**A systems view of spliceosomal assembly and branchpoints with iCLIP**

Michael Briese[1,2]*, Nejc Haberman[3,4]*, Christopher R. Sibley[1,4,5,6]*, Rupert Faraway[3,4], Andrea S. Elser[3,4], Anob M. Chakrabarti[3,7], Zhen Wang[1], Julian König[1,8], David Perera[9], Vihandha O. Wickramasinghe[9,10], Ashok R. Venkitaraman[9], Nicholas M. Luscombe[3,7,11], Luciano Saieva[12,13], Livio Pellizzoni[12], Christopher W.J. Smith[14], Tomaž Curk[15], Jernej Ule[1,3,4]§

[1]MRC Laboratory of Molecular Biology, Cambridge, UK

[2]Institute of Clinical Neurobiology, University of Wuerzburg, Wuerzburg, Germany

[3]The Francis Crick Institute, London, UK

[4]Department of Neuromuscular Disease, UCL Institute of Neurology, London, UK

[5]Division of Brain Sciences, Department of Medicine, Imperial College London, London, UK

[6]Institute of Quantitative Biology, Biochemistry and Biotechnology, Edinburgh University, UK

[7]Department of Genetics, Environment and Evolution, UCL Genetics Institute, London, UK

[8]Institute of Molecular Biology (IMB) GmbH, Mainz, Germany

[9]MRC Cancer Unit at the University of Cambridge, Cambridge, UK

[10]RNA Biology and Cancer Laboratory, Peter MacCallum Cancer Centre, Melbourne, Australia

[11]Okinawa Institute of Science & Technology Graduate University, Okinawa, Japan

[12]Center for Motor Neuron Biology and Disease, Department of Pathology and Cell Biology, Columbia University, New York, NY, USA

[13]Institute of Neuroscience, Newcastle University, Newcastle upon Tyne, UK

[14]Department of Biochemistry, University of Cambridge, Cambridge, UK

[15]Faculty of Computer and Information Science, University of Ljubljana, Ljubljana, Slovenia

**Equal contributions:**
Michael Briese, Nejc Haberman and Christopher R Sibley contributed equally to this work.

**Corresponding author:**
§Jernej Ule:          jernej.ule@crick.ac.uk

**Abstract**

Studies of spliceosomal interactions are challenging due to their dynamic nature. Here we employed spliceosome iCLIP, which immunoprecipitates SmB along with snRNPs and auxiliary RNA binding proteins (RBPs), to map spliceosome engagement with pre-mRNAs in human cell lines. This revealed seven peaks of spliceosomal crosslinking around branchpoints (BPs) and splice sites. We identified RBPs that crosslink to each peak, including known and candidate splicing factors. Moreover, we detected use of over 40,000 BPs with strong sequence consensus and structural accessibility, which align well to nearby crosslinking peaks. We show how the position and strength of BPs affect the crosslinking patterns of spliceosomal factors, which bind more efficiently upstream of strong or proximally located BPs, and downstream of weak or distally located BPs. These insights exemplify spliceosome iCLIP as a broadly applicable method for transcriptomic studies of splicing mechanisms.

## Introduction

Splicing is a multi-step process in which small nuclear ribonucleoprotein particles (snRNPs) and associated splicing factors bind at specific positions around intron boundaries in order to assemble an active spliceosome through a series of remodeling steps. The splicing reactions are coordinated by dynamic pairings between different snRNAs, between snRNAs and pre-mRNA, and by protein-RNA contacts[1]. Spliceosome assembly begins with ATP-independent binding of U1 snRNP at the 5' splice site (ss), and of U2 small nuclear RNA auxiliary factors 1 and 2 (U2AF1 and U2AF2, also known as U2AF35 and U2AF65) to the 3'ss. ATP-dependent remodeling then leads to the formation of complex A in which U2 snRNP contacts the branchpoint (BP), stabilized through interactions with the U2AF and U2 snRNP splicing factor 3 (SF3a and SF3b) complex. Next, U4/U6 and U5 snRNPs are recruited to form complex B. The actions of many RNA helicases and pre-mRNA processing factor 8 (PRPF8) then facilitate rearrangements of snRNP interactions and establishment of the catalytically competent $B_{act}$ and C complexes. These catalyze the two trans-esterification reactions leading to lariat formation, intron removal and exon ligation[2].

Transcriptome-wide studies of splicing reactions are valuable to unravel the multi-component and dynamic assembly of the spliceosome on the pre-mRNA substrate[3-5]. Accordingly, "spliceosome profiling" has been developed through affinity purification of the tagged U2·U5·U6·NTC complex from *Schizosaccharomyces pombe* to monitor its interactions using a RNA footprinting-based strategy[3,4]. However, it is unclear if this method can be applied to mammalian cells which might be more sensitive to introduction of affinity tags into splicing factors. Furthermore, no method has simultaneously monitored the full complexity of the interactions of diverse RBPs on pre-mRNAs from the earliest to the latest stages of spliceosomal assembly.

Here, we have adapted the individual nucleotide resolution UV crosslinking and immunoprecipitation (iCLIP) method[6] to develop spliceosome iCLIP. This approach identifies crosslinks of endogenous, untagged spliceosomal factors on pre-mRNAs at nucleotide resolution. In a previous study, we demonstrated validity of this approach by showing how PRPF8 remodels spliceosomal contacts at 5'ss[5]. Here, we comprehensively characterize spliceosome iCLIP and show that it simultaneously maps the crosslink profiles of core and accessory spliceosomal factors that are known to participate across the diverse stages of the splicing cycle. Due to iCLIP's nucleotide precision, we distinguished 7 binding peaks corresponding to distinct RBPs that differ in their requirement for ATP or the factor PRPF8. Spliceosome iCLIP also purifies intron lariats and identified 132,287 candidate BP positions. Compared to BPs identified in previous RNA-seq studies[7-9], those identified by spliceosome iCLIP contain more canonical sequence and structural features. We further examined the binding profiles of spliceosomal RBPs around the BPs. This demonstrates that assembly of SF3 and associated spliceosomal complexes tends to be determined by a primary BP in most introns, even though alternative BPs are detected by lariat-derived reads in RNA-seq. Moreover, we identify complementary roles of U2AF and SF3 complexes in BP definition. Taken together, these findings demonstrate the value of spliceosome iCLIP for

94　transcriptome-wide studies of BP definition and spliceosomal interactions with pre-
95　mRNAs.

