

1 Pathogenicity and selective constraint on variation near 2 splice sites

3 AUTHORS

4 Jenny Lord¹, Giuseppe Gallone¹, Patrick J. Short¹, Jeremy F. McRae¹, Holly Ironfield¹, Elizabeth H.
5 Wynn¹, Sebastian S. Gerety¹, Liu He¹, Bronwyn Kerr^{2,3}, Diana S. Johnson⁴, Emma McCann⁵, Esther
6 Kinning⁶, Frances Flinter⁷, I. Karen Temple^{8,9}, Jill Clayton-Smith^{2,3}, Meriel McEntagart¹⁰, Sally Ann
7 Lynch¹¹, Shelagh Joss¹², Sofia Douzgou^{2,3}, Tabib Dabir¹³, Virginia Clowes¹⁴, Vivienne P. M.
8 McConnell¹³, Wayne Lam¹⁵, Caroline F. Wright¹⁶, David R. FitzPatrick^{1,15}, Helen V. Firth^{1,17}, Jeffrey
9 C. Barrett¹, Matthew E. Hurles¹, on behalf of the Deciphering Developmental Disorders study

10 AFFILIATIONS

11 ¹Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, CB10 1SA, UK

12 ²Manchester Centre for Genomic Medicine, St Mary's Hospital, Manchester University Hospitals NHS
13 Foundation Trust Manchester Academic Health Sciences Centre

14 ³Division of Evolution and Genomic Sciences School of Biological Sciences University of Manchester

15 ⁴Sheffield Clinical Genetics Service, Sheffield Children's Hospital, OPD2, Northern General Hospital,
16 Herries Road, Sheffield, S5 7AU

17 ⁵Liverpool Women's Hospital Foundation Trust, Crown Street, Liverpool, L8 7SS

18 ⁶West of Scotland Regional Genetics Service, NHS Greater Glasgow and Clyde, Institute of Medical
19 Genetics, Yorkhill Hospital, Glasgow G3 8SJ, UK

20 ⁷South East Thames Regional Genetics Centre, Guy's and St Thomas' NHS Foundation Trust, Guy's
21 Hospital, Great Maze Pond, London SE1 9RT, UK

22 ⁸Faculty of Medicine, University of Southampton, Institute of Developmental Sciences, Tremona
23 Road, Southampton SO16 6YD

1 ⁹Wessex Clinical Genetics Service, University Hospital Southampton, Princess Anne Hospital, Coxford
2 Road, Southampton SO16 5YA, UK

3 ¹⁰South West Thames Regional Genetics Centre, St George's Healthcare NHS Trust, St George's,
4 University of London, Cranmer Terrace, London SW17 0RE, UK

5 ¹¹Temple Street Children's Hospital, Dublin 1, Ireland

6 ¹²West of Scotland Regional Genetics Service, NHS Greater Glasgow & Clyde, Level 2, Laboratory
7 Medicine Building, Queen Elizabeth University Hospital, Glasgow G51 4TF

8 ¹³Northern Ireland Regional Genetics Centre, Belfast Health and Social Care Trust, Belfast City
9 Hospital, Lisburn Road, Belfast BT9 7AB, UK

10 ¹⁴North West Thames Regional Genetics Service, London North West University Healthcare NHS
11 Trust, Northwick Park and St Mark's Hospitals, Watford Road, Harrow HA1 3UJ, UK

12 ¹⁵MRC Human Genetics Unit, MRC IGMM, University of Edinburgh, Western General Hospital,
13 Edinburgh EH4 2XU, UK

14 ¹⁶Institute of Biomedical and Clinical Science, University of Exeter Medical School, RILD Level 4,
15 ED&E, Barrack Road, Exeter, EX2 5DW, UK

16 ¹⁷East Anglian Medical Genetics Service, Box 134, Cambridge University Hospitals NHS foundation
17 Trust, Cambridge Biomedical Campus, Cambridge, CB2 0QQ, UK

18 **CONTACT DETAILS OF CORRESPONDING AUTHOR**

19 Matthew E. Hurles, meh@sanger.ac.uk

20 Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, CB10 1SA, UK

21

22 **Running Title:** Selection and pathogenicity near splice sites

23 **Keywords:** Splicing, selection, constraint, de novo mutation, developmental disorders

1 **Abstract**

2 Mutations which perturb normal pre-mRNA splicing are significant contributors to human disease.
3 We used exome sequencing data from 7,833 probands with developmental disorders (DD) and their
4 unaffected parents, as well as >60,000 aggregated exomes from the Exome Aggregation Consortium,
5 to investigate selection around the splice site, and quantify the contribution of splicing mutations to
6 DDs. Patterns of purifying selection, a deficit of variants in highly constrained genes in healthy
7 subjects and excess *de novo* mutations in patients highlighted particular positions within and around
8 the consensus splice site of greater functional relevance. Using mutational burden analyses in this
9 large cohort of proband-parent trios, we could estimate in an unbiased manner the relative
10 contributions of mutations at canonical dinucleotides (73%) and flanking non-canonical positions
11 (27%), and calculated the positive predictive value of pathogenicity for different classes of
12 mutations. We identified 18 patients with likely diagnostic *de novo* mutations in dominant DD-
13 associated genes at non-canonical positions in splice sites. We estimate 35-40% of pathogenic
14 variants in non-canonical splice site positions are missing from public databases.

15 **Introduction**

16 Pre-mRNA splicing in humans is mediated by the major and minor spliceosomes, highly dynamic,
17 metalloenzyme complexes comprised of five key small nuclear RNAs (snRNA), along with over 100
18 protein components and accessory molecules (Brody and Abelson 1985; Hang et al. 2015; Scotti and
19 Swanson 2016). Accurate recruitment and function of the spliceosome is reliant on a plethora of cis-
20 acting regulatory elements encoded within the pre-mRNA itself. Whilst our understanding of the
21 underlying mechanistic processes regulating splicing has greatly increased in recent years, our ability
22 to predict whether or not a mutation will affect splicing remains limited. However, with estimates
23 that up to 50% of monogenic disease-causing variants may affect splicing (Teraoka et al. 1999; Ars et
24 al. 2000), a better understanding and more coherent approach to interpretation of variants affecting

1 splicing is badly needed(Cartegni et al. 2002; Baralle and Buratti 2017). With a plethora of *in silico*
2 splicing pathogenicity predictors available, there is little consensus on what a “gold standard” for
3 splicing pathogenicity prediction would be(Houdayer et al. 2012; Jian et al. 2014b; Tang et al. 2016).
4 Whilst many of these methods perform well within the canonical splice site dinucleotides (CSS, the
5 two highly-conserved bases flanking the acceptor and donor sites), their utility for other splice
6 relevant regions is less clear(Tang et al. 2016). In the clinical setting, often multiple algorithms and
7 expert judgment are used to predict pathogenicity, while for large scale research projects,
8 classification of variants is often binary, with CSS mutations typically classified as likely splice
9 affecting, whilst mutations in other splicing regulatory components are typically overlooked(Iossifov
10 et al. 2014; Wright et al. 2015; Deciphering Developmental Disorders Study 2017). Previous attempts
11 to estimate the relative contribution of pathogenic variants at non-canonical splice sites were based
12 on collating diverse published datasets of pathogenic variants (Lewandowska 2013) or data
13 submitted to databases of clinically interpreted variation (Krawczak et al. 2007) and are therefore
14 sensitive to the inherent heterogeneity and biases of such data, especially given the inevitable
15 subjectivity involved in clinical interpretation of this class of variation. Both clinical and research
16 interpretation of potential splice-disrupting variants lacks a robust quantitative foundation.

17 Utilising large scale exome sequencing data from 13,750 unaffected parents recruited as part of the
18 Deciphering Developmental Disorders (DDD) project(Wright et al. 2015) and >60,000 aggregated
19 exomes from the Exome Aggregation Consortium (ExAC) (Lek et al. 2016), we explore selective
20 constraint around splice regions across a set of 148,244 stringently defined exons well covered
21 (median coverage >15× at both CSS) across the DDD cohort (see Methods). Since selection is driven
22 by a number of factors, including monogenic developmental disorders (DD), as a complementary,
23 disease-based approach, we analyse enrichment of *de novo* mutations (DNMs) in DDD probands in
24 the same regions. We provide an unbiased, exome wide view of the signatures of selection and the
25 relative contribution of pathogenic splice altering mutations between the CSS and other, near splice
26 positions.

1 Results

2 Signatures of purifying selection around the splice site

3 Since purifying selection acts to keep deleterious alleles rare, population variation data can be used
4 to identify and assess the relative strengths of signals of purifying selection. To assess selective
5 constraint acting on positions around the canonical splice site we used the Mutability Adjusted
6 Proportion of Singletons (MAPS) metric(Lek et al. 2016) (a measure for inferring the degree of
7 selection robust to local variance in mutation rate) in 13,750 unaffected parents enrolled in the DDD
8 study as well as >60,000 aggregated exomes from ExAC (Figure 1a). The canonical splice acceptor
9 and donor dinucleotides show a clear signal of purifying selection in both datasets.

10 Outside of the CSS, other positions clearly show a signal of purifying selection beyond the
11 background level, including the donor site (last base of the exon, which is particularly pronounced
12 when the reference allele is G (Figure 1a)), and the intronic positions proximal to the canonical
13 donor site, peaking at the don+5 position, which exhibits a signal of purifying selection intermediate
14 between missense and stop-gained variants. Although no sites within the polypyrimidine tract
15 (PolyPy) show a signal of purifying selection individually, when these sites are grouped together
16 (Methods) and stratified by changes from a pyrimidine to a purine (PyPu), versus all other changes,
17 there is a clear difference between the two types of variants, with PyPu changes exhibiting an
18 increased signal of purifying selection when compared to non-PyPu changes (bootstrap $p < 0.001$,
19 Figure 1a, Supplemental Fig S1).

20 Deficit of splicing variants in highly constrained genes in healthy individuals

21 We also examined the distribution of variants of different classes among genes that are known to be
22 under different levels of selective constraint. Highly constrained genes should contain fewer
23 deleterious variants than less constrained genes. We investigated the proportion of variants
24 observed in the 13,750 unaffected parents which fell within highly constrained genes (probability of

