

# MRSD: a novel quantitative approach for assessing suitability of RNA-seq in the clinical investigation of mis-splicing in Mendelian disease

Charlie F. Rowlands<sup>1,2</sup>, Algy Taylor<sup>2</sup>, Gillian Rice<sup>1</sup>, Nicola Whiffin<sup>3</sup>, Hildegard Nikki Hall<sup>4</sup>, William G. Newman<sup>1,2</sup>, Graeme C.M. Black<sup>1,2</sup>, kConFab Investigators<sup>5,6</sup>, Raymond T. O’Keefe<sup>1</sup>, Simon Hubbard<sup>1</sup>, Andrew G.L. Douglas<sup>7,8</sup>, Diana Baralle<sup>7,8</sup>, Tracy A. Briggs<sup>1,2</sup>, Jamie M. Ellingford<sup>1,2</sup>

1. Division of Evolution and Genomic Sciences, School of Biological Sciences, Faculty of Biology, Medicine and Health, University of Manchester, Manchester, UK

2. Manchester Centre for Genomic Medicine, St Mary’s Hospital, Manchester University NHS Foundation Trust, Health Innovation Manchester, Manchester, UK

3. Wellcome Centre for Human Genetics, University of Oxford, Oxford, UK

4. MRC Human Genetics Unit, Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh, UK

5. Sir Peter MacCallum Department of Oncology, The University of Melbourne, Parkville, VIC 3010, Australia

6. kConFab, Research Department, Peter MacCallum Cancer Centre, Melbourne, VIC 3000, Australia

7. Wessex Clinical Genetics Service, Princess Anne Hospital, University Hospital Southampton NHS Foundation Trust, Coxford Rd, Southampton, SO16 5YA, UK

8. Faculty of Medicine, University of Southampton, Duthie Building, Southampton General Hospital, Tremona Road, Southampton, SO16 6YD, UK

## Abstract

**Background:** RNA-sequencing of patient biosamples is a promising approach to delineate the impact of genomic variants on splicing, but variable gene expression between tissues complicates selection of appropriate tissues. Relative expression level is often used as a metric to predict RNA-sequencing utility. Here, we describe a gene- and tissue-specific metric to inform the feasibility of RNA-sequencing, overcoming some issues with using expression values alone.

**Results:** We derive a novel metric, *Minimum Required Sequencing Depth* (MRSD), for all genes across three human biosamples (whole blood, lymphoblastoid cell lines (LCLs) and skeletal muscle). MRSD estimates the depth of sequencing required from RNA-sequencing to achieve user-specified sequencing coverage of a gene, transcript or group of genes of interest. MRSD predicts levels of splice junction coverage with high precision (90.1-98.2%) and overcomes transcript region-specific sequencing biases. Applying MRSD scoring to established disease gene panels shows that LCLs are the optimum source of RNA, of the three investigated biosamples, for 69.3% of gene panels. Our approach demonstrates that up to 59.4% of variants of uncertain significance in ClinVar predicted to impact splicing could be functionally assayed by RNA-sequencing in at least one of the investigated biosamples.

**Conclusions:** We demonstrate the power of MRSD as a metric to inform choice of appropriate biosamples for the functional assessment of splicing aberrations. We apply MRSD in the context of Mendelian genetic disorders and illustrate its benefits over expression-based approaches. We anticipate that the integration of MRSD into clinical pipelines will improve variant interpretation and, ultimately, diagnostic yield.

## 50 **Introduction**

51 Pinpointing disease-causing genomic variation informs diagnosis, treatment and  
52 management for a wide range of rare disorders. An underappreciated group of  
53 pathogenic variants is those that lie outside of canonical splice sites but act through  
54 disruption of pre-mRNA splicing, the process whereby introns are removed from  
55 nascent pre-mRNA to produce mature and functional transcripts (Supplementary  
56 Figure 1a). The ways through which genomic variants can disrupt pre-mRNA splicing  
57 are diverse (Supplementary Figures 1b-g), including both protein-coding and intronic  
58 variants that are well described as causes of rare disorders (1-3). However, the  
59 omission of intronic regions in targeted sequencing approaches (4, 5), discordance  
60 between *in silico* variant prioritization tools (6) and the lack of availability of the  
61 appropriate tissue from which to survey RNA for splicing disruption (7, 8) limit  
62 effective identification of pathogenic splice-impacting variants.

63

64 RNA sequencing (RNA-seq) offers a potential route to overcome issues of variant  
65 interpretation (3, 9-12). The complex impacts of variants on splicing can be fully  
66 characterized through RNA-seq. Moreover, aberrant splicing events can be identified  
67 from RNA-seq datasets without prior knowledge of genomic variants driving their  
68 impact. Whilst targeted analyses, such as RT-PCR, also enable detection of splicing  
69 aberrations (3), such approaches are designed to test the presence of specific  
70 disruptions and may not identify the complete spectrum of splicing disruption caused  
71 by a single genomic variant.

72

73 There is growing evidence that RNA-seq can substantially improve diagnostic yield  
74 across a variety of disease subtypes (3, 10, 13-15) through identification of variants

75 impacting splicing, or leading to impairment of transcript expression or stability (16).  
76 However, there remain several hurdles to the effective and routine integration of  
77 RNA-seq into diagnostic pipelines. For example, surveying a whole transcriptome  
78 identifies a large number of aberrant splicing events – in the order of hundreds of  
79 thousands – and there is little consensus regarding the best approach to filter for true  
80 positive and pathogenic events. Furthermore, diagnostic analysis using RNA-seq is  
81 only effective when sufficient levels of sequence coverage of a relevant gene  
82 transcript are present in the sampled tissue.

83

84 In this study, we develop an informatics approach to quantify the likelihood that a  
85 gene/transcript, or a defined set of genes or transcripts, can be appropriately  
86 surveyed using RNA-seq. We name our framework the *Minimum Required*  
87 *Sequencing Depth* (MRSD), which can be utilized in a flexible and customized  
88 manner to assess the suitability of RNA-seq derived from different tissues to identify  
89 pathogenic splicing aberrations in specific genes of interest. MRSD scores (available  
90 at: <https://mcgm-mrds.github.io/>) can be utilized to select the most appropriate  
91 biosample to detect splicing aberrations for a candidate set of genes/transcripts or to  
92 guide the amount of sequencing reads from a specific biosample required to  
93 generate appropriate transcriptomic datasets for a gene of interest. We apply these  
94 techniques to the study of monogenic disease genes, and assess three clinically  
95 accessible biosamples for their appropriateness to survey all known monogenic  
96 disease genes.

97

98

99

## 100 **Results**

### 101 *Minimum Required Sequencing Depth scores differ across biosamples*

102 We first derived MRSD scores, corresponding to the required sequencing depth (in  
103 M uniquely mapping sequencing reads) for a specified level of coverage of a  
104 transcript, for 3112 known multi-exon disease genes in three distinct tissues (blood,  
105 LCLs and skeletal muscle). Three parameters can be altered for the MRSD model;  
106 we observed that MRSDs differed dependent on the values chosen for these  
107 parameters, comprising the number of reads desired to cover each splice junction,  
108 the proportion of splicing junctions for each gene that meet this coverage threshold  
109 (75% or 95%), and the proportion of samples for which the prediction is predicted to  
110 be sufficient (the “confidence level” of either 95% or 99%; Figure 1). For example,  
111 across all three tissues at a specified read coverage level of eight reads per splicing  
112 junction, we observed that increases in the desired proportion of covered splice  
113 junctions from 75-95% was associated with an increase in median MRSD of between  
114 5.4% (in blood) to 61.2% (in LCLs; Figure 1a, top). In general, increasing desired  
115 confidence level for appropriate splice junction coverage from 95% to 99% resulted  
116 in an increase in median MRSD of between 25.8-85.8%. Conversely, for skeletal  
117 muscle samples, when stipulating 95% splice junction coverage, we observed a  
118 decrease of 3.1% in MRSD scores when desired confidence level was increased  
119 from 95% ( $n = 1241$ , median = 41.83) to 99% ( $n = 921$ , median = 40.54); this was  
120 accounted for by an increase in the number of genes that were considered  
121 “unfeasible” for surveillance, i.e. those for which zero reads cover the given  
122 proportion of junctions ( $n$  unfeasible at 95% confidence = 1873,  $n$  unfeasible at 99%  
123 confidence = 2193). This definition of feasibility is limited by the sequencing depth of  
124 the control models on which the predictions are based. For example, no coverage of

125 splice junctions in a particular transcript may have been observed simply due to low  
126 sequencing depth; with ultra-deep sequencing of the same sample, we may have  
127 observed coverage of splice junctions and so have been able to generate a feasible  
128 MRSD prediction.

129

130 Overall, these analyses suggested that, of the three investigated biosamples, LCLs  
131 would enable investigation of the most comprehensive set of genes for aberrant  
132 splicing. This conclusion was supported by LCLs displaying, across all four  
133 parameter combinations, the lowest median MRSDs (range = 12.86-33.77, Figure 1b,  
134 top), and the fewest “unfeasible” genes (43-63%). On the other hand, whole blood  
135 exhibited the highest number of unfeasible genes across the different parameter  
136 combinations (61-84%).