96　**Results**

97　**Spliceosome iCLIP identifies interactions between splicing factors, snRNAs and pre-**
98　**mRNAs**

99　SmB/B' proteins are part of the highly stable Sm core common to all spliceosomal snRNPs
100　except U6[1]. In order to adapt iCLIP for the study of a multi-component machine like the
101　spliceosome, we immunopurified endogenous SmB/B' proteins[10] using a range of
102　conditions with differing stringency of detergents and salt concentrations for the lysis and
103　washing steps (Supplementary Table 1, Fig. 1a and Supplementary Fig. 1a,b). First, to
104　enable denaturing purification, we generated HEK293 cells stably expressing Flag-tagged
105　SmB and employed 6M urea during cell lysis to minimize co-purification of additional
106　proteins[11] ('stringent' purification, Supplementary Table 1), followed by dilution of the
107　lysis buffer (see Methods) to facilitate immunopurification of SmB via the Flag tag. We
108　observed a 25 kDa band corresponding to the molecular weight of SmB-RNA complexes,
109　which was absent when UV light or anti-Flag antibody were omitted, or when cells not
110　expressing Flag-SmB were used (Supplementary Fig. 1c). Next, we employed the
111　standard, non-denaturing iCLIP condition, which uses a high concentration of detergents
112　in the lysis buffer, and wash buffer with 1M NaCl ('medium' purification, Supplementary
113　Table 1). This disrupts most protein-protein interactions but can preserve stable
114　complexes such as snRNPs, as evident by the multiple radioactive bands in addition to the
115　25 kDa SmB-RNA complex upon treatment with low RNase (Fig. 1b). Of note, similar
116　profiles of protein-RNA complexes were obtained when using different monoclonal
117　SmB/B' antibodies (Supplementary Fig. 1d). Last, we further decreased the concentration
118　of detergents in the lysis buffer, used 0.1M NaCl in the washing buffer ('mild' purification,
119　Supplementary Table 1), and employed the low RNase treatment that leaves snRNAs
120　generally intact such that they serve as a scaffold for purifying the multi-protein
121　spliceosomal complexes (Fig. 1a).

122　To produce cDNA libraries with spliceosome iCLIP, we immunoprecipitated SmB/B'
123　under the three different stringency conditions from lysates of UV-crosslinked cells, and
124　isolated a broad size distribution of protein-RNA complexes in order to recover the
125　greatest possible diversity of spliceosomal protein-RNA interactions (Fig. 1b and
126　Supplementary Fig. 1c,d). An antibody against endogenous SmB/B' was used for medium
127　and mild purification from HEK293, K562 and HepG2 cells, and an anti-Flag antibody for
128　stringent purification from HEK293 cells expressing Flag-SmB (Supplementary Table 2
129　and 3). As in previous iCLIP studies[6], the nucleotide preceding each cDNA was used for all
130　analyses. When stringent conditions were used, >75% of iCLIP cDNAs mapped to snRNAs,
131　likely corresponding to the direct binding of Flag-SmB (Fig. 1c). However, the proportion
132　of snRNA crosslinking reduced to ~40-60% under mild and medium conditions, with a
133　corresponding increase of crosslinking to introns and exons that likely reflects binding of
134　snRNP-associated proteins to pre-mRNAs (Fig. 1a,c).

## Spliceosome iCLIP identifies seven crosslinking peaks on pre-mRNAs

Assembly of the spliceosome on pre-mRNA is guided by three main landmarks: the 5'ss, 3'ss and BP. Therefore, we evaluated if spliceosomal crosslinks are located at specific positions relative to splice sites and computationally predicted BPs[12]. For this purpose we performed spliceosome iCLIP from human Cal51 cells, which we have previously used as a model system to study the roles of spliceosomal factors in cell cycle[5]. RNA maps of summarized spliceosomal crosslinking revealed 7 peaks around these landmarks (Fig. 2a). Importantly, similar positional patterns were also seen in HEK293, K562 and HepG2 cell lines (Supplementary Fig. 2a). The centers of the peaks were 15 nt upstream of the 5'ss (peak 1), 10 nt downstream of the 5'ss (peak 2), 31 nt downstream of the 5'ss (peak 3), 26 nt upstream of the BP (peak 4), 20 nt upstream of the BP (peak 5), 11 nt upstream of the 3'ss (peak 6) and 3 nt upstream of the 3'ss (peak 7). We also observed alignment of cDNA starts to the start of the intron and the BPs, which we refer to as positions A and B, respectively (Fig. 2a and Supplementary Fig. 2a). The crosslinking enrichment at most peaks was generally stronger under the mild condition, especially at the 3'ss (Supplementary Fig. 2a). This indicates that spliceosome iCLIP performed under mild conditions is most suitable for investigating spliceosomal assembly on pre-mRNAs.

## Spliceosome iCLIP monitors multiple stages of spliceosomal remodeling

Next, we investigated whether spliceosome iCLIP is able to monitor spliceosome assembly at different stages during the splicing cycle. For this purpose we knocked down (KD) PRPF8 in Cal51 cells (Supplementary Fig. 2b) and performed spliceosome iCLIP under mild conditions. As an integral component of the U4/U6.U5 tri-snRNP, PRPF8 is essential for both catalytic reactions[1]. We previously showed that PRPF8 is required for efficient spliceosomal assembly at 5'ss[5]. Here, we additionally find that PRPF8 is essential for efficient spliceosomal assembly at peaks 4 and 5 (Fig. 2a). Moreover, we also observed a major decrease of reads truncating at the positions A and B, whereas crosslinking at peaks 2 and 6 is increased upon PRPF8 KD.

To further investigate whether spliceosome iCLIP can monitor distinct stages of the splicing reaction, we performed an *in vitro* splicing assay in which an exogenous pre-mRNA splicing substrate was incubated with HeLa nuclear extract in the presence or absence of ATP. ATP is required for the progression of early, ATP-independent, spliceosomal complexes to later assembly stages mediating the catalytic splicing reactions. The RNA substrate was produced by *in vitro* transcription of a minigene construct containing a short intron and flanking exons from the human *C6orf10* gene. Gel electrophoresis analysis confirmed that the minigene RNA was efficiently spliced *in vitro* in an ATP-dependent manner (Supplementary Fig. 2c). We performed spliceosome iCLIP from the splicing reactions using the mild purification condition (Supplementary Fig. 2d). Following sequencing, the reads mapping to the exogenous splicing substrate or spliced product represented ~1%, whereas the remaining reads were derived from endogenous RNAs present in the nuclear extract (Supplementary Table 4). The spliced product was detected with exon-exon junction reads primarily in the presence of ATP (364 reads in +ATP vs. 5 reads in -ATP condition) (Supplementary Fig. 2e and Supplementary Table 4).

177 As expected given that the spliceosome rapidly disassembles upon completion of the
178 splicing reaction, very few reads mapped to the spliced (364 reads) compared to
179 unspliced substrate (48,584 reads) (Supplementary Table 4) in the +ATP condition. It
180 should be considered, however, that some reads from exogenous minigene could
181 represent RNA that did not enter the splicing pathway.

