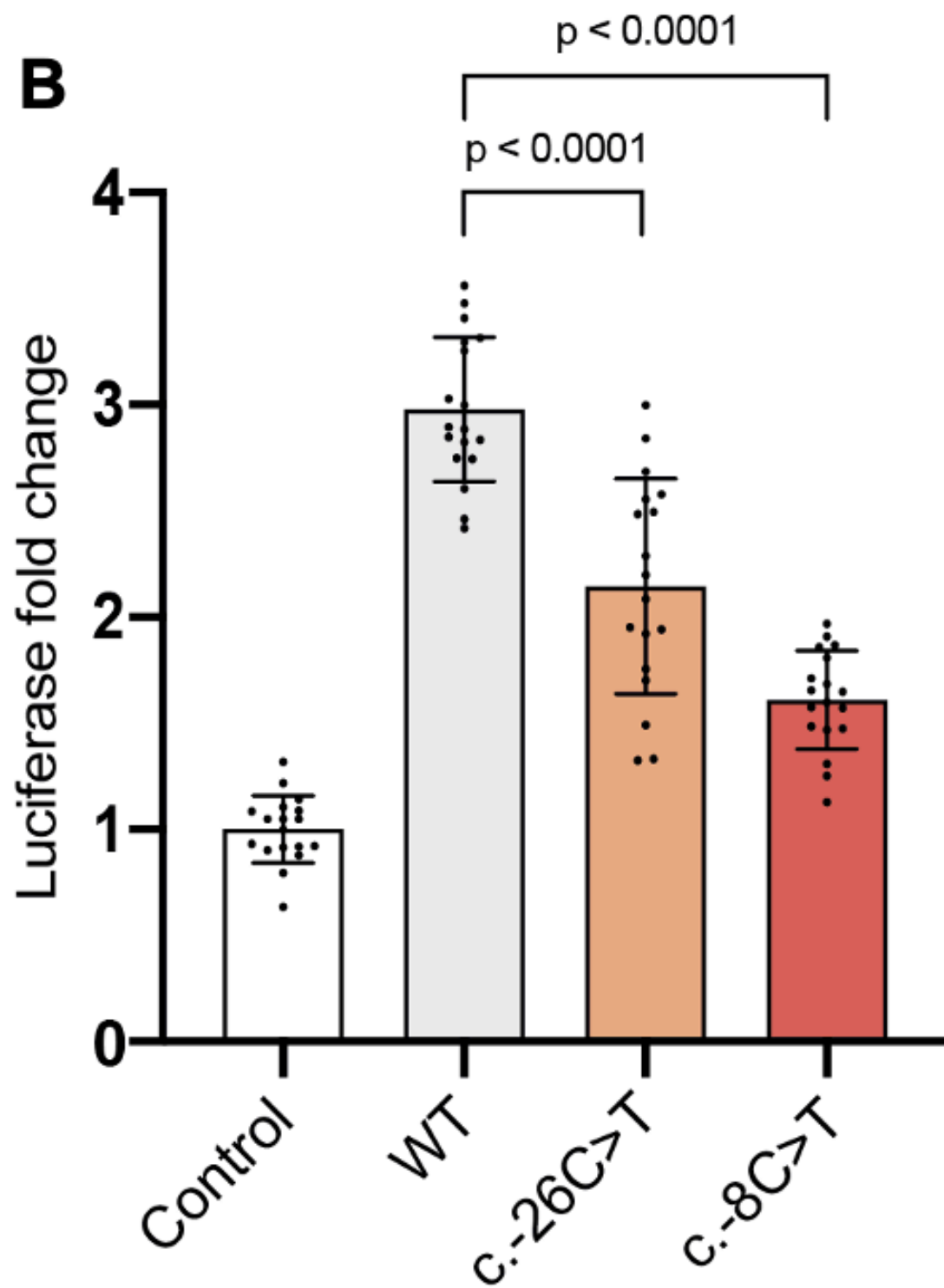
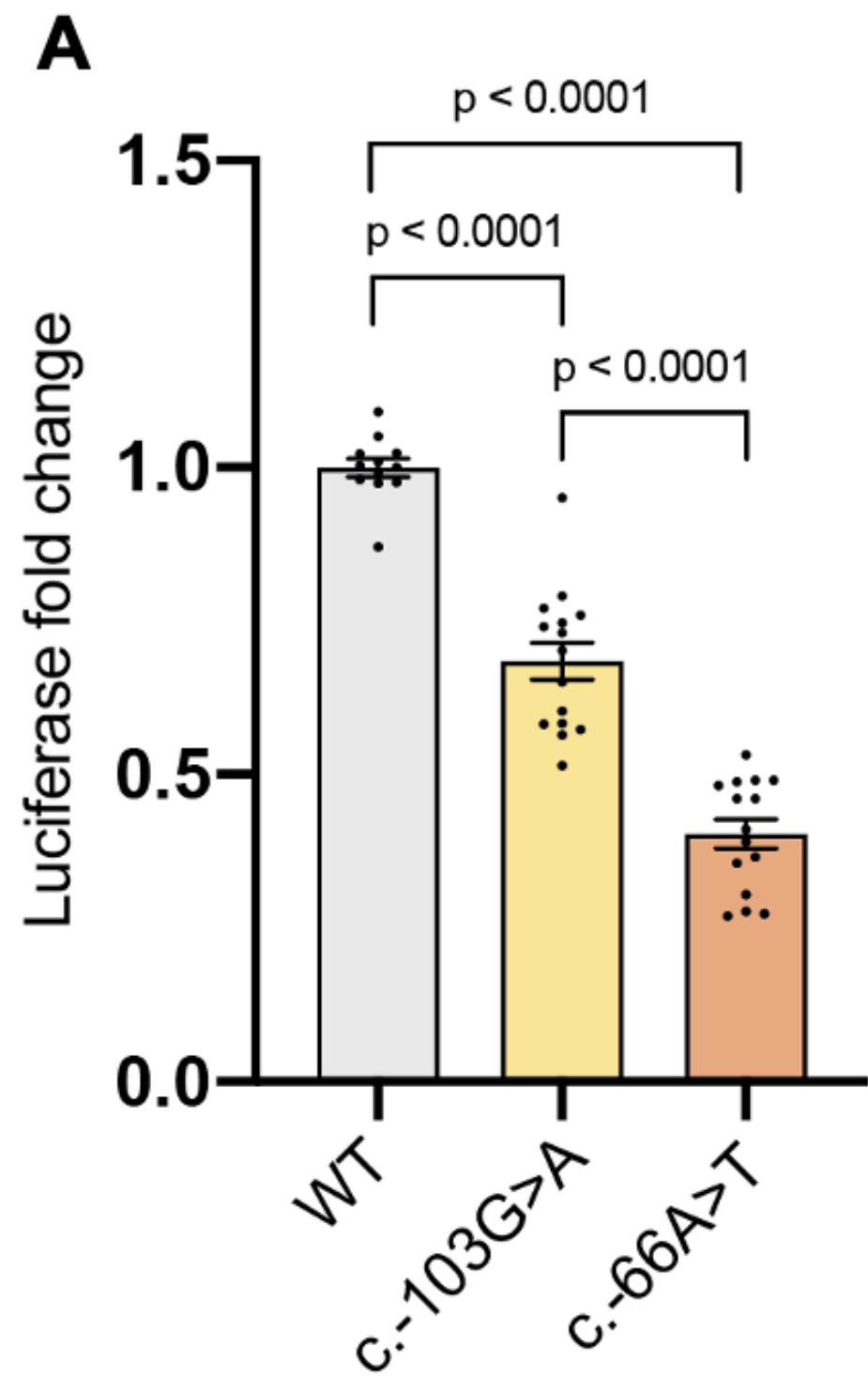
919 **Table 1:** Details of *MEF2C* uAUG-creating and upstream deletion variants discussed in this
 920 work. Shown are the four uAUG SNVs identified in DDD, uAUG SNVs observed in gnomAD
 921 v3.0, and non-coding CNVs found upstream of *MEF2C* in DDD. oORF = overlapping ORF;
 922 uORF = upstream ORF; AC = allele count. Proband IDs refer to those used in Table S2.