137

### 138 *Accuracy of Minimum Required Sequencing Depth calculations*

139 We next obtained RNA-seq datasets for 68 samples from the three investigated  
140 tissues (blood,  $n = 12$ ; LCLs,  $n = 4$ ; muscle,  $n = 52$ ), with a wide range of sequencing  
141 depths (Supplementary Figure 2). We assessed the performance of the MRSD  
142 model against these datasets, defining the positive predictive value (PPV) of MRSD  
143 as the likelihood that appropriate sequencing coverage was obtained given that the  
144 level of sequencing depth exceeded the MRSD prediction. Conversely, the negative  
145 predictive value (NPV) was defined as the likelihood that appropriate sequencing  
146 coverage was not obtained, given that the sample did not meet the specified criteria  
147 of the MRSD prediction. Across all investigated MRSD parameters, we observed 96%  
148 PPV and 79% NPV, on average, for the 68 samples (Figure 2a). We observed a  
149 general trend that the PPV and NPV of MRSD decreased and increased,

150 respectively, as higher levels of required coverage were imposed (Figure 2b-c).  
151 Across all parameter combinations, PPV values ranged from 90.1-98.2%, while NPV  
152 ranged from 56.4-94.7%, suggesting MRSD is a fairly conservative model that  
153 primarily returns positive results with high certainty.

154

155 Interestingly, although MRSD scores were derived from 75 bp paired-end RNA-seq  
156 data, evaluating the ability of the model to predict transcript coverage in 150 bp  
157 paired-end data (LCLs,  $n=20$ ) shows higher PPV than with 75 bp data for half of the  
158 four parameter combinations tested, while NPV was only slightly lower for all  
159 combinations (Supplementary Figure 3). This suggests that, while care must be  
160 taken applying this approach to datasets derived using alternative experimental  
161 approaches, the MRSD model described here may provide a suitable approximation  
162 in the case of alternative sequencing read lengths.

163

#### 164 *Comparison of MRSD and TPM as a guide for appropriate surveillance*

165 We compared MRSD to the use of relative expression level (in transcripts per million,  
166 TPM) as a possible indicator of RNA-seq suitability for the detection of aberrant  
167 splicing events. We compared the expression levels, in TPM, of PanelApp genes  
168 against tissue-specific MRSD predictions, finding a negative correlation between the  
169 level of gene expression and its predicted MRSD across all three tissues ( $r^2 = 0.539$ -  
170  $0.669$ ; Figure 3a-c). This comparison confirms that more highly-expressed genes are  
171 associated with lower MRSD scores. However, we noted significant overlap between  
172 genes grouped into low-MRSD ( $< 100$  M reads) and high-MRSD ( $\geq 100$  M reads)  
173 brackets. For example, among genes considered low-MRSD, TPM values ranged  
174 from 1.25-1390, while genes with high-MRSD values had TPM values between 0-

175 4880 (Figure 3d). We quantified the overlap between these distributions,  
176 demonstrating that 98.6% of high-MRSD genes had higher TPM values than at least  
177 one low-MRSD gene. We calculated the tissue-specific median and the lowest TPM  
178 values within the low-MRSD bracket for the top 95% and 70% percentiles, and  
179 observed higher TPM values in 52.2%, 13.3% and 5.3% of high-MRSD genes,  
180 respectively (Figure 3d). The substantial overlap in the TPM values for low and high  
181 MRSD genes suggests that relative expression does not provide a wholly accurate  
182 representation of transcript coverage in RNA-seq data. Such inconsistencies may  
183 arise from bias in the regions of genes that are sequenced, for example, genes with  
184 high degrees of 3' bias in RNA-seq datasets (Supplementary Figure 4).

185

#### 186 *Traits of pathogenic splicing variation vary widely between genes and events*

187 We identified pathogenic aberrations to splicing in 20 of the 88 samples utilizing a  
188 previously described analysis pipeline (13) with a wide variety of mis-splicing effects  
189 (Supplementary Figure 5), and calculated respective median TPM and MRSD values  
190 (Supplementary Table 1). The method for aberrant splicing detection pooled  
191 evidence for splicing junctions in reference sets to generate tissue-specific models of  
192 “healthy” splicing. We incorporated RNA-seq datasets from relevant samples into the  
193 healthy splicing models (Supplementary Table 1) and collected metrics indicative of  
194 aberrant splicing events (Box 1). We observed high variability in all metrics  
195 associated with pathogenic aberrant splicing events (Table 1). All patients harbored  
196 at least one pathogenic splicing event supported by two reads and with normalized  
197 read counts (NRCs)  $\geq 0.19$ , and 80% of these events had a relative fold change in  
198 NRC  $> 19x$  relative to controls (Table 1). While a blanket set of parameters for all  
199 aberrant splicing events may be unsuitable, our data suggests that 90% of



200 pathogenic events could be retained if filtering for events that were singletons  
201 (evident only in a single sample), or were non-singletons with an NRC > 0.25.

202

203 **Box 1. Metrics collated during splice event analysis**

- 204 - Read count – Number of split reads supporting the existence of a given splice  
205 junction
- 206 - Normalized read count (NRC) – Ratio of reads supporting a given junction  
207 compared to the adjoining canonical junction with the highest read count
- 208 - NRC fold change – fold difference in NRC for a given event between an  
209 individual and the control individual with the next-highest NRC for that event
- 210 - Number of samples – the number of individuals, across both case and  
211 controls, in which an event is present
- 212 - Rank – position of a given event in a list of significant events, when ordered  
213 by decreasing read count (for singleton events) or fold change (for non-  
214 singleton events)

215

216 **Table 1. Range of metrics observed for pathogenic splicing events**

217

| Metric            | Tissue            |                      |                        |
|-------------------|-------------------|----------------------|------------------------|
|                   | Whole blood (n=3) | LCLs (n=7)           | Skeletal muscle (n=10) |
| Read count        | 2-40              | 4-38                 | 2-462                  |
| NRC               | 0.48-1.25         | 0.19-1.52            | 0.34-3.19              |
| NRC fold change   | Singletons        | 3.7-8.2 + singletons | 19.6-442 + singletons  |
| Number of samples | 1                 | 1-48                 | 1-110                  |
| Rank              | 2-5               | 10-232               | 1-342                  |

218



219 *Factors influencing the likelihood of pathogenic splicing variation identification &*  
220 *MRSD predictions*

221 To further define the most informative parameters for use in the MRSD model, we  
222 investigated the impact of a variety of metrics on the capability to identify pathogenic  
223 splicing events, including number of samples within the healthy reference set, the  
224 extent of read support for splicing junctions, and the relative expression of genes of  
225 interest. Overall, our analyses suggested that two supporting reads for an aberrant  
226 splicing event that is novel or has an NRC > 0.25 would reliably highlight pathogenic  
227 aberrations amongst transcriptome-wide splicing variation. These parameters are  
228 conservative and could be relaxed for the targeted investigation of variants of  
229 interest.

230

231 We first identified how the number of control samples used as a reference set for  
232 “healthy splicing” impacted our ability to identify aberrant splicing events. For all  
233 samples within our healthy splicing set, we iteratively selected groups of control  
234 samples at sizes of 30, 60 or 90. We observed that moving from 30 to 60 controls is  
235 associated with a mean reduction in event count of 19.3% (28.1% of non-singleton  
236 events, 17.1% of singleton events) across the three tissues, while increasing the  
237 control size to 90 results in a further reduction of 10.2% of events (16.5% of non-  
238 singleton events, 9.5% of singleton events; Figure 4); this effect was consistent  
239 across tissue types.

240

241 We next investigated how read count filters impacted the number of events observed  
242 for a given individual (Figure 4). Filtering out all splicing events supported by just a  
243 single read against a background of 90 control samples removes, on average, 91.2%

244 of events (60.4% of non-singleton events, 97.3% of singleton events). Increasing  
245 read support thresholds to 10 unique sequencing reads results in a total of 99.4% of  
246 events being excluded on average (96.2% of non-singleton events, 99.99% of  
247 singleton events), while retaining only those events supported by 100 reads or more  
248 removes an average of 99.97% of events (99.8% of non-singleton events, 100.0% of  
249 singleton events). To understand how the level of read support impacted the ability  
250 to identify specific events, we collated 31 aberrant splicing events across 22 muscle-  
251 derived RNA-seq samples, and downsampled reads in the genes containing these  
252 events. We observed that we could identify the same aberrant splicing events at  
253 reduced relative expression levels, and, while read support decreased (Figure 5a),  
254 the ranked position of the event within the rank-ordered output remained  
255 approximately the same in most cases (Figure 5b). However, the weakened read  
256 support increased the risk of eliminating the variant from consideration when read  
257 count filters were applied (Figure 5c). This analysis further emphasized that TPM  
258 values alone may not be a reliable measure of ability to survey all splicing junctions  
259 within a gene; we observed that splice junctions in different samples covered by the  
260 same number of sequencing reads belonged to genes with widely ranging TPM  
261 values (Supplementary Figure 6). For example, splice junctions covered by eight  
262 reads were associated with TPMs ranging between 0.17 and 52.