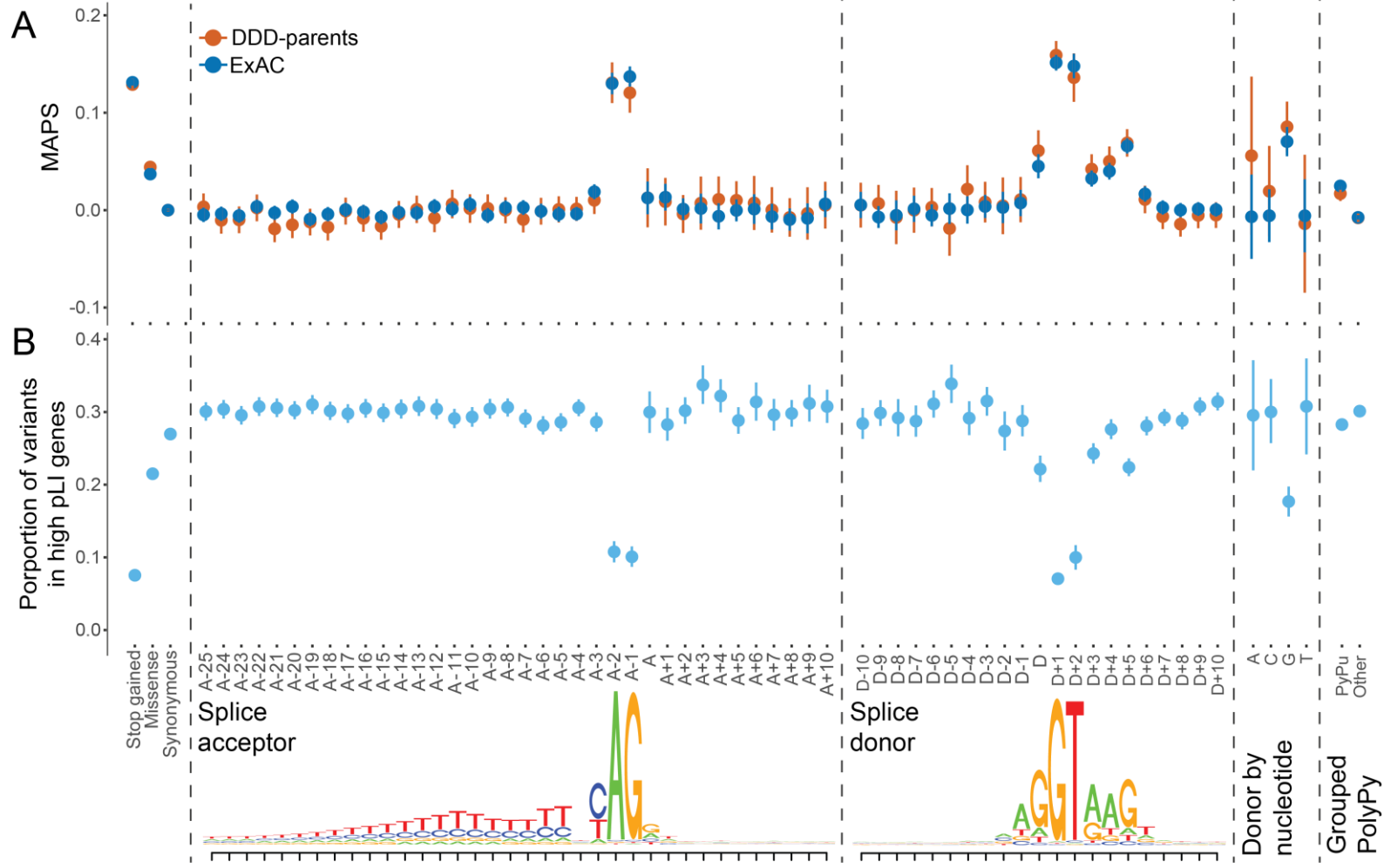
1 loss of function (LoF) intolerance (pLI, (Lek et al. 2016)) > 0.9) in our splicing regions of interest
2 (Figure 1b). In the near splice positions at which the highest MAPS values were seen (CSS, donor,
3 donor+5), we also observed a stronger depletion of variants in high pLI genes within the unaffected
4 parents, again supporting the potential pathogenicity of variants at these positions. The proportion
5 of parental variants in high pLI genes also recapitulates the signals of purifying selection seen in the
6 MAPS analyses with regard to the donor position split by reference allele (Figure 1b) and the PolyPy
7 region (Figure 1b, Supplemental Fig. S1), with the lowest proportions in high pLI genes observed for
8 sites with the highest MAPS values.

1 **Figure 1 – Signals of purifying selection around splice sites**

2 A. Selective constraint across splicing region in 13,750 unaffected parents of DDD probands and >60,000 aggregated exomes from ExAC. Mutability adjusted
3 proportion of singletons (MAPS) with 95% confidence intervals (CI) shown for Ensembl's Variant Effect Predictor (VEP) annotated exonic sites, extended
4 splice acceptor and splice donor regions, the last base of the exon, split by reference nucleotide, and grouped sites in the polypyrimidine tract region
5 (PolyPy), split by changes from a pyrimidine to a purine (PyPu) vs all other changes. B. Proportion of variants with 95% CI in 13,750 unaffected parents of
6 DDD probands which fall within genes with high probability of loss of function intolerance ($pLI > 0.9$) across VEP annotated exonic sites, extended splice
7 acceptor and splice donor regions, the last base of the exon, split by reference nucleotide, and grouped sites in the polypyrimidine tract region, split by
8 changes from a pyrimidine to a purine (PyPu) vs all other changes. Lower panel shows splice acceptor and splice donor consensus sequences, based on our
9 exons of interest.

10

11



1 **Assessing the significance of mutational burden for different classes of splicing mutations**

2 We identified 871 autosomal high confidence DNMs (non-synonymous consequences excluded)
3 within canonical and near-splice regions of interest well covered by exome data in the 7,833
4 probands, allowing us to test for enrichment of DNMs relative to expectations based on a
5 trinucleotide null model of mutation rate (Samocha et al. 2014) across different sets of genes (DD-
6 associated with dominant or recessive mechanisms, and non-DD associated, see Methods). Across
7 recessive DD and non-DD associated genes, no enrichment of DNMs beyond the null expectation
8 was observed (Figure 2a). In dominant DD genes, a significant cumulative excess of DNMs was noted
9 across the full splicing region (Poisson p (false discovery rate (FDR) adjusted) = 1.33×10^{-14} , fold
10 enrichment = 3.47), which remained significant upon exclusion of the canonical dinucleotide
11 positions (Poisson p (FDR adjusted) = 0.0035, fold enrichment = 1.86). Individually, the four canonical
12 splice site positions each showed a significant (10-27 fold) excess of DNMs (Poisson p (FDR adjusted),
13 fold enrichment: acc-2 = 4.22×10^{-12} , 26.6; acc-1 = 3.43×10^{-8} , 16.6; don+1 = 1.33×10^{-14} , 20.1; don+2
14 = 0.004, 10.0), as did the don+5 site (9.7×10^{-5} , 9.29). The similar level of enrichment between don+5
15 and don+2 implies these positions harbour comparable proportions of splice disrupting mutations.
16 No individual positions within the PolyPy region showed an individual excess of DNMs, however,
17 when the positions were considered cumulatively and split between PyPu and non-PyPu changes
18 (Figure 2b), an excess of DNMs was observed in the PyPu group for dominant DD genes (fold
19 enrichment = 3.46), although this was not significant at an FDR of 5% (Poisson p (FDR corrected) =
20 0.086).

21 When the same analysis was performed for dominant genes in subsets of the DDD cohort with
22 (n=1417) and without (n=3364) robust diagnoses from the standard diagnostic protocol, which only
23 assesses splicing mutations at the CSS (Supplemental Fig. S2), the enrichment within the diagnosed
24 subset was confined to the CSS (Poisson p (FDR adjusted), fold enrichment: CSS = 1.33×10^{-14} , 69.74;
25 other positions = 0.658, 1.82), whilst in the undiagnosed subset, the opposite pattern was observed

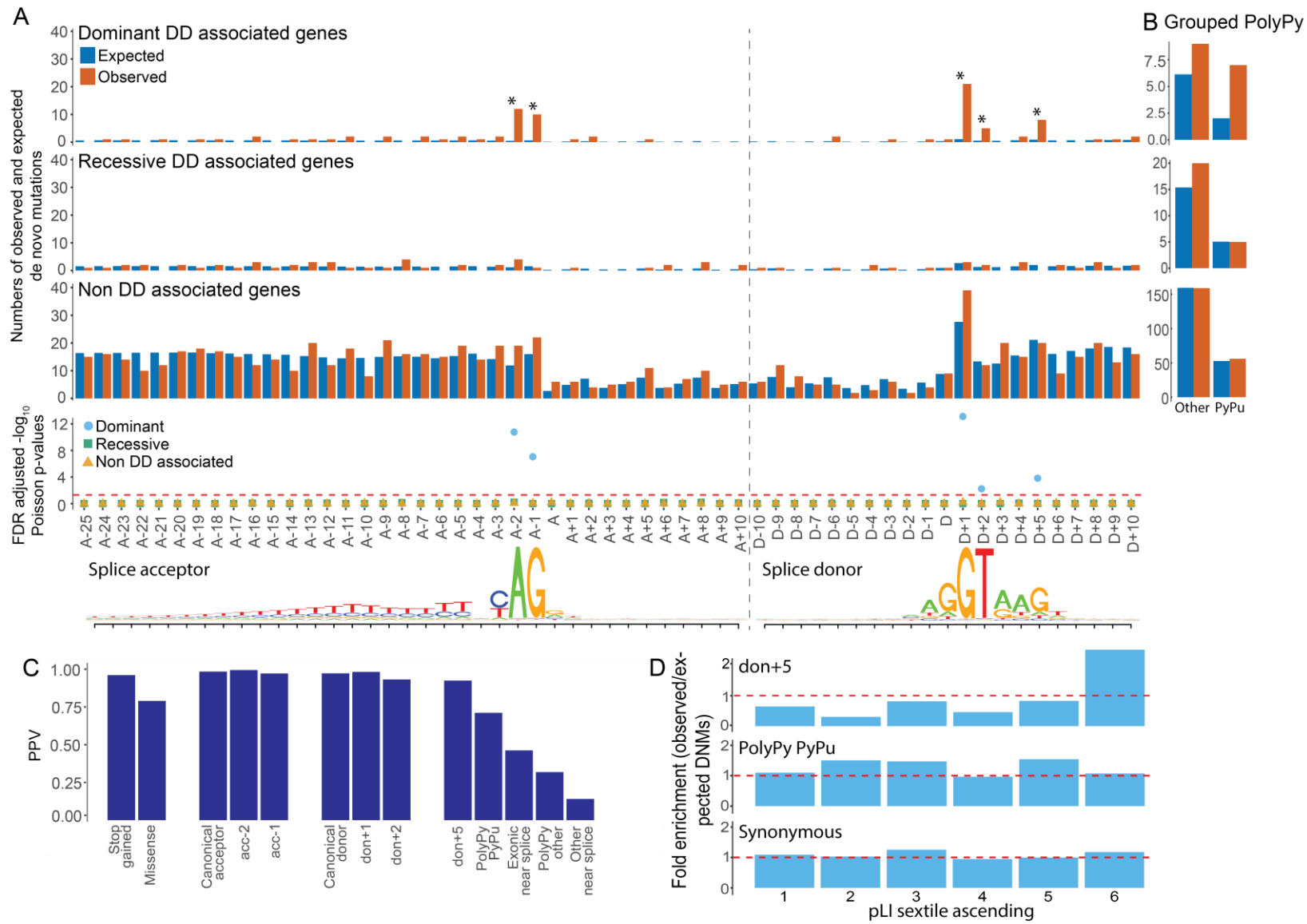
- 1 (Poisson p (FDR adjusted), fold enrichment: CSS = 1, 0; other positions = 0.012, 2.21), with the don+5
- 2 site showing the greatest enrichment (16.18, Poisson p (FDR adjusted) = 5.35×10^{-5}).
- 3 These results are highly concordant with the signatures of purifying selection identified using the
- 4 MAPS metric, and the deficit of parental variants in high pLI genes, providing multiple independent
- 5 lines of evidence that mutations in positions outside of the CSS can disrupt normal splicing.

1 **Figure 2 – *De novo* mutations around splice sites**

2 Enrichment of *de novo* mutations (DNMs) across the splicing region in 7,833 DDD probands A. Numbers of observed and expected DNMs across splicing
3 region, in known dominant and recessive DD genes, as well as non-DD associated genes, with FDR corrected Poisson p-values. Splice acceptor and splice
4 donor consensus sequences shown below, as in Figure 1. B. Aggregation of observed and expected numbers of DNMs in the polypyrimidine tract (PolyPy)
5 region, with changes from a pyrimidine to a purine (“PyPu”) and all other changes shown separately for known dominant and recessive DD genes, as well as
6 non-DD associated genes. C. Positive predictive values (PPVs) for *de novo* mutations in dominant DD-associated genes in positions across the splicing region,
7 as well as VEP annotated stop gained and missense changes, calculated from observed and expected numbers of DNMs. D. Enrichment
8 (observed/expected) of *de novo* mutations (DNMs) by gene probability of loss of function intolerance (pLI), split in to sextiles for donor+5, pyrimidine to
9 purine PolyPy, and synonymous sites. pLI scores encompassed by each sextile: 1 = 5.36E-91 - 0.000000605, 2 = 0.000000609 - 0.000558185, 3 =
10 0.000559475 - 0.027905143, 4 = 0.027908298 - 0.377456159, 5 = 0.377491926 - 0.919495985, 6 = 0.91955878 - 1.