182 We visualized crosslinking on the substrate RNA, and marked positions that correspond
183 to peaks on the transcriptome-wide RNA maps (Fig. 2b). Whilst crosslinking peaks on a
184 metagene plot might not necessarily be representative of individual splicing substrates,
185 we nevertheless observed crosslinking in corresponding regions of the *C6orf10* substrate
186 (comparing Fig. 2a and 2b). When comparing crosslinking in the presence or absence of
187 ATP, an unchanged crosslinking profile was seen in regions of peaks 1, 2, 6 and 7,
188 indicating these are ATP-independent contacts of early spliceosomal factors. In contrast,
189 the presence of ATP led to a ~11 fold increase of crosslinking in the region upstream of
190 the BP where the PRPF8-dependent peaks 4 and 5 are located on endogenous transcripts
191 (Fig. 2b). This indicates that spliceosome iCLIP detects pre-mRNA binding of factors
192 contributing to early, ATP-independent and late, ATP-dependent stages of spliceosomal
193 assembly.

194 Following crosslinking, the peptide that remains bound to the RNA after RBP digestion
195 will normally terminate reverse transcription to produce so-called 'truncated cDNAs'[13-15].
196 Accordingly, analysis of data from iCLIP and derived methods, such as eCLIP[16], generally
197 refer to the nucleotide preceding the iCLIP read on the reference genome as the 'crosslink
198 site'. However, in spliceosome iCLIP we additionally expect cDNAs that truncate at the
199 three-way junction formed by intron lariats, where the 5' end of the intron is linked via a
200 2'-5' phosphodiester bond to the BP (Fig. 2c). Following RNase digestion, such lariat
201 three-way-junction RNAs present two available 3' ends for ligation of adapters, such that
202 cDNAs can truncate at the BP (i.e. position B) or at the start of the intron (i.e. position A).
203 Interestingly, the medium purification condition was optimal to produce cDNAs
204 truncating at positions A and B (Supplementary Fig. 2a), possibly because spliceosomal C
205 complexes containing lariat intermediates are known to be stable under high-salt
206 conditions[17]. Note that peaks A and B are higher in HEK293 compared to HepG2 and K562
207 cells under medium purification conditions, and likely reflect differences in lariat co-
208 purification. Meanwhile, the number of cDNAs truncating at the positions A and B is
209 dramatically decreased under conditions that inhibit splicing progression and lariat
210 formation: PRPF8 KD *in vivo* (2-fold, Fig. 2a), or absence of ATP *in vitro* (≥18-fold, Fig. 2b).
211 This further confirms that spliceosome iCLIP can monitor spliceosome assembly at
212 distinct stages of the splicing cycle.
213
214 **Specific RBPs are enriched at each peak of spliceosomal crosslinking**

215 Next, to identify RBPs that crosslink at peaks identified by spliceosome iCLIP, we
216 examined the eCLIP data for 110 RBPs (from 157 eCLIP samples of 68 RBPs in the HepG2,
217 and 89 RBPs in the K562 cell line) provided by the ENCODE consortium[16]. Of note,
218 comparisons between iCLIP and eCLIP are justified due to their use of identical lysis and
219 wash buffers (analogous to medium stringency from the present study), use of truncated

220  cDNAs to identify crosslink sites and similar RNase digestion conditions, and comparable
221  crosslinking profiles for RBPs such as PTBP1 and U2AF2[15]. Accordingly, we analyzed the
222  eCLIP data to identify RBPs with enriched normalized crosslinking at each spliceosomal
223  iCLIP peak. This identified a specific set of RBPs at each peak, with good overlap between
224  RBPs identified in K562 and HepG2 cells (Fig. 3 and Supplementary Data Set 1). As
225  expected, SF3 components SF3B4, SF3A3 and SF3B1 bind to peaks 4 or 5[18], U2AF2 binds
226  the polypyrimidine (polyY) tract (peak 6), and U2AF1 close to the intron-exon junction
227  (peak 7)[19].

**Spliceosome iCLIP identifies BPs with canonical sequence and structural features**

229  To determine whether spliceosome iCLIP could experimentally identify human BPs, we
230  used spliceosome iCLIP data produced under medium purification from Cal51 cells. Most
231  cDNA starts in spliceosome iCLIP overlap with a uridine-rich motif (Fig. 4a), in agreement
232  with an increased propensity of protein-RNA crosslinking at uridine-rich sites[14]. In
233  contrast, cDNAs ending at the last nucleotide of introns, which are thus likely derived from
234  intron lariats, have starts overlapping the YUNAY motif matching the consensus BP
235  sequence (Fig. 4b). Further, these cDNAs have higher enrichment of mismatches of
236  adenosines at their first nucleotide (Supplementary Fig. 3a), which is consistent with
237  mismatch, insertion and deletion errors during reverse transcription across the three-
238  way junction of the BP[9]. For comparison, reads that start in regions where BPs are
239  typically located, but which do not align with intron ends, have less enrichment of the BP
240  consensus motif at their starts (Supplementary Fig. 3b,c). To identify a confident set of
241  putative BPs in a transcriptome-wide manner, we therefore used the spliceosome iCLIP
242  cDNAs that aligned with the end of introns (Fig. 4b). These cDNAs started at adenines in
243  132,287 intronic positions, which we considered as BP candidates. The 41 read-length
244  limited our analysis to the region where most BPs are located, but more distal BPs cannot
245  be identified by this approach. For further study, we selected BPs with the highest number
246  of truncated cDNAs per intron. This identified candidate BPs in 43,637 introns of 9,565
247  genes.

248  To examine the BPs identified by spliceosome iCLIP ('iCLIP BPs'), we compared them with
249  the 'computational BPs' recently identified with a sequence-based deep learning
250  predictor, LaBranchoR, which predicted BPs for over 90% of 3'ss[12]. We also compared
251  with 'RNA-seq BPs', including the 138,314 BPs from 43,637 introns that were identified
252  by analysis of lariat-spanning reads from 17,164 RNA-seq datasets[8]. Initially, 65% of iCLIP
253  BPs overlapped with the top-scoring computational BPs (Supplementary Fig. 3d).
254  Interestingly, in cases where iCLIP and computational BPs were located <5 nt apart, they
255  frequently occurred within A-rich sequences (Supplementary Fig. 3e). This mismatch
256  could be of technical nature, as truncation of iCLIP cDNAs may not always be precisely
257  aligned to the BPs in case of A-rich sequences. Alternatively, more than one A might be
258  capable of serving as the BP. When allowing a 1 nt shift for comparison between methods,
259  as has been done previously[12], 70% of iCLIP BPs overlapped with the top-scoring
260  computational BPs, whilst 26% overlapped with the RNA-seq BPs (Fig. 4c, Supplementary
261  Data Set 2). If the computational BPs overlapped either with an iCLIP BP and/or RNA-seq
262  BP, it generally had a strong BP consensus motif (o-BP, Fig. 4d).