263

#### 264 *Implications for investigation of variants in known disease-causing genes*

265 We applied our MRSD model to all established disease genes included in the  
266 Genomics England PanelApp repository, encompassing 275 distinct gene panels  
267 and 3199 unique genes. 87 single-exon genes were excluded from analysis, leaving  
268 3112 unique disease genes. Based on our investigations of MRSD, we applied the

269 following parameters: read coverage = 8; proportion of junctions = 75%; confidence  
270 level = 95%. Using this approach (with expected PPV = 0.936-0.974, NPV = 0.776-  
271 0.880 across the three tissues) we observed that 58.0% (1806/3112) of PanelApp  
272 genes were predicted to be low-MRSD (< 100 M reads required) in at least one of  
273 whole blood, LCLs or skeletal muscle (Figure 6a). At the individual tissue level, 27.0%  
274 (841/3112) of PanelApp genes in whole blood, 49.0% (1524/3112) in LCLs and 44.0%  
275 (1369/3112) in skeletal muscle were predicted to be low-MRSD (Figure 6a). Of note,  
276 LCLs were observed to have the highest proportion of low-MRSD panel genes in  
277 190/275 disease-gene panels (69.3%, Figure 6c). Whole blood exhibited the highest  
278 proportion of genes with low MRSDs in just 24/275 disease-gene panels (8.8%).  
279

280 MRSD predictions revealed many use cases for specific tissues: in the familial  
281 rhabdomyosarcoma panel, for example, none of the 11 genes were predicted to be  
282 low-MRSD in blood, while 10/11 were predicted low-MRSD in LCLs (Figure 6c), of  
283 which nine were actually assigned an MRSD < 50 M reads. Results across all 275  
284 panels are shown in Supplementary Figures 8 & 9.

285

286 Overall, this analysis suggests both that whole blood may often represent the  
287 poorest choice of RNA source tissue in terms of disease gene coverage; in contrast,  
288 LCLs appear to show robustly high expression of many disease genes across  
289 diverse disease subtypes, and so may constitute a more reliable source of RNA for  
290 clinical transcriptomic investigations.

291

292

293 *Quantifying the resolving power of RNA-seq for variants of uncertain significance*

294 To analyze the possible impact of diagnostic RNA-seq integration on variant  
295 interpretation, we curated variants of uncertain significance (VUSs) from the ClinVar  
296 variant database (17) that were predicted by SpliceAI (18) to impact splicing (score  $\geq$   
297 0.5; see Materials and Methods). Of a total of 352,011 ClinVar variants, 185,119  
298 (52.6%) were identified as VUSs, and 7,507 (2.1%) were retained after filtering  
299 based on SpliceAI score. Cross-referencing the MRSDs of the genes harboring  
300 SpliceAI prioritized variants across tissues revealed that, depending on model  
301 stringency, between 22.1% and 59.4% of these variants may lie in genes that are  
302 low-MRSD in at least one of the three tissues (Figure 7a). Further, among the 30  
303 genes in which the greatest number of predicted splice-impacting VUSs were  
304 identified, 21 were predicted to be low-MRSD in at least one tissue (Figure 7b).  
305 Similar patterns were observed when using a more relaxed SpliceAI score filter of  
306 0.25 (Supplementary Figure 10). The guided integration of RNA-seq into diagnostic  
307 services alongside predictive bioinformatics tools is therefore likely to provide a  
308 significant improvement to interpretation of VUSs in a variety of disease contexts.

309

## 310 **Discussion**

311 The recent development of machine learning approaches has underpinned  
312 improvements to the prioritization of variants that impact splicing and cause rare  
313 disease (19). Despite these advances, corroboration of the effect of such variants  
314 remains a major obstacle to improving diagnostic yield for Mendelian disorders. This  
315 obstacle is amplified by the unexpected functional impact of some variants on  
316 splicing, which may change the way the variant is classified in accordance with  
317 current guidelines (6). The MRSD-based approach described here allows the  
318 informed selection of biosample(s) for bulk RNA-seq, based on the required number

319 of sequencing reads that need to be generated for appropriate surveillance of genes  
320 of interest. This approach enables the effective identification of patients, disease  
321 groups and genomic variants that are amenable for functional assessment of mis-  
322 splicing through RNA-seq, and may help to improve the efficiency and accuracy of  
323 genomic diagnostic approaches.

324

325 The primary purpose of MRSD is to predict the likelihood of observing pathogenic  
326 splicing defects in a given gene and tissue, and we quantify the utility of three distinct  
327 biosamples in this manner for known monogenic disease genes (Figure 6). Through  
328 this analysis, we are able to highlight biosamples that may be most informative for  
329 RNA-seq based analysis datasets for specific disease subsets. Although our model  
330 is conservative (Figure 2), we demonstrate through MRSD-guided re-inspection of  
331 VUSs in ClinVar that it may be possible to use RNA-seq to clarify the effect of up to  
332 2.4% of variants of uncertain significance (Figure 7a).

333

334 Other approaches to select genes amenable to functional analysis through RNA-seq  
335 include leveraging relative gene expression metrics (14, 20), or tools which assess  
336 the similarity of transcript isoforms between tissues, e.g. MAGIQ-CAT (7). We show  
337 that, whilst TPM values are well correlated with MRSD scores (Figure 3a-c), uneven  
338 sequencing coverage across the length of the transcript may, in some cases, falsely  
339 identify specific genes or splice junctions as being amenable to RNA-seq-based  
340 analysis (Supplementary Figure 5). 3' sequencing bias, which is a known artefact of  
341 poly-A enriched mRNA sequencing (21-23), may elevate the risk of inaccurately  
342 selecting genes that could be surveyed through RNA-seq when considering TPM  
343 alone. Additionally, the normalization against sequencing depth that occurs during

344 the calculation of TPM obscures information about raw read count, which is  
345 important when analyzing the utility of RNA-seq for clinical diagnostics. MRSD  
346 scoring, conversely, leverages variation in sample read depth to provide quantitative  
347 predictions about optimal sequencing depths.

348

349 On the other hand, the recently released tool MAGIQ-CAT (7) assesses the degree  
350 to which transcript isoforms in a sampled tissue accurately resemble those in the  
351 primary disease-affected tissue. However, MAGIQ-CAT primarily captures the  
352 degree of similarity between isoform structure and does not aim to provide a  
353 quantitative readout to guide the diagnostic route. Thus, a proxy tissue may be  
354 described as suitable for RNA-seq-based analysis despite having poor coverage of  
355 splice junctions. We envision that the use of both MAGIQ-CAT and MRSD could  
356 comprehensively capture information about the utility of RNA-seq, both in terms of  
357 similarity of isoform structure relative to the disease-affected tissue and in terms of  
358 the likelihood of observing disruptions to this structure.

359

360 There are several limitations of the current MRSD model, which could be  
361 incorporated into future work. Firstly, the MRSD model cannot directly be extended  
362 to predict the suitability of datasets to detect allele-specific expression biases and  
363 differential gene expression, which have been demonstrated to be evidence of  
364 pathogenic mechanisms in known disease-causing genes (10, 11, 14, 24). Although  
365 further investigations are required to quantify and prove this suitability, it is likely that  
366 genes with low MRSD scores (Figure 3d) are also amenable to investigations of  
367 differential gene expression and isoform imbalance.

368

369 Secondly, further extensions to the model could incorporate genomic background  
370 which influences gene expression profiles. For example, interferonopathies are a  
371 class of genomic immune disorders (25, 26) that are characterized by the aberrant  
372 upregulation of large numbers of transcripts belonging to so-called “interferon-  
373 stimulated genes” (25, 27). As a result of these wide-ranging impacts on their  
374 transcriptomes, MRSD predictions, which ostensibly represent the “normal”  
375 transcriptomic landscape, may not accurately reflect the degree of sequencing  
376 coverage for certain transcripts in patients with interferonopathies, or indeed other  
377 disease groups where disrupted expression of many transcripts is characteristic,  
378 such as disorders where chromatin structure (28, 29) or the function of the  
379 spliceosome (30-32) is disrupted. Moreover, the current MRSD model does not  
380 explicitly account for the presence of expression quantitative trait loci (eQTLs) or  
381 splicing quantitative trait loci (sQTLs) which are known to influence gene expression  
382 profiles (33-35). We have demonstrated that modulation in expression levels may  
383 disrupt our ability to reliably highlight pathogenic splicing events (Figure 5c). As a  
384 greater number of paired transcriptome and genomic datasets become available, we  
385 expect that MRSD scores can be generated in a dynamic manner to account for the  
386 presence of eQTLs, sQTLs or other modifiers of gene expression profiles.