11

12



1 **Estimating positive predictive values for different classes of splice mutation**

2 We used the fold-enrichment of the numbers of observed DNMs in dominant DD genes in the DDD
3 cohort over the number expected under the null mutation model to calculate positive predictive
4 values (PPVs) for groupings of near splice site positions. We compared these with PPVs for other,
5 more commonly disease-associated variant classes within the same exons of the same genes (Figure
6 2c). We observe minor differences in PPV for the individual positions of the canonical acceptor and
7 donor sites, with don+2 showing the lowest PPV at 0.90, which is approximately the same as for the
8 don+5 position (PPV 0.89). Variants within the PolyPy region which change a pyrimidine for a purine
9 have a PPV of 0.71, which is below the PPV for missense mutations (0.79), but still predicts a
10 substantive number of pathogenic mutations arising from disruption of the PolyPy.

11 Despite the modest number of observed DNMs used to make these PPV estimates, we see
12 concordance with the population based metrics described above (MAPS and deficit of splicing
13 variants in high pLI genes in unaffected parents of DDD patients – Supplemental Fig. S3), suggesting
14 these estimates are robust.

15 We looked at the distribution of observed DNMs in genes with respect to their probability of being
16 LoF intolerant (using the pLI metric(Lek et al. 2016), Figure 2d). For synonymous variants, we
17 observed no significant enrichment of DNMs in high pLI genes. For don+5 mutations, there is a clear
18 excess of DNMs in genes most likely to be intolerant to LoF mutations in the DDD cohort, further
19 supporting the likely pathogenicity of mutations in these positions. For the PolyPy PyPu mutations,
20 although there is a nominally significant enrichment of DNMs in general, this does not show a
21 significant skew towards high pLI genes in our cohort.

22

23

24

1 **Identifying diagnostic non-canonical splice mutations**

2 After exclusion of probands with likely diagnostic protein-coding or canonical splice site variants, 38
3 DNMs in our near splice site positions of interest in dominant DD genes were identified. The clinical
4 phenotypes of patients carrying these mutations were reviewed by a consultant clinical geneticist,
5 blinded to the precise mutation and PPVs estimated above, and the patient's recruiting clinician, to
6 assess the phenotypic similarity between the proband and the disorder expected from a LoF
7 mutation in that gene. The 38 variants were classified as likely diagnostic (Table 1), or unlikely
8 diagnostic/unknown (Supplemental Table S1), depending on the strength of phenotypic similarity.
9 Phenotypic information for probands with likely diagnostic variants is given in Supplemental Table
10 S2, and pathogenicity prediction scores for the SNVs in Supplemental Table S3. The clinical review
11 resulted in 18 variants (47%) being classified as likely diagnostic, highly concordant with the number
12 predicted from the overall PPV of non-canonical sites of 46%, moreover, a higher proportion of likely
13 diagnostic variants were classified at sites with higher PPVs (Pearson's correlation coefficient = 0.91,
14 $p = 0.033$) (Figure 3). With 48 CSS DNMs observed within the same exons in our probands, we
15 estimate that 73% (95% CI: 60-82%) of disease causing splice disrupting DNMs occur within the CSS,
16 while 27% (95% CI: 18-39%) are in non-canonical, near-splice positions.

17 Eight DNMs were selected for functional validation via a minigene vector system, including six likely
18 diagnostic PolyPy variants, a PolyPy variant of uncertain clinical significance, and a likely diagnostic
19 don+5 variant, where both the phenotype of the patient and that associated with the gene (*MBD5*)
20 are nonspecific, along with two negative controls (untransmitted variants identified in unaffected
21 parents within the same PolyPys as test variants). For six of the variants selected for validation,
22 differences in splicing between the reference and mutant constructs were observed (Supplemental
23 Fig. S4a and S4b). One of the likely diagnostic PolyPy mutations, the PolyPy mutation of uncertain
24 significance, and both negative controls showed no difference in splicing between the reference and
25 mutant constructs (Supplemental Fig. S4c and S4d).

1 **Table 1 – Diagnostic *de novo* mutations in non-canonical dinucleotide near splice positions**

2 Variant and proband information for 18 *de novo* likely diagnostic splice region variants identified in previously undiagnosed DDD probands in known
 3 dominant DD-associated genes (hg19 coordinates).

chrom:pos_ref/alt	symbol	VEP annotation	Splice annotation	Associated disorder	Clinician classification
7:42063221_G/C	GLI3	intron_variant	acc-14	Greig Cephalopolysyndactyly Syndrome	Likely pathogenic, full contribution
16:3819367_C/T	CREBBP	intron_variant	acc-13	Rubinstein-Taybi Syndrome Type 1	Likely pathogenic, full contribution
22:24143120_T/G	SMARCB1	intron_variant	acc-11	Rhabdoid Predisposition Syndrome 1 / Coffin-Siris Syndrome 3	Likely pathogenic, full contribution
18:52895603_T/C	TCF4	intron_variant	acc-11	Pitt-Hopkins Syndrome	Likely pathogenic, full contribution
5:88025173_A/C	MEF2C	splice_region_variant	acc-9	Mental Retardation-Stereotypic Movements-Epilepsy And/Or Cerebral Malformations	Likely pathogenic, full contribution
9:130988306_G/A	DNM1	splice_region_variant	acc-8	Epileptic Encephalopathy	Likely pathogenic, full contribution
8:61763045_G/A	CHD7	splice_region_variant	acc-7	CHARGE / Kallmann Syndrome Type 5 / Idiopathic Hypogonadotropic Hypogonadism	Definitely pathogenic, full contribution
17:38801875_T/C	SMARCE1	splice_region_variant	acc-4	Coffin-Siris Syndrome 5	Likely pathogenic, full contribution
1:27097607_C/A	ARID1A	splice_region_variant	acc-3	Coffin-Siris Syndrome 2	Likely pathogenic, full contribution
9:140728798_C/G	EHMT1	splice_region_variant	acc-3	9q Subtelomeric Deletion Syndrome / Kleefstra Syndrome 1	Definitely pathogenic, full contribution
2:223160248_T/C	PAX3	splice_region_variant	don-1	Waardenburg Syndrome, Type 1 / Craniofacial-Deafness-Hand Syndrome	Likely pathogenic, partial contribution
2:166229861_A/G	SCN2A	splice_region_variant	don+4	Nonspecific Severe Id / Benign Familial Neonatal Infantile Seizures / Infantile Epileptic Encephalopathy	Likely pathogenic, full contribution
9:130422391_A/G	STXBP1	splice_region_variant	don+4	Angelman/Pitt Hopkins Syndrome-Like Disorder / Epileptic Encephalopathy Early Infantile Type 4	Likely pathogenic, full contribution
22:41556731_G/A	EP300	splice_region_variant	don+5	Rubinstein-Taybi Syndrome Type 2	Likely pathogenic, full contribution
2:149221493_G/C	MBD5	splice_region_variant	don+5	Ehmt1-Like Intellectual Disability	Likely pathogenic, full contribution
9:130427615_G/C	STXBP1	splice_region_variant	don+5	Angelman/Pitt Hopkins Syndrome-Like Disorder / Epileptic Encephalopathy Early Infantile Type 4	Likely pathogenic, full contribution
17:42956919_C/T	EFTUD2	splice_region_variant	don+5	Mandibulofacial Dysostosis With Microcephaly	Definitely pathogenic, full contribution
20:61452890_C/G	COL9A3	splice_region_variant	don+8	Multiple Epiphyseal Dysplasia Type 3	Likely pathogenic, partial contribution

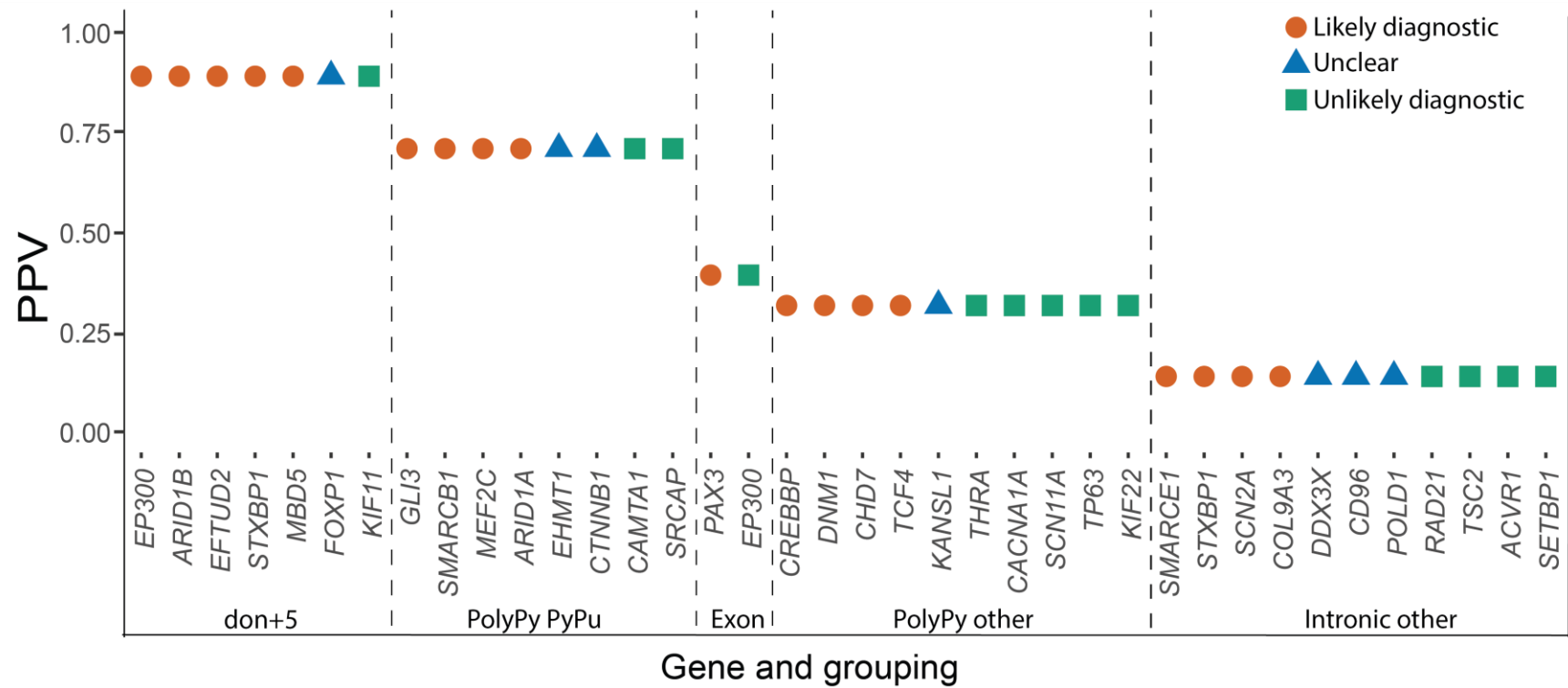
4

5

1 **Figure 3 – Clinical classifications of non-canonical near splice *de novo* mutations**

2 Relationship between clinical classifications of 38 splice region *de novo* mutations (DNMs) in undiagnosed DDD probands and positive predictive values
 3 (PPVs) calculated using observed and expected numbers of DNMs in 7,833 probands.