263  To gain insight into the differences between the methods, we focused on BPs that were
264  identified by a single method and located >5 nt away from BPs identified by other
265  methods. Notably, the computational- or iCLIP-specific BPs have a strong enrichment of
266  the consensus YUNAY motif (c-BP, i-BP, Fig. 4e,f,h,i). In contrast, RNA-seq-specific BPs
267  contain a larger proportion of non-canonical BP motifs, which agrees with previous
268  observations[7,9,12] (Fig. 4g,j). To evaluate further, we compared iCLIP BPs with two studies
269  that identified 59,359 BPs by exoribonuclease digestion and targeted RNA-sequencing[9],
270  and 36,078 BPs by lariat-spanning reads refined by U2 snRNP/pre-mRNA base-pairing
271  models[7]. Considering the introns that contained BPs defined both by RNA-seq and iCLIP,
272  we found 57% and 47% overlapping BPs (Supplementary Fig. 3f-i). Again, the iCLIP-
273  specific BPs were more strongly enriched in the consensus YUNAY motif compared to BPs
274  specifically identified by either RNA-seq method (Supplementary Fig. 3j-o). We also
275  examined the local RNA structure around each category of BPs. Overlapping, iCLIP-
276  specific and computational-specific BPs had a decreased pairing probability at the
277  position of the BP, which was not seen for the RNA-seq-specific BPs (Fig. 4k,l). The
278  difference in RNA-seq BPs derives from the presence of non-canonical, non-A branched
279  BPs, which have a generally increased pairing probability (Supplementary Fig. 3p,q). This
280  indicates that the non-A BPs might be structurally less accessible for pairing with U2
281  snRNP.

**Alignment of RBP binding profiles signifies the functionality of BPs**

283  Peaks 4, 5 and position B align to BP position, and therefore we could evaluate how the
284  crosslinking profiles of RBPs binding at these peaks align to the different classes of BPs.
285  First, we examined the crosslinking of SF3B4, which binds in the region of peak 4 as part
286  of the U2 snRNP complex that recognises the BP[1]. Analysis of the overlapping BPs (o-BP)
287  defines the peak of SF3B4 crosslinking at the $25_{th}$ nt upstream of BPs (Fig. 5 and
288  Supplementary Fig. 4a,b). However, the peak of SF3B4 crosslinking is shifted from this
289  $25_{th}$ position for the non-overlapping, method-specific BPs; it is generally closer than 25
290  nt to the BPs located upstream of another BP (up BP), and further than 25 nt away from
291  BPs located downstream of another BP (down BP) (Fig. 5). The shift from the expected
292  position is greatest for RNA-seq-specific BPs (R-BP), and smallest for computationally
293  predicted BPs, as evident by eCLIP data from two cell lines (Fig. 5a,b). Moreover, the same
294  result is seen with U2AF2, where the strongest shift away from expected positions is seen
295  for RNA-seq BPs, and weakest for computational BPs (Supplementary Fig. 4c,d). The
296  cDNA starts from PRPF8 eCLIP are highly enriched at position B, corresponding to the
297  lariat-derived cDNAs that truncate at BPs (Fig. 3). Interestingly, the PRPF8 cDNA starts
298  had the strongest peak at the overlapping BPs, but also peaked at all the remaining classes
299  of BPs (Supplementary Fig. 4e,f). This indicates that all classes of BPs contribute to lariat
300  formation, and that the non-overlapping BPs most likely act as alternative BPs within the
301  introns.

**Effects of BP position on spliceosomal assembly**

303  To assess how BP positioning determines spliceosome assembly, we evaluated binding
304  profiles of the RBPs that are enriched at peaks 4-7 and at positions A and B (Fig. 3). We

divided BPs based on their distance from 3'ss, and normalized RBP binding profiles within each subclass of BP. This showed that crosslinking of U2AF1 and U2AF2 aligns to the region between the BPs and 3'ss, which is covered by the polyY tract (Supplementary Fig. 5 and 6). Whilst SF3B4 is the primary RBP crosslinking at peak 4, and SF3A3 at peak 5, binding of SMNDC1, SF3B1, EFTUD2, BUD13, GPKOW and XRN2 to peaks 4 and 5 was also evident (Supplementary Fig. 5, 6 and Fig. 3). PRPF8, RBM22 and SUPV3L1 have their cDNA starts truncating at positions A and B (Supplementary Fig. 5 and 6), corresponding to the three-way junction formed by intron lariats (Fig. 2c). This is in agreement with the association of PRPF8 and RBM22 with intron lariats as part of the human catalytic step I spliceosome[1]. The positions of SF3B4 and SF3A3 crosslinking peaks also agree with CryoEM studies of the human spliceosome that show closer pre-mRNA binding of SF3A3 (also referred to as SF3a60) to the BP compared to SF3B4 (also referred to as SF3b49)[20].

In order to quantify how BP positioning affects the intensity of RBP binding, we divided BPs into 10 equally sized groups based on the distance from 3'ss. We then normalized the relative binding intensity of each RBP at each position on the RNA maps across the ten groups, and revealed strong relationships between BP position and binding intensity of certain RBPs (Fig. 6a, Supplementary Fig. 7a). For example, if a BP is located distally from the 3'ss, then U2AF components bind stronger to peaks 6 and 7. In contrast, if a BP is located proximally to the 3'ss, then EFTUD2, SF3 components and several other RBPs bind stronger to the peaks 4 or 5 (Fig. 6b). Notably, increased BP distance causes increased binding of BUD13 and GPKOW at peaks 6 or 7 and decreased binding at peaks 4 and 5. The more efficient recruitment of U2AF and associated factors to peaks 6 and 7 could be explained by the long polyY-tracts at distal BPs (Supplementary Fig. 5), while their decreased binding at proximal BPs appears to be compensated by increased binding of SF3 and other U2 snRNP-associated factors at peaks 4 and 5.

In contrast to effects on individual splicing factors, we did not observe any effect of BP distance on the relative intensity of spliceosome iCLIP crosslinking in peaks 4 and 5 compared to 6 and 7 (Fig. 6c). This indicates that the effects may be masked during later stages of spliceosome assembly. To ask if this is the case, we turned to PRPF8, a protein that is essential for later stages of spliceosomal assembly, a role it plays together with EFTUD2 and BRR2 as part of U5 snRNP[1]. PRPF8 KD leads to decreased spliceosomal binding at peaks 4 and 5, and this effect is stronger at distal compared to proximal BPs (Fig. 6c). In conclusion, our results reveal differences in the binding profiles of splicing factors in relation to BP distance, but these differences are neutralized upon full spliceosome assembly in a manner that requires the presence of PRPF8.