387

388 Thirdly, our approach is built for a specific set of RNA-seq-based analyses; namely,  
389 the analysis of a selection of tissues by bulk short-read poly-A enrichment RNA-seq,  
390 followed by a specific bioinformatics analysis pipeline (13). This experimental RNA-  
391 seq approach currently remains widespread (3, 10, 13-15); however, our model may  
392 be readily applicable to RNA-seq generated using alternative methodologies, such  
393 as increased read length, with only minor variations in model performance



394 (Supplementary Figure 3). As other technologies, such as long-read (36-38), single-  
395 cell (39, 40) and spatially resolved RNA-seq (41-44), become more prevalent in a  
396 clinical setting, appropriate control datasets must be generated to develop  
397 corresponding MRSD models. Similarly, recent research has shown noticeable  
398 improvements to diagnostic yield for neuromuscular disorders by conducting RNA-  
399 seq on *in vitro* myofibrils generated by a fibroblast-to-myofibril transdifferentiation  
400 protocol (45). Such patient-derived cell line approaches represent a promising  
401 avenue to scrutinize transcripts not otherwise observable in proxy tissues (31, 46).  
402 As these protocols gain wider use, generation of control RNA-seq data from healthy  
403 individuals using these approaches will be vital both to allow the generation of MRSD  
404 scores and to accurately assess pathogenicity of any identified mis-splicing events.

405

## 406 **Conclusions**

407 In summary, the novel MRSD model presented here offers a gene-specific readout  
408 to predict the most suitable biosample for interrogation of splicing disruption at the  
409 transcript level. This may uncover previously unintuitive choices of biosample, as  
410 discussed above in the case of familial rhabdomyosarcoma (Figure 6c). The use of  
411 different biosamples is associated with different costs: while whole blood is routinely  
412 taken in the clinic, cell-based RNA-seq requires harvesting and culturing of patient  
413 cells, and muscle biopsy is an invasive procedure that is generally only undertaken if  
414 deemed necessary. Our tool may allow clinical staff to make informed decisions  
415 about the likely cost-benefit balance of RNA-seq analysis to ensure such costs are  
416 not incurred unnecessarily. We expect that the use of MRSD will allow effective and  
417 appropriate integration of RNA-seq into diagnostic genomic services, and ultimately  
418 improve variant interpretation and diagnostic yield.

419

## 420 **Methods**

### 421 *Minimum required sequencing depth (MRSD) score*

422 We generated a collated map of splice junction coverage for GTEx samples from  
423 three tissues (peripheral blood:  $n = 151$ ; LCLs:  $n = 91$ ; skeletal muscle:  $n = 184$ ; see  
424 *RNA-seq data acquisition*, below), using established methods (Cummings et al.,  
425 2017). These samples were designated as *reference sets*. Our model considers the  
426 level of sequencing coverage for splice junctions in each tissue-specific reference  
427 set and calculates the minimum required sequencing depth (MRSD), in millions of  
428 uniquely mapping 75 bp reads, that would be required for the desired proportion of  
429 splice junctions in a given gene to be covered by a desired number of sequencing  
430 reads. Our model is dynamic, and can be adjusted by the user to account for  
431 customized levels of desired sequencing coverage per splicing junction, the  
432 proportion of splicing junctions covered, and the confidence level with which MRSD  
433 will generate datasets with the specified level of coverage (suggested usage of 95 or  
434 99%).

435

436 MRSD is defined for a given gene in a given sample as:

437

$$MRSD = r / \left( \frac{R_{1-p}}{d} \cdot 10^6 \right)$$

438

439 Where  $r$  is the desired level of read coverage across desired proportion  $p$  of splice  
440 junctions,  $R$  is the set of read counts supporting all junctions in the transcript of  
441 interest, and  $d$  is the total number of sequencing reads in the RNA-seq sample (by  
442 default, the number of uniquely mapping sequencing reads). The term  $R_{1-p}$

443 corresponds to the number of reads covering the junction with the “ $1 - p$ ”-th-highest  
444 read count across all splice junctions in the transcript of interest.

445

446 MRSD scores have been generated for specified transcripts across all samples  
447 within the reference set in the three tissues of interest. The score at the  $X^{\text{th}}$  percentile  
448 position in the reference set list is returned as the MRSD, where  $X$  is termed the  
449 “confidence level” and is customizable by the user (default = 95%, Supplementary  
450 Methods 1).

451

#### 452 *Transcript selection*

453 MRSD can be calculated for any transcript sets of interest. Here, we utilized a  
454 hierarchy for transcript selection for all genes present in the GENCODE v19 human  
455 genome annotation (Supplementary Methods 2). We prioritized transcripts in the  
456 MANE v0.7 curated transcript list, providing that all splicing junctions were supported  
457 in the GENCODE v19 annotation. Genes without MANE transcripts were assigned  
458 composite transcripts, consisting of the union of all junctions found in transcripts in  
459 NCBI RefSeq transcripts. For genes that matched neither criteria, the union of all  
460 junctions present in all GENCODE v19-listed transcripts for that gene were used as  
461 the transcript model.

462

#### 463 *Control RNA-seq data acquisition*

464 FASTQs were downloaded from the Database of Genotypes and Phenotypes  
465 (dbGaP) under the project accessions phs000424.v8.p2 and phs000655.v3.p1.c1 for  
466 GTEx control individuals and neuromuscular disease patients, respectively. GTEx  
467 controls were selected for LCLs ( $n = 91$ ), skeletal muscle ( $n = 184$ ) and whole blood

468 ( $n = 151$ ) according to tissue-specific criteria (Supplementary Methods 3) to ensure  
469 use of only high-quality samples in generating control splicing datasets.

470

#### 471 *In-house RNA-seq generation*

472 RNA-seq datasets used to evaluate model performance were accessed from  
473 previously published datasets (13), under dbGaP project accession  
474 phs000655.v3.p1.c1, through international consortia (47), or for individuals in whom  
475 written informed consent was obtained and ethical approval for the study granted by  
476 Scotland A (refs: 06/MRE00/76 and 16/SS/0201), South Central-Hampshire A (ref:  
477 17/SC/0026), South Central-Oxford B (ref:11/SC/0269) or South Manchester (ref:  
478 11/H10003/3).

479

480 For in-house peripheral blood samples, RNA was extracted from PAXgene Blood  
481 RNA Kits and underwent poly-A enrichment library preparation using the TruSeq  
482 Stranded mRNA assay (Illumina) followed by 76 bp paired end sequencing using an  
483 Illumina HiSeq 4000 sequencing platform. For in-house LCL samples, RNA was  
484 extracted from pelleted LCLs thawed directly into TRIzol reagent (Invitrogen, 15596-  
485 026) using chloroform, and treated with TURBO DNase (Invitrogen, AM1907), both  
486 following the manufacturers' instructions. RNA was prepared using the NEBNext  
487 Ultra II Directional RNA Library Prep kit (NEB #7760) with the Poly-A mRNA  
488 magnetic isolation module (NEB #E7490), according to manufacturer's instructions,  
489 and 75bp paired end sequencing was performed using the Illumina NextSeq 550  
490 sequencing platform. Ribosomal RNA depleted datasets were generated using RNA  
491 extracted via the PAXgene Blood RNA system, and 150bp paired end sequencing  
492 performed via Novogene (Hong Kong) using the NEBNext Globin and rRNA

493 Depletion and NEBNext Ultra Directional RNA Library Prep Kits on a HiSeq 2000  
494 instrument (Illumina). RNA samples from 20 LCLs were obtained from the kConFab  
495 consortium. Poly(A)-selected RNA was generated using the TruSeq Stranded mRNA  
496 Library Prep Kit (Illumina), and 150bp paired end reads created using the NextSeq  
497 500 instrument (Illumina).

498

#### 499 *Splice event identification*

500 All FASTQs were aligned and processed as previously described (Cummings et al.,  
501 2017). Briefly, this analysis consisted of two-pass alignment using the STAR v2.4.2  
502 aligner, marking of suspected PCR duplicates, and processing of the resultant  
503 alignments to generate tissue-by-tissue lists of splice junctions present within the  
504 cohort. Metrics for each splicing event were collected (Box 1), and splicing junctions  
505 were filtered to retain only those events that were unique to single samples  
506 (singletons) or that were present in multiple samples (non-singletons) but with an  
507 increased usage in the sample of interest, that is, with a higher normalized read  
508 count (NRC), than any control. The resulting list was ranked according to NRC fold  
509 change, with singletons with high read counts considered the most significant events.  
510 The resulting junctions were considered “events of interest”.

511

#### 512 *Factors influencing the likelihood of aberrant splicing identification*

513 To calculate how the level of background splicing aberrations was altered by sample  
514 size, each individual in the three control splicing datasets was processed using the  
515 above pipeline (13) and compared against 2000 bootstraps of 30, 60 and 90 controls  
516 each from their respective control tissue dataset with replacement. Events were then  
517 filtered to retain only those events for which the NRC was higher in the given

518 individual than in any controls, and then counted for each bootstrap. Median counts  
519 for singleton and non-singleton events were then collated for each control group size.  
520 We selected 32 aberrant splicing events identified in neuromuscular patient RNA-seq  
521 data. From the genes in which we identified these variants, samtools was used to  
522 remove random subsets of reads in 10% intervals from each of these events to  
523 simulate variability in the number of reads generated for the gene of interest. The  
524 resulting datasets, exhibiting variable expression of a single gene, were then rerun  
525 through the splice analysis pipeline and the above metrics gathered for these  
526 simulated datasets.