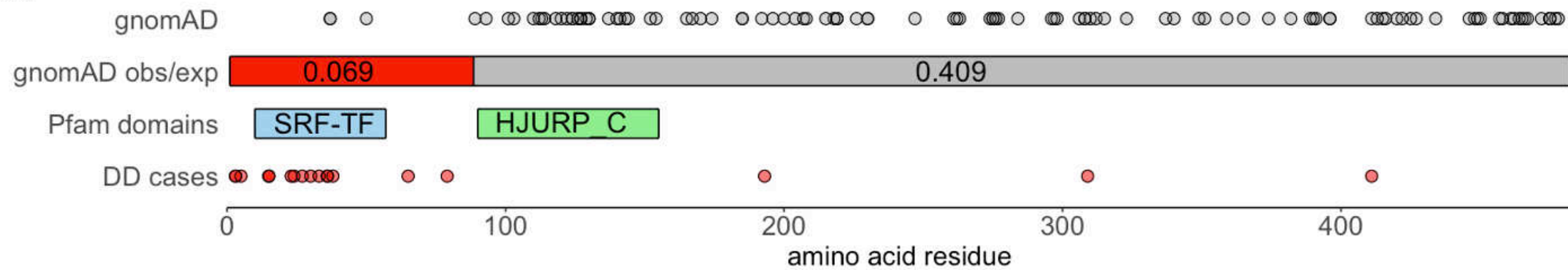
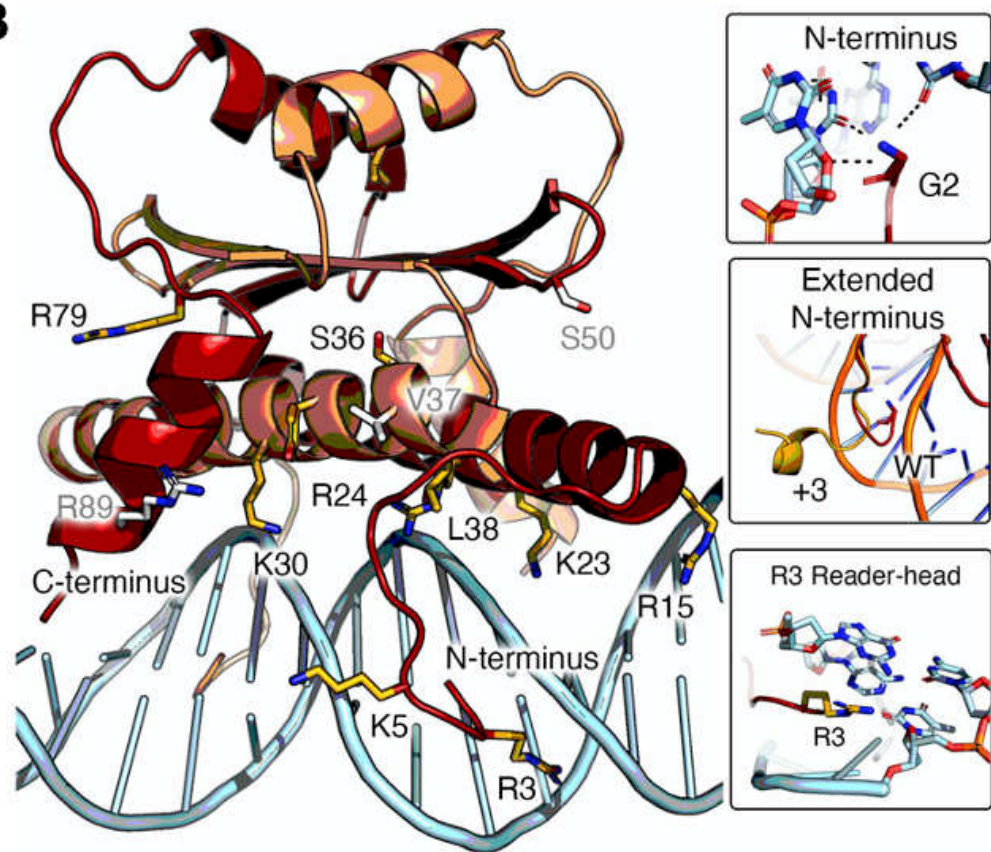
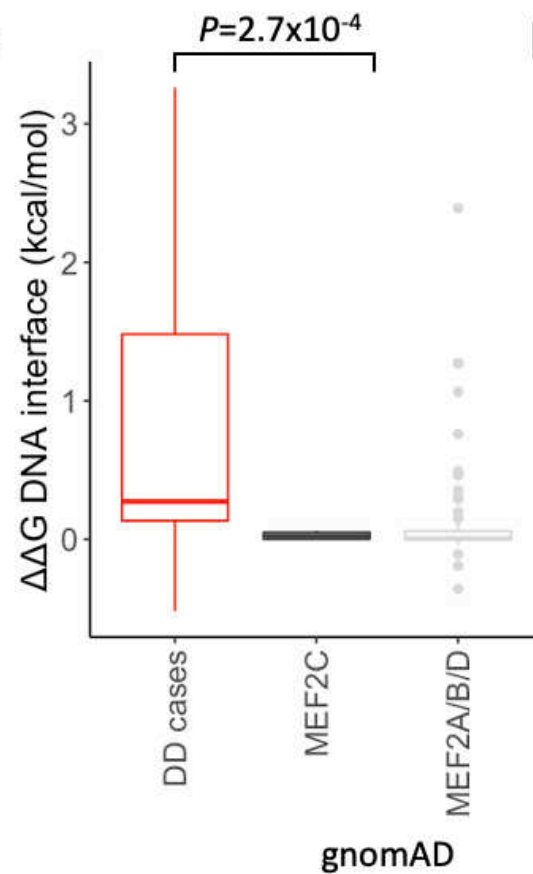
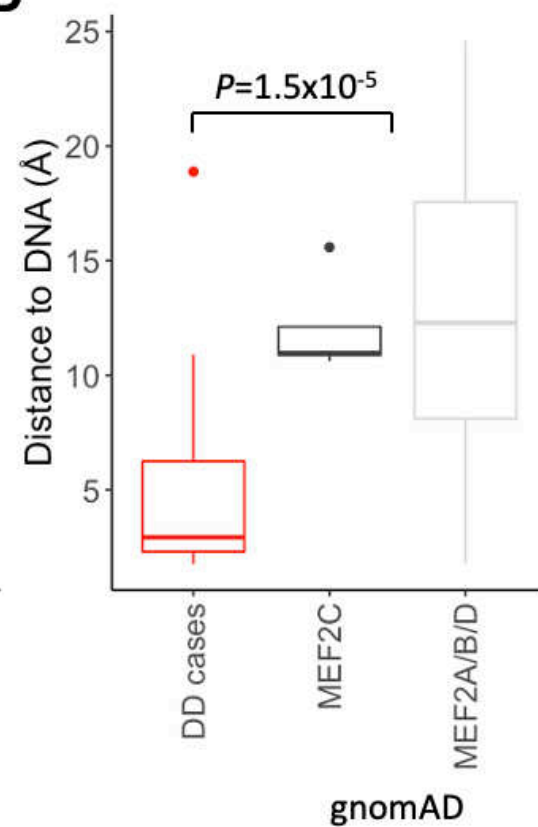
923