**Effects of BP strength on spliceosomal assembly**

To examine how BP strength affects spliceosomal assembly we focused on BPs that have been identified both by spliceosome iCLIP and computational modelling, and which are located at 23-28 nt upstream of the 3'ss. Of note, this is the most common position of BPs (Supplementary Data Set 3). As an estimate of BP strength we used the BP score, which was determined with a deep-learning model[12]. This showed strong correlation between BP strength and RBP binding intensities, such that most RBPs have increased crosslinking

347 at peaks 4 and 5 at BPs with very high scores, and, conversely, increased crosslinking at
348 peaks 6 and 7 at BPs with very low scores (Fig. 7a,b, Supplementary Fig. 7b). Since SF3
349 components primarily bind at peaks 4 and 5, and U2AF components at peaks 6 and 7, an
350 over 4-fold change is seen in the ratio of crosslinking when comparing the extreme deciles
351 of BP strength (Supplementary Fig. 7c). We did not observe any correlation between the
352 polyY tract coverage and BP score (Supplementary Fig. 7d), which indicates that BP
353 strength directly affects the RBP binding profiles.

354 Similar to the effects on individual splicing factors, the relative intensity of spliceosome
355 iCLIP crosslinking in peaks 4 and 5 was increased with increasing BP strength (Fig. 7c,
356 compare blue lines on the left and right graphs). PRPF8 KD decreased spliceosomal
357 binding at peaks 4 and 5 of both classes of BPs, and this led to stronger crosslinking at
358 peaks 6 and 7 relative to peaks 4 and 5 at weak BPs, even though the peaks 4 and 5 are
359 usually stronger. The signal at position B of weak BPs is almost completely lost upon
360 PRPF8 KD, which likely reflects the absence of intron lariats due to perturbed splicing of
361 introns with weak BPs (Fig. 7c). In conclusion, our results suggest that the assembly
362 efficiency of spliceosomal factors at peaks 4 and 5 closely correlates with BP strength,
363 which indicates that recognition of weak BPs might be more sensitive to perturbed
364 spliceosome function.

365
366 **Discussion**

367 Here we established spliceosome iCLIP to study the interactions of endogenous snRNPs
368 and accessory splicing factors on pre-mRNAs. We identified peaks of spliceosomal
369 protein-pre-mRNA interactions, which precisely overlap with crosslinking profiles of 15
370 splicing factors. Interestingly, the contacts of RBPs in peaks 4 and 5 don't overlap with
371 any sequence motif, and thus the constrained conformation of the larger spliceosomal
372 complex appears to act as a molecular ruler that positions each associated RBP on pre-
373 mRNAs at a specific distance from BPs. Moreover, the presence of lariat-derived reads in
374 spliceosome iCLIP identified >40,000 BPs that have canonical sequence and structural
375 features. Due to the precise alignment of splicing factors relative to the positions of BPs,
376 we could use their binding profiles to show that the assembly of U2 snRNP is primarily
377 coordinated by the computationally predicted BPs, whilst alternative BPs, identified only
378 by iCLIP or RNA-seq, are more rarely used. Finally, we reveal the major effect of the
379 position and strength of BPs on spliceosomal assembly, which can explain why distally
380 located or weak BPs are particularly sensitive to perturbed spliceosome function upon
381 PRPF8 KD. These findings demonstrate the broad utility of spliceosome iCLIP for
382 simultaneous and transcriptome-wide analysis of the assembly of diverse spliceosomal
383 components.

384 **The value of spliceosome iCLIP for identifying BPs**

385 Both RNA-seq and iCLIP identify BPs by analyzing cDNAs derived from intron lariats.
386 Thus, the efficiency of these methods depends on the abundance of intron lariats, which
387 depends on the kinetics of lariat debranching. Several studies demonstrated that lariats

388　formed at non-canonical BPs are less efficiently debranched[21-23], and therefore these non-
389　canonical BPs are expected to be more efficiently detected. This is especially true for RNA-
390　seq-based methods, because they monitor steady state RNA levels. In contrast, iCLIP only
391　captures lariats in complex with spliceosomes, thus minimizing bias for lariats that are
392　stable after their release from the spliceosome. This could explain why the BPs identified
393　by iCLIP contain a stronger consensus sequence than BPs identified from lariat-spanning
394　reads in RNA-seq. The further value of spliceosome iCLIP is that, in addition to
395　experiments under the medium condition that permit BP identification through lariat-
396　derived cDNAs, experiments under the mild condition identify the SF3 complex and other
397　U2 snRNP-associated RBPs that crosslink at peaks 4 and 5. These can crucially be used to
398　independently validate the functional role of BPs in the assembly of U2 snRNP. Thus, use
399　of spliceosome iCLIP under both conditions, combined with computational modelling of
400　BPs[12], is well suited to studying the functionality of BPs.

401　**The role of BP position and strength in spliceosomal assembly**

402　We show that BP position and the computationally defined strength of BPs correlate with
403　the relative binding of splicing factors around BPs. This is exemplified by strong binding
404　of SF3 components at strong BPs, or BPs located close to 3'ss, whilst U2AF components
405　bind stronger to weak BPs, or BPs located further from 3'ss (Fig. 7d). In the cases of SF3B1,
406　BUD13 and GPKOW, we observed enriched binding at peaks 4 and 5 as well as 6 and 7,
407　with reciprocal changes between the two peak regions dependent on BP features (Fig. 6
408　and 7). These RBPs are not known to bind at peaks 6 or 7, and it is plausible that the signal
409　at some peaks represents binding of U2AF or other spliceosomal factors that are co-
410　purified during eCLIP. It is presently not possible to fully distinguish between direct and
411　indirect binding from eCLIP data, because purified protein-RNA complexes have not been
412　visualized after their separation on SDS-PAGE gels in eCLIP[13]. Nevertheless, it is clear that
413　BP characteristics determine the balance between binding of SF3 and associated factors
414　at peaks 4 and 5 and of U2AF and associated factors at peaks 6 and 7. This suggests further
415　study of RBP binding profiles around BPs could unravel a BP 'code' that facilitates specific
416　stages of BP recognition and function.

