



Long-read transcriptome sequencing of CLL and MDS patients uncovers molecular effects of *SF3B1* mutations

Alicja Pacholewska, Matthias Lienhard, Mirko Brueggemann, et al.

Genome Res. published online September 13, 2024
Access the most recent version at doi:[10.1101/gr.279327.124](https://doi.org/10.1101/gr.279327.124)

P<P	Published online September 13, 2024 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **Long-read transcriptome sequencing of CLL and MDS patients uncovers**
2 **molecular effects of *SF3B1* mutations**

3 Alicja Pacholewska^{1,2,*}, Matthias Lienhard^{3,*}, Mirko Brüggemann^{4,*}, Heike Hänel⁵, Lorina Bilalli¹,
4 Anja Königs^{1,2}, Felix Heß^{1,2}, Kerstin Becker^{6,7}, Karl Köhrer⁶, Jesko Kaiser⁸, Holger Gohlke^{8,9},
5 Norbert Gattermann¹⁰, Michael Hallek^{2,11}, Carmen D. Herling^{11,12}, Julian König⁵, Christina
6 Grimm^{1,2,\$}, Ralf Herwig^{3,\$,+}, Kathi Zarnack^{4,\$,+}, Michal R. Schweiger^{1,2,\$,+}

7 ¹ Institute for Translational Epigenetics, Faculty of Medicine, University of Cologne, 50931 Cologne,
8 Germany.

9 ² Center for Molecular Medicine Cologne (CMMC), Faculty of Medicine, University of Cologne, 50931
10 Cologne, Germany.

11 ³ Department of Computational Biology, Max Planck Institute (MPI) for Molecular Genetics, 14195 Berlin,
12 Germany.

13 ⁴ Buchmann Institute for Molecular Life Sciences and Institute of Molecular Biosciences, Goethe University
14 Frankfurt, 60438 Frankfurt, Germany.

15 ⁵ Institute of Molecular Biology, 55128 Mainz, Germany.

16 ⁶ Genomics & Transcriptomics Laboratory, Biological and Medical Research Center, Heinrich-Heine-
17 University, and West German Genome Center, 40225 Düsseldorf, Germany.

18 ⁷ Cologne Center for Genomics (CCG), University of Cologne, 50931 Cologne, Germany.

19 ⁸ Institute for Pharmaceutical and Medicinal Chemistry, Heinrich Heine University Düsseldorf, 40225
20 Düsseldorf, Germany

21 ⁹ Institute of Bio- and Geosciences (IBG-4: Bioinformatics), Forschungszentrum Jülich, 52428 Jülich,
22 Germany

23 ¹⁰ Department of Haematology, Oncology and Clinical Immunology, University Hospital Düsseldorf, 40225
24 Düsseldorf, Germany.

25 ¹¹ Department I of Internal Medicine, Center for Integrated Oncology Aachen-Bonn-Cologne-Düsseldorf,
26 University Hospital Cologne, 50937 Cologne, Germany.

27 ¹² Medical Clinic and Polyclinic 1, Hematology, Cellular Therapy, Hemostaseology, and Infectious Diseases,
28 University of Leipzig Medical Center, 04103 Leipzig, Germany

29 ^{*, \$} Equal contribution

30 ⁺ Corresponding authors

31

32 Running Title: Long-read transcriptomics in CLL and MDS

33 **Abstract**

34 Mutations in splicing factor 3B subunit 1 (*SF3B1*) frequently occur in patients with chronic
35 lymphocytic leukemia (CLL) and myelodysplastic syndromes (MDS). These mutations have different
36 effects on the disease prognosis with beneficial effect in MDS and worse prognosis in CLL patients.
37 A full-length transcriptome approach can expand our knowledge on *SF3B1* mutation effects on RNA
38 splicing and its contribution to patient survival and treatment options. We applied long-read
39 transcriptome sequencing (LRTS) to 44 MDS and CLL patients, as well as two pairs of isogenic cell
40 lines with and without *SF3B1* mutations, and found > 60% of novel isoforms. Splicing alterations
41 were largely shared between cancer types and specifically affected the usage of introns and 3' splice
42 sites. Our data highlighted a constrained window at canonical 3' splice sites in which dynamic splice
43 site switches occurred in *SF3B1*-mutated patients. Using transcriptome-wide RNA binding maps and
44 molecular dynamics simulations, we showed multimodal SF3B1 binding at 3' splice sites and
45 predicted reduced RNA binding at the second binding pocket of SF3B1^{K700E}. Our work presents the
46 hitherto most complete LRTS study of the *SF3B1* mutation in CLL and MDS and provides a resource
47 to study aberrant splicing in cancer. Moreover, we showed that different disease prognosis results
48 most likely from the different cell types expanded during carcinogenesis rather than different
49 mechanisms of action of the mutated SF3B1. These results have important implications for
50 understanding the role of *SF3B1* mutations in hematological malignancies and other related diseases.

51

52 **Keywords:**

53 Long-read transcriptome sequencing; LRTS, chronic lymphocytic leukemia, CLL, myelodysplastic
54 syndrome, MDS, Iso-Seq, differential splicing, iCLIP

55

56 **Introduction**

57 Splicing is a fundamental step in eukaryotic gene expression in which non-coding introns are removed
58 from pre-messenger RNA (pre-mRNA) transcripts and exons are joined to form mature mRNAs. This
59 intricate process is often disrupted in cancer, either by mutations in spliceosomal genes or by other
60 mechanisms that affect normal splicing function (Yang et al. 2021; Seiler et al. 2018; Shiozawa et al.
61 2018; Quesada et al. 2012; Bradley and Anczuków 2023). In turn, aberrant splicing can lead to
62 changes in the composition of expressed isoforms and the formation of new isoforms that alter the
63 encoded proteins and can have far-reaching consequences for cellular function. One striking example
64 of splicing alterations in cancer are mutations in the gene encoding the splicing factor 3B subunit 1
65 (*SF3B1*) that have divergent ramifications for treatment efficiency and prognosis (Rossi et al. 2011;
66 Papaemmanuil et al. 2011). Somatic *SF3B1* mutations are frequently found in myelodysplastic
67 syndrome (MDS, 20%), chronic lymphocytic leukemia (CLL, 15%), acute myeloid leukemia (3%),
68 uveal melanoma (20%), cutaneous melanoma (4%), prostate cancer (1%) and in 2% of all breast,
69 pancreatic and lung cancers (Bland et al. 2023).

70 In CLL patients, *SF3B1* mutations are typically subclonal and have been linked to disease progression
71 and shorter survival (Wan and Wu 2013; Landau et al. 2015). On the other hand, *SF3B1* mutations in
72 MDS patients have been associated with specific disease phenotypes that show erythroid dysplasia
73 with ring sideroblasts and ineffective erythropoiesis (Malcovati et al. 2015). Unlike in CLL patients, a
74 positive effect of the *SF3B1* mutation on survival has been observed in almost all groups of MDS
75 patients, except those with excess blasts, for whom no significant effect has been observed (Malcovati
76 et al. 2020). However, so far there is no explanation for the divergent ramifications of *SF3B1*
77 mutations in CLL and MDS pathology.

78 Pre-mRNA splicing is a highly dynamic process, and the spliceosome undergoes several structural
79 changes from E to A, B, and C complex during splicing. SF3B1 is part of the spliceosome and plays a
80 critical role in 3' splice site usage. As a subunit of the U2 small nuclear ribonucleoprotein complex
81 (snRNP), SF3B1 is UV-crosslinked with the pre-mRNA on both sides of the branch point (BP)

82 adenosine in the A-complex, at nucleotide positions -6 and +5 (Gozani et al. 1996, 1998). The most
83 common mutations in *SF3B1* accumulate in the HEAT (Huntington, Elongation Factor 3, PR65/A,
84 TOR) domain at its C-terminus (Supplemental Fig. S1). The HEAT domain consists of 20 non-
85 identical HEAT repeats that form the RNA binding interface (Cretu et al. 2016). These mutations are
86 predicted to impact the N-terminal domain involved in complex formation with other splicing factors
87 (Canbezdi et al. 2021). Mutations in SURP and G-patch domain containing 1 (*SUGPI*) mimic the
88 splice alterations of mutant SF3B1 (Liu et al. 2020; Alsafadi et al. 2021) and mutations in *DHX15*
89 partially recapitulate the splicing alterations of mutant SF3B1 (Zhang et al. 2022). Both proteins have
90 been shown to bind less to mutated SF3B1 (Zhang et al. 2019).

91 Previous studies based on short-read RNA sequencing (RNA-seq) have reported alternative 3' splice
92 site usage (3'AS) and intron retention (IR) as the most prominent splicing alterations in CLL and
93 MDS patients with mutated *SF3B1* (Shiozawa et al. 2018; Wang et al. 2016; Tang et al. 2020;
94 DeBoever et al. 2015; Kesarwani et al. 2017). The alternative 3' splice sites (referred to as AG') that
95 were preferably used upon *SF3B1* mutation, are enriched at approximately 20 nucleotides (nt)
96 upstream of the canonical splice sites (AG) (Obeng et al. 2016; Wang et al. 2016; Tang et al. 2020;
97 DeBoever et al. 2015). This strong positional constraint suggested that the mutations impacted SF3B1
98 binding and BP recognition upstream of the 3' splice sites. Additionally, it was proposed that the
99 mutation promotes the usage of otherwise inaccessible AG' within the RNA secondary structure
100 (Kesarwani et al. 2017). Despite these hypotheses, the exact mechanism of the effect of mutations in
101 *SF3B1* is still not resolved.

102 Here, we aimed to comprehensively characterize the effects of *SF3B1* mutations in cancer using long-
103 read transcriptome sequencing (LRTS) and combined complementary data derived from MDS and
104 CLL patients with isogenic cell lines.

105 **Results**

106 **Long-read RNA sequencing expands patient transcriptome landscapes in** 107 **divergent biological contexts**

108 To investigate the effect of *SF3B1* mutations on splicing, we characterized the transcriptomes of three
109 datasets: CLL patients, MDS patients with ring sideroblasts, and isogenic cell lines with or without
110 somatic *SF3B1* mutations (Supplemental Fig. S2). In brief, we collected CLL cells or whole-blood
111 samples from 19 CLL and 25 MDS patients, including eight CLL patients and 14 MDS patients with
112 mutations in the SF3B1 HEAT repeat domain (Fig. 1A). These were complemented by two isogenic
113 leukemia cell line pairs (K562 and Nalm6), both with *SF3B1*^{wt/wt} and *SF3B1*^{mut/wt}. K562 cells
114 originated from a patient with chronic myeloid leukemia (CML) and the *SF3B1*^{mut/wt} cells carried
115 K700E the mutation, whereas the Nalm6 cells originated from a patient with B cell acute
116 lymphoblastic leukemia (B-ALL) and the *SF3B1*^{mut/wt} cells carried the H662Q mutation. As controls,
117 we further included B cells from six healthy donors (Supplemental Table S1). The RNA expression
118 level of mutated *SF3B1* ranged from 14% to 52% (43% on average) in patients, 43% in K562, and
119 29% in the Nalm6 cell line (Supplemental Table S1). To detect complete transcript isoforms, we
120 performed LRTS using Iso-Seq[®] (Pacific Bioscience). We reached a mean of 582,135 full-length non-
121 chimeric reads per sample, cumulating in a total of 33,763,806 reads with an average length of 2,721
122 bp (Supplemental Table S1). Only 9% of the reads were potentially affected by technology-specific
123 technical artefacts (Cocquet et al. 2006) (Supplemental Fig. S3–S4).

124 In total, we identified 89,659 substantially expressed transcripts that contributed to at least 1% to a
125 gene's total expression and were covered by at least five full-length reads (Supplemental Fig. S5).
126 Almost one third of these reads (28,261; 31.5%) were classified as full splice matches (FSM) to
127 annotated isoforms. Moreover, 58,168 (64.9%) represented novel isoforms that only partially
128 overlapped with gene annotations, and 3,230 (3.6%) reads originated from non-annotated, novel genes
129 (Supplemental Fig. S6). Even for transcripts expressed at ≥ 1 transcript per million (TPM), the novel
130 isoforms consisted of more than half of all transcripts detected (Fig. 1B). A large fraction of isoforms

131 was shared by the different patient cohorts and isogenic cell lines, with a larger overlap of expressed
132 isoforms between $SF3B1^{mut/wt}$ and $SF3B1^{wt/wt}$ of the same dataset, than between datasets (Fig. 1C,
133 Supplemental Fig. S7).

134 The *SF3B1* gene has multiple isoforms annotated and was indeed expressed in several isoforms in
135 both, samples with or without *SF3B1* mutations (Fig. 1D, Supplemental Fig. S8). Although the most
136 frequently expressed *SF3B1*-FSM isoform (around 70% of *SF3B1* transcripts) fully corresponded to
137 the annotated isoform, two shorter novel isoforms showed a disease-specific expression almost
138 exclusively in either CLL or MDS patients. These contributed approximately 10% each to the gene's
139 overall expression, irrespective of the *SF3B1* mutational status (Fig. 1E). The CLL-specific isoform
140 (*SF3B1*-CLL) showed retention of the fourth intron which introduced a premature termination codon
141 and likely targeted the isoform for nonsense-mediated mRNA decay (NMD). In the MDS-specific
142 isoform (*SF3B1*-MDS), the penultimate exon was skipped and induced a frameshift that probably
143 resulted in an NMD-resistant isoform that encoded for a C-terminally truncated protein devoid of
144 HEAT repeats 18–20 and the anchor domain. In addition to the divergent splicing pattern, *SF3B1* also
145 showed three times higher expression in CLL compared to MDS patients, whereas its levels were
146 reduced in MDS patients when compared to B cells from healthy donors (Supplemental Fig. S8).

147 Overall, our results demonstrated a high transcriptome information content in the patient cohorts,
148 which was dominated by a large number of novel transcripts.

149 **Patients and cell lines with *SF3B1* mutations show similar splicing defects**

150 In order to investigate the transcriptome diversity at the splice-site level, we used the recently
151 developed IsoTools software (Lienhard et al. 2023) to identify alternative splicing events (ASEs) in
152 the transcripts expressed. IsoTools distinguishes exon skipping (ES), intron retention (IR), mutually
153 exclusive exons (ME), and 5' and 3' alternative splice sites (5'AS and 3'AS) events, as well as
154 alternative first and last junctions (AFJ and ALJ). Using a cut-off of ≥ 100 reads, ASEs were
155 quantified as the proportion of reads supporting the ASE in relation to the sum of reads for all
156 transcript isoforms, referred to as percent spliced-in index (PSI). This threshold is motivated from

157 extensive testing to optimize the detection of true splicing events while minimizing false positives and
158 also ensures a sufficient number of individual samples (≥ 5) supporting the vast majority of ASEs
159 (Supplementary Fig. S9). Across all samples, we discovered 80,995 ASEs in 9,746 genes, for which
160 the less expressed ASE made up for at least 10% of the reads. For 75% of these events, at least one of
161 the alternatives was not annotated (novel event) (Fig. 2A).

162 Next, we used IsoTools (Lienhard et al. 2023) to detect significant differences in splicing associated
163 with *SF3B1* mutations. Because *SF3B1* mutations have been reported to convey either beneficial
164 (MDS) or disadvantageous (CLL) effects on patient survival (Rossi et al. 2011; Papaemmanuil et al.
165 2011), we first tested for differential splicing in *SF3B1*^{mut/wt} vs. *SF3B1*^{wt/wt} samples, separately in each
166 dataset. We detected 82, 288, and 219 ASEs in the isogenic cell lines, CLL, and MDS patients,
167 respectively (adjusted p-value [q-value] with false discovery rate, FDR (Benjamini and Hochberg
168 1995) $< 10\%$, Supplemental Table S2). Although we observed only a moderate overlap of the
169 identified events between the datasets (Fig. 2B, Supplemental Fig. S10), to our surprise, the
170 correlation of the PSI changes for the ASEs identified was high (Fig. 2C), indicating a common
171 mutational effect. We found that the genes altered by the disease-specific ASE were generally
172 expressed significantly higher in the corresponding group of patients (Fig. 2D, Supplemental Fig.
173 S11) and out of the union of 531 genes with an ASE, 149 were relatively higher in MDS and 155 were
174 higher in the CLL samples by at least 2-fold (FDR $< 1\%$) based on the Iso-Seq data. Overall, about
175 two thirds of ASEs called in the MDS or CLL dataset separately were significantly differentially
176 expressed (FDR $< 5\%$) when comparing MDS to CLL (Supplemental Fig. S12, Supplemental Table
177 S3).

178 Our findings suggested that while *SF3B1* mutations introduced shared splicing effects in both CLL
179 and MDS patients, the divergence in the disease outcome could be attributed to the differential
180 transcriptomic profiles. Specifically, the mutation seemed to exert its most potent effects on genes that
181 were already dominantly expressed in each disease.