4



5

1 **Assessing splicing pathogenicity prediction tools**

2 The population genetic metrics of purifying selection and mutation enrichment metric for
3 pathogenicity that we have derived provide an orthogonal approach to assessing the accuracy of
4 splicing pathogenicity prediction tools, compared to the standard approach of assessing
5 classification accuracy for clinically interpreted variants. We assessed four splicing pathogenicity
6 prediction tools: two recently published genome-wide ensemble learning methods: AdaBoost and
7 RandomForest, Spidex (based on deep learning trained on RNA sequencing data), and the longer
8 standing, widely used MaxEntScan (MES) (Yeo and Burge 2004; Jian et al. 2014a; Xiong et al. 2015).

9 We divided the scores from each prediction tool, plus CADD(Kircher et al. 2014), into 20 equal-sized
10 bins to facilitate cross-method comparability. We calculated the MAPS for each bin of each of the
11 scoring metrics for the splicing variants observed in the 13,750 DDD unaffected parents, and saw a
12 strong positive correlation between pathogenicity metric and MAPS for all tools (Figure 4). AdaBoost
13 had the highest absolute MAPS value for the top scoring bin, suggesting that it is best able to identify
14 variants under the strongest purifying selection. The proportion of variants in the unaffected parents
15 falling in genes with pLI > 0.9 broadly recapitulates this pattern, with fewer variants in high pLI genes
16 in the highest scoring brackets for all metrics (Supplemental Fig. S5). We then looked at the
17 distribution of scores for each tool for the 83 splicing DNMs observed in DDD probands in autosomal
18 dominant DD-associated genes which were covered by all five scoring systems to compare
19 performance of the metrics on mutations more likely to have a deleterious impact on splicing, with
20 the expectation that these potentially damaging variants would be scored highly by the metrics,
21 giving high values of area under the curve (AUC, Figure 5). Again, all metrics performed well, with
22 the majority of DNMs being classified in the most deleterious score brackets. Here AdaBoost gave
23 the highest AUC value, suggesting it weighted these likely damaging variants as more deleterious
24 than the other metrics comparatively. When CSS positions were removed from the analysis,
25 AdaBoost remained the tool with the highest AUC. The largest reduction in the AUC metric was seen

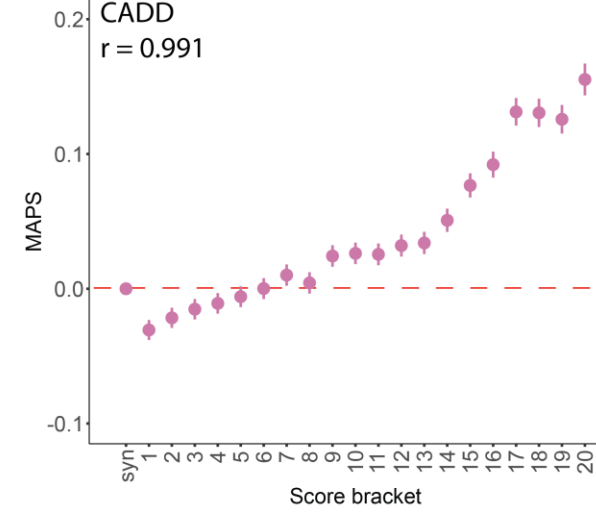
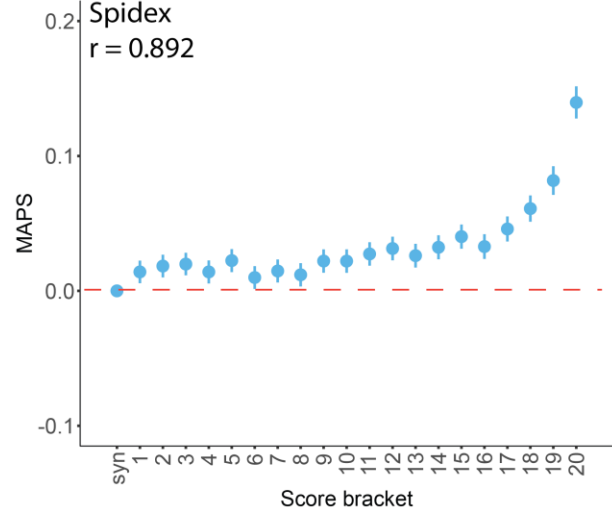
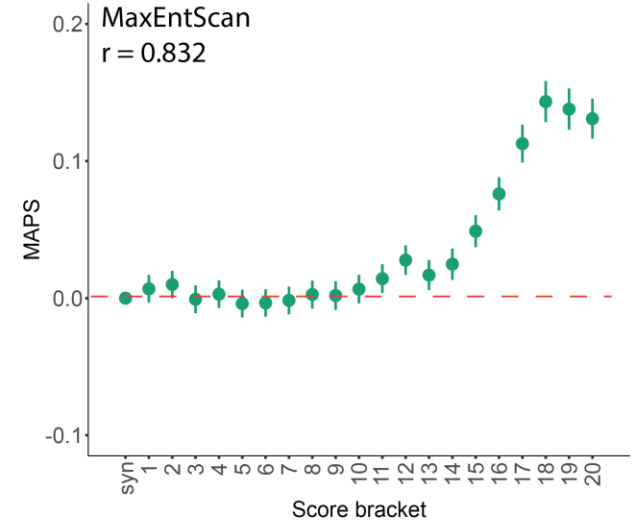
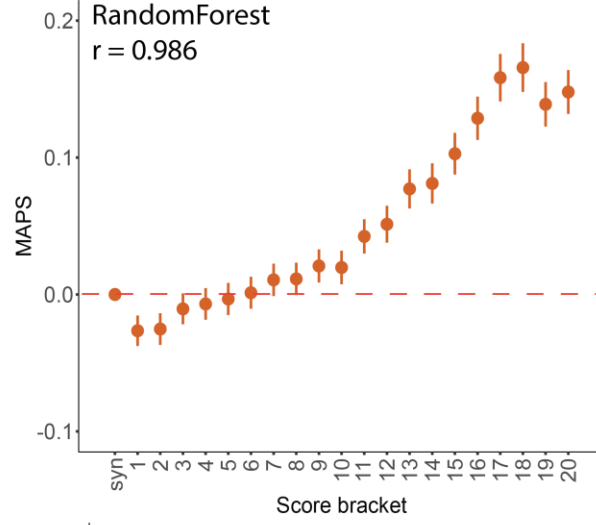
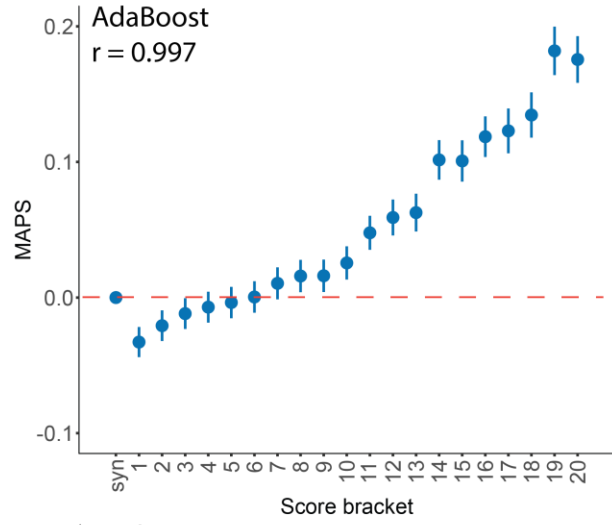
1 for Spidex and CADD, indicating these tools may be least informative for positions outside of the CSS.
2 Upon removal of the CSS positions from the analyses of MAPS and deficit of parental variants in high
3 pLI genes, similar results were revealed, with the highest AdaBoost scores retaining strong signals of
4 purifying selection but a marked reduction in signal from the highest Spidex scores (Supplemental
5 Fig. S6 and Supplemental Fig. S7).

6 Taken together, these data show a strong relationship between the considered splicing
7 pathogenicity scoring systems and the general landscape of purifying selection on splicing control,
8 but demonstrate that the utility of these systems in identifying likely diagnostic variants is limited
9 outside of the CSS.

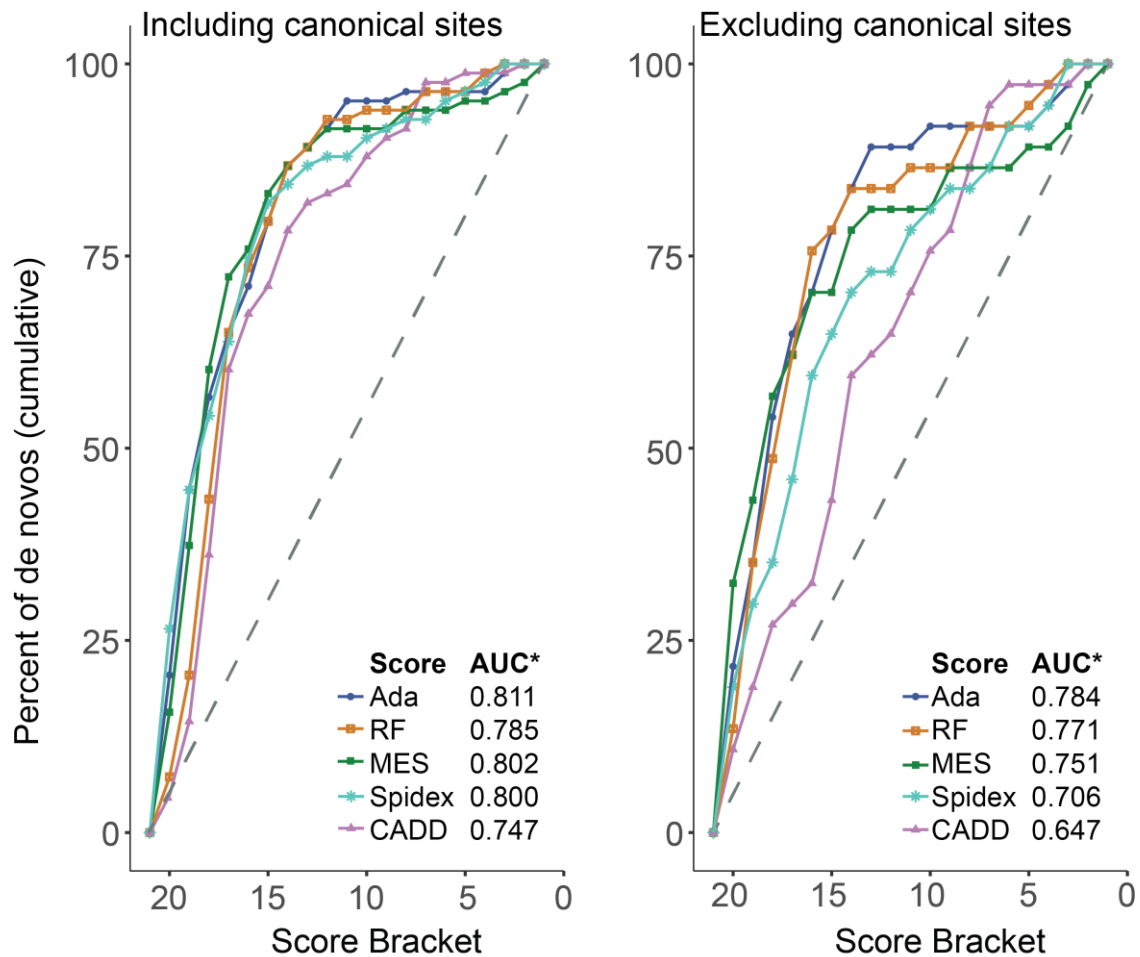
1 **Figure 4 – Selective constraint and pathogenicity scores**

2 Mutability adjusted proportion of singletons (MAPS), with 95% CI, calculated for pathogenicity score brackets (least to most severe) in 13,750 unaffected
3 parents from the DDD project, with Spearman's rank correlation coefficient.