417　In conclusion, spliceosome iCLIP monitors concerted pre-mRNA binding of many types of
418　spliceosomal complexes with nucleotide resolution, allowing their simultaneous study
419　due to the distinct position-dependent binding pattern of components acting at multiple
420　stages of the splicing cycle. The method can now be used to study the endogenous
421　spliceosome and BPs across tissues, species and stages of development without need for
422　the protein tagging used in yeast[3,4]. Further, several spliceosomal components, including
423　U2AF1, SF3B1 and PRPF8, are targets for mutations in myeloid neoplasms, retinitis
424　pigmentosa and other diseases[24]. Spliceosome iCLIP could now be used to monitor global
425　impacts of these mutations on spliceosome assembly in human cells. More generally, our
426　study demonstrates the value of iCLIP for monitoring the position-dependent assembly
427　and dynamics of multi-protein complexes on endogenous transcripts.

428

429

## Acknowledgements

## Author contributions

M.B., C.R.S. and J.U. conceived the project, designed the experiments and wrote the manuscript, with assistance of all co-authors. M.B., C.R.S., Z.W., R.F. and A.S.E. performed experiments, with assistance from J.U., J.K. and C.W.S.. N.H. performed most computational analyses, with assistance from C.R.S., T.C., R.F., A.M.C. and N.M.L.. V.O.W., D.P. and A.R.V. provided crosslinked pellets from wild-type and PRPF8-depleted Cal51 cells. L.S. and L.P. developed and characterized the monoclonal antibody 18F6.

## Competing interests

The authors declare no competing interests.

**References**

459

460    1.    Fica, S.M. & Nagai, K. Cryo-electron microscopy snapshots of the
461          spliceosome: structural insights into a dynamic ribonucleoprotein
462          machine. *Nat Struct Mol Biol* **24**, 791-799 (2017).
463    2.    Wahl, M.C., Will, C.L. & Lührmann, R. The spliceosome: design principles of
464          a dynamic RNP machine. *Cell* **136**, 701-18 (2009).
465    3.    Chen, W. et al. Transcriptome-wide Interrogation of the Functional
466          Intronome by Spliceosome Profiling. *Cell* **173**, 1031-1044 e13 (2018).
467    4.    Burke, J.E. et al. Spliceosome Profiling Visualizes Operations of a Dynamic
468          RNP at Nucleotide Resolution. *Cell* **173**, 1014-1030 e17 (2018).
469    5.    Wickramasinghe, V.O. et al. Regulation of constitutive and alternative
470          mRNA splicing across the human transcriptome by PRPF8 is determined
471          by 5' splice site strength. *Genome Biol* **16**, 201 (2015).
472    6.    König, J. et al. iCLIP reveals the function of hnRNP particles in splicing at
473          individual nucleotide resolution. *Nat Struct Mol Biol* **17**, 909-15 (2010).
474    7.    Taggart, A.J. et al. Large-scale analysis of branchpoint usage across species
475          and cell lines. *Genome Res* **27**, 639-649 (2017).
476    8.    Pineda, J.M.B. & Bradley, R.K. Most human introns are recognized via
477          multiple and tissue-specific branchpoints. *Genes Dev* **32**, 577-591 (2018).
478    9.    Mercer, T.R. et al. Genome-wide discovery of human splicing
479          branchpoints. *Genome Res* **25**, 290-303 (2015).
480    10.   Carissimi, C., Saieva, L., Gabanella, F. & Pellizzoni, L. Gemin8 is required
481          for the architecture and function of the survival motor neuron complex. *J
482          Biol Chem* **281**, 37009-16 (2006).
483    11.   Huppertz, I. et al. iCLIP: protein-RNA interactions at nucleotide resolution.
484          *Methods* **65**, 274-87 (2014).
485    12.   Paggi, J.M. & Bejerano, G. A sequence-based, deep learning model
486          accurately predicts RNA splicing branchpoints. *RNA* **24**, 1647-1658
487          (2018).
488    13.   Lee, F.C.Y. & Ule, J. Advances in CLIP Technologies for Studies of Protein-
489          RNA Interactions. *Mol Cell* **69**, 354-369 (2018).
490    14.   Sugimoto, Y. et al. Analysis of CLIP and iCLIP methods for nucleotide-
491          resolution studies of protein-RNA interactions. *Genome biology* **13**, R67
492          (2012).
493    15.   Haberman, N. et al. Insights into the design and interpretation of iCLIP
494          experiments. *Genome Biol* **18**, 7 (2017).
495    16.   Van Nostrand, E.L. et al. A Large-Scale Binding and Functional Map of
496          Human RNA Binding Proteins. *bioRxiv* (2017).
497    17.   Bessonov, S., Anokhina, M., Will, C.L., Urlaub, H. & Luhrmann, R. Isolation
498          of an active step I spliceosome and composition of its RNP core. *Nature*
499          **452**, 846-50 (2008).
500    18.   Gozani, O., Feld, R. & Reed, R. Evidence that sequence-independent
501          binding of highly conserved U2 snRNP proteins upstream of the branch
502          site is required for assembly of spliceosomal complex A. *Genes Dev* **10**,
503          233-43 (1996).
504    19.   Zarnack, K. et al. Direct Competition between hnRNP C and U2AF65
505          Protects the Transcriptome from the Exonization of Alu Elements. *Cell*
506          **152**, 453-66 (2013).

507 20. Zhang, X. et al. Structure of the human activated spliceosome in three
508     conformational states. *Cell Res* **28**, 307-322 (2018).
509 21. Jacquier, A. & Rosbash, M. RNA splicing and intron turnover are greatly
510     diminished by a mutant yeast branch point. *Proc Natl Acad Sci U S A* **83**,
511     5835-9 (1986).
512 22. Hesselberth, J.R. Lives that introns lead after splicing. *Wiley Interdiscip*
513     *Rev RNA* **4**, 677-91 (2013).
514 23. Talhouarne, G.J.S. & Gall, J.G. Lariat intronic RNAs in the cytoplasm of
515     vertebrate cells. *Proc Natl Acad Sci U S A* **115**, E7970-E7977 (2018).
516 24. Scotti, M.M. & Swanson, M.S. RNA mis-splicing in disease. *Nat Rev Genet*
517     **17**, 19-32 (2016).
518 25. Lorenz, R. et al. ViennaRNA Package 2.0. *Algorithms Mol Biol* **6**, 26 (2011).
519
520

**Figure legends**

**Fig. 1 | Spliceosome iCLIP identifies protein interactions with snRNAs and splicing**
**substrates.**

524 (a) Schematic representation of the spliceosome iCLIP method performed under
525 conditions of varying purification stringency.

526 (b) Autoradiogram of crosslinked RNPs immunopurified from HeLa cells under medium
527 conditions by a SmB/B' antibody following digestion with high (++) or low (+) amounts
528 of RNase I. The dotted line depicts the region typically excised from the nitrocellulose
529 membrane for spliceosome iCLIP. As control, the antibody (Ab) was omitted during
530 immunopurification.