182

183 ***SF3B1* mutations affect alternative 3' splice site usage and intron retention**

184 Since the individual analyses for patients and isogenic cell lines indicated a common effect of *SF3B1*
185 mutations, we combined all *SF3B1*^{mut/wt} and *SF3B1*^{wt/wt} samples to increase the statistical power for an
186 overall estimation of the *SF3B1* mutational impact. In total, we identified 775 differential splicing
187 events in 530 different genes (Fig. 3A, Supplemental Table S4). As reported previously (Shiozawa et
188 al. 2018; Wang et al. 2016; Tang et al. 2020; DeBoever et al. 2015; Kesarwani et al. 2017; Darman et
189 al. 2015), the splicing changes upon *SF3B1* mutation were strongly enriched for 3'AS (326, 42%) and
190 IR (213, 27%) events which together accounted for more than two thirds of the significant changes.
191 The majority of IR events showed decreased IR (89%), whereas the majority of 3'AS events (78%)
192 showed higher PSI values, corresponding to longer exons in *SF3B1*^{mut/wt} (Fig. 3B–C, Supplemental
193 Fig. S13, Supplemental Table S4). Consistent with the common mutational effect, the regulated 3'AS
194 events showed a uniform response across cohorts and allowed to cluster *SF3B1*^{mut/wt} and *SF3B1*^{wt/wt}
195 samples (Fig. 3B).

196 Upon closer inspection, we found that two CLL-*SF3B1*^{mut/wt} patient samples with mutations in the
197 HEAT domain clustered with the *SF3B1*^{wt/wt} samples when a subset of ASEs was used (208 ASEs
198 with q-value < 0.01). One of these patients carried a rare *SF3B1* mutation, Q699E, that had so far
199 been reported only once, in a single patient with bladder urothelial carcinoma included in the TCGA
200 Pan-Cancer Atlas (according to cBioPortal accessed on 2023.10.16). Moreover, overexpression of the
201 SF3B1-Q699H construct in HEK293FT cells did not lead to any aberrant splicing, suggesting that this
202 mutation is weakly pathogenic (Darman et al. 2015). The second patient carried two mutations:
203 D894G in the HEAT repeat 11 (allele frequency (AF) = 51 %), and I704N in the HEAT repeat 6 (AF
204 = 18%, Supplemental Table S1). Although the mutations in these two patients might have a weaker
205 effect on the global splicing, we decided to keep the samples in our analysis to allow for more
206 biological variance in the statistical analysis and detect stronger signals. For all other mutations, the
207 position within the SF3B1 HEAT domain had no discernible influence on the clustering, indicating
208 that the different mutations impair SF3B1 similarly (Fig. 3B). In fact, using general splicing
209 information (PSI values), irrespective of regulation, we found that 3'AS events, but no other type of

210 ASEs, clearly differentiated the samples based on the *SF3B1* mutational status in an unsupervised
211 principal component analysis performed on each dataset (Fig. 3C, Supplemental Fig. S14–S15),
212 underlining the predominant effect of *SF3B1* mutation on 3'AS events.

213 In addition to many new differential ASEs, previously published discoveries could be confirmed,
214 including e.g., five from 35 3'AS events (in *SEPTIN2*, *ERGIC3*, *RHNO1*, *FDPS*, and *SNRPN*)
215 identified in CLL based on Nanopore sequencing (Tang et al. 2020), a 3'AS in *SEPTIN6* reported in
216 MDS-*SF3B1*^{mut/wt} patients (Dolatshad et al. 2016) (Supplemental Fig. S16), as well as thirteen 3'AS
217 events (*BCL2L1*, *COASY*, *DPH5*, *DYNLL1*, *EI24*, *ERGIC3*, *MED6*, *METTL5*, *SERBP1*, *SKIV2L*,
218 *TMEM14C*, *ZBED5*, and *ZDHHC16* that were consistently found in CLL, MDS, and uveal melanoma
219 patients (Inoue et al. 2019; Pellagatti et al. 2018). Moreover, we confirmed 13 from 32 (Zhou et al.
220 2020) and 8 from 11 genes (Liu et al. 2020) reported as aberrantly spliced in either MDS or CLL
221 patients.

222 The LRTS data opened the possibility to assess the splicing alterations in the context of complete
223 transcript isoforms. Generally, we found that the effect of the *SF3B1* mutation on splicing did not
224 influence the choice of transcript start or end sites, nor the probability of other splicing events of the
225 same gene. This means that the effect was local, and the resulting alternative transcript corresponded
226 to the canonical transcript, except for the single alternative event. This is exemplified by the
227 *SF3B1*^{mut/wt}-induced inclusion of the poison cassette exon (PCE) in the *BRD9* gene that introduces a
228 premature termination codon (PTC) (Inoue et al. 2019) (Supplemental Fig. S17–S18). Of note, the
229 full-length reads allowed us not only to confirm differential splicing of this PCE, but also to locate it
230 to a specific isoform that has not been annotated yet (NIC class). This long-read-derived novel
231 isoform otherwise resembled the canonical *BRD9* isoform, whereas the annotated PTC isoforms were
232 presumed to also harbor an alternative first exon and additional splicing alterations. Such incomplete
233 and incorrect isoform annotations are likely to cause problems in quantifying transcriptome changes,
234 especially when using short-read sequencing data.

235 Notably, we found multiple splicing factors among the genes affected by *SF3B1* mutations,
236 Overrepresentation analysis revealed 18 from the spliceosome pathway (KEGG: 03040,

237 q-value 2.8×10^{-3}) (Supplemental Table S4). This indicated a broad impact of *SF3B1* mutations on
238 the general splicing machinery, which could potentially lead to secondary effects on splicing. Taking
239 a closer look at the 208 highly significant 3'AS events (q-value, < 0.01), we identified four clusters
240 based on PSI values detected in *SF3B1*^{mut/wt} and *SF3B1*^{wt/wt} (Fig. 3B, Supplemental Fig. S19). The
241 cluster II with low PSI values in *SF3B1*^{mut/wt} and no expression in *SF3B1*^{wt/wt} was enriched in
242 spliceosome (q-value = 0.0106) and cell cycle genes (GO: 0007049 q-value = 2.8×10^{-4} ,
243 Supplemental Fig. S20).

244 To independently validate the detected splicing changes, we performed short-read RNA-seq on the
245 isogenic cell line pairs as well as on 27 CLL patient samples, which included the same 19 CLL
246 samples used for Iso-Seq, and collected publicly available RNA-seq data from 398 MDS patients
247 (Supplemental Table S1, Supplemental Fig. S21). Although we detected more events using rMATS
248 (Shen et al. 2014) (Supplemental Fig. S22 bottom, Supplemental Table S4) we observed a high and
249 significant (p-values < 0.001) correlation of PSI values in the ASEs detected with IsoTools (Lienhard
250 et al. 2023) and rMATS (Shen et al. 2014) on the same cell lines (Pearson correlation coefficient $R =$
251 0.840 , p-value = 2.47×10^{-17}). The same held true for the samples from CLL and MDS patients ($R =$
252 0.794 with p-value = 2.95×10^{-36} and $R = 0.800$ with p-value = 3.60×10^{-24} , respectively,
253 Supplemental Fig. S23).

254 As an orthogonal approach, we employed semi-quantitative reverse transcription PCR (RT-PCR) to
255 test 15 differential 3'AS events in the isogenic cell line pairs (Supplemental Fig. S24, Supplemental
256 Table S5). From these 15 tested, 12 (80%) clearly showed an increase in alternatively spliced isoform
257 expression in the *SF3B1*^{mut/wt} conditions. The strongest effects were observed for 3'AS events in
258 *MAP3K7*, *SEPTIN6*, and *SETD4*, which showed an almost complete switch to the alternative splicing
259 variant (Fig. 3D–E). We additionally performed splicing reporter assays using minigenes for six 3'AS
260 events in HEK293T cells (Fig. 3F). Indeed, we observed differential 3'AS usage upon ectopic
261 *SF3B1*^{K700E} expression in four out of six minigenes tested (*SETD4*, *PRPF38A*, *THOC1*, and *SEPTIN2*,
262 Fig. 3G), supporting that the effect of the *SF3B1* mutation persist in an unrelated cell line. In the two

263 remaining cases (*TPP2* and *BRCA1*) the usage of the upstream AG' was already low in the validation
264 assays using the K562 and Nalm6 cell line pairs (Supplemental Fig. S24).

265 Overall, these results supported a common effect of *SF3B1* mutations in different biological
266 backgrounds and confirmed their predominant impact on alternative 3' splice site usage and intron
267 retention.

268 **Computational prediction of protein function and stability of individual** 269 **splicing isoforms.**

270 Further on, we used the LRTS information on full-length transcripts to predict the potential coding
271 sequences of the transcripts (Supplemental [Methods](#)). Among the 28,261 known transcripts (FSM),
272 we found that 73.4% had a matching reference coding sequence (CDS) (Fig. 4A). In contrast, the
273 majority of the 58,168 novel transcripts (89.7%) did not match a CDS and either lacked an open
274 reading frame (ORF) (20.1%), initiated from an unannotated start codon (13.2%), or began at an
275 annotated initiation site but deviated from the reference CDS due to alternative splicing (56.4%).
276 Additionally, we observed that 35.1% of the expressed novel transcripts were likely to be targeted by
277 NMD, compared to only 8.2% among the known transcripts (Fig. 4A).

278 We next determined how *SF3B1* mutation-induced alternative splicing impacts the coding potential
279 and the function of the proteins. To this end, we classified ASEs into categories based on their relative
280 location and the impact on the coding sequence: 5'UTR, disrupted start codon, in-frame, frameshift,
281 disrupted stop codon, and 3'UTR. Of the 326 events featuring 3'AS, the majority led to either
282 frameshift modifications (121 events) or in-frame changes (94 events) in the CDS (Fig. 4B). In total,
283 we identified 274 ASEs predicted to yield stable alternative proteins which are not predicted to
284 undergo NMD.

285 We further examined more thoroughly the functional consequences of the two novel *SF3B1* isoforms,
286 one predominantly expressed in CLL and one in MDS (Fig. 1D-E). In the *SF3B1*-CLL transcript, the
287 fourth intron was retained which contained a PTC that shortened the CDS to 465 nt. Our prediction
288 showed a strong signal for NMD due to the presence of 19 downstream exon-exon junctions. In

289 contrast, in the *SF3B1*-MDS transcript, the penultimate exon was skipped, a frameshift was
290 introduced and, subsequently, a PTC. However, because this PTC was located within the last exon,
291 the *SF3B1*-MDS transcript was unlikely to be targeted by NMD and should result in a protein product
292 missing its C-terminal section i.e., HEAT domains 18–20 and the terminal anchor domain (Fig. 4C).

293 To further investigate the functional outcomes of the altered proteins, we examined the impact on
294 protein domain levels, by aligning Pfam (Mistry et al. 2021) domains to the predicted protein
295 sequences. For 57 of the ASEs, we found at least one Pfam domain that overlapped the divergent part
296 of the protein sequence, indicating partially altered protein functions (Fig. 4D, Supplemental Table
297 S4) and we found evidence that some aberrantly spliced transcripts may induce a translations re-
298 initiation event (Fig. 4E, Supplemental Fig. S25). We show this as an example for *MAP3K7*, a
299 frequently described gene with an alternative splicing event in *SF3B1*^{mut/wt} (Shiozawa et al. 2018;
300 Wang et al. 2016; DeBoever et al. 2015).

301 **The effect of *SF3B1* mutations depends on the distance and sequence** 302 **context of alternative 3' AS splice sites**

303 When we plotted the fraction of differential 3'AS against the splice site differences, we noticed,
304 consistent with previous findings, that the alternative splice sites of differential 3'AS events were
305 enriched within 12–21 nt upstream of the canonical splice site (AG) mainly used in the *SF3B1*^{wt/wt}
306 samples (Fig. 5A) (Obeng et al. 2016; Wang et al. 2016; Tang et al. 2020; DeBoever et al. 2015;
307 Kesarwani et al. 2017). Within this range, 30.8% of 3'AS were significantly differentially used in
308 *SF3B1*^{mut/wt} compared to 1.8% outside this range. To investigate the significance of the AG'–AG
309 distance and sequence context for the differential splicing event we constructed a minigene assay with
310 a part of the *THOC1* transcript harboring the significantly affected 3'AS. The insert was then
311 modified by replacing the 21 nt fragment between the AG' and AG of *THOC1* with 45–50 nt long
312 AG'–AG fragments from alternatively but non-differentially spliced introns (*PABCL1*, *USP1*,
313 *ZNF124*, Supplemental Fig. S26). As a control, we mutated AG' to GG', to disrupt any alternative
314 splicing. (Fig. 5B). These experiments suggested that increasing the AG'–AG distance was sufficient
315 to remove the 3'AS from SF3B1 regulation. We also found that specific sequences, such as the branch

316 point region upstream of AG', were responsible for the differential splicing between *SF3B1*^{wt} and
317 *SF3B1*^{K700E} expressing cells and confirmed that a strong AG is required for AG' usage (Darman et al.
318 2015) (Fig. [5C](#), for details please see Supplemental [Notes](#)).

319 Our results confirmed that mutations in *SF3B1* primarily affected proximal AG', but the AG'–AG
320 distance did not seem to be the sole factor required for the usage of AG'. Moreover, we did not find
321 any motif enriched at this position that could indicate a binding of another protein potentially
322 disrupting or competing with SF3B1. We therefore speculated that SF3B1 binding at the sites with
323 ASE may be altered in patients carrying *SF3B1* mutation.

324 **K700E mutation may lead to destabilization of SF3B1-mRNA binding**

325 To scrutinize the effect of the most common *SF3B1* mutation, K700E, we performed molecular
326 dynamics (MD) simulations of 20 transcripts with a 3'AS within 50 nt distance, including 14
327 transcripts which were differentially spliced between *SF3B1*^{mut/wt} and *SF3B1*^{wt/wt} and 6 non-
328 differentially spliced transcripts (Supplemental [Methods](#)). For each transcript, we performed four
329 replicas of 200 nt long MD simulations of the mRNA with: i) the BP of the downstream AG (BP)
330 bound to SF3B1^{wt}; ii) the BP of the upstream AG (BP') bound to SF3B1^{wt}; iii) the BP bound to the
331 SF3B1^{K700E} mutant; and iv) the BP' bound to the SF3B1^{K700E} mutant. There were no differences in BP
332 binding between all transcript–protein combinations in the first binding pocket (Supplemental Fig.
333 [S27A–C](#), whereas the frequency of the interaction in the second binding pocket decreased in
334 SF3B1^{K700E} due to repulsive interactions of like-charged atoms (Supplemental [Methods](#), Supplemental
335 Fig. [S27A, D–E](#)).

336 To investigate whether this decrease in interactions is associated with an increase in mRNA mobility,
337 we computed the per-residue root mean square fluctuations (RMSF) of the mRNA nucleobases.
338 Indeed, the mobility of nucleobases significantly increased for nucleobases in the 3' direction after the
339 K700 binding site for SF3B1^{K700E} compared to SF3B1^{wt}, when the BP' was bound to SF3B1
340 (Supplemental Fig. [S27F](#) and Supplemental Fig. [S28–S32](#)).

341 Taken together, these results indicated that the K700E mutation did not lead to differences in the BP
342 recognition. However, the mutation led to significant differences in SF3B1-mRNA contacts within the
343 second mRNA-binding pocket at least for the 20 mRNAs tested with an upstream alternative AG
344 within 50 nt.

345 **SF3B1 shows multimodal binding at 3' splice sites**

346 Our splicing analyses showed that mutations in *SF3B1* resulted in the activation of alternative splice
347 sites 12 to 21 nt upstream from the canonical AG (Fig. 5A, 6A). To understand how SF3B1
348 recognizes these sites, we performed individual-nucleotide resolution UV crosslinking and
349 immunoprecipitation (iCLIP) to map SF3B1 binding sites throughout the transcriptome (König et al.
350 2010). After UV crosslinking, we immunoprecipitated SF3B1 from K562-SF3B1^{wt/wt} and K562-
351 SF3B1^{K700E/wt} cells, yielding together more than 100 million SF3B1 crosslink events (Fig. 6B,
352 Supplemental Table S1). To facilitate direct comparisons, we randomly subsampled the sequencing
353 reads to adjust the library size of the replicates (see Methods). Based on the merged iCLIP data we
354 identified 96,852 SF3B1 reproducible binding sites with an optimal width of 5 nt (Fig. 6C–D). The
355 binding sites occurred in 8,127 genes, with the vast majority being protein-coding genes (93%). As
356 expected, within the protein coding transcripts, SF3B1 mostly bound introns (94%, Fig. 6E–F).