527

#### 528 *Genomics England PanelApp data collection*

529 Tabulated versions of 284 gene panels were downloaded from the Genomics  
530 England PanelApp repository. Each panel was filtered to retain only genes assigned  
531 a “green” classification for that panel, representing the highest level of confidence of  
532 a real genotype-phenotype association.

533

#### 534 *Curation of ClinVar variants of uncertain significance*

535 A tabulated version of the comprehensive ClinVar variant listing (17) for January  
536 2021 was downloaded and filtered to retain only those variants that were annotated  
537 as either “Uncertain significance” or “Conflicting interpretations of pathogenicity”.  
538 SpliceAI scores (v1.2.1; (18)) were generated for these variants and those with a  
539 score of 0.5 or greater retained for downstream analysis.

540

#### 541 **Declarations**

542 *Ethics approval and consent to participate*

543 External datasets utilized in this study were accessed under dbGaP project  
544 accessions phs000655.v3.p1.c1 and phs000424.v8.p2. Informed written consent  
545 was obtained for all inhouse analyses, with ethical and study approval from South  
546 Central-Hampshire A (ref: 17/SC/0026), South Central-Oxford B (ref:11/SC/0269),  
547 South Manchester (ref:11/H10003/3) and Scotland A (refs: 06/MRE00/76 and  
548 16/SS/0201) Research Ethics Committees.

549

#### 550 *Consent for publication*

551 No identifiable patient information is reported in this study.

552

#### 553 *Availability of data and materials*

554 The control datasets used to generate the MRSD model are available through the  
555 dbGaP repository as part of the GTEx v8 release (accession phs000424.v8.p2).  
556 Publicly available muscle-derived RNA-seq datasets to test the model are available  
557 at dbGaP (accession phs000655.v3.p1.c1). Source code will be made available  
558 upon publication. All MRSD scores are available at <http://mcgm-mrds.github.io/>.

559

#### 560 *Competing interests*

561 The authors declare no competing interests.

562

#### 563 *Funding*

564 C.F.R. is funded by the Medical Research Council (MRC; 1926882) as part of a  
565 CASE studentship with QIAGEN. The Baralle lab is supported by an NIHR Research  
566 Professorship to D.B. (RP-2016-07-011). W.G.N. is supported by the NIHR  
567 Manchester Biomedical Research Centre (IS-BRC-1215-20007). We acknowledge



568 funding from the Wellcome Trust Transforming Genomic Medicine Initiative  
569 (200990/Z/16/Z) and the Medical Research Foundation. J.M.E is funded by a  
570 postdoctoral research fellowship from the Health Education England Genomics  
571 Education Programme (HEE GEP). The views expressed in this publication are  
572 those of the authors and not necessarily those of the HEE GEP.

573

#### 574 *Authors' contributions*

575 The study was designed and coordinated by C.F.R., G.C.M.B., R.T.O, S.H., T.A.B. &  
576 J.M.E. All authors contributed genetic or phenotypic data. C.F.R. and J.E. wrote the  
577 manuscript. A.T. designed and implemented the MRSD web portal. All authors  
578 contributed to the editing and revision of the manuscript.

579

#### 580 *Acknowledgements*

581 We wish to thank Heather Thorne, Eveline Niedermayr, all the kConFab research  
582 nurses and staff, the heads and staff of the Family Cancer Clinics, and the Clinical  
583 Follow Up Study (which has received funding from the NHMRC, the National Breast  
584 Cancer Foundation, Cancer Australia and the National Institute of Health (USA)) for  
585 their contributions to this resource, and the many families who contribute  
586 to kConFab. kConFab is supported by a grant from the National Breast Cancer  
587 Foundation, and previously by the National Health and Medical Research Council  
588 (NHMRC), the Queensland Cancer Fund, the Cancer Councils of New South Wales,  
589 Victoria, Tasmania and South Australia, and the Cancer Foundation of Western  
590 Australia. We also wish to thank members of the Wessex Investigational Sciences  
591 Hub (WISH) Laboratory, Southampton, UK, for their help in facilitating RNA-seq of  
592 kConFab LCL samples (particularly Christopher Mattocks, Daniel Ward and Jade

593 Forster), as well as the work of the University of Manchester Genomics Core  
594 Technology and Bioinformatics Facilities for their assistance in sample processing.

595

## 596 **References**

- 597 1. Anna A, Monika G. Splicing mutations in human genetic disorders: examples,  
598 detection, and confirmation. *J Appl Genet.* 2018;59(3):253-68.
- 599 2. Scotti MM, Swanson MS. RNA mis-splicing in disease. *Nat Rev Genet.*  
600 2016;17(1):19-32.
- 601 3. Wai HA, Lord J, Lyon M, Gunning A, Kelly H, Cibir P, et al. Blood RNA  
602 analysis can increase clinical diagnostic rate and resolve variants of uncertain  
603 significance. *Genet Med.* 2020;22(6):1005-14.
- 604 4. Sangermano R, Garanto A, Khan M, Runhart EH, Bauwens M, Bax NM, et al.  
605 Deep-intronic ABCA4 variants explain missing heritability in Stargardt disease and  
606 allow correction of splice defects by antisense oligonucleotides. *Genet Med.*  
607 2019;21(8):1751-60.
- 608 5. Khan M, Cornelis SS, Pozo-Valero MD, Whelan L, Runhart EH, Mishra K, et  
609 al. Resolving the dark matter of ABCA4 for 1054 Stargardt disease probands through  
610 integrated genomics and transcriptomics. *Genet Med.* 2020;22(7):1235-46.
- 611 6. Rowlands CF, Thomas H, Lord J, Wai H, Arno G, Beaman G, et al.  
612 Comparison of in silico strategies to prioritize rare genomic variants impacting RNA  
613 splicing for the diagnosis of genomic disorders. *Authorea [Internet].* 2020.
- 614 7. Aicher JK, Jewell P, Vaquero-Garcia J, Barash Y, Bhoj EJ. Mapping RNA  
615 splicing variations in clinically accessible and nonaccessible tissues to facilitate  
616 Mendelian disease diagnosis using RNA-seq. *Genet Med.* 2020;22(7):1181-90.
- 617 8. Marston S, Copeland O, Jacques A, Livesey K, Tsang V, McKenna WJ, et al.  
618 Evidence from human myectomy samples that MYBPC3 mutations cause  
619 hypertrophic cardiomyopathy through haploinsufficiency. *Circ Res.* 2009;105(3):219-  
620 22.
- 621 9. Mertes C, Scheller IF, Yépez VA, Çelik MH, Liang Y, Kremer LS, et al.  
622 Detection of aberrant splicing events in RNA-seq data using FRASER. *Nat Commun.*  
623 2021;12(1):529.
- 624 10. Kremer LS, Bader DM, Mertes C, Kopajtich R, Pichler G, Iuso A, et al.  
625 Genetic diagnosis of Mendelian disorders via RNA sequencing. *Nat Commun.*  
626 2017;8:15824.
- 627 11. Byron SA, Van Keuren-Jensen KR, Engelthaler DM, Carpten JD, Craig DW.  
628 Translating RNA sequencing into clinical diagnostics: opportunities and challenges.  
629 *Nat Rev Genet.* 2016;17(5):257-71.
- 630 12. Marco-Puche G, Lois S, Benítez J, Trivino JC. RNA-Seq Perspectives to  
631 Improve Clinical Diagnosis. *Front Genet.* 2019;10:1152.
- 632 13. Cummings BB, Marshall JL, Tukiainen T, Lek M, Donkervoort S, Foley AR, et  
633 al. Improving genetic diagnosis in Mendelian disease with transcriptome sequencing.  
634 *Sci Transl Med.* 2017;9(386).
- 635 14. Frésard L, Smail C, Ferraro NM, Teran NA, Li X, Smith KS, et al. Identification  
636 of rare-disease genes using blood transcriptome sequencing and large control  
637 cohorts. *Nat Med.* 2019;25(6):911-9.