531 (c) Genomic distribution of spliceosome iCLIP cDNAs produced under stringent, medium
532 and mild conditions from HEK293 cells. Data was mapped first to snRNAs, allowing
533 multiple mapping reads, and then to the genome, allowing only uniquely mapped reads.
534 Proportions of cDNAs mapping to snRNAs, introns, coding sequence of mRNAs (CDS),
535 untranslated regions of mRNAs (UTR) and long non-coding RNAs (lncRNAs) are shown
536 (but not the intergenic reads and other types of RNAs). Data are shown as mean±s.e.m
537 from three independent experiments for the medium and mild purification condition and
538 two independent experiments for the stringent purification condition. Source data for
539 panel c are available online.

540

541 **Fig. 2 | Analysis of spliceosomal interactions with pre-mRNAs *in vitro* and *in vivo*.**

542 (a) Metagene plots of spliceosome iCLIP from Cal51 cells. Plots are depicted as RNA maps
543 of summarized crosslinking at all exon-intron and intron-exon boundaries, and around
544 BPs to identify major binding peaks, and to monitor changes between control and PRPF8
545 knockdown (KD) cells. Crosslinking is regionally normalized to its average crosslinking
546 across the -100..50 nt region relative to splice sites or BPs depending on the RNA map in
547 order to focus the comparison on the relative positions of peaks.

548 (b) Normalized spliceosome iCLIP cDNA counts on the *C6orf10 in vitro* splicing substrate.
549 Exons are marked by grey boxes, intron by a line, and the BP by a green dot. The positions
550 of crosslinking peaks are marked by numbers and letters corresponding to the peaks in
551 Figure 2a.

552 (c) Schematic description of the three-way junctions of intron lariats. The three-way
553 junction is produced after limited RNase I digestion of intron lariats. This can lead to
554 cDNAs that don't truncate at sites of protein-RNA crosslinking, but rather at the three-
555 way junction of intron lariats. These cDNAs initiate from the end of the intron and
556 truncate at the BP (position B), or initiate downstream of the 5'ss and truncate at the first
557 nucleotide of the intron (position A).

558

559 **Fig. 3 | Identification of RBPs overlapping with spliceosomal peaks at BPs and 3'ss.**

560 Enrichment of eCLIP crosslinking within each of the spliceosome iCLIP peaks, which are
561 defined by the positions marked in the figure. We first regionally normalized the
562 crosslinking of each RBP to its average crosslinking over -100..50 nt region relative to 3'ss,

563 which generates the RNA maps as shown in Supplementary Fig. 5 and 6. We then ranked
564 the RBPs according to the average normalized crosslinking across the nucleotides within
565 each peak. We analyzed peaks 4-7 and positions A and B, as marked on the top of each
566 plot. The top-ranking RBPs in each peak are shown on the left plot, and the full
567 distribution of RBP enrichments is shown on the right plot.

568

569 **Fig. 4 | Comparison of BPs identified by spliceosome iCLIP, RNA-seq lariat reads or**
570 **computational prediction.**

571 (a) Weblogo around the nucleotide preceding all spliceosome iCLIP reads.

572 (b) Weblogo around the nucleotide preceding only those spliceosome iCLIP reads that
573 align with ends of introns.

574 (c) Introns that contain at least one BP identified either by published RNA-seq[8] or by
575 spliceosome iCLIP are used to examine the overlap between the top BPs identified by
576 RNA-seq (i.e., the BP with most lariat-spanning reads in each intron), iCLIP (BP with
577 most cDNA starts) or computational predictions (highest scoring BP)[12]. BPs that are 0 or
578 1 nt apart are considered as overlapping. At the right, BP categories that are used for all
579 subsequent analyses are defined, along with their acronyms. If a BP defined by one
580 method is >5 nt upstream of a BP defined by another method, then 'up' is added to its
581 acronym, and if it is >5 nt downstream, 'down' is added.

582 (d) Weblogo of o-BP category of BPs.

583 (e) Weblogo of C-BPup category of BPs.

584 (f) Weblogo of i-BPup category of BPs.

585 (g) Weblogo of R-BPup category of BPs.

586 (h) Weblogo of C-BPdown category of BPs.

587 (i) Weblogo of i-BPdown category of BPs.

588 (j) Weblogo of R-BPdown category of BPs.

589 (k, l) The 100 nt RNA region centered on the BP was used to calculate pairing probability
590 with the RNAfold program using default parameters[25], and the average pairing
591 probability of each nucleotide around BPs is shown for the 40 nt region around method-
592 specific BPs located upstream (k) or downstream (l).

593

594 **Fig. 5 | Spliceosome assembly at BPs identified by spliceosome iCLIP, RNA-seq lariat**
595 **reads or computational prediction.**

596 Violin plots depicting the positioning of SF3B4 cDNA starts relative to the indicated BP
597 categories. SF3B4 eCLIP data were from K562 (a) and HepG2 (b) cells. Box-plot elements
598 are defined by center line, median; box limits, upper and lower quartiles; and whiskers,
599 1.5× interquartile range. Each data point corresponds to an eCLIP crosslink event, and the
600 total number of eCLIP crosslinks that map in the area analysed around each set of BPs
601 (sample size) is shown under the plot.

602

**Fig. 5 | Spliceosome assembly at BPs identified by spliceosome iCLIP, RNA-seq lariat reads or computational prediction.**

Violin plots depicting the positioning of SF3B4 cDNA starts relative to the indicated BP categories. SF3B4 eCLIP data were from K562 (a) and HepG2 (b) cells. Box-plot elements are defined by center line, median; box limits, upper and lower quartiles; and whiskers, 1.5× interquartile range.


**Fig. 6 | BP position defines the binding patterns of splicing factors at 3'ss.**

(a) Heatmaps depicting the normalized crosslinking of RBPs in peak regions around 10 groups of BPs that were categorized according to the distance of the BP from 3'ss. Crosslinks were derived as cDNA starts from eCLIP of HepG2 cells.

(b) RNA maps showing normalized crosslinking profiles of selected RBPs relative to BPs and 3'ss for the two deciles of BPs that are located most proximal (interrupted light lines) or most distal (solid dark lines) from 3'ss.

(c) RNA maps showing crosslinking profile of spliceosome iCLIP from control and PRPF8 KD Cal51 cells in the same format as panel b.


**Fig. 7 | BP strength correlates with the binding of splicing factors.**

(a) Heatmaps depicting the normalized crosslinking of RBPs in peak regions around 10 groups of BPs that were categorized according to the computational scores that define BP strength. Crosslinks were derived as cDNA starts from eCLIP of HepG2 cells.