357 Since K562-SF3B1^{K700E/wt} cells expressed both wild-type and mutated SF3B1 protein and both
358 variants were recognized by the anti-SF3B1 antibody that specifically binds to the SF3B1 N-terminus,
359 we tested for differences in the SF3B1 binding between K562-SF3B1^{K700E/wt} and K562-SF3B1^{wt/wt}.
360 Consistent with a recent study (Porter et al. 2021), the K700E mutation did not generally impair RNA
361 binding (Supplemental Fig. S33). Moreover, at the level of binding sites, we detected only minor
362 differences between K562-SF3B1^{K700E/wt} and K562-SF3B1^{wt/wt} (Fig. 6G). Thus, although subtle
363 differences may have been masked by the overlay of both protein variants in the heterozygous cells,
364 the K700E mutation does not obviously change the global RNA binding behavior of SF3B1.
365 However, local changes as predicted with the molecular dynamics simulations might be too dynamic
366 to be caught by the global iCLIP analysis.

367 Next, we examined SF3B1 binding at 3' splice sites. Using meta-profiles, we detected two prominent
368 peaks of SF3B1 binding (Fig. 6H). The two peaks were centered at about -50 nt and -10 nt upstream
369 of the canonical 3' splice site and surrounded the BP (Fig. 6H). Visual inspection indicated multiple
370 SF3B1 binding sites within each peak (Supplemental Fig. S34). When centering the metaprofiles to
371 the predicted branch point adenosine, SF3B1 binds about 25 nt upstream and directly downstream of
372 the predicted branch point adenosine (Fig. 6H). The binding peak at -10 nt of the canonical 3' splice
373 site, coincided with the Py-tract region bound by U2AF2 (Zarnack et al. 2013). To test this, we
374 performed iCLIP experiments with U2AF2 in K562-SF3B1^{wt/wt} cells which confirmed that the -10 nt
375 SF3B1 peak overlapped with U2AF2 binding (Fig. 6H). Together, these observations indicated that
376 SF3B1 binds at both sites of the BP and that the binding site directly downstream of the BP
377 encompasses the Py-tract.

378 Consistent with the two peaks of SF3B1 binding at 3' splice sites, we found that SF3B1 binding sites
379 frequently occurred at distances of ~30 nt to each other (Fig. 6I). To globally classify the SF3B1
380 binding pattern, we merged adjacent binding sites into equal-sized binding regions (80 nt, 56,224
381 regions, Supplemental Table S6) and performed unsupervised uniform manifold approximation and
382 projection (UMAP) followed by density-based of applications with noise (DBSCAN). This yielded
383 three distinct clusters: Cluster C1 harbored mostly isolated binding sites (35,907 regions), cluster C2
384 included two closely spaced binding sites (5,635 regions), whereas cluster C3 showed a more
385 complex arrangement of 3–4 binding sites with wider spacing (13,847 regions) (Fig. 6J–K,
386 Supplemental Fig. S35). The latter were located closest to 3' splice sites (Fig. 6L), suggesting that
387 multiple SF3B1 binding sites assemble into complex binding patterns at 3' splice sites. We then
388 investigated the differences in binding between K562-SF3B1^{K700E/wt} and K562-SF3B1^{wt/wt} solely in the
389 C3 cluster. We noticed that K562-SF3B1^{K700E/wt} showed a slight decrease in the left peak, that was
390 more distal to the canonical AG (Supplemental Fig. S36). Within cluster C3, we noticed a slight
391 decrease in binding of K562-SF3B1^{K700E/wt} compared to K562-SF3B1^{wt/wt} specifically in the left peak,
392 that was more proximal to the canonical AG (Supplemental Fig. S36). This suggested that the K700E

393 mutation led to change in the complex arrangement of SF3B1 binding sites at 3' splice sites,
394 preferentially affecting AG-proximal binding.

395 Altogether, our SF3B1 iCLIP data showed that SF3B1 adopted a multimodal mode of binding at 3'
396 splice sites, with two major peaks of SF3B1 binding that surround the BP. The peaks include multiple
397 binding sites, which may reflect the dynamic binding rearrangements during the splicing process. The
398 strong enrichment of this binding pattern at 3' splice sites suggested that the defined arrangement of
399 binding is required for SF3B1's function in splicing.

400 **SF3B1 alternates within a small window of alternative splice site distances** 401 **and the K700E mutation leads to the use of the proximal upstream AG**

402 We and others (Alsafadi et al. 2016; DeBoever et al. 2015; Darman et al. 2015; Zhang et al. 2019)
403 found that 3' splice sites are particularly sensitive to *SF3B1* mutations when they are directly
404 preceded by an alternative 3' splice site. To test how this relates to binding, we overlaid the iCLIP
405 data with the splicing quantifications from the same isogenic cell lines (K562-SF3B1^{wt/wt} and K562-
406 SF3B1^{K700E/wt}). As shown above, the *SF3B1* mutations showed most prominent effects on alternative
407 3' splice sites 12–21 nt upstream of the canonical 3' splice site (Fig. 5A). At this distance, the
408 upstream AG was at the border of the proximal peak of SF3B1 binding (Fig. 7A). Moving further
409 away, the effects subsided drastically as soon as the upstream AG emerged from the proximal peak
410 (Fig. 7B).

411 The 3'AS events within the critical window showed a distinct behavior already in wild-type K562-
412 SF3B1^{wt/wt} cells. If the distance between 3' potential splice sites was either < 12 nt or
413 > 21 nt, splicing predominantly occurred at the upstream AG in the vast majority of cases, indicating
414 that the spliceosome generally favored upstream AG usage. However, within the critical window, this
415 pattern was inverted, such that upstream AG at these distances were generally outdone by the
416 canonical downstream AG. This indicated that in presence of two AGs within 12–21 nt the
417 downstream AG is predominately used for splicing. In line with this notion, we found that upstream
418 AGs were generally depleted from this region (Fig. 7C, top panel). Moreover, upstream AG still

419 falling into this window showed basal usage already in wild-type cells in more than 50% of cases and
420 were predominantly activated upon *SF3B1* mutation (Fig. 7C, middle panel), which could be
421 facilitated by a less stable binding of SF3B1^{K700E} to the pre-mRNA as suggested by the molecular
422 dynamics model.

423 Based on these observations, we hypothesized that diminished SF3B1^{K700E/wt} binding to the pre-
424 mRNA results in increased upstream AG usage within a critical window that may particularly impair
425 splicing fidelity. Indeed, *SF3B1* mutation preferentially affected the first (nearest) upstream AG (Fig.
426 7D) which lay significantly closer to the 3' splice site compared to non-differential 3'AS (Fig. 7E). In
427 contrast, the BP of the differential spliced 3'AS events were neither moved closer to the 3' splice site
428 nor differed in their predicted strengths (Fig. 7E, Supplemental Fig. S37), indicating that the
429 positioning of the upstream AG rather than the BP was the primary determinant for the observed
430 effects. The differentially used upstream AG had slightly lower predicted splice site strengths
431 (MaxEnt score (Yeo and Burge 2003)) than non-differential 3'AS. The downstream AGs were
432 considerably stronger than unused upstream AGs (Fig. 7C, lower panel). This was accompanied by
433 higher splice site strengths of the canonical 3' splice sites, indicating that a strong canonical 3' splice
434 site was required to support differentially spliced 3'AS events (Fig. 7F).

435 Taken together, we propose that SF3B1 binding often directly overlaps with alternative AG' in a
436 constrained window upstream of 3' splice sites (> 20 % Fig. 5D), thereby using downstream 3' splice
437 sites in wild-type conditions. In the presence of SF3B1 mutations, SF3B1 cannot properly bind the
438 Py-tract of the downstream canonical AG when a strong upstream AG' lies within a short distance,
439 leading to increased alternative splicing with the use of an alternative branch point (Supplemental Fig.
440 S38). This is partly explained by changes in the second binding pocket of SF3B1, as predicted by
441 molecular dynamics simulations, that destabilize the SF3B1-mRNA interaction. Thus, although
442 mutated SF3B1 still binds to both sides of the BP in this scenario, the changed SF3B1 protein
443 structure reduces usage of the downstream AG, resulting in widespread splicing defects.

444 **Discussion**

445 Alternative splicing plays a critical role in generating transcriptome diversity, and its dysregulation
446 has been linked to various diseases, including cancer. Yet, complex transcriptome studies are
447 challenging to perform with classical RNA-seq due to the frequent ambiguity in mapping short reads
448 and the difficulties in identifying novel isoforms (Hooper 2014; Rehrauer et al. 2013; Zhang et al.
449 2017). Of note, our capacity to study transcriptomes and alternative splicing events has expanded with
450 the recent advancements in long-read sequencing technologies. To the best of our knowledge, up to
451 now only few studies applied long-read Oxford Nanopore sequencing to analyze splice site alterations
452 mediated by mutated *SF3B1* in CLL (bulk (Tang et al. 2020) or single cell (Peng et al. 2024)) and
453 MDS (single-cell (Cortés-López et al. 2023)).

454 Here, we used long-read Iso-Seq (PacBio) sequencing of 44 patients to investigate the impact of
455 *SF3B1* mutations on alternative splicing. Our LRTS analysis revealed a wide variety of transcripts,
456 with more than two thirds of yet unannotated, or even over 3,000 transcripts from novel genes. This
457 highlighted the importance of long-read sequencing for comprehensive transcriptome profiling,
458 particularly for detecting novel splice variants and alternative splicing events. Our results supported
459 previous findings that *SF3B1* mutations specifically alter the usage of 3' alternative splice sites and
460 intron retention (DeBoever et al. 2015; Darman et al. 2015; Alsafadi et al. 2016). Importantly, using a
461 comprehensive setup of two cohorts of CLL and MDS patients, complemented by two isogenic cell
462 line pairs we were able to substantially expand the catalog of differentially spliced 3'AS to a total of
463 326 3'AS events in 266 genes (Supplemental Table S4). We observed similar effects in both patient
464 cohorts and the cell lines studied, indicating a common effect of *SF3B1* mutation on splicing. The
465 clinical differences and prognosis of the CLL and MDS patients with *SF3B1* mutations most likely
466 depend on the specific gene expression profiles and thus their different relevance of splicing
467 alterations.

468 Our results revealed that *SF3B1* mutations affect genes involved in the major mRNA splicing
469 pathway, indicating a broader impact on the splicing machinery. This may trigger secondary effects

470 on splicing, potentially leading to altered transcriptomes and disease phenotypes. In particular and in
471 concordance with previous studies, we observed an enrichment of the 3'AS events 12–21 nt upstream
472 of the canonical 3' splice sites (DeBoever et al. 2015; Darman et al. 2015; Alsafadi et al. 2016). This
473 region is typically depleted of alternative AG dinucleotides. Why AGs are depleted in the critical
474 region remains unclear. One possibility is that they are removed by purifying selection and might
475 have an evolutionary advantage. However, a subset of introns still contains AGs within this critical
476 region leading to alternative 3'AS. Why some AGs remain present in the critical region is also
477 unclear. We and others (Darman et al. 2015) observed that in these cases, the canonical AG is often
478 stronger than at other 3' splice sites, indicating that a strong 3' splice site may be required to tolerate
479 an upstream AG within the critical window.

480 Our minigene assays did not confirm that these 3'AS events depend solely on the distance between
481 canonical and alternative AGs. Coinciding with the critical range of AG'–AG distances of 12–21 nt,
482 we observed a bimodal binding of SF3B1 surrounding the BP, whereby the BP often coincided or was
483 near to the upstream AG'. Thereby, SF3B1 binding may shield the upstream AG from recognition
484 during splicing as was proposed by (Kesarwani et al. 2017). We did not detect obvious global changes
485 in SF3B1^{K700E} binding, which might be due to a temporary effect and/or binding affinity of the
486 SF3B1^{K700E} that might not be caught by iCLIP analyses. However, our MD simulations predict that
487 the number of contacts between the mRNA and residue 700 of SF3B1 resulted in increased mobility
488 of the mRNA at both, canonical and alternative 3' splice sites in SF3B1^{K700E}.

489 We confirmed with minigene assays that SF3B1^{mut} uses an alternative BP that leads to 3'AS usage in
490 SF3B1^{mut/wt} cells (Darman et al. 2015; Alsafadi et al. 2016). However, transcriptome-wide studies
491 revealed that about one third of all human exons have multiple branch points (Pineda and Bradley
492 2018; Mercer et al. 2015), which argues against the hypothesis that mutated SF3B1 always prefers
493 usage of an alternative BP. Consistently, we did not observe any strong alterations in the binding of
494 SF3B1^{K700E/wt} to mRNAs, although a slight increase of the peak at the Py-tract directly upstream of the
495 3' splice site was observed (Supplemental Fig. S39). Whereas we cannot exclude that we partially co-
496 precipitated U2AF2, we and others did not find obvious changes in U2AF2 binding to SF3B1 in

497 immunoprecipitations (Alsafadi et al. 2016; Cretu et al. 2016). Furthermore, U2AF2 binds to the Py-
498 tract only during early stages of the splicing process and is released during transition to the activated
499 B complex (Agafonov et al. 2011). In contrast, SF3B1 is assembled into the spliceosome later, and
500 subsequently replaces U2AF2 in the activated B complex and binds to the Py-tract with its binding
501 pocket consisting of HEAT domains 3–7 that harbor most of the mutational hotspots (Schmitzová et
502 al. 2023). Therefore, the downstream peak observed in our SF3B1 iCLIP experiments most likely
503 corresponds to SF3B1 binding.

504 The differential splicing observed in *SF3B1*^{mut/wt} may result from a weakened binding of SF3B1^{mut}-
505 containing spliceosomes. There might be an additional effect of decreased binding of DDX46/PRP5, a
506 kinase involved in proof reading of the pre-mRNA branch site (Carrocci et al. 2017; Tang et al. 2016;
507 Zhao et al. 2022; Zhang et al. 2021). Other splicing proteins that have been shown to bind less to
508 SF3B1^{mut}, are DDX42 (Zhao et al. 2022), DHX15 (Zhang et al. 2024) and SUGP1 (Zhang et al. 2019,
509 2023). DDX42 and DDX46 have been shown to sequentially occupy the RNA binding pocket
510 consisting of HEAT repeats 3–7 during early steps of the splicing process (Zhang et al. 2021, 2024;
511 Yang et al. 2023).

512 Besides gaining additional insight into the SF3B1 splicing mechanism, we also explored the splicing
513 and expression alterations identified through the Iso-Seq sequencing approach in CLL and MDS.
514 Apart from identification of large numbers of new differentially spliced genes, we were able to
515 specifically map the toxic exon of *BRD9* to its isoform and predict its amino acid composition. We
516 also identified *SF3B1* isoforms specifically more present in MDS or CLL, albeit their expression
517 levels were at approximately 10%. The *SF3B1*-CLL transcript is predicted to undergo NMD, whereas
518 the *SF3B1*-MDS is predicted to miss its C-terminal part. This shortened protein might rescue part of
519 the detrimental effect of the mutated *SF3B1* leading to a slightly favorable prognosis. *SF3B1*
520 overexpression in CLL B cells with respect to normal B cells has been reported before (Wan and Wu
521 2013). Possible reasons for this increase might be: i) regulatory mechanisms, such as alterations in
522 transcription factors or epigenetic changes; ii) oncogenic pathways involving growth factors or
523 cytokines; or iii) feedback mechanisms (Huang et al. 2011). We would speculate that due to the

524 NMD-sensitive *SF3B1*-CLL transcript the cell increases the transcription of the regular version of the
525 *SF3B1* transcript to compensate for this loss. This is not necessary in the case of MDS, because the
526 *SF3B1*-MDS transcript is functional and only misses its C-terminal part. This would explain how the
527 *SF3B1*-CLL isoform under NMD is related to the overexpression of *SF3B1* in CLL.

528 This overexpression can now be brought in context with the worse prognosis of mutated SF3B1 in
529 CLL, because in the presence of the mutation this overexpression leads to a more massive
530 dysregulation of alternative splicing caused by SF3B1 in CLL (compared to MDS) which in turn
531 disrupts multiple pathways and lowers survival of the patients. Indeed, if we further investigate the
532 disease specific ASEs in mutated vs wild-type *SF3B1* patients (Fig. 2B; Supplemental Table S2) we
533 observed that the overall number of ASEs is fairly similar (CLL: 288, MDS: 219) with 1.31-fold
534 increase in ASEs in CLL, but there was a drastic ($\times 2.75$) increase of CLL-specific IR events (CLL:
535 77, MDS: 28) (Supplemental Fig. S40A). The ASEs break down to 69 (CLL) and 27 (MDS) unique
536 genes and were largely different with only six genes in common (Supplemental Fig. S40B). The
537 pathways affected by CLL-IRs were related to mRNA splicing machinery, oncogenic signaling
538 pathways, as well as immune pathways that might affect the survival of the patients (Supplemental
539 Fig. S40C). Thus, we would argue that the overexpression of *SF3B1* in CLL compared to MDS leads
540 to elevated splicing effects, in particular introns, that target a different, more signaling-related panel of
541 genes with multiple cellular functions that promote tumorigenesis and reduce survival. However,
542 further functional analyses will show the impact of these altered SF3B1 proteins and if they influence
543 MDS and CLL pathomechanisms.