4



1 **Figure 5 – Pathogenicity scores for observed near splice site *de novo* mutations**
 2 Cumulative percentage of *de novo* mutations (DNMs) in known dominant DD genes with decreasing
 3 pathogenicity score bracket, shown with canonical splice site positions included (left) and excluded
 4 (right). * AUC = area under curve



5

6 Discussion

7 Our study represents a large, unbiased exploration of the perturbation of splicing by genetic
 8 variation in near-splice regions, with complementary signals of selection observed through two
 9 different population-based analyses. Selection can be driven by many factors, including monogenic
 10 disease resulting in reduced reproductive fitness. DDs represent the largest single class of
 11 monogenic diseases. With a clear enrichment of DNMs demonstrated by previous studies of DDs
 12 (Deciphering Developmental Disorders Study 2017), analysis of such DNMs provided an

1 independent, disease-based approach, complementary to our population-based analyses.
2 Concordance in results from the different approaches indicates the robustness of our conclusions.
3 Our analyses, taken together, suggest the pathogenic contribution of non-canonical splice positions
4 has been under-appreciated. We estimate that around 27% (95% CI: 18-39%) of splice disrupting
5 pathogenic mutations within the DDD cohort are in non-canonical positions. In sites with pathogenic
6 or likely pathogenic clinical significance in ClinVar(Landrum et al. 2016) overlapping with our splicing
7 positions of interest (and with non-synonymous consequences removed), we found 83.5% of
8 variants fell within canonical positions, with just 16.5% in non-canonical positions. When adjusted
9 for number of submissions as a proxy for allele count, this figure was 17.5%, perhaps indicative that
10 recurrence strengthens evidence of pathogenicity. Both of these values are significantly below our
11 estimate of 27% (Fisher's Exact $p = 1.22 \times 10^{-15}$ and $p < 2.2 \times 10^{-16}$ respectively), suggesting under-
12 ascertainment of non-canonical splicing variants by around 35-40% in clinical databases, despite a
13 growing understanding of the importance of such sites in splicing regulation(Kircher et al. 2014;
14 Ferreira et al. 2016; Soukarieh et al. 2016; Cummings et al. 2017; Ito et al. 2017; Soemedi et al. 2017;
15 Ke et al. 2018).

16 Estimates of the relative contribution of canonical and non-canonical splice site mutations are sparse
17 in the literature. These estimates are also typically based on clinically interpreted variants and so are
18 likely to be biased by the accuracy of current clinical practices. When comparing canonical and non-
19 canonical mutations within the Human Gene Mutation Database (HGMD, based on variation
20 described in publications, Krawczak *et al.*(Krawczak et al. 2007) stated canonical mutations
21 accounted for 64% of mutations at donor sites and 77.4% of mutations at acceptor sites, giving an
22 estimated non-canonical contribution of ~30% overall (consistent with our data), while data taken
23 from Caminsky *et al.*(Caminsky et al. 2014) put this estimate at around 43% (above our upper
24 bound). These values are closer to our 27% estimate than to the ClinVar proportion of ~17%, despite
25 our approach focussing on DNMs and dominant disorders, whereas the other two studies did not

1 discriminate on mode of inheritance and included recessive disorders, which can also be caused by
2 non-canonical splicing mutations(Basel-Vanagaite et al. 2013; Brunham et al. 2015) and exonic
3 variants. Our findings highlight the complementarity of assessing the clinical importance of non-
4 canonical splice variants both through the traditional approach based on clinically-interpreted
5 variation accrued through diagnostic practice and through unbiased approaches that leverage
6 population variation and unbiased models of germline mutation.

7 Our analysis of non-canonical splice position mutations did not include exonic missense
8 variants(Teraoka et al. 1999; Ars et al. 2000; Krawczak et al. 2007), nor did it explicitly include
9 branchpoints(Maslen et al. 1997; Di Leo et al. 2004; Crotti et al. 2009; Aten et al. 2013), splicing
10 enhancers and suppressors(Lorson et al. 1999; Liu et al. 2001) or deep intronic mutations(Cummings
11 et al. 2017; Vaz-Drago et al. 2017). Detecting splice disrupting variants at these sites is even more
12 challenging, as despite recent efforts(Corvelo et al. 2010; Wang and Wang 2014; Mercer et al. 2015;
13 Badr et al. 2016; Taggart et al. 2017), comprehensive catalogues of all branchpoints and exonic and
14 intronic splicing enhancers and silencers are currently unavailable, algorithms that predict the
15 impact of mutations at such sites are not highly accurate, and some of these sites are not covered by
16 exome sequencing (the greater utilisation of whole genome sequencing will allow greater
17 opportunity to find and assess the contributions of more distal splice disrupting variants). As such,
18 our estimate of the contribution of non-canonical splicing position mutations is likely to be a lower
19 bound. Thus our estimate of 35-40% under-ascertainment in clinical databases may be conservative,
20 and the true extent of missed diagnoses may be even higher.

21 The size of the available datasets determines our power to detect signals of selection and
22 enrichment of DNMs, so although we could demonstrate signal at splice-important non-canonical
23 positions, there may be other positions with more subtle signals of selection which we lacked the
24 power to detect. Another limiting factor for such analyses is the specificity with which we can
25 identify splicing related sequences. For splice sites themselves, well curated resources of intron-exon

1 junctions exist (e.g. from GENCODE), giving a high degree of confidence that what we are assessing is
2 indeed near-splice sequence. For exonic and intronic splicing enhancers and suppressors, although
3 attempts at comprehensive identification have been made (Fairbrother et al. 2002; Zhang and
4 Chasin 2004; Goren et al. 2006; Ke et al. 2011), there is little concordance between available
5 resources (Caceres and Hurst 2013), meaning non-enhancer/suppressor sequence would almost
6 certainly be included in analyses, limiting power to detect any signal. Collation of yet larger datasets
7 and a greater understanding of other splicing elements will help to identify these sites, using the
8 same methodology applied here.

9 The nature of many developmental disorders makes obtaining RNA samples from relevant tissues of
10 patients (i.e. neural tissue) acutely problematic, so we investigated the effects on splicing of several
11 of the potentially diagnostic DNMs using a minigene vector system. We were able to demonstrate
12 changes to splicing for five out of six likely diagnostic PolyPy variants as well as the likely diagnostic
13 don+5 variant, supporting the clinical interpretation based on clinical phenotype. We did not
14 observe an effect on splicing for one likely diagnostic PolyPy variant, and one PolyPy variant of
15 uncertain clinical significance. Although the accuracy of minigene assays when compared with
16 patient RNA is generally high (Bonnet et al. 2008; They et al. 2011; van der Klift et al. 2015), known
17 limitations of the system (e.g. lack of full endogenous genetic context (Baralle et al. 2006;
18 Sangermano et al. 2017), and sensitivity to cell type utilised (Lastella et al. 2006)) mean we cannot
19 definitively state that the effects seen in the minigene assay would be the same in the full genetic,
20 developmental and cellular context within the patient.

21 We envisage that greater appreciation of the importance of near splice site mutations will increase
22 diagnostic yields, as well as providing increased power for the detection of new genetic associations,
23 both within the field of rare disease and beyond. We highlight two challenges to improving detection
24 of pathogenic non-canonical splice site mutations.

1 First, many commonly used *in silico* tools for annotating the likely functional impact of variants do
2 not discriminate between different non-CSS positions with very different likelihoods of being
3 pathogenic. Moreover, commonly used annotation tools differ in the ways in which variants are
4 annotated, with splicing variants displaying the highest level of disagreement between
5 tools(McCarthy et al. 2014). This highlights the need for a more consistent and evidence based
6 annotation of splicing variants. Of the positions shown in our analyses to be most damaging, don+5
7 sites are annotated by VEP(McLaren et al. 2016) and SnpEff(Cingolani et al. 2012) as
8 “splice_region_variant”, while most positions of the PolyPy are annotated as intronic, so are
9 potentially easily overlooked. With Annovar’s(Wang et al. 2010) default settings, only the CSSs are
10 flagged as splicing variants, although with both Annovar and SnpEff, the user can optionally extend
11 the region to be annotated as splice variants. We note that Ensembl have recently implemented a
12 VEP plugin which allows greater granularity in splice region annotation
13 (https://github.com/Ensembl/VEP_plugins/blob/release/88/SpliceRegion.pm), including annotating
14 the don+5 and other near-donor positions, as well as the PolyPy region. This type of increased
15 granularity of splicing annotation should facilitate consideration of these variants in future studies.

16 Second, current tools that predict the pathogenicity of non-CSS mutations have limited accuracy,
17 and it is not clear how to translate the scores that they output into a likelihood of pathogenicity. The
18 quantitative framework that we introduced here of estimating PPVs for different classes of
19 mutations by comparing the number of observed mutations to the number expected under a well-
20 calibrated null model of germline mutation has much more direct relevance to clinical
21 interpretation, although the interpretation of specific DNMs still proves problematic, particularly for
22 DNMs in sites of intermediate PPV. We propose that the scores generated by such splicing
23 prediction tools could be calibrated by performing analogous analyses of mutation enrichment to
24 estimate PPVs for different bins of scores. As the size of trio-based cohorts increases, the accuracy of
25 calibration will improve.

1 In summary, our results demonstrate a significant contribution of non-canonical splicing mutations
2 to the genetic landscape of DDs, a finding which is highly likely to be recapitulated across other
3 monogenic disorders and contexts. We demonstrate the importance of non-canonical positions
4 (particularly the don+5 site and pyrimidine-removing mutations in the PolyPy region). These
5 inferences are supported by both population genetic investigations of purifying selection, as well as
6 a disease based approach, considering the burden of DNMs in ~8,000 children with severe DDs.
7 Mutations at some non-canonical splicing positions convey a risk of disease similar to that of protein
8 truncating and missense mutations, but are under-represented in existing databases of disease-
9 causing variants.

10 **Materials and methods**

11 **Cohort and sequencing**

12 For full description of cohort and analytical methodology, see previous DDD publications
13 (Deciphering Developmental Disorders Study 2015; Deciphering Developmental Disorders Study
14 2017). Briefly, 7,833 patients with severe, undiagnosed developmental disorders were recruited to
15 the DDD study from 24 clinical genetics centres from across the UK and Ireland. Whole exome
16 sequencing was conducted on the proband and both parents, with exome capture using SureSelect
17 RNA baits (Agilent Human All-Exon V3 Plus with custom ELID C0338371 and Agilent Human All-Exon
18 V5 Plus with custom ELID C0338371) and sequencing using 75 base paired-end reads using Illumina's
19 HiSeq. Mapping was conducted to GRCh37 using the Burrows-Wheeler Aligner (BWA, v0.59(Li and
20 Durbin 2009)) and variant identification was conducted using the Genome Analysis Toolkit (GATK,
21 v3.5.0(McKenna et al. 2010)). Realignment to GRCh38 should not affect the conclusions of this work, as
22 only high confidence intron-exon boundaries were used in the analyses. These were taken from
23 GENCODE v19 (GRCh37) but filtered to exclude a small subset of exons which no longer met our
24 stringent criteria in GENCODE v22 (GRCh38), as described below. Variant annotation was conducted
25 with Ensembl's VEP (<https://www.ensembl.org/info/docs/tools/vep/index.html>), using Ensembl

1 gene build 76(McLaren et al. 2016). DNMs were identified using DeNovoGear (v0.54)(Ramu et al.
2 2013), and filtered using an in house pipeline - denovoFilter - developed by Jeremy F.
3 McRae(Deciphering Developmental Disorders Study 2017)
4 (<https://github.com/jeremymcrae/denovoFilter>). Exome sequencing and phenotype data are
5 accessible via the European Genome-phenome Archive (EGA) under accession number
6 EGAS00001000775 (<https://www.ebi.ac.uk/ega/studies/EGAS00001000775>).