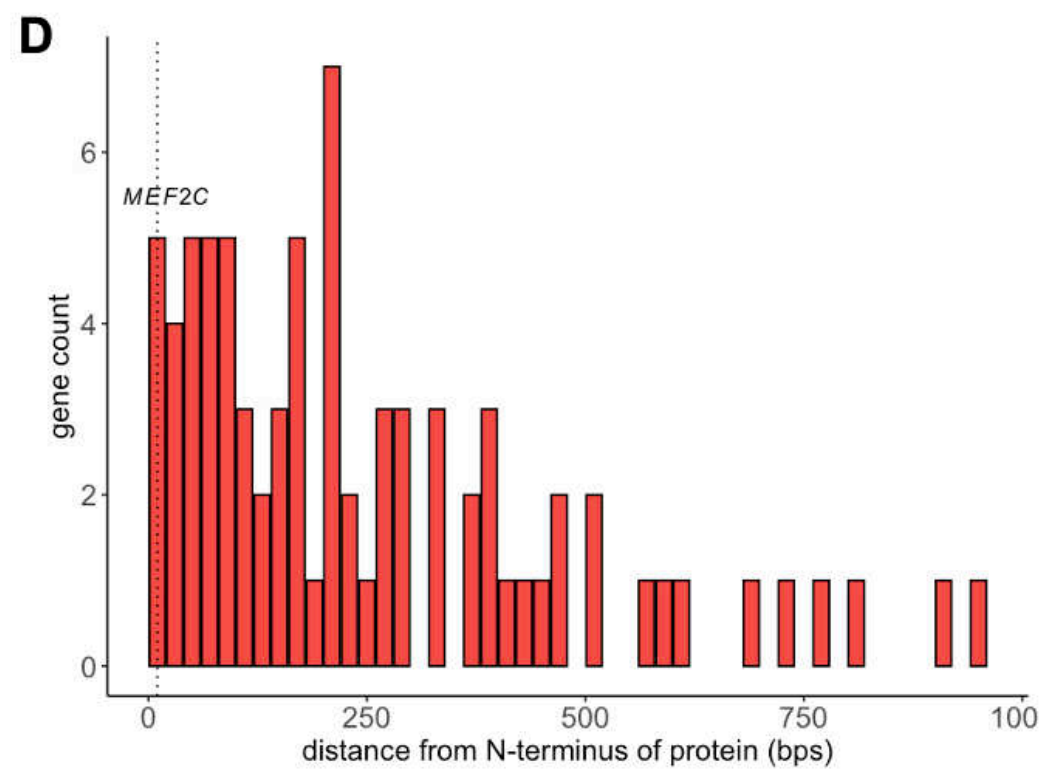
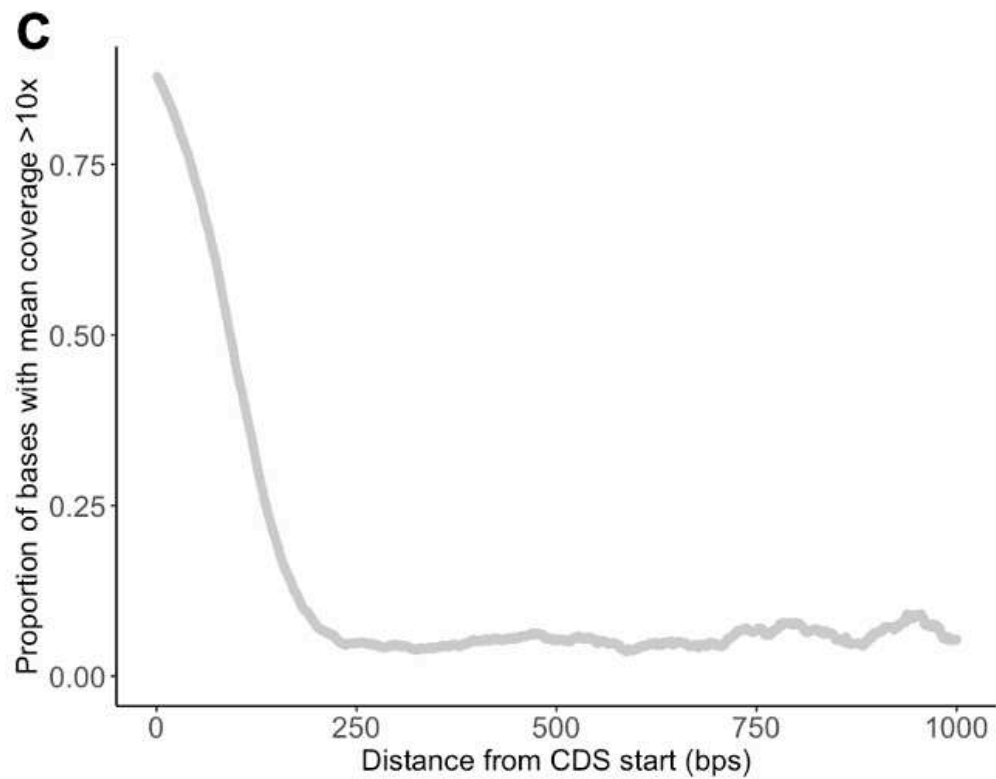