544 Another example, which is frequently reported to be differentially spliced within *SF3B1* mutated
545 cancers is *MAP3K7*. We were able to show that mutations in *SF3B1* led to reduced expression of
546 longer isoforms and increased expression of isoforms with shortened protein kinase domain, likely
547 impacting its function. Thus, with this data at hand it is possible to not only identify splicing events
548 but to also map them to their cognate isoform and thus to provide information on the resulting protein
549 composition. This is a fundamental information for understanding splicing data and to gain insight
550 into pathomechanisms underlying CLL and MDS.

551 Our study provides new insights into the mechanism by which *SF3B1* mutations affect splicing
552 regulation, and the potential consequences on protein function. Our findings highlight the importance
553 of long-read sequencing for investigating differential alternative splicing usage and splicing factor
554 function. These results have implications for understanding the role of *SF3B1* mutations in
555 hematological malignancies and other diseases, and may be used in the future to predict new
556 approaches for targeted therapies for these conditions.

557 **Methods**

558 **Ethics approval**

559 The study was approved by the Ethics Committee of the University of Cologne (Ethikvotum 11-319
560 from 11th December 2011, with an amendment from 7th June 2016) and the Ethics Committee of the
561 University of Düsseldorf (Ethikvotum 3768, amendment from 24th October 2018). Informed consent
562 has been obtained from all patients involved.

563 **Cell lines and patients' samples**

564 The isogenic cell line pairs, K562-SF3B1^{K700E/wt} and its parental K562-SF3B1^{wt/wt}
565 (RRID:CVCL_0004), as well as Nalm6-SF3B1^{H662Q/wt} and its parental Nalm6-SF3B1^{wt/wt} were
566 obtained from Horizon Discovery (HD181-012, HD115-110). Since homozygous *SF3B1* mutations
567 were reported to be lethal (Lee et al. 2016), we used heterozygous cell lines. The K700E mutation is
568 the most frequent *SF3B1* mutation reported in CLL and MDS, and H662Q mutation is also frequently
569 reported (Wan and Wu 2013; Rossi et al. 2011; Quesada et al. 2012). The *SF3B1* mutated cell lines
570 were described previously (Darman et al. 2015). HEK293-FT (RRID: CVCL_6911) was purchased
571 from Thermo Fisher Scientific (#R70007). Information on cell line authentication and cell growth
572 conditions are provided in Supplemental [Methods](#).

573 CLL and B cell samples were obtained from the CLL-Biobank Cologne. *IGHV* mutational status was
574 determined, as previously described (Rosenquist et al. 2017). Peripheral blood B cells were isolated

575 via negative selection using RosetteSep immunodensity-based cell separation (Stemcell Technologies,
576 Vancouver, BC, Canada). The purity of CLL/ B cells was analyzed by flow cytometry and revealed
577 that $\geq 90\%$ cells co-expressed CD5/CD19.

578 Specimens from MDS with ring sideroblast (MDS-RS) patients were obtained from the MDS Biobank
579 of the University Clinic Düsseldorf. Either RNA or cells were obtained from the Biobank. If cells
580 were obtained, RNA was isolated using the Nucleospin RNA kit (Maceray Nagel). RNA quality was
581 assessed by RNA ScreenTape analysis (Agilent) or a Bioanalyzer (Agilent).

582 Clinical information on the patients is summarized in Supplemental Table S1.

583 **Plasmids**

584 Plasmids pCMV-3Tag-1A-SF3B1^{wt} and pCMV-3Tag-1A-SF3B1^{K700E} (Alsafadi et al. 2016) were
585 designed by Angelos Constantinou (Department of Molecular Bases of Human Diseases, CNRS UPR
586 1142, IGH-Institute of Human Genetics, Montpellier 34090, France) and kindly provided by Marc-
587 Henri Stern, Institut Curie, Paris, France. Plasmids pcDNA3.1-FLAG-SF3B1-WT and pcDNA3.1-
588 FLAG-hSF3B1-K700E (Kesarwani et al. 2017) were obtained from Addgene (#82576 and #82577).
589 The human full-length *SF3B1* sequence has been previously reported to be impossible to clone into
590 bacteria (Yokoi et al. 2011; Wang et al. 1998). Therefore, the plasmids consisted of synthetic
591 sequences, codon-optimized for expression in bacteria (Alsafadi et al. 2016; Kesarwani et al. 2017).
592 For the minigene constructs the intron and parts/complete adjacent upstream and downstream exons
593 were PCR-amplified from K562 genomic DNA using Phusion Hot Start Flex polymerase (New
594 England Biolabs) and cloned by the Hot Fusion (Fu et al. 2014) method into the BamHI restriction
595 site of pcDNA3 (Invitrogen, <https://www.addgene.org/vector-database/2092/>). The open reading
596 frame of the exons was left intact. The oligonucleotides used for cloning of the constructs are listed in
597 Supplemental Table S5. Mutations and insertions were introduced by site directed mutagenesis using
598 the Q5 site-directed mutagenesis kit (New England Biolabs). Oligonucleotides for site-directed
599 mutagenesis were designed using the NEBaseChanger version 1.3.3 (New England Biolabs) and are

600 listed in Supplemental Table [S5](#). All constructs were verified by Sanger sequencing (Microsynth
601 Seqlab AG, Göttingen, Germany).

602 **cDNA synthesis and validation of the splicing alterations**

603 Transfections and RNA isolation were performed following standard procedures with PEI Max
604 (PolyScience, #24765 1) and NucleoSpin RNA Mini kit (Macherey Nagel, #740955.250). An amount
605 of 500 ng total RNA was reverse transcribed using SuperScript™ II (Thermo Fisher Scientific,
606 #18064014) and hexamer oligonucleotides for the cDNA synthesis from K562 and Nalm6 RNA. For
607 the minigene assays 500 ng RNA was reverse transcribed using SuperScript™ IV Reverse
608 Transcriptase (Thermo Fisher Scientific, #18090010) with the plasmid-specific BGH-rev oligo in 20
609 µl. Subsequently, RNA in DNA-RNA hybrids was digested by RNase H incubation. For RT-PCR we
610 used 1 µl of cDNA, Taq DNA polymerase, recombinant (Thermo Fisher Scientific, #10342020), and
611 specific oligonucleotides (Supplemental Table [S5](#)) in a volume of 25 µl. The PCR ran for 35 PCR
612 cycles. PCR products were separated on a 3–4% TAE-agarose gel.

613 **PacBio Iso-Seq library preparation and sequencing**

614 For the cDNA synthesis we used oligo(dT) oligonucleotides and the TeloPrime v2 Kit (Lexogen) to
615 ensure the amplification of full-length mRNAs that contained a Cap-structure. Barcoded primers were
616 used in the cDNA amplification step to enable multiplexing before library preparation. To enrich for
617 slightly larger cDNAs, we adjusted the magnetic bead concentration in the bead clean-up after cDNA
618 amplification. Subsequent library preparation was performed with the Express Template Kit 2.0
619 (PacBio).

620 In total, 58 libraries were sequenced on the PacBio Sequel II platform, by multiplexing 4 samples per
621 8M SMRT cell at the Genomics and Transcriptomics Laboratory, the production site of the West
622 German Genome Center in Düsseldorf (Heinrich-Heine Universität, Düsseldorf, Germany)
623 (Supplemental Table [S1](#)).

624 **Processing of PacBio Iso-Seq data**

625 Preprocessing of raw Iso-Seq sequencing data was performed with *Iso-Seq* (PacBio 2020 2021)
626 software version 3.4, using the recommended parameters. In brief, we used the *ccs* tool to call circular
627 consensus sequences by clustering and collapsing steps, *lima* to remove primers and adapters, and
628 *isoseq refine* to demultiplex samples and filter out reads not featuring poly(A) sequences. This
629 resulted more than 33 million aligned full-length non-chimeric (flnc) polyA HiFi reads (100,302 –
630 1,258,653 per library, average 582,135, Supplemental Table S1), with an average length of 2,721 bp.
631 At this sequencing depth, transcript isoforms expressed at one transcript per million (TPM) are
632 expected to be sequenced by at least 25 reads, and two TPM isoforms by at least 50 reads, with more
633 than 95% probability (Supplemental Fig. S1).

634 According to the base quality values, 58 samples had an error rate of less than 1% in at least 99.7% of
635 the HiFi reads and only four samples had higher error rates with 96.7% to 96.9% reads with < 1%
636 error rate. Overall, the quality of the reads was high, with only 9% of reads potentially affected by
637 technology-specific technical artefacts (Cocquet et al. 2006) (Supplemental Fig. S2).

638 The flnc reads were converted to FASTQ using SAMtools v.1.18 (Li and Durbin 2009), and, without
639 an additional clustering step, aligned to the human genome GRCh38.p13 using minimap2 (Li 2018)
640 version 2.22 with the preset parameters for high quality spliced reads (-ax splice:hq). For each sample,
641 at least 99.85% of the reads were mapped to the genome, and at least 91.9% were uniquely mapped.
642 Samples for which we sequenced more than one library were merged using SAMtools v.1.18 (Li and
643 Durbin 2009) after the mapping step. Sequencing and mapping statistics per sample are detailed in
644 Supplemental Table S1.

645 *SF3B1* mutation calling was done with BCFtools v.1.13 (Danecek et al. 2021) mpileup and SnpEff
646 v.5.1d (Cingolani et al. 2012).

647 Further analysis of Iso-Seq data was performed in Python v3, using IsoTools (Lienhard et al. 2023)
648 version 0.2.8. In brief, aligned reads were imported and compared to the human reference annotation

649 version 36 from GENCODE (Frankish et al. 2021), to call, annotate, classify, and quantify transcripts,
650 using IsoTools' `add_sample_from_bam` function.

651 For exploratory analysis, alternative splicing events were detected using IsoTools'
652 `alternative_splicing_events` function. For each sample, individual events were quantified by percent
653 spliced index (PSI) values, i.e., the number of reads supporting transcripts that include additional
654 exonic sequence over all transcripts spanning that event. Based on these PSI values, principal
655 component analysis (PCA) plots for different alternative splicing categories were computed using
656 IsoTools' `plot_embedding` function.

657 Differential splicing events between $SF3B1^{mut/wt}$ and $SF3B1^{wt/wt}$ CLL, MDS and cell line samples were
658 computed with the IsoTools `altsplice_test` function, using the betabinomial likelihood ratio test. This
659 test models the variability within the tested groups with a beta-binomial mixture distribution, a
660 binomial distribution where the probability parameter, p , of the binomial distribution $B(n,p)$ follows a
661 beta distribution, $Beta(a, b)$. The test compares the group-wise coverage of the splicing event with the
662 total coverage.

663 $\Lambda = -2(l_0 - l_1)$, where:

$$664 \quad l_1 = \log\left(BB(k_1|\hat{\alpha}_1, \hat{\beta}_1, n_1)\right) + \log\left(BB(k_2|\hat{\alpha}_2, \hat{\beta}_2, n_2)\right) \text{ and}$$

$$665 \quad l_0 = \log\left(BB(k_1 + k_2|\hat{\alpha}, \hat{\beta}, n_1 + n_2)\right).$$

666 Here, $BB(k|\alpha, \beta, n)$ is the probability mass function of the beta-binomial distribution and $\hat{\alpha}, \hat{\beta}$ are
667 maximum likelihood estimates for the parameters. The maximum log-likelihood parameters are
668 determined numerically by a quasi-Newton optimization method (LM-BFGS from SciPy (Virtanen et
669 al. 2020)). Under the null hypothesis (i.e. no differential splicing) the test statistic is X^2 distributed
670 with two degrees of freedom.

671 This formulation allows for considering within-group variability in a similar manner as tests based on
672 negative binomial distribution for RNA-seq data, which is crucial for heterogeneous samples such as
673 individual cancer patients. To be tested, we required the events to be covered by at least 10 reads in at

674 least 4 samples per group, and the minor alternative to be covered by at least 5% of the total reads
675 (test = 'betabinom_lr', min_n = 10, min_sa = 4, min_alt_fraction = 0.05). Due to the limited number of
676 samples, we did not include any covariates in the model analysis. We did not observe any bias
677 towards highly expressed genes (Supplemental Fig. S41).

678 All 3' alternative splicing events (including those that were not differentially expressed between
679 *SF3B1*^{mut/wt} and *SF3B1*^{wt/wt}) were exported for further analysis (Supplemental Table S4).

680 Differential expression analysis on the Iso-Seq read counts was performed with DESeq2 (Love et al.
681 2014).

682 **Estimation of branch point position and splice-site strength score**

683 For all 3' alternative splicing events, we used the R (R Core Team 2021) Bioconductor package
684 branchpointR (Signal et al. 2018) to predict branchpoint probabilities for both, canonical and
685 alternative splice sites. We used the position with the highest branchpoint probability as the predicted
686 branch point.

687 The strength of the 3' splice sites for the minigene constructs was calculated with SpliceRover
688 (Zuallaert et al. 2018) (<http://bioit2.irc.ugent.be/rover/splicerover>, accessed on 23.06.2023 using the
689 model "human acceptors").

690 **Transcript coding potential**

691 For a detailed description of the methods, please see Supplemental Notes.

692 **Illumina RNA-seq library preparation and sequencing**

693 RNA from K562 and Nalm6 cells was isolated using the NucleoSpin RNA Mini kit (Macherey Nagel)
694 followed by DNase-digestion with the DNase I set (Zymo Research, #E1010) and clean-up using the
695 NucleoSpin RNA Clean-up Mini Kit (Macherey-Nagel, #740948.50). RNA quality was surveyed
696 using RNA ScreenTapes (Agilent) and the RNA integrity number (RIN) was 10 for all samples. CLL

697 cells from 19 CLL patients used for Iso-Seq and additional eight patients (including four patients with
698 *SF3B1* mutation) were stored in RNA later and RNA was isolated using the RNeasy Mini Kit
699 (Qiagen) followed by DNase digestion using the DNase I Amplification Grade Kit (Invitrogen) and a
700 clean-up using RNeasy Min Elute columns (Qiagen). RNA-seq libraries for the cell lines and the CLL
701 RNA were prepared using the TruSeq-stranded-total RNA SamplePrep kit (Illumina) according to the
702 manufacturer's protocol. In brief, 2 µg of total RNA were depleted for ribosomal RNA using a Ribo-
703 Zero-rRNA removal kit (Illumina), followed by random primed cDNA synthesis. Sequencing libraries
704 were run at the sequencing core unit of the Max-Planck Institute for Molecular Genetics on a HiSeq
705 2500 (Illumina) using 50 bp paired-end reads.

706 **Processing of Illumina RNA-seq data**

707 We collected publicly available MDS RNA-seq data from the Short Read Archive using the following
708 criteria: "MDS" in the study description, paired-end Illumina reads, available FASTQ files, human
709 blood cell samples, and at least 5 samples per study. After merging technical replicates, our dataset
710 consisted of 263 samples from *SF3B1*^{wt/wt} MDS patients and 135 *SF3B1*^{mut/wt} patients. All study and
711 sample IDs are listed in Supplementary Table S1.

712 Illumina RNA-seq reads were aligned to the human reference genome GRCh38.p13 using STAR
713 aligner version 2.7.6a (Dobin et al. 2013), with provided GFF annotation from GENCODE release 36
714 including annotation of non-chromosomal scaffolds. Alternative splicing events were called and
715 quantified using rMATS (v4.1.1) (Shen et al. 2014).

716 Mutations in *SF3B1* were called as for Iso-Seq described above.

717 **iCLIP experiments**

718 iCLIP of K562-SF3B1^{wt/wt} (three replicates) and K562-SF3B1^{K700E/wt} (two replicates) was performed
719 as described previously (Sutandy et al. 2016). To this end, exponentially growing K562 cells were
720 pelleted, washed once with PBS and 10×10^6 were subjected to UV-cross linking at 400 mJ/cm² at

721 254 nm in 6 ml PBS in a 10 cm petri dish on ice. Crosslinked cells were scraped from the dish,
722 collected by centrifugation for 2 min at $500 \times g$ at 4 °C, snap frozen in liquid nitrogen and stored at -
723 80 °C. About 3×10^6 cells were immunoprecipitated using 10 μg of a monoclonal SF3B1 antibody
724 (clone 16, MBL D221-3) or 10 μg of the monoclonal U2AF2 antibody (U4759, SIGMA).