7 **Defining exons of interest**

8 We took exons from GENCODE v19 (<https://www.encodegenes.org/releases/19.html>) which met
9 the following criteria: annotation_type = "exon", gene_type = "protein_coding", gene_status =
10 "KNOWN", transcript_type = "protein_coding", transcript_status = "KNOWN", annotation != "level
11 3" (automated annotation), and tag = "CCDS", "appris_principal", "appris_candidate_longest",
12 "appris_candidate", or "exp_conf" (n = 255,812 exons)(Harrow et al. 2012). We removed a small
13 subset of exons which no longer met these criteria in the more recent, GRCh38 based GENCODE v22
14 release (leaving 253,275 exons). We removed any exons where the median coverage at the
15 canonical acceptor or donor positions was $<15\times$ in two sets of DDD data which used different exon
16 capture methods (Agilent Human All-Exon V3 Plus with custom ELID C0338371 and Agilent Human
17 All-Exon V5 Plus with custom ELID C0338371). 148,244 exons passed these criteria.

18 We annotated individual genomic positions relative to the acceptor and donor sites, removing any
19 exons $<14\text{bp}$, and any positions which had multiple potential annotations. At the acceptor end, we
20 considered 25bp of intronic sequence (acc-25 to acc-1) and 11bp exonic sequence (acc to acc+10). At
21 the donor end, we considered 10bp of intronic sequence (don+1 to don+10) and 11bp exonic
22 sequence (don to don-10). This yielded ~ 6.9 million near-splice positions of interest.

23 The reference nucleotide composition at each position of the splicing region of interest was
24 calculated using all sites and a weighted position weight matrix graph was generated using the

1 seqLogo package via Bioconductor(Huber et al. 2015)
2 (<https://bioconductor.org/packages/release/bioc/html/seqLogo.html>) in R (R Core Team
3 2018)(version 3.1.3).

4 We define the PolyPy region as acc-3, and acc-5 to acc-17, based on pyrimidine content > 70% in our
5 exons of interest. We assess changes from a pyrimidine to a purine (PyPu) adjusting for the strand
6 the exon is on.

7 **Mutability adjusted proportion of singletons (MAPS)**

8 In 13,750 unaffected parents enrolled as part of the DDD study, as well as >60,000 aggregated
9 exomes from ExAC v0.3.1 (<http://exac.broadinstitute.org/>), we calculated the MAPS metric(Lek et al.
10 2016) using code developed in house by Patrick J. Short (Short et al. 2018)
11 (<https://github.com/pjshort/dddMAPS>). The MAPS metric is based on the principle that negative
12 selection acts to keep deleterious variation rare at a population level, but more mutable sequence
13 contexts can contain variants that appear more common because of recent recurrent mutational
14 events, so the metric corrects frequencies based on local sequence context using synonymous
15 mutations. Only relevant ExAC sites with "PASS" in the VCF "FILTER" column were counted, and ExAC
16 and DDD variants were filtered for FisherStrand (FS) <10. MAPS was calculated for all SNVs
17 overlapping out splice positions of interest (201,587 near splice variants for DDD, and 678,241 for
18 ExAC), the last base of the exon split by reference nucleotide (2109 variants for DDD, 6325 for ExAC),
19 and the PolyPy, split by PyPu (15,847 variants for DDD, 58,762 for ExAC) vs all other PolyPy changes
20 (52,300 variants for DDD, 175,287 for ExAC), as well as VEP(McLaren et al. 2016) ascertained
21 synonymous (580,066 variants for DDD, 1,513,758 for ExAC), missense (1,125,167 variants for DDD,
22 2,786,533 for ExAC) and stop-gained (25,863 variants for DDD, 78,496 for ExAC) sites across
23 autosomal regions. To establish whether the MAPS metric was significantly different between PolyPy
24 PyPu vs all other PolyPy changes, a bootstrap resampling method was run with 1000 iterations.

1 **Parental variants in high pLI genes**

2 We annotated all variant sites used in the MAPS calculations above in the 13,750 DDD parents with
3 the gene in which the variant fell, and the pLI score of that gene, and calculated the proportion of
4 these sites which fell within genes with high pLI scores(Lek et al. 2016) (> 0.9).

5 ***De novo mutations***

6 DNMs were identified using DeNovoGear(Ramu et al. 2013) as described in McRae *et al*,
7 2017(Deciphering Developmental Disorders Study 2017), and a stringent confidence threshold
8 (posterior probability > 0.8) was applied. We used triplet-based mutation rates(Samocho et al. 2014)
9 for each potential single nucleotide change across our splicing regions of interest to calculate the
10 expected number of DNMs across autosomes in the 7,833 probands. Expected values were adjusted
11 for depth of sequencing coverage < 50 to account for poorer ascertainment of variants in low
12 coverage regions (exon depth <1, $\text{exp} \times 0.119$; exon depth >1 and <50, $\text{exp} \times$
13 $(0.119 + 0.204 \times \log(\text{depth}))$). The values used for this correction are based on the relationship
14 between observed and expected synonymous DNMs at different levels of coverage. We stratified
15 this analysis into known dominant, known recessive and non-DD associated genes using the DDG2P
16 gene list (<http://www.ebi.ac.uk/gene2phenotype>), downloaded in January 2016. Genes with
17 recessive and dominant modes of inheritance were restricted to the recessive list (see Supplemental
18 Table S4). Observed and expected numbers of DNMs were also calculated in subsets of the DDD
19 probands with confident diagnoses (individuals with a reported variant classed as pathogenic or
20 likely pathogenic by the referring clinician) and those lacking a potential diagnosis (diagnosed
21 $n=1417$, undiagnosed $n=3364$, with the remainder of the cohort having possible or uncertain
22 diagnostic states, as of January 2018). We used the Poisson test (using R's `poisson.test`, with two
23 sided alternative hypothesis) to examine differences in the observed and expected values, and a 5%
24 FDR correction to control for multiple testing using the `p.adjust` R package (method=fdr) across all
25 tests (R Core Team 2018) (R v3.1.3).

1 PPVs were calculated $((\text{observed} - \text{expected}) / \text{observed})$ for CSS positions, combined and
2 individually, don+5 sites, PolyPy PyPu, PolyPy other, other near splice exonic and intronic variants, as
3 well as VEP defined missense and stop gained mutations.

4 We divided our exons into sextiles based on the pLI (Lek et al. 2016) of the gene to which they
5 belong, and calculated the observed and expected number of DNMs in each sextile for don+5,
6 PolyPy PyPu and synonymous variants (as above) to see if the enrichment of don+5 and PolyPy PyPu
7 changes was concentrated in genes more likely to be intolerant of loss of function (LoF) mutations.

8 **Potential diagnostic variants**

9 DNMs overlapping with our near-splice positions of interest within dominant DDG2P genes were
10 identified in DDD probands lacking a potential explanatory variant (Dec 2016, n = 5907). The Human
11 Phenotype Ontology (HPO, <http://compbio.charite.de/hpweb/>) encoded (Kohler et al. 2017)
12 phenotypes of the probands were assessed by consultant clinical geneticist Helen V. Firth, along with
13 the patient's recruiting clinician, and compared to the known clinical presentation of individuals with
14 LoF mutations within those genes, classifying each variant as likely diagnostic, unlikely diagnostic, or
15 unsure, depending on the strength of similarity between the proband and the disorder, and the
16 specificity of the phenotype. The relationship between our PPVs and the proportion of clinical
17 diagnoses in each class of near splice mutation was assessed using Pearson's product-moment
18 correlation using the `cor.test` function in R (R Core Team 2018) (version 3.4.4).

19 The proportion of CSS to non-CSS splicing diagnoses was calculated, along with 95% CIs, based on 18
20 non-CSS diagnoses and 48 CSS diagnoses in the same regions using the `prop.test` package in R (R
21 Core Team 2018) (version 3.4.4).

22 **Validation of putative splicing variants**

23 Eight variants were selected for validation via a minigene vector system. These comprised six likely
24 diagnostic variants from the PolyPy, a PolyPy variant of uncertain clinical significance, and a likely

1 diagnostic don+5 variant. Additionally, two untransmitted variants identified in unaffected parents
2 within the same PolyPys as test variants were selected as negative controls. Details of all variants
3 selected for validation are shown in Supplemental Table S5.

4 **Cloning splicing vectors**

5 The minigene splice assay vector was adapted from that used in Singh *et al.* (Singh et al. 2016), by
6 replacing intron 1 with the first intron from the rat insulin 2 gene (Ins-2)(Rnor_6.0 Chr1:215857148-
7 215857695). To generate individual assay vectors, either the 5' most 231bp (for the don+5 variant)
8 or the 3' most 274bp (for PolyPy variants) of this vector was replaced with the appropriate
9 endogenous intronic sequence encompassing the DNM of interest (Figure S4a and S4b), as described
10 below. Between 114 and 202bp flanking endogenous intronic sequence was included, along with
11 6bp local exonic sequence from the gene of interest.

12 First, proband genotypes (Supplemental Table S5) were verified by capillary sequencing of genomic
13 PCR products (Supplemental Table S6). Genomic regions containing the reference and alternate
14 sequences were then either amplified by nested PCR, generated by site directed mutagenesis, or
15 generated using gene synthesis (IDT). These fragments were sub-cloned, by Gibson Assembly (NEB),
16 into our minigene vector (Supplemental Table S7, Supplemental Table S8). The regions assayed in
17 our vectors are detailed by genomic coordinates in Supplemental Table S5.

18 ***In vitro* splicing assay**

19 HeLa cells were seeded into 12-well plates at a density of 160,000 cells per well, grown for 24 hours
20 and transfected with 1 microgram of plasmid vector using Lipofectamine 3000 (Invitrogen). All
21 transfections were carried out in duplicate and cultured for 48 hours. HeLa cells were cultured in
22 DMEM (10% FCS + 1% pen/strep) at 37°C in a humidified incubator. Total RNA was extracted using a
23 Micro RNeasy Qiagen kit and mRNA converted into cDNA using SuperScript IV (Invitrogen). RT-PCR
24 was carried out using primers designed to span from exon 1 to exon 2, exon 2 to exon 3 and exon 1

1 to exon 3 and amplified on a thermocycler for either 25 or 35 cycles (Supplemental Table S9).
2 Amplicons were capillary sequenced (GATC). For amplicons showing more than one splice variant
3 (mixed capillary traces, for *CHD7*-Alt and *MBD5*-Alt), we cloned the PCR amplicons (Zero Blunt PCR
4 cloning kit, Invitrogen) and sequenced individual colonies by capillary sequencing to identify the
5 splice variants present (Supplemental Table S10).

6 Chromatograms were generated in R (R Core Team 2018) from .ab1 files using the sangerseqR(Hill et
7 al. 2014) package via Bioconductor (Huber et al. 2015)
8 (<http://bioconductor.org/packages/release/bioc/html/sangerseqR.html>, R v3.1.3), and likely
9 consequences on the protein primary structure were generated using reference and alternative RNA
10 sequences with the ExPASy Nucleotide Sequence Translation tool(Artimo et al. 2012)
11 (<https://web.expasy.org/translate/>).