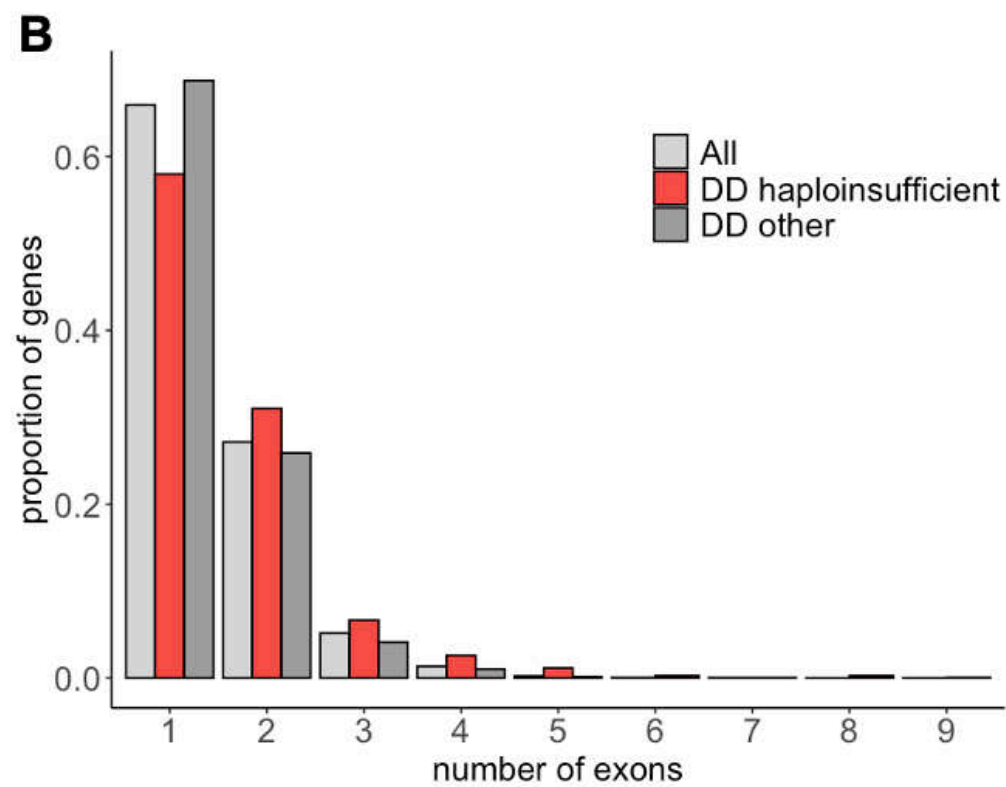
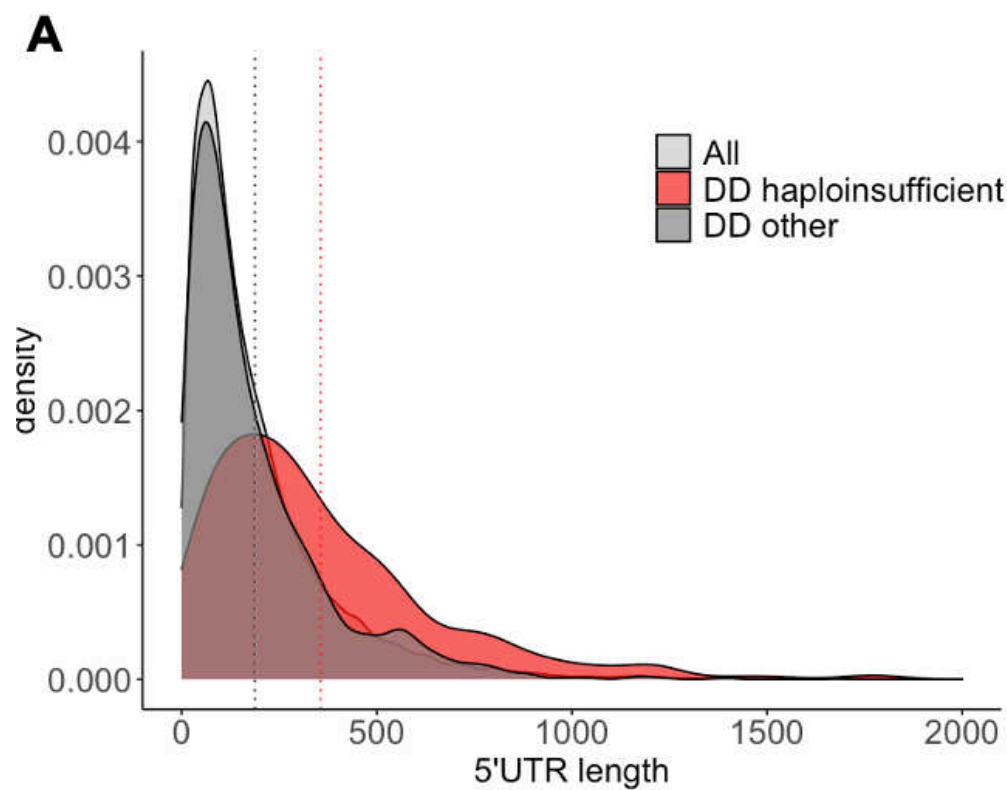
924