725 **iCLIP data processing**

726 Initial quality control was done using FastQC (Andrews 2010) before and after quality filtering). All
727 reads having at least one position with a sequencing quality < 10 in the barcode area (positions 1–3,
728 4–7, 8–9) were removed. De-multiplexing and adapter trimming were done on quality filtered data
729 using Flexbar (Roehr et al. 2017). No mismatches were allowed during de-multiplexing, while an
730 error rate of 0.1 was accepted when trimming the adapter at the right end of the reads. Furthermore, a
731 minimal overlap of 1 bp between reads and adapter was required and only trimmed reads with a
732 minimal length of 15 bp (24 bp including the barcode) were kept for further analysis. Barcodes of
733 remaining reads were trimmed (but kept as additional information with the reads). Trimmed reads
734 were then mapped to genome assembly version GRCh38 using STAR (v. 2.6.1b) (Dobin et al. 2013)
735 with 4% mismatched bases allowed and turned off soft-clipping on the 5-prime end, and GENCODE
736 gene annotation v31 (Frankish et al. 2019). While technical duplicates were removed with UMI-tools
737 (Smith et al. 2017) with *unique* method, all real duplicate reads were kept. Then, we checked the
738 crosslink quality with iCLIPPro (Hauer et al. 2015).

739 For facilitate comparisons, the crosslink events, i.e., reads after duplicate removal, of the replicates
740 were randomly subsampled to the size of the smallest replicate (n = 13,892,358) (replicate 2 of
741 SF3B1^{K700E/wt}; Supplemental Table S1).

742 **Definition and classification of SF3B1 binding sites**

743 Binding sites were identified from the merged crosslink events of SF3B^{wt/wt} and SF3B1^{K700E/wt} as
744 described in Busch *et al.* (Busch et al. 2020) (Supplemental Fig. S29). For this, the crosslink events of
745 all five replicates were combined and subjected to peak calling with PureCLIP (version 1.3.1)
746 (Krakau et al. 2017) with default parameters. The PureCLIP-called sites (Psites) were filtered by first
747 removing 5% of the Psites with the lowest score associated and then keeping only the top 20% of
748 Psites within each gene annotated (GENCODE release 36, GRCh38; only annotations with a gene
749 support level of 1 or 2 and transcript support level from 1 to 3). The Psites were then merged into

750 binding sites using the R/Bioconductor package BindingSiteFinder (version 1.0.3) (Brüggemann and
751 Zarnack), using the following options: width of 5 nt (bsSize =5); ≥ 2 Psites (minWidth = 2,
752 minCISites = 1) and ≥ 1 crosslink positions within each binding site (minCrosslinks = 1). In brief,
753 Psites closer than 5 nt were merged into regions, and isolated Psites were discarded. Within each
754 region, binding site centers were iteratively placed at the position with most crosslink events and
755 extended by 2 nt on both sides. Binding site centers were required to harbor the maximum crosslink
756 signal within the binding site. The optimal binding site width of 5 nt was determined by an evaluation
757 of the ratio of crosslink events within binding sites of increasing width over the mean background
758 signal in flanking windows of the same size (Supplemental Fig. S30). Next, binding sites that were
759 not supported by all replicates in at least one condition (SF3B1^{wt/wt} or SF3B1^{K700E/wt}) were filtered out
760 (Supplemental Fig. S31). The threshold for sufficient coverage in a replicate was determined using the
761 5th percentile and a lower boundary of two crosslink events as described in Busch *et al.* (Busch et al.
762 2020). Finally, binding sites were assigned to target genes using GENCODE annotation (release 36,
763 GRCh38; filtered as above) as described in Busch *et al.* (Busch et al. 2020). In total, this procedure
764 identified 96,852 SF3B1 binding sites in 8,127 genes.

765 To classify distinct SF3B1 binding patterns in introns, bound regions were defined by merging
766 intronic binding sites within a distance < 55 nt and resizing the obtained regions to 81 nt around the
767 center, resulting in 56,224 regions harboring 87,199 binding sites. Following the approaches
768 suggested in Heyl *et al.* (Heyl and Backofen 2021), we used unsupervised clustering to separate the
769 crosslink patterns in the bound regions. For this, the crosslink coverage (sum of all replicates) was
770 subjected to min-max normalization (Tarantola 2008) within each window (i.e., scaling such that the
771 lowest and highest number of crosslink events are set to 0 and 1, respectively), followed by spline-
772 smoothing using the smooth.spline function (R package stats, version 4.1.0) with lambda 0.2 (spar =
773 0.2) and inflated dimensions (dim = 150). This changed the shape of the matrix from $A \times B$ ($56,224 \times$
774 81) to $A \times B'$ ($56,244 \times 150$), where A is the number of regions, B are the nucleotide positions and B'
775 are the inflated nucleotide positions. The matrix $A \times B'$ of normalized and smoothed crosslink
776 coverages was then subjected to dimension reduction using uniform manifold approximation and

777 projection (UMAP) (McInnes et al. 2018) with the *umap* function (package *umap*, version 0.2.7) with
778 parameters `n_epochs = 5000`, `n_components = 2`, `min_dist = 0.01` and `n_neighbors = 5` (Supplemental
779 Fig. S34A). The UMAP results were assigned to clusters using density-based clustering of
780 applications with noise (DBSCAN) (Ester et al. 1996) with the *dbscan* function (R package *dbscan*
781 (Hahsler et al. 2019), version 1.1, `eps = 0.3`), with a minimum number of 150 points per cluster
782 (`MinPts = 150`), yielding three clusters: C1 (`n = 35,907` regions), C2 (`n = 5,635`) and C3 (`n = 13,847`).
783 Bound regions in cluster C0 (`n = 835`) were deemed as outliers that could not be assigned to any of the
784 fitted density centers and excluded from further analysis. Bound regions in cluster C3 (wide pattern)
785 were smoothed more finely (`spar = 0.1`, `dim = 500`) and then subjected to a second round of UMAP
786 dimension reduction (parameters as above) and DBSCAN clustering (`MinPts = 60`, `eps = 0.23`),
787 yielding subclusters #0–33 (Supplemental Fig. S34B). Cluster numbering is based on the increasing
788 distance between the two modes in the arrangement of binding sites, calculated on the summed and
789 smoothed coverages within each cluster using the *locmodes* function (R package *multimode*, version
790 1.5 (Ameijeiras-Alonso et al. 2021)) (Supplemental Fig. S34C).

791 **Data Access**

792 Cell lines' and patients' transcriptome raw FASTQ files or PacBio CCS unaligned BAM files have
793 been submitted to the NCBI BioProject database (<https://www.ncbi.nlm.nih.gov/bioproject/>) under
794 accession number PRJNA1037338 or to the European Genome-Phenome Archive ([https://ega-
795 archive.org/](https://ega-archive.org/)) under accession number EGAS50000000053, respectively. iCLIP raw FASTQ data and
796 processed files have been submitted to the NCBI Gene Expression Omnibus (GEO;
797 <https://www.ncbi.nlm.nih.gov/geo/>) under accession number GSE247658. Patient data access will be
798 controlled by the Data Access Committee at the Institute for Translational Epigenetics at the
799 University Hospital Cologne, University of Cologne, Cologne, Germany.

800 The complete Iso-Seq/iCLIP data analysis code is available as a Jupyter Notebook in the GitHub
801 repository https://github.com/ZarnackGroup/go_long2023 as well at Zenodo at the following link
802 <https://zenodo.org/doi/10.5281/zenodo.12597945>.

803 **Competing Interest Statement**

804 MH: Honoraria: (speaker's bureau and/or advisory board): Roche, Janssen, Abbvie; Research support:
805 Roche, Janssen, Abbvie, Astra Zeneca, Beigene.

806 **Acknowledgements**

807 The study was funded by the German Research Foundation: KFO286-RP8/SCHW1605/1-1,
808 SCHW1605/4-1 (GO-LONG), SFB1399 and SFB1530 to M.R.S., KFO-286-RP6 to M. H., KFO-286-
809 CP to C.D.H., SFB1530 to M.H., the Volkswagen Stiftung Lichtenberg program to M.R.S. and the
810 Center for Molecular Medicine Cologne, CMMC (A12 to M.R.S.).

811 We acknowledge IMB Genomics Core Facility and its NextSeq 500 sequencer [funded by the
812 Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) INST 247/870-1 FUGG].

813 We would like to thank Angelos Constantinou (IGH-Institute of Human Genetics, France) and Marc-
814 Henri Stern (Institut Curie, Paris, France) for sharing the *SF3B1* plasmids, Anke Busch (IMB Mainz,
815 Germany) for iCLIP data pre-processing, Elena Wasserburger-Zichel (University Hospital Cologne,
816 Germany), and Bernd Timmermann (Sequencing Core Unit, Max Planck Institute for Molecular
817 Genetics, Berlin, Germany) for their technical assistance. We furthermore thank the Regional
818 Computing Center of the University of Cologne (RRZK) for providing computing time on the DFG-
819 funded (Funding number: INST 216/512/1FUGG) High Performance Computing (HPC) system
820 CHEOPS as well as IT support. In addition, we acknowledge computational support of the Centre for
821 Information and Media Technology, especially the HPC team at the Heinrich-Heine University, as
822 well as the computing time provided by the John von Neumann Institute for Computing on the

823 supercomputer JUWELS at Jülich Supercomputing Centre (user IDs: VSK33). This work was
824 supported by the DFG Research Infrastructure West German Genome Center (407493903) as part of
825 the Next Generation Sequencing Competence Network (project 423957469). High-throughput
826 sequencing was carried out at the West German Genome Center, and production sites in Cologne and
827 Düsseldorf.

828 **Authors' contributions**

829 M.R.S., C.G., R.H., K.Z., H.G., N.G., M.H., and J.K. designed the study. H.H., L.B., A.K., K.B., K.K.
830 and C.D.H. acquired the data. A.P., M.L., M.B., J.K., H.G., C.G., R.H., K.Z., and M.R.S. analyzed
831 and interpreted the data. A.P., M.L., C.G., R.H., K.Z., and M.R.S. drafted and wrote the manuscript.
832 All authors have read and approved the final manuscript.

833 **Figure legends**

834 **Figure 1. Long-read sequencing of CLL and MDS patient samples discovers novel isoforms. (A)**
835 Distribution of *SF3B1* mutations in CLL and MDS patient samples used for Iso-Seq: each dot
836 represents a mutated sample. One CLL patient is marked twice due to two mutations (I704N and
837 D894G). Note that the A284T mutation is outside the HEAT repeat domains, and thus was grouped as
838 a wild-type sample, also according to further analysis described below. *SF3B1* is shown as the major
839 isoform expressed, with the full splice match to annotated isoform 201. **(B)** The number of substantial
840 transcripts identified in each group and expressed at the level of at least 1 transcript per million (TPM)
841 colored by the category of isoform novelty: full splice match (FSM), with incomplete splice matches
842 (ISM), with combinations of annotated splice junctions (novel in catalog, NIC), with at least one
843 novel splice site (novel not in catalog, NNC), or from novel genes (NOVEL) (Tardaguila et al. 2018).
844 **(C)** Venn diagrams showing the overlap between isoforms from **(B)** expressed at ≥ 1 TPM in each
845 group. **(D)** *SF3B1* isoforms expressed at $> 10\%$ relative expression level. **(E)** Relative expression
846 levels of *SF3B1* isoforms from **(D)**.

847 **Figure 2. *SF3B1* mutation effect is independent of the biological background, but its**
 848 **manifestation depends on the transcriptomic profile. (A)** The number of alternative splicing events
 849 identified with Iso-Seq separated by splicing event type in all groups investigated, differentiated by
 850 the novelty class. **(B)** Overlap between significantly altered alternative splicing events (ASEs) in
 851 samples with *SF3B1* mutation identified in the three datasets used (cell lines, CLL patients, or MDS
 852 patients). **(C)** Correlation of isoform usage measured by the difference in PSI from all events listed in
 853 (B) i.e., events called significant in at least one of the three datasets are shown. The colors of the dots
 854 correspond to significance reached only in one set (blue: X axis only; red: Y axis only; light blue:
 855 both; grey: none, i.e. called significant in a dataset absent from the graph). Pearson correlation
 856 coefficient (*R*) and associated p-value (*p*) are given. **(D)** Violin plots with boxplots show the
 857 distribution of expression values of the genes with dataset-specific ASEs from (C). Significant
 858 differences are marked with *** for paired, two-tailed Student's *t*-test p-value < 0.001, ** p-value <
 859 0.01, * p-value < 0.05, N.S. – not significant with p-value ≥ 0.05.

860 **Figure 3. *SF3B1* mutation increases 3' alternative splice sites usage and decreases intron**
 861 **retention. (A)** Number of differential alternative splicing events with shorter or longer variant
 862 expression (percent spliced index, PSI) in *SF3B1*^{mut/wt} vs. *SF3B1*^{wt/wt} samples. **(B)** Highly significantly
 863 altered (q-value < 0.01) 3' alternative splice sites clearly separate samples by *SF3B1* mutations in
 864 leukemia cell lines as well as CLL and MDS patients based on the longer variant PSI values showing
 865 four clusters described more in Supplementary Fig. S19. **(C)** Principal component (PC) analysis,
 866 based on the isoform usage of 3' alternative splice sites, clearly separates CLL and MDS patients, as
 867 well as cell lines, according to the *SF3B1* mutational status. **(D)** Swarm plots showing the distribution
 868 of the isoform usage (PSI) among groups with or without *SF3B1* mutation. **(E)** Validation of the
 869 differential splicing associated with *SF3B1* mutation with RT-PCR experiment in isogenic K562 and
 870 Nalm6 cell lines. **(F)** Minigene assays workflow. HEK293T cells were co-transfected with minigenes
 871 and either *SF3B1*^{wt} or *SF3B1*^{K700E} for 48 h. RNA was extracted and used for amplification of splicing
 872 products with minigene-specific primers. **(G)** The results from the minigene assay from (F). The

873 lower band in the agarose gel corresponds to the usage of the canonical AG, the upper band to the
874 upstream AG'.

875 **Figure 4. *SF3B1* mutation results in altered mRNAs potentially translated into modified**
876 **proteins. (A)** Coding potential of the known and novel isoforms identified divided by CDS similarity
877 to annotated isoforms and NMD prediction. **(B)** Effect of *SF3B1*^{mut}-associated ASEs on the protein-
878 coding potential. **(C)** *SF3B1* isoforms detected in this study with CLL- or MDS-specific isoforms. **(D)**
879 The *PRPF38A* isoform expression levels (left) and structure with Pfam domains indicated (right).
880 Highlighted in yellow is the *SF3B1*^{mut}-associated 3'AS that may influence the protein function. **(E)**
881 Schematic of major *MAP3K7* isoforms (left) with the protein kinase domain showed as green boxes.
882 The *SF3B1*^{mut}-associated 3'AS is highlighted in yellow and ORF start/end indicated by the triangles.
883 Highlighted in light green are predicted upstream ORFs (uORFs). Red box highlights the additional
884 exon 12 in the isoform 202, which is absent in the isoform 205. Expression of each isoform is shown
885 on the right. The expression of isoforms with 3'AS event is shown as striped bars.

886 **Figure 5. *SF3B1* mutations promote upstream alternative 3' splicing sites and partially**
887 **dependent on the sequence context. (A)** 3' alternative splice site distance distribution. Negative
888 distances indicate an alternative was located upstream and positive values indicate alternative located
889 downstream leading to a shorter exon. Blue line represents proportion and green total number of
890 3'alternative splicing events (3'AS). Dotted vertical lines indicate the enriched region of 12–21 nt
891 upstream of the canonical AG. **(B)** Minigene assays with long (45–50 bases) AG–AG' inserts,
892 shortened inserts containing about 20 nt directly upstream the AG including the polypyrimidine (Py)
893 tract, and short (15–20 nt) AG–AG' inserts from non-differentially alternatively spliced 3'AS events.
894 The chosen events without differential splicing detected with *SF3B1* mutation were from *PAPCLI*,
895 *USP1* and *ZNF124* (AG'–AG distance > 50 nt) as well as *GPR98B*, *UROD* and *CELF2* (AG'–AG
896 distance < 20 nt). **(C)** Table showing splice site strength for AG and AG' calculated with SpliceRover
897 (Zuallaert et al., 2018).

898 **Figure 6. Multimodal SF3B1 binding.** (A) SF3B1 choice of AG within a narrow window of 12–21
 899 nt is strongly affected by the alternative splice site distance. The AG usage is shown as percent
 900 spliced index (PSI) as a function of the splice site distance in SF3B1^{wt/wt}. Rolling mean across 20 nt
 901 and smoothed (loess method) trend line are shown for upstream (violet) and downstream (pink) AG.
 902 (B) Schematic workflow of processing iCLIP reads and calling SF3B1 binding sites. (C) Defining
 903 optimal site binding width. A binding site width of 5 nt optimally captures the SF3B1 crosslink
 904 events. Dot plot shows average ratio of crosslink events within binding sites of increasing widths (x-
 905 axis) over the mean background signal in flanking windows of the same size, indicating how much
 906 more signal occurs within the binding sites compared to their immediate surrounding. (D) SF3B1
 907 binding sites reproducibility across replicates. Upper panel shows overlaps of supported binding sites
 908 in the replicates, with threshold for sufficient coverage individually adjusted to the signal depth in
 909 each replicate (Busch et al. 2020). (E) Gene classes targeted by SF3B1 based on iCLIP. (F)
 910 Transcript regions of protein-coding genes targeted by the SF3B1 based on iCLIP. (G) Differential
 911 SF3B1 binding sites in K562-SF3B1^{K700E/wt} vs. K562-SF3B1^{wt/wt}. (H) Meta-profiles of SF3B1 (top)
 912 and U2AF2 (bottom) binding centered at branch point adenosine (left) and 3' splice site AG (right).
 913 (I) Distribution of distances between neighboring binding sites. (J) Density plot and heat map
 914 showing SF3B1 patterns in regions from three cluster types identified in Supplemental Fig. S35. (K)
 915 Distribution of the number of binding sites per region for each of the three clusters. (L) Distribution of
 916 the distance between SF3B1 binding region and closest splice site.