12 **Splicing pathogenicity scores**

13 Since our region of interest spanned >6 million individual positions, each with three potential single
14 nucleotide changes, we were restricted in the choice of splicing pathogenicity prediction tools we
15 could utilise, as many function primarily through a low throughput web interface model. We
16 identified three resources recently published which provide “genome wide” splicing pathogenicity
17 scores. Two methods, dbSCSNV’s AdaBoost and RandomForest are based on ensemble learning
18 combining predictions from multiple other splice prediction tools as well as conservation and CADD
19 scores(Jian et al. 2014a). The targeted region at the acceptor end spans 14 bases (12 intronic, 2
20 exonic) and at the donor end spans 11 bases (8 intronic, 3 exonic). Spidex utilises deep learning
21 methods trained on RNA sequencing data to estimate the consequence of variants on the “percent
22 spliced in” of an exon, relative to the reference sequence(Xiong et al. 2015). Spidex scores positions
23 up to 300bp from intron/exon boundaries, so provides greater coverage of our splicing region of
24 interest. We also utilised the longer standing, and widely used MaxEntScan (MES)(Yeo and Burge
25 2004), for which Perl scripts were available, allowing the tool to be run locally for all alternative

1 alleles of all positions of interest. The metric used for MES was the percent difference between the
2 scores for the reference and alternative alleles, with the greatest reduction in score classed as most
3 pathogenic. All sites were also scored with CADD(Kircher et al. 2014).

4 To allow cross-tool comparison, we ordered positions by increasing pathogenicity from each metric,
5 and split positions into 20 brackets, such that the cumulative triplet based mutation rate for all
6 variants in each bracket was equal, and the 20th bracket contained the positions with the most
7 pathogenic scores. We calculated MAPS and the proportion of parental variants falling in high pLI
8 genes for each bracket for all five metrics, as above, and looked at the number of DNMs in known
9 dominant genes which fell in each bracket for the five metrics. Each of these analyses was conducted
10 including and excluding CSS dinucleotide positions.

11 **Splice region variants in the ClinVar database**

12 We extracted all ClinVar(Landrum et al. 2016) (<https://www.ncbi.nlm.nih.gov/clinvar/>) variants using
13 the UCSC Table Browser(Karolchik et al. 2004) on 02.05.2017 and matched these against our splicing
14 positions of interest, removing exonic sites with non-synonymous consequences. This resulted in
15 3603 positions with clinical significance recorded as pathogenic or likely pathogenic. We calculated
16 the ratio of canonical to non-canonical splice positions within this data. Since each variant is present
17 in this data only once, we used number of submissions as a proxy for allele count, and calculated the
18 ratio of canonical to non-canonical variants adjusting for this. Differences between these observed
19 values and our expectations, based on 27% of splice affecting mutations being in non-canonical
20 positions, were assessed using Fisher's exact test (R Core Team 2018)(R v3.1.3).

21 **Software availability**

22 Code and data to reproduce the analyses within this paper are available in Supplementary code and
23 figures, as well as on GitHub (https://github.com/JLord86/DDD_Splicing).

1 **Acknowledgments**

2 We thank the families for their participation and patience. We are grateful to the Exome Aggregation
3 Consortium for making their data and code available. We thank the Sanger Human Genome
4 Informatics and DNA pipelines teams for their support in generating and processing the data. We are
5 grateful to Adam Frankish for advice selecting an appropriate exon set, and to Sarah Hunt and Fiona
6 Cunningham for help and advice regarding splice annotation, and the development of the VEP
7 SpliceRegion.pm plugin. Thanks also go to Alan Donaldson, Alex Henderson, Anand Sagggar, Diana
8 Baralle, Elisabeth Rosser, Elizabeth Jones, Emma Wakeling, Fleur van Dijk, Joan Paterson, Joanna
9 Jarvis, Kate Chandler, Katherine Lachlan, Miranda Splitt, Neeti Ghali, Rachel Harrison, Sahar
10 Mansour, Shane Mckee, Susan Tomkins and Victoria McKay for providing phenotypic information
11 and insight on the probands with potential splice disrupting variants. The DDD study presents
12 independent research commissioned by the Health Innovation Challenge Fund (grant
13 HICF110091003), a parallel funding partnership between the Wellcome Trust and the UK
14 Department of Health, and the Wellcome Trust Sanger Institute (grant WT098051). The views
15 expressed in this publication are those of the author(s) and not necessarily those of the Wellcome
16 Trust or the UK Department of Health. The study has UK Research Ethics Committee approval
17 (10/H0305/83, granted by the Cambridge South Research Ethics Committee and GEN/284/12,
18 granted by the Republic of Ireland Research Ethics Committee).

19 **Disclosure Declaration**

20 M.E.H. is a co-founder of, consultant to, and holds shares in, Congenica Ltd, a genetics diagnostic
21 company.

1 References

- 2 Ars E, Serra E, Garcia J, Kruyer H, Gaona A, Lazaro C, Estivill X. 2000. Mutations affecting mRNA
3 splicing are the most common molecular defects in patients with neurofibromatosis type 1.
4 *Hum Mol Genet* **9**: 237-247.
- 5 Artimo P, Jonnalagedda M, Arnold K, Baratin D, Csardi G, de Castro E, Duvaud S, Flegel V, Fortier A,
6 Gasteiger E et al. 2012. ExpASy: SIB bioinformatics resource portal. *Nucleic Acids Res* **40**:
7 W597-603.
- 8 Aten E, Sun Y, Almomani R, Santen GW, Messemaker T, Maas SM, Breuning MH, den Dunnen JT.
9 2013. Exome sequencing identifies a branch point variant in Aarskog-Scott syndrome. *Hum*
10 *Mutat* **34**: 430-434.
- 11 Badr E, ElHefnawi M, Heath LS. 2016. Computational Identification of Tissue-Specific Splicing
12 Regulatory Elements in Human Genes from RNA-Seq Data. *PLoS One* **11**: e0166978.
- 13 Baralle D, Buratti E. 2017. RNA splicing in human disease and in the clinic. *Clin Sci (Lond)* **131**: 355-
14 368.
- 15 Baralle M, Skoko N, Knezevich A, De Conti L, Motti D, Bhuvanagiri M, Baralle D, Buratti E, Baralle FE.
16 2006. NF1 mRNA biogenesis: effect of the genomic milieu in splicing regulation of the NF1
17 exon 37 region. *FEBS Lett* **580**: 4449-4456.
- 18 Basel-Vanagaite L, Hershkovitz T, Heyman E, Raspall-Chaure M, Kakar N, Smirin-Yosef P, Vila-Pueyo
19 M, Kornreich L, Thiele H, Bode H et al. 2013. Biallelic SZT2 mutations cause infantile
20 encephalopathy with epilepsy and dysmorphic corpus callosum. *Am J Hum Genet* **93**: 524-
21 529.
- 22 Bonnet C, Krieger S, Vezain M, Rousselin A, Tournier I, Martins A, Berthet P, Chevrier A, Dugast C,
23 Layet V et al. 2008. Screening BRCA1 and BRCA2 unclassified variants for splicing mutations
24 using reverse transcription PCR on patient RNA and an ex vivo assay based on a splicing
25 reporter minigene. *J Med Genet* **45**: 438-446.
- 26 Brody E, Abelson J. 1985. The "spliceosome": yeast pre-messenger RNA associates with a 40S
27 complex in a splicing-dependent reaction. *Science* **228**: 963-967.
- 28 Brunham LR, Kang MH, Van Karnebeek C, Sadananda SN, Collins JA, Zhang LH, Sayson B, Miao F,
29 Stockler S, Frohlich J et al. 2015. Clinical, Biochemical, and Molecular Characterization of
30 Novel Mutations in ABCA1 in Families with Tangier Disease. *JIMD Rep* **18**: 51-62.
- 31 Caceres EF, Hurst LD. 2013. The evolution, impact and properties of exonic splice enhancers.
32 *Genome Biol* **14**: R143.
- 33 Caminsky N, Mucaki EJ, Rogan PK. 2014. Interpretation of mRNA splicing mutations in genetic
34 disease: review of the literature and guidelines for information-theoretical analysis.
35 *F1000Res* **3**: 282.
- 36 Cartegni L, Chew SL, Krainer AR. 2002. Listening to silence and understanding nonsense: exonic
37 mutations that affect splicing. *Nat Rev Genet* **3**: 285-298.
- 38 Cingolani P, Platts A, Wang le L, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM. 2012. A
39 program for annotating and predicting the effects of single nucleotide polymorphisms,
40 SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly*
41 (*Austin*) **6**: 80-92.
- 42 Corvelo A, Hallegger M, Smith CW, Eyra E. 2010. Genome-wide association between branch point
43 properties and alternative splicing. *PLoS Comput Biol* **6**: e1001016.
- 44 Crotti L, Lewandowska MA, Schwartz PJ, Insolia R, Pedrazzini M, Bussani E, Dagradi F, George AL, Jr.,
45 Pagni F. 2009. A KCNH2 branch point mutation causing aberrant splicing contributes to an
46 explanation of genotype-negative long QT syndrome. *Heart Rhythm* **6**: 212-218.
- 47 Cummings BB, Marshall JL, Tukiainen T, Lek M, Donkervoort S, Foley AR, Bolduc V, Waddell LB,
48 Sandaradura SA, O'Grady GL et al. 2017. Improving genetic diagnosis in Mendelian disease
49 with transcriptome sequencing. *Sci Transl Med* **9**.

1 Deciphering Developmental Disorders Study. 2015. Large-scale discovery of novel genetic causes of
2 developmental disorders. *Nature* **519**: 223-228.

3 Deciphering Developmental Disorders Study. 2017. Prevalence and architecture of de novo
4 mutations in developmental disorders. *Nature* **542**: 433-438.

5 Di Leo E, Panico F, Tarugi P, Battisti C, Federico A, Calandra S. 2004. A point mutation in the lariat
6 branch point of intron 6 of NPC1 as the cause of abnormal pre-mRNA splicing in Niemann-
7 Pick type C disease. *Hum Mutat* **24**: 440.

8 Fairbrother WG, Yeh RF, Sharp PA, Burge CB. 2002. Predictive identification of exonic splicing
9 enhancers in human genes. *Science* **297**: 1007-1013.

10 Ferreira PG, Oti M, Barann M, Wieland T, Ezquina S, Friedlander MR, Rivas MA, Esteve-Codina A,
11 Consortium G, Rosenstiel P et al. 2016. Sequence variation between 462 human individuals
12 fine-tunes functional sites of RNA processing. *Sci Rep* **6**: 32406.

13 Goren A, Ram O, Amit M, Keren H, Lev-Maor G, Vig I, Pupko T, Ast G. 2006. Comparative analysis
14 identifies exonic splicing regulatory sequences--The complex definition of enhancers and
15 silencers. *Mol Cell* **22**: 769-781.

16 Hang J, Wan R, Yan C, Shi Y. 2015. Structural basis of pre-mRNA splicing. *Science* **349**: 1191-1198.

17 Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, Aken BL, Barrell D, Zadissa
18 A, Searle S et al. 2012. GENCODE: the reference human genome annotation for The ENCODE
19 Project. *Genome Res* **22**: 1760-1774.

20 Hill JT, Demarest BL, Bisgrove BW, Su YC, Smith M, Yost HJ. 2014. Poly peak parser: Method and
21 software for identification of unknown indels using sanger sequencing of polymerase chain
22 reaction products. *Dev Dyn* **243**: 1632-1636.