925





A**B****C****D**



Supplementary Data: Non-coding variants upstream of *MEF2C* cause severe developmental disorder through three distinct loss-of-function mechanisms

Figure S1: Two non-coding deletions remove the distal 5'UTR exon of *MEF2C* and the entire promoter sequence. Coding exons are shown in black with UTRs in red. The dotted line indicates the start of the coding sequence. The five deletions identified in DDD in the *MEF2C* region are shown as orange bars, the top two of which are entirely non-coding. A representative H3K4me3 dataset from ENCODE is plotted in blue across the top (GN12878) to show active promoter regions.

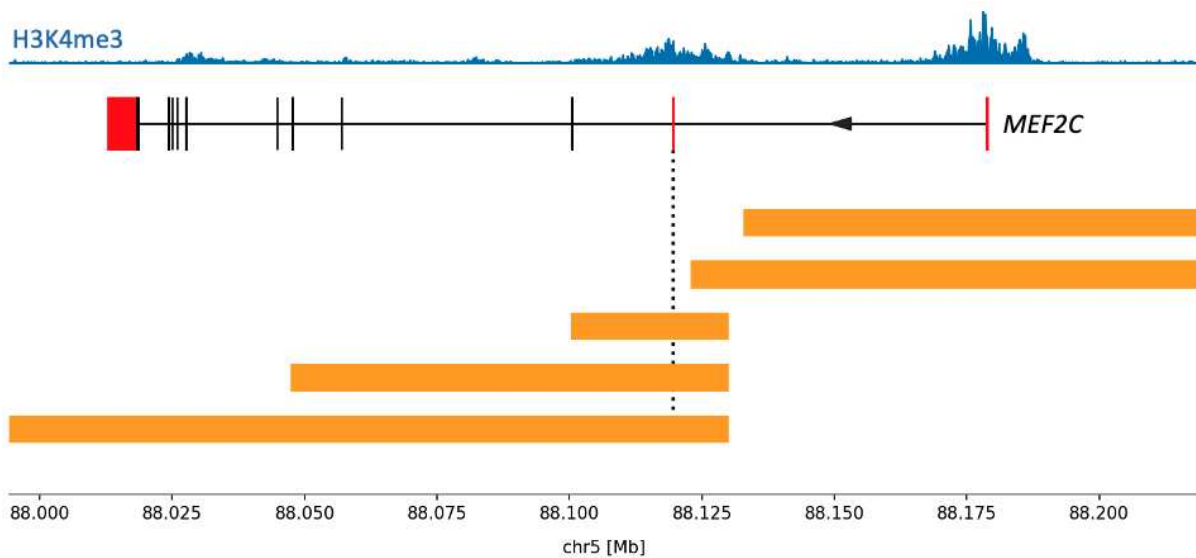


Figure S2: uAUG-creating variants do not alter RNA or protein levels. (A) Relative Gaussia luciferase (GLuc) RNA levels remain unchanged with each out-of-frame oORF-creating variant when normalised to RNA of secreted alkaline phosphatase (SEAP) transfection control. (B and C) The decreases in transactivation seen for the CDS-elongating variants c.-8C>T and c.-26C>T are not accompanied by a significant change in protein levels. For (A) and (C) bars are coloured by Kozak consensus: yellow = weak; orange = moderate; red = strong. ns = not significant.

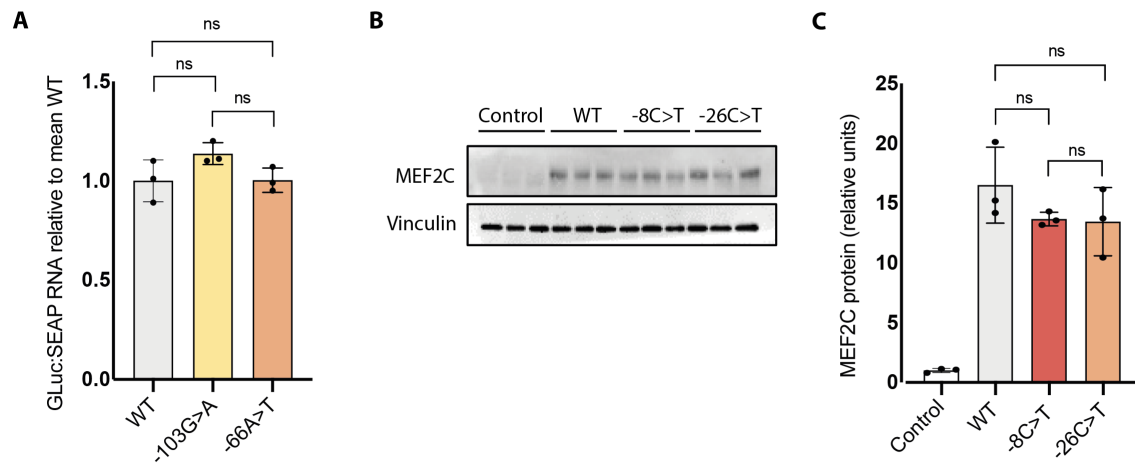


Figure S3: A single base mutation in the context surrounding the c.-103G>A variant which changes a weak Kozak consensus into a moderate consensus significantly reduces translational efficiency. (A) oORF-creating variants c.-103G>A and c.-66A>T reduce downstream luciferase expression relative to wild-type (WT) 5' UTR in a translation reporter assay. Reduction is stronger for c.-66A>T (moderate Kozak context) than for c.-103G>A (weak Kozak context). Modifying the context surrounding the c.-103G>A variant into a moderate Kozak context (as shown in B) reduces downstream luciferase expression compared to the unmodified vector. The translational efficiency of the modified vector is equivalent to the c.-66A>T variant which also has a moderate Kozak consensus. ns = not significant.