917 **Figure 7. SF3B1 promotes alternative proximal AG' usage.** (A) U2AF2 and SF3B1 binding to pre-
 918 mRNA based on the iCLIP signal. U2AF2 iCLIP was performed using K562-SF3B1^{wt/wt} cells and
 919 SF3B1 iCLIP was performed on K562-SF3B1^{wt/wt} and K562-SF3B1^{K700E/wt}. On the left panel the
 920 splice site distance is shown, followed by the iCLIP signal aligned to the more downstream AG used.
 921 Significance of the differential 3'AS usage between K562-SF3B1^{wt/wt} and K562-SF3B1^{K700E/wt} is
 922 shown as annotation bar on the right side of the iCLIP signal heatmap. (B) For every 3'AS from (A)
 923 the upstream AG PSI value is denoted for SF3B1^{wt/wt} (blue, left), SF3B1^{mut/wt} (red, center), and the
 924 difference between SF3B1^{mut/wt} and SF3B1^{wt/wt} (black, right). (C) Distribution of AG occurrence (top),

925 proportion of significantly alternatively used AG' (middle) and AG' scores among introns with
 926 significant (pink) or non-significant (violet) difference in usage between *SF3B1*^{mut/wt} and *SF3B1*^{wt/wt}.
 927 **(D)** Number of alternative AGs (AG') among regions with multiple AG's. AG's within 6 nt distance
 928 from AG were removed to avoid NAGNAG acceptor sites (Hiller et al. 2004). **(E)** Distance of AG's
 929 and BPs from AG among non-significant (grey) and significant (red) differentially splicing events.
 930 **(F)** Splice site strength of canonical and alternative AGs as calculated with MaxEnt score (Yeo and
 931 Burge 2003). For figures (A-F) only 3'ASs with the following features were used: i) placed
 932 chromosome scaffold; ii) classical AGs; iii) an intron became shorter in the mutant (the use of
 933 upstream alternative AG'); and iv) overlap with an iCLIP crosslink. The more used AG in *SF3B1*^{wt/wt}
 934 was set as canonical AG.

935 **References**

- 936 Agafonov DE, Deckert J, Wolf E, Odenwalder P, Bessonov S, Will CL, Urlaub H, Luhrmann R.
 937 2011. Semiquantitative Proteomic Analysis of the Human Spliceosome via a Novel Two-
 938 Dimensional Gel Electrophoresis Method. *Mol Cell Biol* **31**: 2667–2682
 939
- 940 Alsafadi S, Dayot S, Tarin M, Houy A, Bellanger D, Cornella M, Wassef M, Waterfall JJ, Lehnert E,
 941 Roman-Roman S, et al. 2021. Genetic alterations of SUGP1 mimic mutant-SF3B1 splice pattern
 942 in lung adenocarcinoma and other cancers. *Oncogene* **40**: 85–96.
 943
- 944 Alsafadi S, Houy A, Battistella A, Popova T, Wassef M, Henry E, Tirode F, Constantinou A, Piperno-
 945 Neumann S, Roman-Roman S, et al. 2016. Cancer-associated SF3B1 mutations affect alternative
 946 splicing by promoting alternative branchpoint usage. *Nat Commun* **7**: 10615.
 947
- 948 Ameijeiras-Alonso J, Crujeiras RM, Rodriguez-Casal A. 2021. multimode: An R Package for Mode
 949 Assessment. *J Stat Softw* **97**: 1–32.
 950
- 951 Andrews S. 2010. FASTQC - a quality control tool for high throughput sequence data.
 952 www.bioinformatics.babraham.ac.uk/projects/fastqc.
 953
- 954 Benjamini Y, Hochberg Y. 1995. Controlling the False Discovery Rate: A Practical and Powerful
 955 Approach to Multiple Testing. *Journal of the Royal Statistical Society Series B (Methodological)*
 956 **57**: 289–300.
 957

- 958 Bland P, Saville H, Wai PT, Curnow L, Muirhead G, Nieminuszczy J, Ravindran N, John MB,
 959 Hedayat S, Barker HE, et al. 2023. SF3B1 hotspot mutations confer sensitivity to PARP
 960 inhibition by eliciting a defective replication stress response. *Nat Genet* **55**: 1311–1323.
 961
- 962 Bradley RK, Anczuków O. 2023. RNA splicing dysregulation and the hallmarks of cancer. *Nat Rev*
 963 *Cancer* **23**: 135–155.
 964
- 965 Brüggemann M, Zarnack K. Binding site definition based on iCLIP data.
 966 <https://bioconductor.org/packages/release/bioc/html/BindingSiteFinder.html>.
 967
- 968 Busch A, Brüggemann M, Ebersberger S, Zarnack K. 2020. iCLIP data analysis: A complete pipeline
 969 from sequencing reads to RBP binding sites. *Methods* **178**: 49–62.
 970
- 971 Canbezdi C, Tarin M, Houy A, Bellanger D, Popova T, Stern M-H, Roman-Roman S, Alsafadi S.
 972 2021. Functional and conformational impact of cancer-associated SF3B1 mutations depends on
 973 the position and the charge of amino acid substitution. *Comput Struct Biotechnol J* **19**: 1361–
 974 1370.
 975
- 976 Carrocci TJ, Zoerner DM, Paulson JC, Hoskins AA. 2017. SF3b1 mutations associated with
 977 myelodysplastic syndromes alter the fidelity of branchsite selection in yeast. *Nucleic Acids Res*
 978 **45**: 4837–4852.
 979
- 980 Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM. 2012. A
 981 program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff.
 982 *Fly (Austin)* **6**: 80–92.
 983
- 984 Cocquet J, Chong A, Zhang G, Veitia RA. 2006. Reverse transcriptase template switching and false
 985 alternative transcripts. *Genomics* **88**: 127–131.
 986
- 987 Cortés-López M, Chamely P, Hawkins AG, Stanley RF, Swett AD, Ganesan S, Mouhieddine TH, Dai
 988 X, Kluegel L, Chen C, et al. 2023. Single-cell multi-omics defines the cell-type-specific impact
 989 of splicing aberrations in human hematopoietic clonal outgrowths. *Cell Stem Cell*.
 990
- 991 Cretu C, Schmitzová J, Ponce-Salvatierra A, Dybkov O, De Laurentiis EI, Sharma K, Will CL,
 992 Urlaub H, Lührmann R, Pena V. 2016. Molecular Architecture of SF3b and Structural
 993 Consequences of Its Cancer-Related Mutations. *Mol Cell* **64**: 307–319.
 994
- 995 Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane T,
 996 McCarthy SA, Davies RM, et al. 2021. Twelve years of SAMtools and BCFtools. *Gigascience*
 997 **10**: giab008.
 998
- 999 Darman RB, Seiler M, Agrawal AA, Lim KH, Peng S, Aird D, Bailey SL, Bhavsar EB, Chan B, Colla
 1000 S, et al. 2015. Cancer-Associated SF3B1 Hotspot Mutations Induce Cryptic 3' Splice Site
 1001 Selection through Use of a Different Branch Point. *Cell Rep* **13**: 1033–1045.

- 1002
- 1003 DeBoever C, Ghia EM, Shepard PJ, Rassenti L, Barrett CL, Jepsen K, Jamieson CHM, Carson D,
1004 Kipps TJ, Frazer KA. 2015. Transcriptome sequencing reveals potential mechanism of cryptic 3'
1005 splice site selection in SF3B1-mutated cancers. *PLoS Comput Biol* **11**: e1004105.
- 1006
- 1007 Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras TR.
1008 2013. STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics* **29**: 15–21.
- 1009
- 1010 Dolatshad H, Pellagatti A, Liberante FG, Llorian M, Repapi E, Steeples V, Roy S, Scifo L, Armstrong
1011 RN, Shaw J, et al. 2016. Cryptic splicing events in the iron transporter ABCB7 and other key
1012 target genes in SF3B1-mutant myelodysplastic syndromes. *Leukemia* **30**: 2322–2331.
- 1013
- 1014 Ester M, Kriegel H-P, Sander J, Xu X. 1996. A density-based algorithm for discovering clusters in
1015 large spatial databases with noise. pp. 226–231, AAAI Press.
- 1016
- 1017 Frankish A, Diekhans M, Ferreira A-M, Johnson R, Jungreis I, Loveland J, Mudge JM, Sisu C,
1018 Wright J, Armstrong J, et al. 2019. GENCODE reference annotation for the human and mouse
1019 genomes. *Nucleic Acids Res* **47**: D766–D773.
- 1020
- 1021 Frankish A, Diekhans M, Jungreis I, Lagarde J, Loveland JE, Mudge JM, Sisu C, Wright JC,
1022 Armstrong J, Barnes I, et al. 2021. GENCODE 2021. *Nucleic Acids Res* **49**: D916–D923.
- 1023
- 1024 Fu C, Donovan WP, Shikapwashya-Hasser O, Ye X, Cole RH. 2014. Hot Fusion: an efficient method
1025 to clone multiple DNA fragments as well as inverted repeats without ligase. *PLoS One* **9**:
1026 e115318.
- 1027
- 1028 Gozani O, Feld R, Reed R. 1996. Evidence that sequence-independent binding of highly conserved
1029 U2 snRNP proteins upstream of the branch site is required for assembly of spliceosomal
1030 complex A. *Genes Dev* **10**: 233–243.
- 1031
- 1032 Gozani O, Potashkin J, Reed R. 1998. A Potential Role for U2AF-SAP 155 Interactions in Recruiting
1033 U2 snRNP to the Branch Site. *Mol Cell Biol* **18**: 4752–4760.
- 1034
- 1035 Hahsler M, Piekenbrock M, Doran D. 2019. dbscan: Fast Density-Based Clustering with R. *J Stat*
1036 *Softw* **91**: 1–30.
- 1037
- 1038 Hauer C, Curk T, Anders S, Schwarzl T, Alleaume A-M, Sieber J, Hollerer I, Bhuvanagiri M, Huber
1039 W, Hentze MW, et al. 2015. Improved binding site assignment by high-resolution mapping of
1040 RNA–protein interactions using iCLIP. *Nat Commun* **6**: 7921.
- 1041
- 1042 Heyl F, Backofen R. 2021. StoatyDive: Evaluation and classification of peak profiles for sequencing
1043 data. *Gigascience* **10**: giab045.
- 1044

- 1045 Hiller M, Huse K, Szafranski K, Jahn N, Hampe J, Schreiber S, Backofen R, Platzer M. 2004.
1046 Widespread occurrence of alternative splicing at NAGNAG acceptors contributes to proteome
1047 plasticity. *Nat Genet* **36**: 1255–1257.
- 1048
- 1049 Hooper JE. 2014. A survey of software for genome-wide discovery of differential splicing in RNA-
1050 Seq data. *Hum Genomics* **8**: 3.
- 1051
- 1052 Huang L, Lou C-H, Chan W, Shum EY, Shao A, Stone E, Karam R, Song H-W, Wilkinson MF. 2011.
1053 RNA Homeostasis Governed by Cell Type-Specific and Branched Feedback Loops Acting on
1054 NMD. *Mol Cell* **43**: 950–961.
- 1055
- 1056 Inoue D, Chew G-L, Liu B, Michel BC, Pangallo J, D’Avino AR, Hitchman T, North K, Lee SC-W,
1057 Bitner L, et al. 2019. Spliceosomal disruption of the non-canonical BAF complex in cancer.
1058 *Nature* **574**: 432–436.
- 1059
- 1060 Kesarwani AK, Ramirez O, Gupta AK, Yang X, Murthy T, Minella AC, Pillai MM. 2017. Cancer-
1061 associated SF3B1 mutants recognize otherwise inaccessible cryptic 3’ splice sites within RNA
1062 secondary structures. *Oncogene* **36**: 1123–1133.
- 1063
- 1064 König J, Zarnack K, Rot G, Curk T, Kayikci M, Zupan B, Turner DJ, Luscombe NM, Ule J. 2010.
1065 iCLIP reveals the function of hnRNP particles in splicing at individual nucleotide resolution.
1066 *Nat Struct Mol Biol* **17**: 909–915.
- 1067
- 1068 Krakau S, Richard H, Marsico A. 2017. PureCLIP: capturing target-specific protein–RNA interaction
1069 footprints from single-nucleotide CLIP-seq data. *Genome Biol* **18**: 240.
- 1070
- 1071 Landau DA, Tausch E, Taylor-Weiner AN, Stewart C, Reiter JG, Bahlo J, Kluth S, Bozic I, Lawrence
1072 M, Böttcher S, et al. 2015. Mutations driving CLL and their evolution in progression and
1073 relapse. *Nature*.
- 1074
- 1075 Lee SC-W, Dvinge H, Kim E, Cho H, Micol J-B, Chung YR, Durham BH, Yoshimi A, Kim YJ,
1076 Thomas M, et al. 2016. Modulation of splicing catalysis for therapeutic targeting of leukemia
1077 with mutations in genes encoding spliceosomal proteins. *Nat Med* **22**: 672–678.
- 1078
- 1079 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**: 3094–3100.
- 1080
- 1081 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
1082 *Bioinformatics* **25**: 1754–1760.
- 1083
- 1084 Lienhard M, van den Beucken T, Timmermann B, Hochradel M, Börno S, Caiment F, Vingron M,
1085 Herwig R. 2023. IsoTools: a flexible workflow for long-read transcriptome sequencing analysis.
1086 *Bioinformatics* **39**: btad364.
- 1087