- 638 15. Lee H, Huang AY, Wang LK, Yoon AJ, Renteria G, Eskin A, et al. Diagnostic  
639 utility of transcriptome sequencing for rare Mendelian diseases. *Genet Med.*  
640 2020;22(3):490-9.
- 641 16. Abdrabo LS, Watkins D, Wang SR, Lafond-Lapalme J, Riviere JB, Rosenblatt  
642 DS. Genome and RNA sequencing in patients with methylmalonic aciduria of  
643 unknown cause. *Genet Med.* 2020;22(2):432-6.
- 644 17. Landrum MJ, Lee JM, Benson M, Brown GR, Chao C, Chitipiralla S, et al.  
645 ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic*  
646 *Acids Res.* 2018;46(D1):D1062-D7.
- 647 18. Jaganathan K, Kyriazopoulou Panagiotopoulou S, McRae JF, Darbandi SF,  
648 Knowles D, Li YI, et al. Predicting Splicing from Primary Sequence with Deep  
649 Learning. *Cell.* 2019;176(3):535-48.e24.
- 650 19. Rowlands CF, Baralle D, Ellingford JM. Machine Learning Approaches for the  
651 Prioritization of Genomic Variants Impacting Pre-mRNA Splicing. *Cells.* 2019;8(12).
- 652 20. Murdock DR, Dai H, Burrage LC, Rosenfeld JA, Ketkar S, Müller MF, et al.  
653 Transcriptome-directed analysis for Mendelian disease diagnosis overcomes  
654 limitations of conventional genomic testing. *J Clin Invest.* 2021;131(1).
- 655 21. Finotello F, Lavezzo E, Bianco L, Barzon L, Mazzon P, Fontana P, et al.  
656 Reducing bias in RNA sequencing data: a novel approach to compute counts. *BMC*  
657 *Bioinformatics.* 2014;15 Suppl 1:S7.
- 658 22. Nagalakshmi U, Wang Z, Waern K, Shou C, Raha D, Gerstein M, et al. The  
659 transcriptional landscape of the yeast genome defined by RNA sequencing. *Science.*  
660 2008;320(5881):1344-9.
- 661 23. Wang Z, Gerstein M, Snyder M. RNA-Seq: a revolutionary tool for  
662 transcriptomics. *Nat Rev Genet.* 2009;10(1):57-63.
- 663 24. Kukurba KR, Zhang R, Li X, Smith KS, Knowles DA, How Tan M, et al. Allelic  
664 expression of deleterious protein-coding variants across human tissues. *PLoS Genet.*  
665 2014;10(5):e1004304.
- 666 25. Rodero MP, Crow YJ. Type I interferon-mediated monogenic  
667 autoinflammation: The type I interferonopathies, a conceptual overview. *J Exp Med.*  
668 2016;213(12):2527-38.
- 669 26. Volpi S, Picco P, Caorsi R, Candotti F, Gattorno M. Type I interferonopathies  
670 in pediatric rheumatology. *Pediatr Rheumatol Online J.* 2016;14(1):35.
- 671 27. Schneider WM, Chevillotte MD, Rice CM. Interferon-stimulated genes: a  
672 complex web of host defenses. *Annu Rev Immunol.* 2014;32:513-45.
- 673 28. Bélanger C, Bérubé-Simard FA, Leduc E, Bernas G, Campeau PM, Lalani SR,  
674 et al. Dysregulation of cotranscriptional alternative splicing underlies CHARGE  
675 syndrome. *Proc Natl Acad Sci U S A.* 2018;115(4):E620-E9.
- 676 29. Liu J, Zhang Z, Bando M, Itoh T, Deardorff MA, Clark D, et al. Transcriptional  
677 dysregulation in NIPBL and cohesin mutant human cells. *PLoS Biol.*  
678 2009;7(5):e1000119.
- 679 30. Wood KA, Rowlands CF, Qureshi WMS, Thomas HB, Buczek WA, Briggs TA,  
680 et al. Disease modeling of core pre-mRNA splicing factor haploinsufficiency. *Hum*  
681 *Mol Genet.* 2019;28(22):3704-23.
- 682 31. Wood KA, Rowlands CF, Thomas HB, Woods S, O'Flaherty J, Douzgou S, et  
683 al. Modelling the developmental spliceosomal craniofacial disorder Burn-McKeown  
684 syndrome using induced pluripotent stem cells. *PLoS One.* 2020;15(7):e0233582.
- 685 32. Buskin A, Zhu L, Chichagova V, Basu B, Mozaffari-Jovin S, Dolan D, et al.  
686 Disrupted alternative splicing for genes implicated in splicing and ciliogenesis causes  
687 PRPF31 retinitis pigmentosa. *Nat Commun.* 2018;9(1):4234.

- 688 33. Richards AL, Jones L, Moskvina V, Kirov G, Gejman PV, Levinson DF, et al.  
689 Schizophrenia susceptibility alleles are enriched for alleles that affect gene  
690 expression in adult human brain. *Mol Psychiatry*. 2012;17(2):193-201.
- 691 34. Takata A, Matsumoto N, Kato T. Genome-wide identification of splicing QTLs  
692 in the human brain and their enrichment among schizophrenia-associated loci. *Nat*  
693 *Commun*. 2017;8:14519.
- 694 35. Westra HJ, Franke L. From genome to function by studying eQTLs. *Biochim*  
695 *Biophys Acta*. 2014;1842(10):1896-902.
- 696 36. Mantere T, Kersten S, Hoischen A. Long-Read Sequencing Emerging in  
697 Medical Genetics. *Front Genet*. 2019;10:426.
- 698 37. Merker JD, Wenger AM, Sneddon T, Grove M, Zappala Z, Fresard L, et al.  
699 Long-read genome sequencing identifies causal structural variation in a Mendelian  
700 disease. *Genet Med*. 2018;20(1):159-63.
- 701 38. Pauper M, Kucuk E, Wenger AM, Chakraborty S, Baybayan P, Kwint M, et al.  
702 Long-read trio sequencing of individuals with unsolved intellectual disability. *Eur J*  
703 *Hum Genet*. 2020.
- 704 39. Del-Aguila JL, Li Z, Dube U, Mihindukulasuriya KA, Budde JP, Fernandez MV,  
705 et al. A single-nuclei RNA sequencing study of Mendelian and sporadic AD in the  
706 human brain. *Alzheimers Res Ther*. 2019;11(1):71.
- 707 40. Nomura S. Single-cell genomics to understand disease pathogenesis. *J Hum*  
708 *Genet*. 2021;66(1):75-84.
- 709 41. Crosetto N, Bienko M, van Oudenaarden A. Spatially resolved transcriptomics  
710 and beyond. *Nat Rev Genet*. 2015;16(1):57-66.
- 711 42. Larsson L, Frisén J, Lundeberg J. Spatially resolved transcriptomics adds a  
712 new dimension to genomics. *Nat Methods*. 2021;18(1):15-8.
- 713 43. Marx V. Method of the Year: spatially resolved transcriptomics. *Nat Methods*.  
714 2021;18(1):9-14.
- 715 44. Navarro JF, Croteau DL, Jurek A, Andrusivova Z, Yang B, Wang Y, et al.  
716 Spatial Transcriptomics Reveals Genes Associated with Dysregulated Mitochondrial  
717 Functions and Stress Signaling in Alzheimer Disease. *iScience*. 2020;23(10):101556.
- 718 45. Gonorazky HD, Naumenko S, Ramani AK, Nelakuditi V, Mashouri P, Wang P,  
719 et al. Expanding the Boundaries of RNA Sequencing as a Diagnostic Tool for Rare  
720 Mendelian Disease. *Am J Hum Genet*. 2019;104(5):1007.
- 721 46. Lin M, Pedrosa E, Shah A, Hrabovsky A, Maqbool S, Zheng D, et al. RNA-  
722 Seq of human neurons derived from iPS cells reveals candidate long non-coding  
723 RNAs involved in neurogenesis and neuropsychiatric disorders. *PLoS One*.  
724 2011;6(9):e23356.
- 725 47. Osborne RH, Hopper JL, Kirk JA, Chenevix-Trench G, Thorne HJ, Sambrook  
726 JF. kConFab: a research resource of Australasian breast cancer families. Kathleen  
727 Cuningham Foundation Consortium for Research into Familial Breast Cancer. *Med J*  
728 *Aust*. 2000;172(9):463-4.
- 729

730 **Figure 1.** *Minimum required sequencing depth (MRSD) predictions vary with*  
731 *changes in model parameters and across tissues. (a)* When all other parameters are  
732 constant (default parameters used here), increasing the desired level of read  
733 coverage of a gene results in a proportional increase in MRSD. The distribution of

734 MRSD scores for 3112 PanelApp genes in lymphoblastoid cell lines (LCLs) appears  
735 to be the lowest of the 3 tissues (median = 14.89 M at 10 reads), while whole blood  
736 exhibits the highest overall MRSD scores (median = 45.91 M at 10 reads),  
737 suggesting coverage of disease genes is generally poorer in blood. **(b, top)** In most  
738 cases, for a given level of splice junction (SJ) coverage, increasing the desired  
739 confidence level (the proportion of RNA-seq runs for which the MRSD prediction is  
740 expected to be sufficient) results in an increase in median MRSD score. **(b, bottom)**  
741 The number of genes for which no amount of sequencing is predicted to yield the  
742 specified level of coverage increases gradually as parameter stringency increases.  
743 At the highest level of stringency, the specified coverage was predicted unfeasible  
744 for between 63.1% (1964/3112, in LCLs) and 84.1% (2616/3112, in blood) of  
745 PanelApp genes.