23 Houdayer C, Caux-Moncoutier V, Krieger S, Barrois M, Bonnet F, Bourdon V, Bronner M, Buisson M,
24 Coulet F, Gaildrat P et al. 2012. Guidelines for splicing analysis in molecular diagnosis derived
25 from a set of 327 combined in silico/in vitro studies on BRCA1 and BRCA2 variants. *Hum*
26 *Mutat* **33**: 1228-1238.

27 Huber W, Carey VJ, Gentleman R, Anders S, Carlson M, Carvalho BS, Bravo HC, Davis S, Gatto L, Girke
28 T et al. 2015. Orchestrating high-throughput genomic analysis with Bioconductor. *Nat*
29 *Methods* **12**: 115-121.

30 Iossifov I, O'Roak BJ, Sanders SJ, Ronemus M, Krumm N, Levy D, Stessman HA, Witherspoon KT, Vives
31 L, Patterson KE et al. 2014. The contribution of de novo coding mutations to autism
32 spectrum disorder. *Nature* **515**: 216-221.

33 Ito K, Patel PN, Gorham JM, McDonough B, DePalma SR, Adler EE, Lam L, MacRae CA, Mohiuddin
34 SM, Fatkin D et al. 2017. Identification of pathogenic gene mutations in LMNA and MYBPC3
35 that alter RNA splicing. *Proc Natl Acad Sci U S A* **114**: 7689-7694.

36 Jian X, Boerwinkle E, Liu X. 2014a. In silico prediction of splice-altering single nucleotide variants in
37 the human genome. *Nucleic Acids Res* **42**: 13534-13544.

38 Jian X, Boerwinkle E, Liu X. 2014b. In silico tools for splicing defect prediction: a survey from the
39 viewpoint of end users. *Genet Med* **16**: 497-503.

40 Karolchik D, Hinrichs AS, Furey TS, Roskin KM, Sugnet CW, Haussler D, Kent WJ. 2004. The UCSC
41 Table Browser data retrieval tool. *Nucleic Acids Res* **32**: D493-496.

42 Ke S, Anquetil V, Zamalloa JR, Maity A, Yang A, Arias MA, Kalachikov S, Russo JJ, Ju J, Chasin LA. 2018.
43 Saturation mutagenesis reveals manifold determinants of exon definition. *Genome Res* **28**:
44 11-24.

45 Ke S, Shang S, Kalachikov SM, Morozova I, Yu L, Russo JJ, Ju J, Chasin LA. 2011. Quantitative
46 evaluation of all hexamers as exonic splicing elements. *Genome Res* **21**: 1360-1374.

47 Kircher M, Witten DM, Jain P, O'Roak BJ, Cooper GM, Shendure J. 2014. A general framework for
48 estimating the relative pathogenicity of human genetic variants. *Nat Genet* **46**: 310-315.

49 Kohler S, Vasilevsky NA, Engelstad M, Foster E, McMurry J, Ayme S, Baynam G, Bello SM, Boerkoel
50 CF, Boycott KM et al. 2017. The Human Phenotype Ontology in 2017. *Nucleic Acids Res* **45**:
51 D865-D876.

1 Krawczak M, Thomas NS, Hundrieser B, Mort M, Wittig M, Hampe J, Cooper DN. 2007. Single base-
2 pair substitutions in exon-intron junctions of human genes: nature, distribution, and
3 consequences for mRNA splicing. *Hum Mutat* **28**: 150-158.

4 Landrum MJ, Lee JM, Benson M, Brown G, Chao C, Chitipiralla S, Gu B, Hart J, Hoffman D, Hoover J et
5 al. 2016. ClinVar: public archive of interpretations of clinically relevant variants. *Nucleic Acids*
6 *Res* **44**: D862-868.

7 Lastella P, Surdo NC, Resta N, Guanti G, Stella A. 2006. In silico and in vivo splicing analysis of MLH1
8 and MSH2 missense mutations shows exon- and tissue-specific effects. *BMC Genomics* **7**:
9 243.

10 Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, O'Donnell-Luria AH, Ware JS, Hill
11 AJ, Cummings BB et al. 2016. Analysis of protein-coding genetic variation in 60,706 humans.
12 *Nature* **536**: 285-291.

13 Lewandowska MA. 2013. The missing puzzle piece: splicing mutations. *Int J Clin Exp Pathol* **6**: 2675-
14 2682.

15 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
16 *Bioinformatics* **25**: 1754-1760.

17 Liu HX, Cartegni L, Zhang MQ, Krainer AR. 2001. A mechanism for exon skipping caused by nonsense
18 or missense mutations in BRCA1 and other genes. *Nat Genet* **27**: 55-58.

19 Lorson CL, Hahnen E, Androphy EJ, Wirth B. 1999. A single nucleotide in the SMN gene regulates
20 splicing and is responsible for spinal muscular atrophy. *Proc Natl Acad Sci U S A* **96**: 6307-
21 6311.

22 Maslen C, Babcock D, Raghunath M, Steinmann B. 1997. A rare branch-point mutation is associated
23 with missplicing of fibrillin-2 in a large family with congenital contractural arachnodactyly.
24 *Am J Hum Genet* **60**: 1389-1398.

25 McCarthy DJ, Humburg P, Kanapin A, Rivas MA, Gaulton K, Cazier JB, Donnelly P. 2014. Choice of
26 transcripts and software has a large effect on variant annotation. *Genome Med* **6**: 26.

27 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytzky A, Garimella K, Altshuler D,
28 Gabriel S, Daly M et al. 2010. The Genome Analysis Toolkit: a MapReduce framework for
29 analyzing next-generation DNA sequencing data. *Genome Res* **20**: 1297-1303.

30 McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, Flicek P, Cunningham F. 2016. The
31 Ensembl Variant Effect Predictor. *Genome Biol* **17**: 122.

32 Mercer TR, Clark MB, Andersen SB, Brunck ME, Haerty W, Crawford J, Taft RJ, Nielsen LK, Dinger ME,
33 Mattick JS. 2015. Genome-wide discovery of human splicing branchpoints. *Genome Res* **25**:
34 290-303.

35 R Core Team. 2018. A language and environment for statistical computing.

36 Ramu A, Noordam MJ, Schwartz RS, Wuster A, Hurles ME, Cartwright RA, Conrad DF. 2013.
37 DeNovoGear: de novo indel and point mutation discovery and phasing. *Nat Methods* **10**:
38 985-987.

39 Samocha KE, Robinson EB, Sanders SJ, Stevens C, Sabo A, McGrath LM, Kosmicki JA, Rehnstrom K,
40 Mallick S, Kirby A et al. 2014. A framework for the interpretation of de novo mutation in
41 human disease. *Nat Genet* **46**: 944-950.

42 Sangermano R, Khan M, Cornelis SS, Richelle V, Albert S, Elmelik D, Garanto A, Qamar R, Lugtenberg
43 D, van den Born LI et al. 2017. ABCA4 midigenes reveal the full splice spectrum of all
44 reported non-canonical splice site variants in Stargardt disease. *Genome Res*
45 doi:10.1101/gr.226621.117.

46 Scotti MM, Swanson MS. 2016. RNA mis-splicing in disease. *Nat Rev Genet* **17**: 19-32.

47 Short PJ, McRae JF, Gallone G, Sifrim A, Won H, Geschwind DH, Wright CF, Firth HV, FitzPatrick DR,
48 Barrett JC et al. 2018. De novo mutations in regulatory elements in neurodevelopmental
49 disorders. *Nature* **555**: 611-616.

- 1 Singh T, Kurki MI, Curtis D, Purcell SM, Crooks L, McRae J, Suvisaari J, Chheda H, Blackwood D, Breen
2 G et al. 2016. Rare loss-of-function variants in SETD1A are associated with schizophrenia and
3 developmental disorders. *Nat Neurosci* **19**: 571-577.
- 4 Soemedi R, Cygan KJ, Rhine CL, Wang J, Bulacan C, Yang J, Bayrak-Toydemir P, McDonald J,
5 Fairbrother WG. 2017. Pathogenic variants that alter protein code often disrupt splicing. *Nat*
6 *Genet* **49**: 848-855.
- 7 Soukarieh O, Gaildrat P, Hamieh M, Drouet A, Baert-Desurmont S, Frebourg T, Tosi M, Martins A.
8 2016. Exonic Splicing Mutations Are More Prevalent than Currently Estimated and Can Be
9 Predicted by Using In Silico Tools. *PLoS Genet* **12**: e1005756.
- 10 Taggart AJ, Lin CL, Shrestha B, Heintzelman C, Kim S, Fairbrother WG. 2017. Large-scale analysis of
11 branchpoint usage across species and cell lines. *Genome Res* **27**: 639-649.
- 12 Tang R, Prosser DO, Love DR. 2016. Evaluation of Bioinformatic Programmes for the Analysis of
13 Variants within Splice Site Consensus Regions. *Adv Bioinformatics* **2016**: 5614058.
- 14 Teraoka SN, Telatar M, Becker-Catania S, Liang T, Onengut S, Tolun A, Chessa L, Sanal O,
15 Bernatowska E, Gatti RA et al. 1999. Splicing defects in the ataxia-telangiectasia gene, ATM:
16 underlying mutations and consequences. *Am J Hum Genet* **64**: 1617-1631.
- 17 Thery JC, Krieger S, Gaildrat P, Revillion F, Buisine MP, Killian A, Duponchel C, Rousselin A, Vaur D,
18 Peyrat JP et al. 2011. Contribution of bioinformatics predictions and functional splicing
19 assays to the interpretation of unclassified variants of the BRCA genes. *Eur J Hum Genet* **19**:
20 1052-1058.
- 21 van der Klift HM, Jansen AM, van der Steenstraten N, Bik EC, Tops CM, Devilee P, Wijnen JT. 2015.
22 Splicing analysis for exonic and intronic mismatch repair gene variants associated with Lynch
23 syndrome confirms high concordance between minigene assays and patient RNA analyses.
24 *Mol Genet Genomic Med* **3**: 327-345.
- 25 Vaz-Drago R, Custodio N, Carmo-Fonseca M. 2017. Deep intronic mutations and human disease.
26 *Hum Genet* doi:10.1007/s00439-017-1809-4.
- 27 Wang K, Li M, Hakonarson H. 2010. ANNOVAR: functional annotation of genetic variants from high-
28 throughput sequencing data. *Nucleic Acids Res* **38**: e164.
- 29 Wang Y, Wang Z. 2014. Systematical identification of splicing regulatory cis-elements and cognate
30 trans-factors. *Methods* **65**: 350-358.
- 31 Wright CF, Fitzgerald TW, Jones WD, Clayton S, McRae JF, van Kogelenberg M, King DA, Ambridge K,
32 Barrett DM, Bayzatinova T et al. 2015. Genetic diagnosis of developmental disorders in the
33 DDD study: a scalable analysis of genome-wide research data. *Lancet* **385**: 1305-1314.
- 34 Xiong HY, Alipanahi B, Lee LJ, Bretschneider H, Merico D, Yuen RK, Hua Y, Gueroussov S, Najafabadi
35 HS, Hughes TR et al. 2015. RNA splicing. The human splicing code reveals new insights into
36 the genetic determinants of disease. *Science* **347**: 1254806.
- 37 Yeo G, Burge CB. 2004. Maximum entropy modeling of short sequence motifs with applications to
38 RNA splicing signals. *J Comput Biol* **11**: 377-394.
- 39 Zhang XH, Chasin LA. 2004. Computational definition of sequence motifs governing constitutive exon
40 splicing. *Genes Dev* **18**: 1241-1250.

41

42