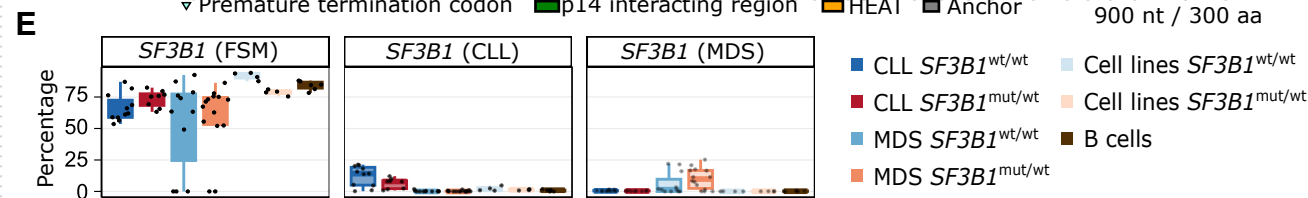
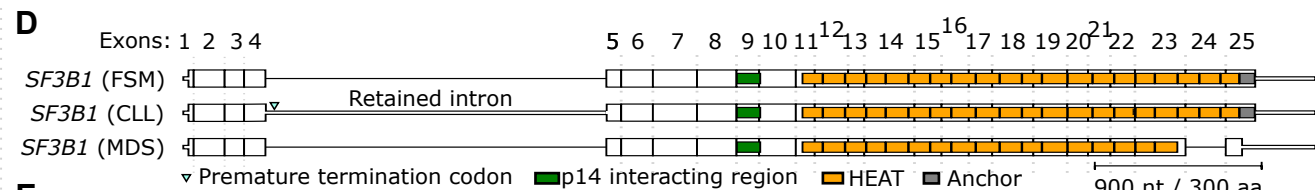
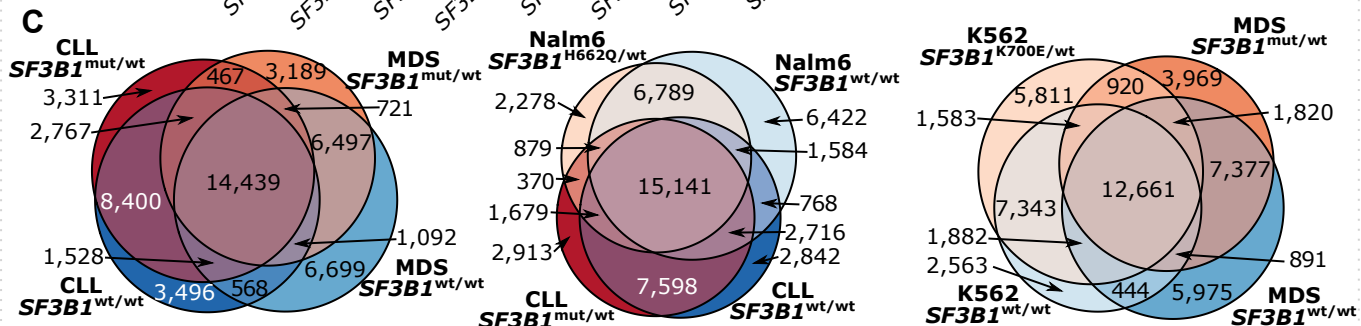
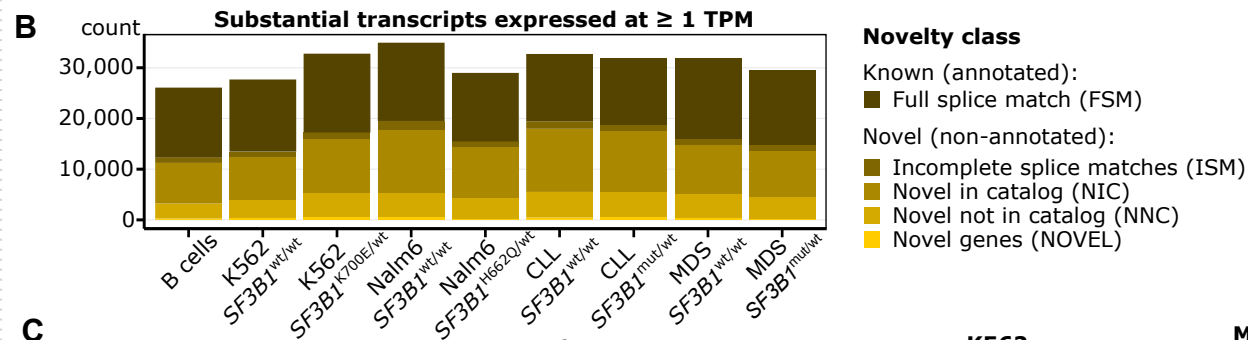
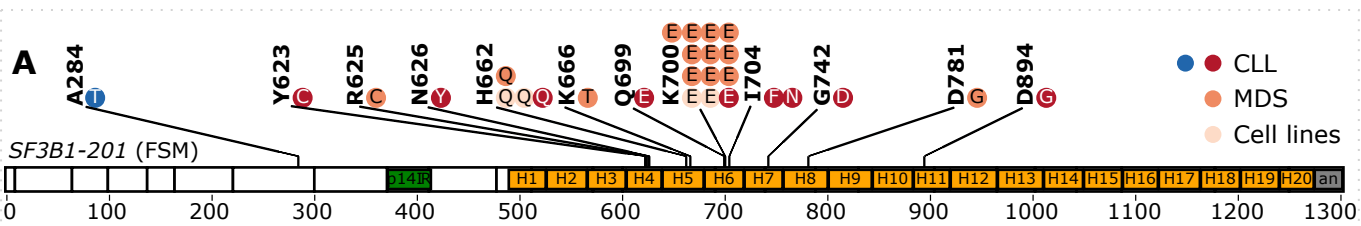
- 1088 Liu Z, Yoshimi A, Wang J, Cho H, Chun-Wei Lee S, Ki M, Bitner L, Chu T, Shah H, Liu B, et al.
1089 2020. Mutations in the RNA Splicing Factor SF3B1 Promote Tumorigenesis through MYC
1090 Stabilization. *Cancer Discov* **10**: 806–821.
- 1091
- 1092 Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for RNA-
1093 seq data with DESeq2. *Genome Biol* **15**: 550.
- 1094
- 1095 Malcovati L, Karimi M, Papaemmanuil E, Ambaglio I, Jädersten M, Jansson M, Elena C, Galli A,
1096 Walldin G, Della Porta MG, et al. 2015. SF3B1 mutation identifies a distinct subset of
1097 myelodysplastic syndrome with ring sideroblasts. *Blood* **126**: 233–241.
- 1098
- 1099 Malcovati L, Stevenson K, Papaemmanuil E, Neuberg D, Bejar R, Boultonwood J, Bowen DT,
1100 Campbell PJ, Ebert BL, Fenaux P, et al. 2020. SF3B1-mutant MDS as a distinct disease subtype:
1101 a proposal from the International Working Group for the Prognosis of MDS. *Blood* **136**: 157–
1102 170.
- 1103
- 1104 McInnes L, Healy J, Melville J. 2018. UMAP: Uniform Manifold Approximation and Projection for
1105 Dimension Reduction. *arXiv* 1802.03426.
- 1106
- 1107 Mercer TR, Clark MB, Andersen SB, Brunck ME, Haerty W, Crawford J, Taft RJ, Nielsen LK,
1108 Dinger ME, Mattick JS. 2015. Genome-wide discovery of human splicing branchpoints.
1109 *Genome Res* **25**: 290–303.
- 1110
- 1111 Mistry J, Chuguransky S, Williams L, Qureshi M, Salazar GA, Sonnhammer ELL, Tosatto SCE,
1112 Paladin L, Raj S, Richardson LJ, et al. 2021. Pfam: The protein families database in 2021.
1113 *Nucleic Acids Res* **49**: D412–D419.
- 1114
- 1115 Obeng EA, Chappell RJ, Seiler M, Chen MC, Campagna DR, Schmidt PJ, Schneider RK, Lord AM,
1116 Wang L, Gambe RG, et al. 2016. Physiologic Expression of Sf3b1(K700E) Causes Impaired
1117 Erythropoiesis, Aberrant Splicing, and Sensitivity to Therapeutic Spliceosome Modulation.
1118 *Cancer Cell* **30**: 404–417.
- 1119
- 1120 PacBio 2020. 2021. IsoSeq v3: Scalable De Novo Isoform Discovery [Online].
1121 <https://github.com/PacificBiosciences/IsoSeq>.
- 1122
- 1123 Papaemmanuil E, Cazzola M, Boultonwood J, Malcovati L, Vyas P, Bowen D, Pellagatti A, Wainscoat
1124 JS, Hellstrom-Lindberg E, Gambacorti-Passerini C, et al. 2011. Somatic SF3B1 Mutation in
1125 Myelodysplasia with Ring Sideroblasts. *New England Journal of Medicine* **365**: 1384–1395.
- 1126
- 1127 Pellagatti A, Armstrong RN, Steeples V, Sharma E, Repapi E, Singh S, Sanchi A, Radujkovic A,
1128 Horn P, Dolatshad H, et al. 2018. Impact of spliceosome mutations on RNA splicing in
1129 myelodysplasia: dysregulated genes/pathways and clinical associations. *Blood* **132**: 1225–1240.
- 1130

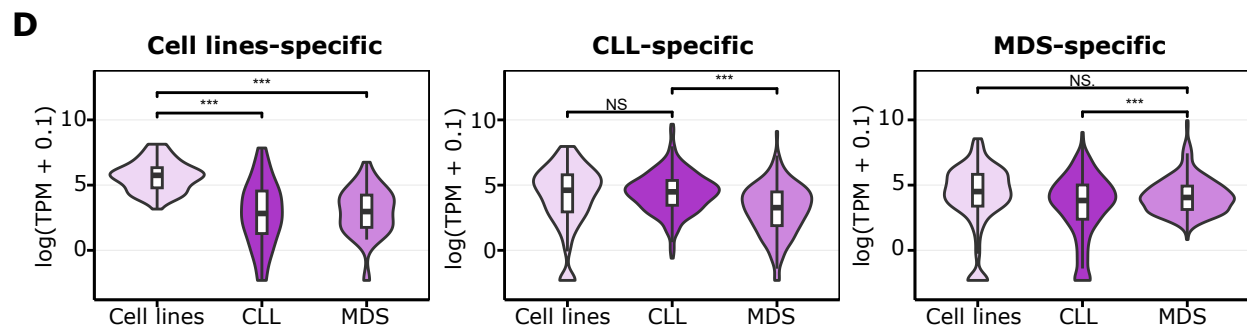
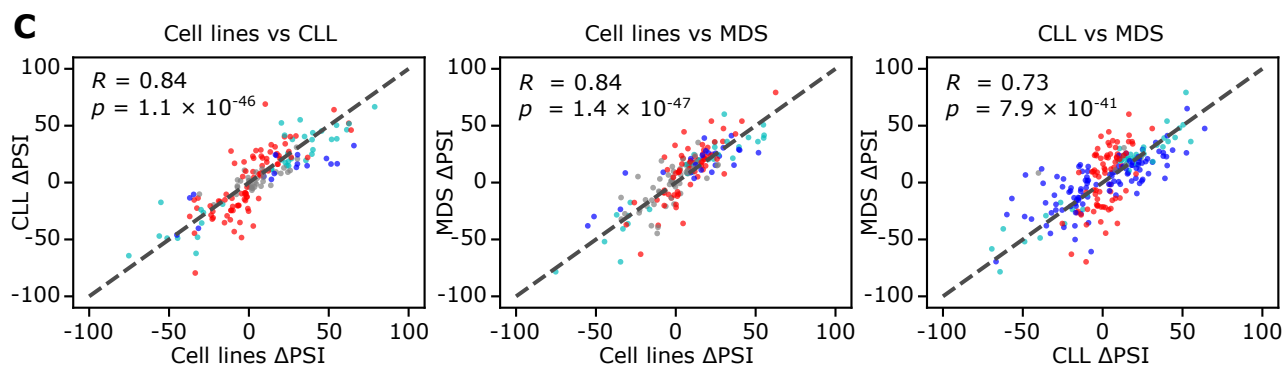
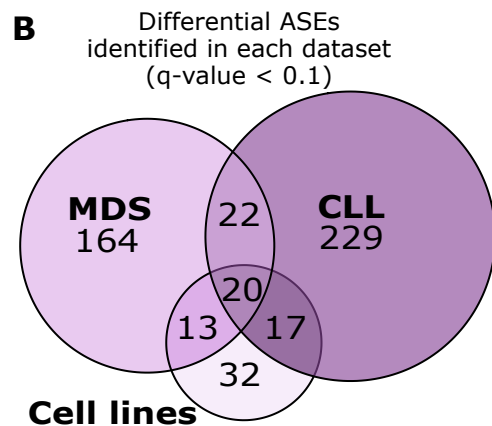
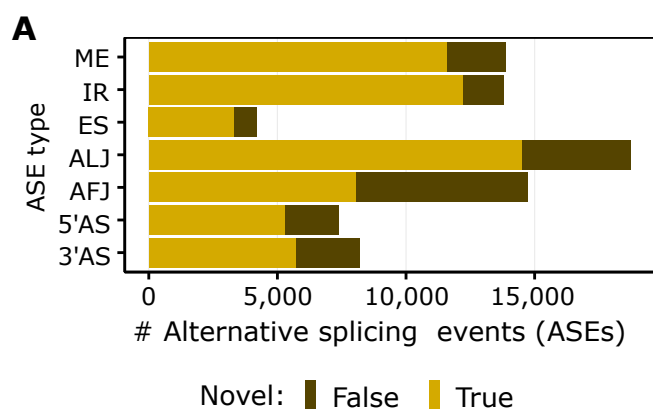
- 1131 Peng H, Jabbari JS, Tian L, Chua CC, Anstee NS, Amin N, Wei AH, Davidson NM, Roberts AW,
 1132 Huang DCS, et al. 2024. Single-cell Rapid Capture Hybridization sequencing (scRaCH-seq) to
 1133 reliably detect isoform usage and coding mutations in targeted genes at a single-cell level.
 1134 *bioRxiv* 2024.01.30.577942.
- 1135
- 1136 Pineda JMB, Bradley RK. 2018. Most human introns are recognized via multiple and tissue-specific
 1137 branchpoints. *Genes Dev* **32**: 577–591.
- 1138
- 1139 Porter DF, Miao W, Yang X, Goda GA, Ji AL, Donohue LKH, Aleman MM, Dominguez D, Khavari
 1140 PA. 2021. easyCLIP analysis of RNA-protein interactions incorporating absolute quantification.
 1141 *Nat Commun* **12**: 1569.
- 1142
- 1143 Quesada V, Conde L, Villamor N, Ordóñez GR, Jares P, Bassaganyas L, Ramsay AJ, Beà S, Pinyol
 1144 M, Martínez-Trillos A, et al. 2012. Exome sequencing identifies recurrent mutations of the
 1145 splicing factor *SF3B1* gene in chronic lymphocytic leukemia. *Nat Genet*.
- 1146
- 1147 R Core Team. 2017. R: A Language and Environment for Statistical Computing. *R Foundation for*
 1148 *Statistical Computing, Vienna, Austria* **0**: {ISBN} 3-900051-07-0.
- 1149
- 1150 Rehrauer H, Opitz L, Tan G, Sieverling L, Schlapbach R. 2013. Blind spots of quantitative RNA-seq:
 1151 the limits for assessing abundance, differential expression, and isoform switching. *BMC*
 1152 *Bioinformatics* **14**: 370.
- 1153
- 1154 Roehr JT, Dieterich C, Reinert K. 2017. Flexbar 3.0 – SIMD and multicore parallelization.
 1155 *Bioinformatics* **33**: 2941–2942.
- 1156
- 1157 Rosenquist R, Ghia P, Hadzidimitriou A, Sutton LA, Agathangelidis A, Baliakas P, Darzentas N,
 1158 Giudicelli V, Lefranc MP, Langerak AW, et al. 2017. Immunoglobulin gene sequence analysis
 1159 in chronic lymphocytic leukemia: Updated ERIC recommendations. *Leukemia*.
- 1160
- 1161 Rossi D, Bruscazzin A, Spina V, Rasi S, Khiabani H, Messina M, Fangazio M, Vaisitti T, Monti S,
 1162 Chiaretti S, et al. 2011. Mutations of the *SF3B1* splicing factor in chronic lymphocytic leukemia:
 1163 Association with progression and fludarabine-refractoriness. *Blood*.
- 1164
- 1165 Schmitzová J, Cretu C, Dienemann C, Urlaub H, Pena V. 2023. Structural basis of catalytic activation
 1166 in human splicing. *Nature* **617**: 842–850.
- 1167
- 1168 Seiler M, Peng S, Agrawal AA, Palacino J, Teng T, Zhu P, Smith PG, Buonamici S, Yu L. 2018.
 1169 Somatic Mutational Landscape of Splicing Factor Genes and Their Functional Consequences
 1170 across 33 Cancer Types. *Cell Rep* **23**: 282-296.e4.
- 1171
- 1172 Shen S, Park JW, Lu Z, Lin L, Henry MD, Wu YN, Zhou Q, Xing Y. 2014. rMATS: Robust and
 1173 flexible detection of differential alternative splicing from replicate RNA-Seq data. *Proceedings*
 1174 *of the National Academy of Sciences* **111**: E5593 LP-E5601.