746

747 **Figure 2.** *Performance metrics of the MRSD model.* The ability of MRSD to  
748 accurately predict levels of PanelApp disease gene coverage based on sequencing  
749 depth was tested on unseen RNA-seq datasets from blood ( $n = 12$ ), LCLs ( $n = 4$ )  
750 and muscle ( $n = 52$ ). **(a)** The mean positive predictive values (PPVs) and negative  
751 predictive values (NPVs) averaged across all parameter combinations for each RNA-  
752 seq dataset show that the median PPV is slightly lower, and the median NPV slightly  
753 higher, for whole blood than for LCLs and skeletal muscle. Breakdown of **(b)** PPVs  
754 and **(c)** NPVs for the MRSD model by parameters shows that specifying an  
755 increasing required read coverage results in a gradual decrease in PPV and  
756 increase in NPV across all tissues and parameter combinations. Dependent on  
757 parameter stringency, and limiting analysis to a maximum specification of 20-read  
758 coverage, PPV predictions range from 90.1-98.2%, while NPV ranges from 56.4-

759 94.7%. Overall, the model is fairly conservative and returns positive predictions only  
760 when they are deemed likely to be true.

761

762 **Figure 3.** *Comparison of MRSD and transcripts per million (TPM) predictions.* MRSD  
763 and TPM predictions for 3112 genes present in the Genomics PanelApp repository  
764 are inversely correlated in (a) whole blood ( $r^2 = 0.549$ ), (b) LCLs ( $r^2 = 0.539$ ) and (c)  
765 skeletal muscle ( $r^2 = 0.669$ ), as might be expected; however, the correlation is broad  
766 and there is high variation in the TPMs both of genes considered low- and high-  
767 MRSD (MRSD  $\leq$  or  $>$  100 M reads, respectively, dotted line). (d) Bracketing  
768 PanelApp genes by MRSD range shows that there is substantial overlap in the TPMs  
769 of genes across different MRSD predictions, to the extent that sufficient coverage of  
770 genes with TPMs up to 2796.5 is predicted unfeasible in some cases. This suggests  
771 relative expression level alone is not an adequate proxy for transcript coverage. The  
772 y-axis is limited to 100 TPM in (a-c) for ease of visualization. Log transformation in (d)  
773 excludes 491 entries with TPMs of 0. Default MRSD parameters (8-read coverage of  
774 75% of splice junctions, confidence level of 95%) used throughout.

775

776 **Figure 4.** *Expanding control datasets and enforcing read count thresholds improves*  
777 *filtering power when analyzing mis-splicing events.* Counting the significant events  
778 identified in each individual in a control splicing dataset when analysed against 2000  
779 bootstraps each of 30, 60 and 90 other individuals from within the control dataset for  
780 the same tissue reveals a small decrease in the number of total events identified as  
781 control dataset size increases, predominantly from non-singleton events. Enforcing a  
782 read coverage threshold has a more significant effect on event counts, particularly  
783 for singleton events, where filtering out events supported by a single read removes



784 up to 95% of singletons. LCLs appear to exhibit the greatest number of splicing  
785 events regardless of filter, although this may be due to differences in sequencing  
786 depth between tissues.

787 **Figure 5.** *Variability in expression level influences the capacity to identify mis-*  
788 *splicing events.* Genes harboring a selection of 31 splicing events that were  
789 identified during analysis of 52 muscle-based RNA-seq datasets (and which would  
790 be identified as events of interest using a filter of normalized read count (NRC) >  
791 0.19) were artificially downsampled to simulate variation in expression. **(a)** Reduction  
792 in expression leads to an intuitive and proportional reduction in the number of reads  
793 supporting each mis-splicing event. **(b)** The rank position – where the event appears  
794 in a list of all splicing events in its respective sample, ordered by decreasing NRC  
795 fold change relative to controls, and – is generally consistent as expression of the  
796 gene decreases; however, for a subset of events, reduction in expression is  
797 sufficient to cause stochastic changes in the NRC value, and so cause movement of  
798 the event down the prioritized list. **(c)** Variation in expression impacts our ability to  
799 identify events of interest when filters of read count supporting the events are  
800 enforced. When the 31 events experience a 50% reduction in expression, for  
801 instance, the application of a minimum 15-read filter leads to the exclusion of 41.9%  
802 (13/31) of events. For ease of visualization, the y-axis in **(a)** is limited to 50 reads,  
803 resulting in the truncation of some data series on the graph.

804

805 **Figure 6.** *Application of MRSD scores to disease genes listed in the Genomics*  
806 *England PanelApp repository.* **(a)** Comparison of PanelApp panel gene MRSD  
807 predictions between tissues shows blood to exhibit markedly poorer coverage of  
808 disease genes than do LCLs or skeletal muscle. **(b)** When comparing MRSD

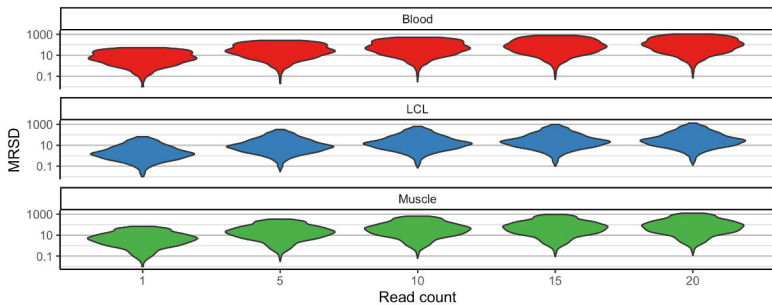
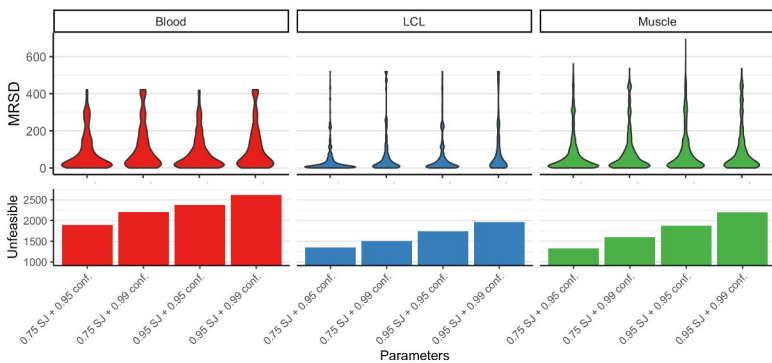


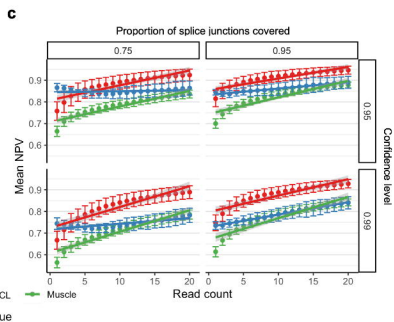
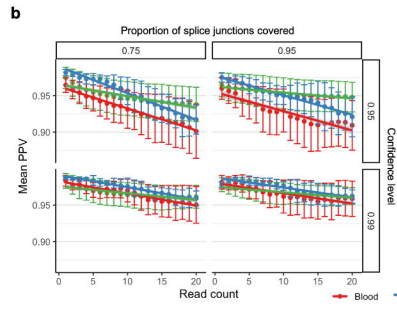
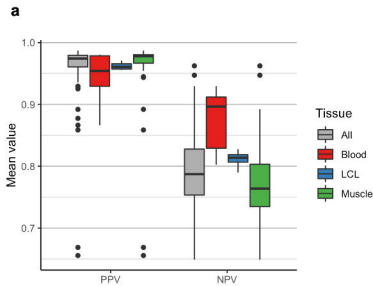
809 predictions for genes in blood and LCLs, 1522 genes are considered "high-MRSD"  
810 (i.e. have an MRSD > 100 M reads) in both tissues (grey). Genes which are  
811 exclusively low-MRSD (i.e. MRSD  $\leq$  100 M) in blood are far fewer in number (with 66  
812 genes, red box), while the remainder are almost evenly split between those that are  
813 low-MRSD in both (775 genes, purple box) and low-MRSD in LCLs only (749 genes,  
814 blue box). **(c)** Comparison of PanelApp panel gene MRSDs between tissues shows  
815 many panel genes have substantially greater coverage in LCLs than blood and, to a  
816 lesser extent, skeletal muscle over a variety of disease subtypes. Panels where  
817 skeletal muscle shows the best coverage of panel genes intuitively correspond to  
818 phenotypes such as neuromuscular disorders and distal myopathies. 40 exemplar  
819 panels shown here, to see results for all 275 panels, see Supp. Figs. 8 & 9. **(d)** Top  
820 10 panels with most significant difference between low- and high-MRSD gene counts  
821 between blood and LCLs (chi-squared test). **(e)** Venn diagrams showing number of  
822 low-MRSD genes predicted in blood and LCLs for (top) the paediatric disorder panel,  
823 the most significantly divergent between the two tissues, and (bottom) the bleeding  
824 and platelet disorders panel, which did not reach statistical significance in the  
825 aforementioned chi-squared analysis.