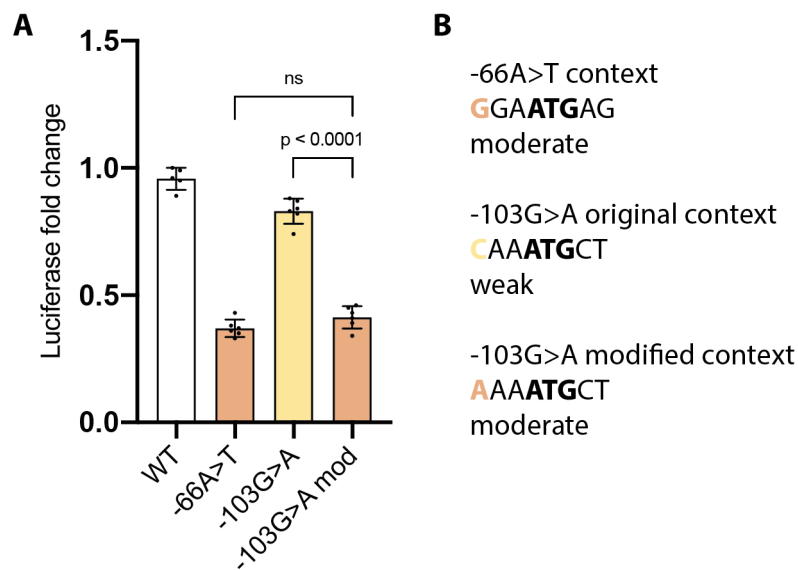


Figure S4: Protein sequence alignment of the four human myocyte enhancer factor 2 proteins (MEF2A-D). The Clustal-Omega default alignment function in UniProt for the first 92 N-terminal residues was used. Coloured by similarity; * = identical amino acids in all 4 proteins; : = similar amino acids in all 4 proteins.

```

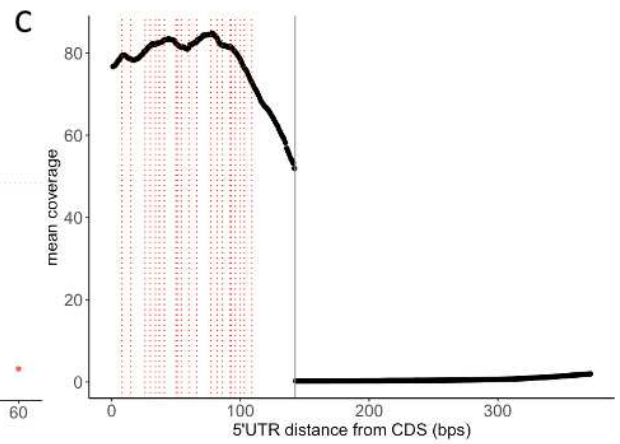
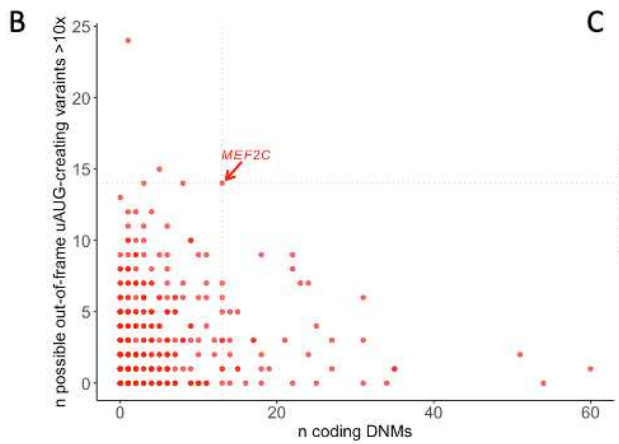
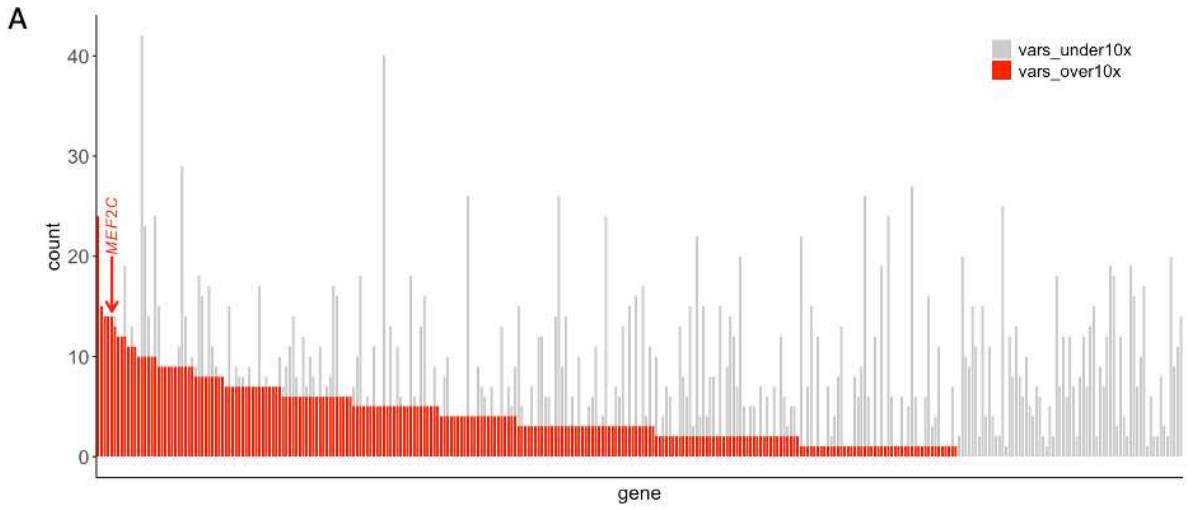
Q06413 MEF2C_HUMAN      1  MGRKKIQITRIMDERNRQVTFTKRKFGLMKKAYELSVLCDCEIALIIFNS*TKLFOYAST 60
Q02078 MEF2A_HUMAN      1  MGRKKIQITRIMDERNRQVTFTKRKFGLMKKAYELSVLCDCEIALIIFNSSNKLFOYAST 60
Q02080 MEF2B_HUMAN      1  MGRKKIQISRIIDORNROVTFTKRKFGLMKKAYELSVLCDCEIALIIFNSANR*LFQYAST 60
Q14814 MEF2D_HUMAN      1  MGRKKIQIQRIIDERNRQVTFTKRKFGLMKKAYELSVLCDCEIALIIFNHSNKLFOYAST 60
***** * * :*****

Q06413 MEF2C_HUMAN      61  DMDK*VLLKYTEYN*EPHESRTNSDIVE*TLRKKG 92
Q02078 MEF2A_HUMAN      61  DMDK*VLLKYTEYN*EPHESRTNSDIVE*ALNKKE 92
Q02080 MEF2B_HUMAN      61  DMDR*VLLKYTEYSEPHESRTNDILE*TLKRRG 92
Q14814 MEF2D_HUMAN      61  DMDK*VLLKYTEYN*EPHESRTNADI*LET*LRKKG 92
***:*****:*****:***:*:

```

Figure S5: Coverage of uAUG-creating sites of DD haploinsufficient genes and the *MEF2C* 5'UTR. (A) Stacked bar chart showing the count of all possible uAUG-creating variants that would create out-of-frame overlapping ORFs that are covered at mean $>10x$ (red), or $\leq 10x$ (grey) per gene. *MEF2C* has a high number of possible variants ($n=14$), all of which are well covered. (B) The number of well-covered uAUG-creating variants that would create out-of-frame overlapping ORFs plotted against the number of coding missense and protein-truncating *de novo* mutations (DNMs) per gene. *MEF2C* has both a high number of well-covered sites and a high diagnostic yield. (C) The mean coverage across the *MEF2C* 5'UTR. All possible uAUG-creating variants that would create either out-of-frame overlapping ORFs or CDS-elongations are plotted as dotted lines. The 5'UTR exon that is adjacent to the CDS is very well covered (mean $>50x$).

NB: (A) and (B) do not include CDS-elongating variants as these would not be predicted to cause loss-of-function unless there is an important N-terminal structure or functional domain.



Supplementary Tables

Table S1: List of haploinsufficient developmental disorder genes and their MANEv0.91 transcripts used for analysis.

Table S2: Clinical details for patients with non-coding MEF2C variants.

Table S3: List of missense variants identified in DD cases. ClinVar variants are filtered to only those identified as de novo or with experimental evidence. Protein changes are with respect to the Ensembl canonical transcript ENST00000340208.5.

Table S4: List of gnomAD v2.1.1 missense variants in MEF2 genes used from protein modelling. Protein changes are with respect to the Ensembl canonical transcript ENST00000340208.5.

Table S5: Residues in the structure of MEF2A and their direction with respect to the bound DNA.

Table S6: Comparing the proportion of DD and gnomAD variants that are in contact/pointing towards DNA to those that are distal or pointing away from the DNA-binding interface.

Table S7: Change in Gibbs free energy ($\Delta\Delta G$) of protein-DNA interaction and complex stability associated with missense variants in MEF2C.

Genomics England Research Consortium

John C. Ambrose¹; Prabhu Arumugam¹; Emma L. Baple¹; Marta Bleda¹; Freya Boardman-Pretty^{1,2}; Jeanne M. Boissiere¹; Christopher R. Boustred¹; Helen Brittain¹; Mark J. Caulfield^{1,2}; Georgia C. Chan¹; Clare E. H. Craig¹; Louise C. Daugherty¹; Anna de Burca¹; Andrew Devereau¹; Greg Elgar^{1,2}; Rebecca E. Foulger¹; Tom Fowler¹; Pedro Furió-Tarí¹; Adam Giess¹; Joanne M. Hackett¹; Dina Halai¹; Angela Hamblin¹; Shirley Henderson^{1,2}; James E. Holman¹; Tim J. P. Hubbard¹; Kristina ibáñez^{1,2}; Rob Jackson¹; Louise J. Jones^{1,2}; Dalia Kasperaviciute¹; Melis Kayikci¹; Athanasios Kousathanas¹; Lea Lahnstein¹; Kay Lawson¹; Sarah E. A. Leigh¹; Ivonne U. S. Leong¹; Javier F. Lopez¹; Fiona Maleady-Crowe¹; Joanne Mason¹; Ellen M. McDonagh^{1,2}; Loukas Moutsianas^{1,2}; Michael Mueller^{1,2}; Nirupa Murugaesu¹; Anna C. Need^{1,2}; Peter O'Donovan¹; Chris A. Odhams¹; Andrea Orioli¹; Christine Patch^{1,2}; Mariana Buongiorno Pereira¹; Daniel Perez-Gil¹; Dimitris Polychronopoulos¹; John Pullinger¹; Tahrima Rahim¹; Augusto Rendon¹; Pablo Riesgo-Ferreiro¹; Tim Rogers¹; Mina Rytén¹; Kevin Savage¹; Kushmita Sawant¹; Richard H. Scott¹; Afshan Siddiq¹; Alexander Sieghart¹; Damian Smedley^{1,2}; Katherine R. Smith^{1,2}; Samuel C. Smith¹; Alona Sosinsky^{1,2}; William Spooner¹; Helen E. Stevens¹; Alexander Stuckey¹; Razvan Sultana¹; Mélanie Tanguy¹; Ellen R. A. Thomas^{1,2}; Simon R. Thompson¹; Carolyn Tregidgo¹; Arianna Tucci^{1,2}; Emma Walsh¹; Sarah A. Watters¹; Matthew J. Welland¹; Eleanor Williams¹; Katarzyna Witkowska^{1,2}; Suzanne M. Wood^{1,2}; Magdalena Zarowiecki¹.

¹Genomics England, London, UK

²William Harvey Research Institute, Queen Mary University of London, London, EC1M 6BQ, UK.