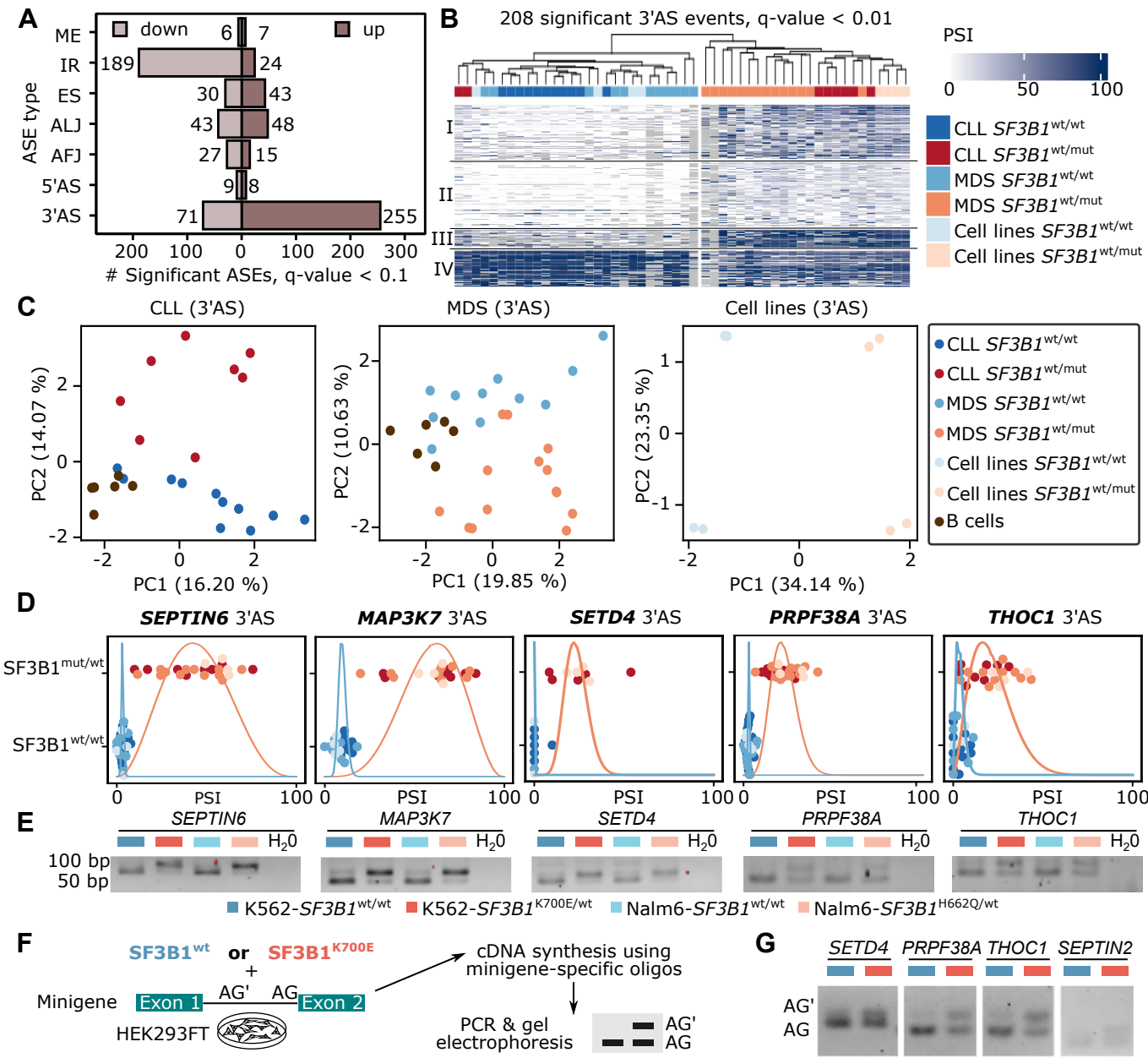
- 1175
- 1176 Shiozawa Y, Malcovati L, Galli A, Sato-Otsubo A, Kataoka K, Sato Y, Watatani Y, Suzuki H,
1177 Yoshizato T, Yoshida K, et al. 2018. Aberrant splicing and defective mRNA production induced
1178 by somatic spliceosome mutations in myelodysplasia. *Nat Commun*.
- 1179
- 1180 Signal B, Gloss BS, Dinger ME, Mercer TR. 2018. Machine learning annotation of human
1181 branchpoints. *Bioinformatics* **34**: 920–927.
- 1182
- 1183 Smith T, Heger A, Sudbery I. 2017. UMI-tools: modeling sequencing errors in Unique Molecular
1184 Identifiers to improve quantification accuracy. *Genome Res* **27**: 491–499.
- 1185
- 1186 Sutandy FXR, Hildebrandt A, König J. 2016. Profiling the Binding Sites of RNA-Binding Proteins
1187 with Nucleotide Resolution Using iCLIP. In *Post-Transcriptional Gene Regulation* (ed. E.
1188 Dassi), pp. 175–195, Springer New York, New York, NY.
- 1189
- 1190 Tang AD, Soulette CM, van Baren MJ, Hart K, Hrabeta-Robinson E, Wu CJ, Brooks AN. 2020. Full-
1191 length transcript characterization of SF3B1 mutation in chronic lymphocytic leukemia reveals
1192 downregulation of retained introns. *Nat Commun*.
- 1193
- 1194 Tang Q, Rodriguez-Santiago S, Wang J, Pu J, Yuste A, Gupta V, Moldón A, Xu Y-Z, Query CC.
1195 2016. SF3B1/Hsh155 HEAT motif mutations affect interaction with the spliceosomal ATPase
1196 Prp5, resulting in altered branch site selectivity in pre-mRNA splicing. *Genes Dev* **30**: 2710–
1197 2723.
- 1198
- 1199 Tarantola S. 2008. European innovation scoreboard: Strategies to measure country progress over time.
1200 *JRC scientific and Technical Reports, EUR* **23526**.
- 1201
- 1202 Tardaguila M, de la Fuente L, Marti C, Pereira C, Pardo-Palacios FJ, del Risco H, Ferrell M, Mellado
1203 M, Macchietto M, Verheggen K, et al. 2018. SQANTI: extensive characterization of long-read
1204 transcript sequences for quality control in full-length transcriptome identification and
1205 quantification. *Genome Res*.
- 1206
- 1207 Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, Burovski E, Peterson
1208 P, Weckesser W, Bright J, et al. 2020. SciPy 1.0: fundamental algorithms for scientific
1209 computing in Python. *Nat Methods* **17**: 261–272.
- 1210
- 1211 Wan Y, Wu CJ. 2013. SF3B1 mutations in chronic lymphocytic leukemia. *Blood*.
- 1212
- 1213 Wang C, Chua K, Seghezzi W, Lees E, Gozani O, Reed R. 1998. Phosphorylation of spliceosomal
1214 protein SAP 155 coupled with splicing catalysis. *Genes Dev* **12**: 1409–1414.
- 1215
- 1216 Wang L, Brooks AN, Fan J, Wan Y, Gambe R, Li S, Hergert S, Yin S, Freeman SS, Levin JZ, et al.
1217 2016. Transcriptomic Characterization of SF3B1 Mutation Reveals Its Pleiotropic Effects in
1218 Chronic Lymphocytic Leukemia. *Cancer Cell* **30**: 750–763.

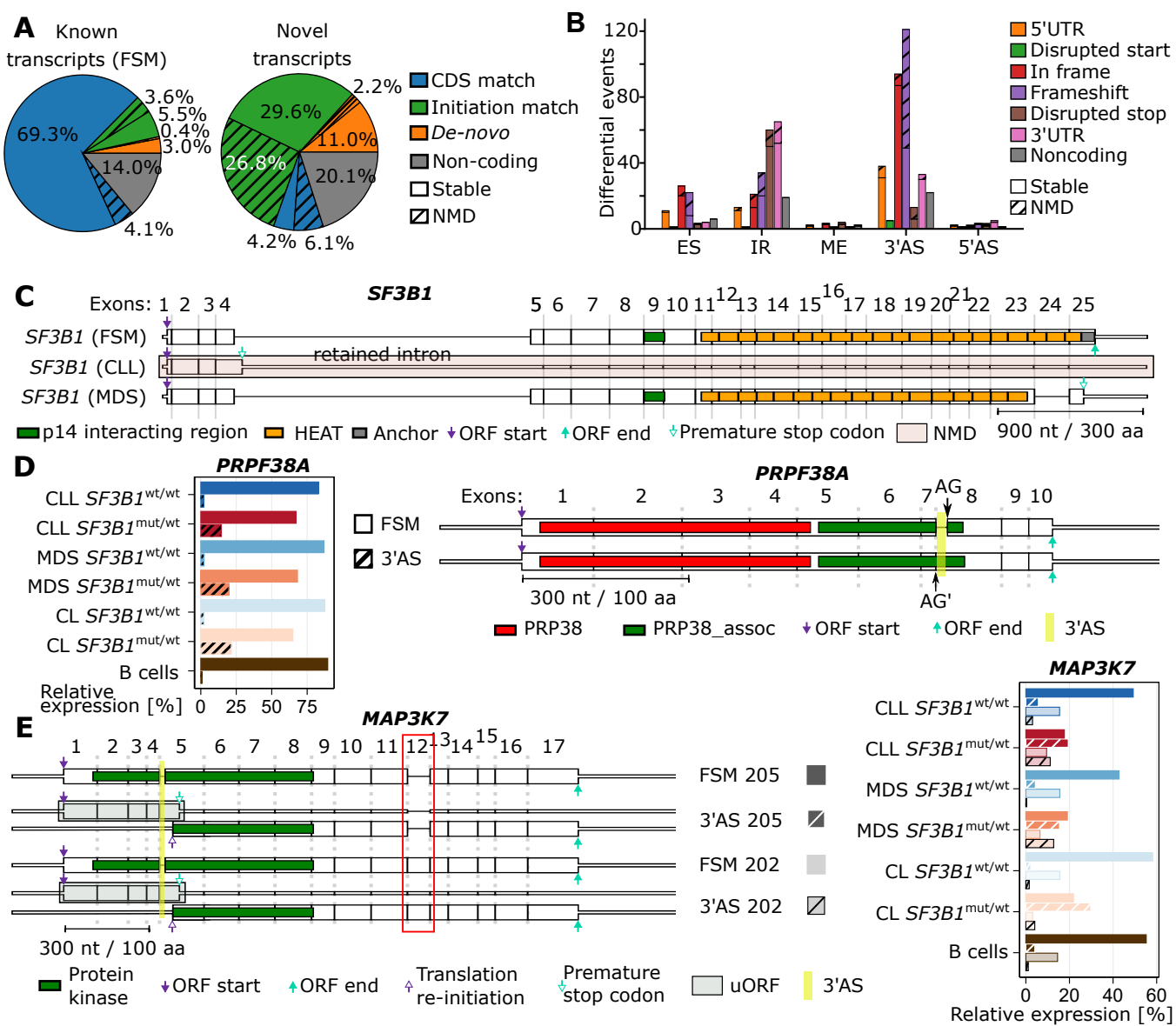
- 1219
- 1220 Yang F, Bian T, Zhan X, Chen Z, Xing Z, Larsen NA, Zhang X, Shi Y. 2023. Mechanisms of the
1221 RNA helicases DDX42 and DDX46 in human U2 snRNP assembly. *Nat Commun* **14**: 897.
- 1222
- 1223 Yang H, Beutler B, Zhang D. 2021. Emerging roles of spliceosome in cancer and immunity. *Protein*
1224 *Cell*.
- 1225
- 1226 Yeo G, Burge CB. 2003. Maximum Entropy Modeling of Short Sequence Motifs with Applications to
1227 RNA Splicing Signals. In *Proceedings of the Seventh Annual International Conference on*
1228 *Research in Computational Molecular Biology, RECOMB '03*, pp. 322–331, Association for
1229 Computing Machinery, New York, NY, USA.
- 1230
- 1231 Yokoi A, Kotake Y, Takahashi K, Kadowaki T, Matsumoto Y, Minoshima Y, Sugi NH, Sagane K,
1232 Hamaguchi M, Iwata M, et al. 2011. Biological validation that SF3b is a target of the antitumor
1233 macrolide pladienolide. *FEBS J* **278**: 4870–4880.
- 1234
- 1235 Zarnack K, König J, Tajnik M, Martincorena I, Eustermann S, Stévant I, Reyes A, Anders S,
1236 Luscombe NM, Ule J. 2013. Direct competition between hnRNP C and U2AF65 protects the
1237 transcriptome from the exonization of Alu elements. *Cell* **152**: 453–466.
- 1238
- 1239 Zhang C, Zhang B, Lin L-L, Zhao S. 2017. Evaluation and comparison of computational tools for
1240 RNA-seq isoform quantification. *BMC Genomics* **18**: 583.
- 1241
- 1242 Zhang J, Ali AM, Lieu YK, Liu Z, Gao J, Rabadan R, Raza A, Mukherjee S, Manley JL. 2019.
1243 Disease-Causing Mutations in SF3B1 Alter Splicing by Disrupting Interaction with SUGP1. *Mol*
1244 *Cell* **76**: 82-95.e7.
- 1245
- 1246 Zhang J, Huang J, Xu K, Xing P, Huang Y, Liu Z, Tong L, Manley JL. 2022. DHX15 is involved in
1247 SUGP1-mediated RNA missplicing by mutant SF3B1 in cancer. *Proceedings of the National*
1248 *Academy of Sciences* **119**: e2216712119.
- 1249
- 1250 Zhang J, Xie J, Huang J, Liu X, Xu R, Tholen J, Galej WP, Tong L, Manley JL, Liu Z. 2023.
1251 Characterization of the SF3B1–SUGP1 interface reveals how numerous cancer mutations cause
1252 mRNA missplicing. *Genes Dev*.
- 1253
- 1254 Zhang X, Zhan X, Bian T, Yang F, Li P, Lu Y, Xing Z, Fan R, Zhang QC, Shi Y. 2024. Structural
1255 insights into branch site proofreading by human spliceosome. *Nat Struct Mol Biol* **31**: 835–845.
- 1256
- 1257 Zhang Z, Rigo N, Dybkov O, Fourmann J-B, Will CL, Kumar V, Urlaub H, Stark H, Lührmann R.
1258 2021. Structural insights into how Prp5 proofreads the pre-mRNA branch site. *Nature* **596**: 296–
1259 300.
- 1260

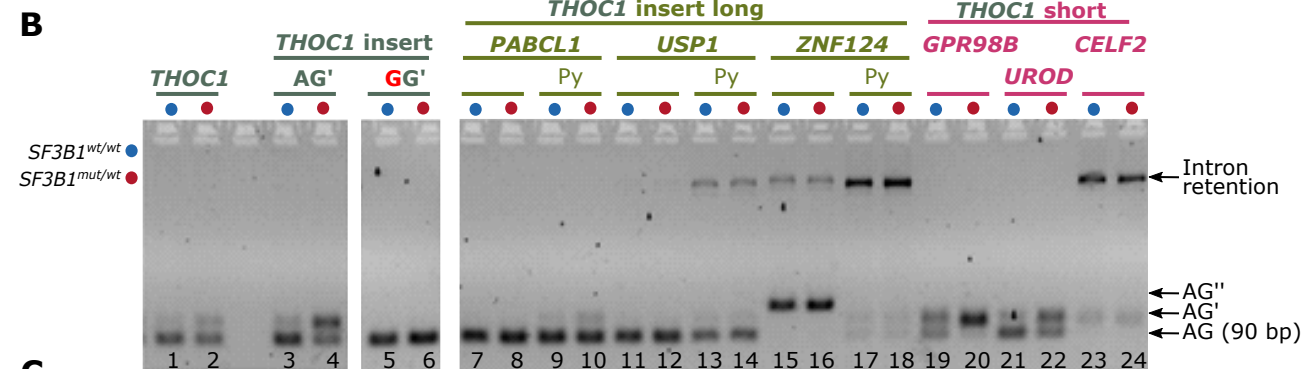
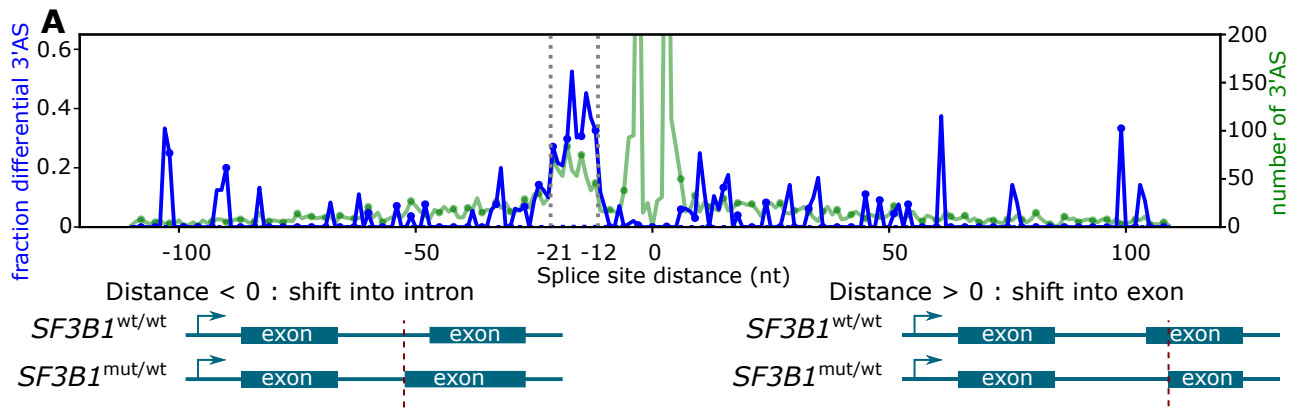
- 1261 Zhao B, Li Z, Qian R, Liu G, Fan M, Liang Z, Hu X, Wan Y. 2022. Cancer-associated mutations in
1262 SF3B1 disrupt the interaction between SF3B1 and DDX42. *The Journal of Biochemistry* **172**:
1263 117–126.
- 1264
- 1265 Zhou Z, Gong Q, Wang Y, Li M, Wang L, Ding H, Li P. 2020. The biological function and clinical
1266 significance of SF3B1 mutations in cancer. *Biomark Res* **8**: 38.
- 1267
- 1268 Zuallaert J, Godin F, Kim M, Soete A, Saeys Y, De Neve W. 2018. SpliceRover: interpretable
1269 convolutional neural networks for improved splice site prediction. *Bioinformatics* **34**: 4180–
1270 4188.
- 1271











C

	<i>THOC1</i> : insertions between AG and AG'										
	<i>THOC1</i>	<i>THOC1</i> AG'>GG	<i>PABCL1</i>	<i>PABCL1</i> (Py)	<i>USP1</i>	<i>USP1</i> (Py)	<i>ZNF124</i>	<i>ZNF124</i> (Py)	<i>GPR98B</i>	<i>UROD</i>	<i>CELF2</i>
AG score	0.988	0.996	0.986	0.995	0.987	0.979	0.719	0.498	0.992	0.995	0.301
AG' score	0.704	/	0.848	0.507	0.210	0.660	0.952	0.928	0.665	0.857	0.902
AG-AG' distance	21	/	58	24	58	23	54	22	19	19	18