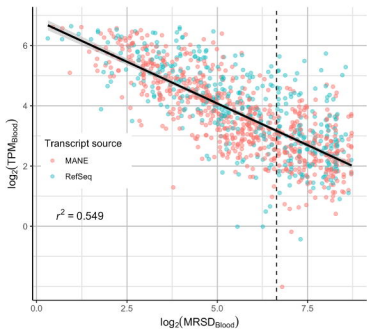
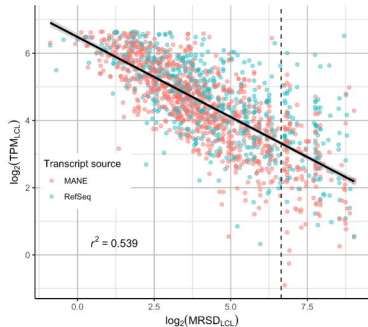
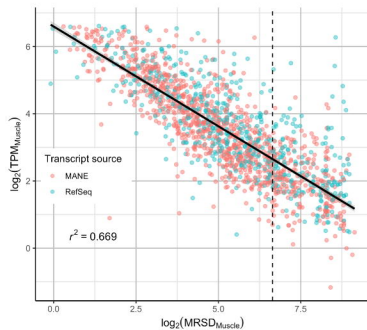
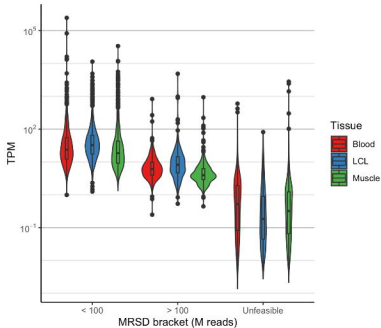
826

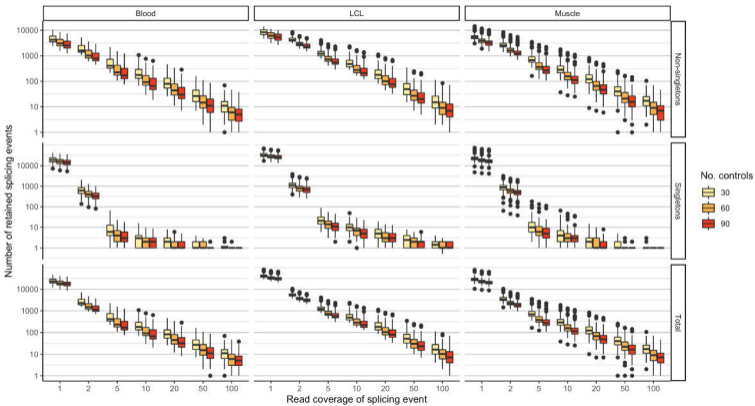
827 **Figure 7.** *The scope for resolution of variants of uncertain significance (VUSs) using*  
828 *RNA-seq-based analysis.* MRSD scores were derived for the genes harbouring  
829 VUSs present in ClinVar if the variants were predicted by the predictive tool SpliceAI  
830 to impact splicing (score  $\geq$  0.5; Jaganathan et al., 2019) **(a)** Depending on the  
831 stringency of the MRSD model parameters, between 22.1% (1663/7507) and 59.4%  
832 (4462/7507) of variants predicted to impact splicing are expected to be adequately  
833 covered by 100 M uniquely mapping reads or fewer in at least one of the 3 tissues

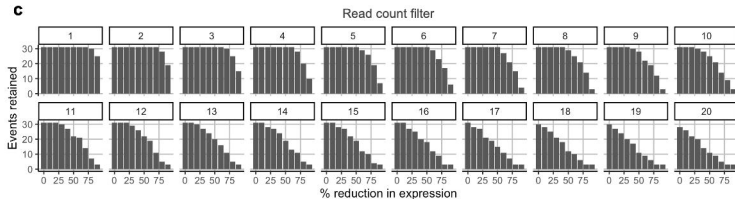
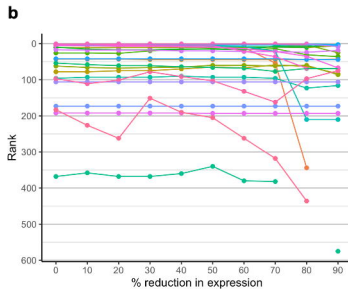
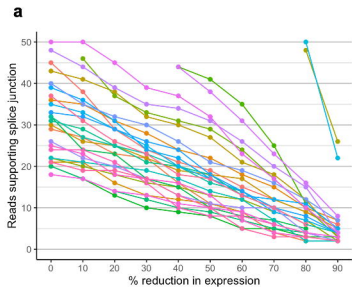
834 (whole blood, LCLs and skeletal muscle). Variants were most likely to be found to be  
835 in low-MRSD genes ( $MRSD \leq 100$  M) in LCLs, irrespective of model parameters. **(b)**  
836 Among the 30 genes with the greatest number of predicted splice-impacting VUSs,  
837 21 were predicted to be adequately covered (using default parameters) with 100 M  
838 uniquely mapping reads or fewer in at least one of the 3 tissues. An 8-read junction  
839 support parameter was used throughout.

**Figure 1****a****b**

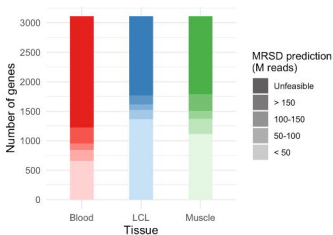
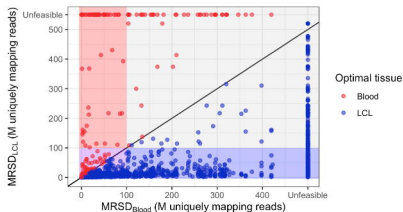
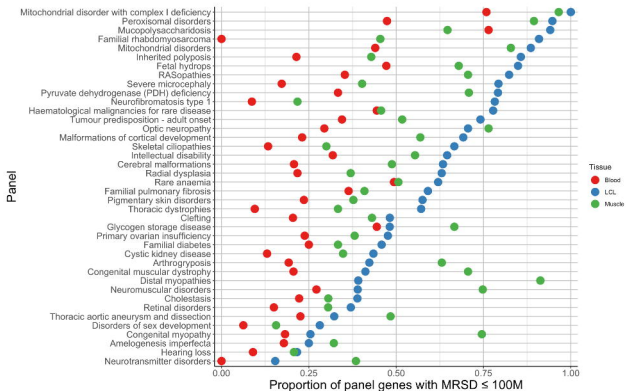


**Figure 3****a Whole blood****b Lymphoblastoid cell lines (LCLs)****c Skeletal muscle****d**

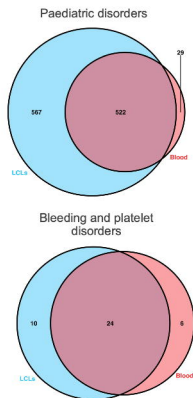
**Figure 4**

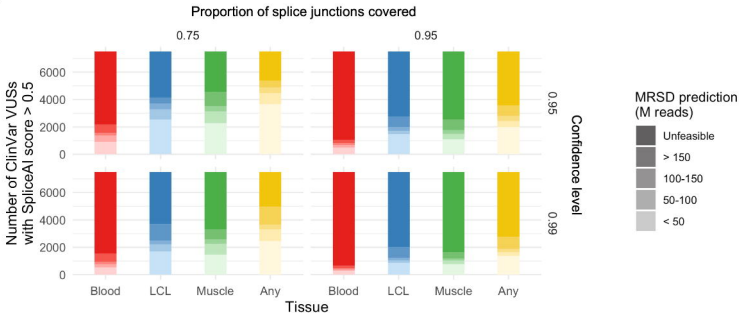
**Figure 5**



**Figure 6****a****b****c****d**

| Panel  | Panel size | $\chi^2$ p-value |
|--|------------|------------------|
| Paediatric disorders   | 3719       | 2.80E-68         |
| White matter disorders – childhood onset                     | 2025       | 2.68E-61         |
| Hypotonic infant   | 1972       | 1.40E-58         |
| Intellectual disability                                      | 1065       | 1.88E-51         |
| DDG2P  | 1167       | 3.14E-42         |
| Fetal anomalies  | 947        | 9.78E-35         |
| Inborn errors of metabolism                                  | 653        | 3.88E-30         |
| Undiagnosed metabolic disorders                              | 602        | 1.20E-26         |
| Possible mitochondrial disorder – nuclear genes              | 214        | 1.21E-19         |
| Mitochondrial disorders                                      | 175        | 3.14E-18         |
| Severe microcephaly  | 87         | 8.94E-16         |
| Skeletal dysplasia   | 351        | 2.31E-13         |
| Tumour predisposition – childhood onset                      | 77         | 2.53E-12         |
| Hereditary ataxia and cerebellar anomalies – childhood onset | 252        | 5.47E-12         |
| Genetic epilepsy syndromes                                   | 402        | 7.74E-12         |

**e**

**Figure 7****a****b**