(b) RNA maps showing normalized crosslinking profiles of selected RBPs relative to 3'ss for the two deciles of BPs that are lowest scoring (interrupted light lines) or highest scoring (solid dark lines).

 (c) RNA maps showing crosslinking profile of spliceosome iCLIP from control and PRPF8 KD Cal51 cells in the same format as panel b.

(d) Schematic representation of the effects that BP position and score have on the assembly of SF3 and U2AF complexes around BPs.




**Online Methods**

**Cell culture**

Flp-In HEK293 T-REx cells were from ThermoFisher (R78007), K562, HepG2 and standard HEK293 cells were obtained from the Francis Crick Cell Services Science Technology Platform,  and Cal51 breast adenocarcinoma cells were obtained from DSMZ (reference 14563). All cell lines tested negative for Mycoplasma contamination. HEK293

640  and HepG2 were cultured in DMEM with 10% FBS (ThermoFisher) and 1× penicillin-
641  streptomycin (ThermoFisher). K562 cells were cultured in RPMI 1640 (IMDM, ATCC)
642  with 10% FBS and 1× penicillin-streptomycin. Cal51 cells were cultured in DMEM
643  (ThermoFisher) with 10% fetal calf serum (FCS, ThermoFisher) and 1× penicillin-
644  streptomycin (ThermoFisher).

645  To generate a plasmid encoding 3×Flag epitope-tagged SmB, the SmB cDNA was amplified
646  using Phusion High-Fidelity DNA polymerase (NEB) with primers carrying the KpnI and
647  NotI restriction enzymes sites and cloned using Rapid DNA Ligation Kit (Thermo Fisher
648  Scientific) into a pcDNA5/FRT/TO vector modified to encode 3×Flag peptide upstream of
649  the multiple cloning site. To produce stable cell lines expressing this construct, the
650  pcDNA5/FRT/TO plasmid with 3×Flag epitope-tagged SmB was co-transfected with
651  pOG44 plasmid into Flp-In HEK293 T-REx cells (ThermoFisher, R78007). Cells stably
652  expressing these proteins were selected by culturing in Dulbecco's Modified Eagle
653  Medium (DMEM, Thermofisher) containing 10% fetal bovine serum (FBS), 3 µg/ml
654  Blasticidine S HCl, 200 µg/ml Hygromycine (InvivoGen). Flp-In 293 T-REx cells (Life
655  Technologies) were cultured in DMEM with 10% FBS, 3 µg/ml Blasticidin S HCl (Life
656  Technologies), 50 µg/ml Zeocin (Life Technologies). Doxycycline was added to media 24
657  hours prior to sample preparation in order to induce construct expression.

658  Cal51 breast adenocarcinoma cells were prepared as described previously[5]. For siRNA-
659  mediated depletion of PRPF8, Cal51 cells were transfected using DharmaFECT1
660  (Dharmafect) with 25 nM siRNA targeting human *PRPF8*. Transfected cells were
661  harvested 54 hrs later, exposed to UV-C light and used for iCLIP as described below. For
662  collection of samples from different stages of the cell cycle, Cal51 cells were synchronized
663  in G1/S by standard double thymidine block. Briefly, cells were treated with 1.5 mM
664  thymidine for 8 hrs, washed and released for 8 hrs, then treated again with thymidine for
665  a further 8 hrs. Cells were also collected 3 hrs (S-phase) and 7 hrs (G2) after release from
666  the thymidine block.

**Antibody production**

668  For production of the anti-SmB/B' monoclonal antibody 18F6, Balb/c females were
669  primed with Immuneasy adjuvant (Qiagen) and 25 mg of 6×His-SmB purified
670  recombinant proteins. Following two boosts at two-week intervals, SP2 myeloma cells
671  were fused with mouse splenocytes and hybridoma supernatants were analyzed onto
672  antigen-coated aminosilane modified slides using a LS400 Scanner (Tecan) and the
673  GenePix Pro 4.1 software as described previously[10]. Hybridoma cells were subcloned by
674  limiting dilution and further screened by ELISA, Western blot and immunofluorescence
675  analysis of HeLa cells.

*In vitro* **splicing**

677  For *in vitro* splicing reactions, a *C6orf10* minigene construct containing exon 8 and 9 and
678  150 nt of the intron around both splice sites was produced (Fig. 2b). The minigene

679    plasmid was linearized and transcribed *in vitro* using T7 polymerase with $_{32}$P-UTP. The
680    transcribed RNA was then subjected to *in vitro* splicing reactions using HeLa nuclear
681    extract. HeLa nuclear extract was depleted of endogenous ATP by pre-incubation and, for
682    each reaction, 10 ng of RNA was incubated with 60% HeLa nuclear extract at 30°C with
683    or without additional 0.5 mM ATP for 1 h in a 20 µl reaction. Afterwards, the reaction
684    mixture was UV-crosslinked at 100 mJ/cm$_2$ and stored at -80°C until further use. To
685    visualize the splicing reaction products, proteinase K was added to the reaction mixture
686    for 30 min at 37°C. The resulting RNA was phenol-extracted, precipitated and subjected
687    to gel electrophoresis on a 5% polyacrylamide-urea gel.

688    **Spliceosome iCLIP protocol**

689    For each experiment, three biological replicate samples of cDNA libraries were prepared
690    (Supplementary Tables 2 and 3). The iCLIP method was done as previously described[11],
691    with the following modifications. Crosslinked cells or tissue were dissociated in the lysis
692    buffer according to the stringency conditions (stringent, medium, mild; Supplementary
693    Table 1) followed by sonication, low RNase I (AM2295, 100 U/µl, ThermoFisher)
694    digestion and centrifugation. RNase at low concentration ensured that cDNAs are of
695    optimal size for comprehensive crosslink determination[15]. For denaturing, high-
696    stringency experiment[11], M2 anti-Flag antibody (Sigma) was used against the 3×Flag-SmB
697    protein that had been stably integrated into HEK-293 FlpIn cells (Supplementary Fig. 1c).
698    6M Urea buffer was first used to lyse cell pellets, before being diluted down 1:9 with a
699    Tween-20-containing IP buffer to allow for immunopurification without denaturing of the
700    M2 anti-Flag antibody, and then proceeded as described previously[15].

701    Standard iCLIP protocol[11] was used for Cal51 cells under mild and medium stringency
702    conditions, and for the *in vitro* splicing reactions under mild conditions, whilst an updated
703    protocol was used for HEK293, HepG2 and K562 cells[26]. For SmB/B' immunopurification
704    anti-SmB/B' antibodies 12F5 (sc-130670, Santa Cruz Biotechnology for Cal51 cells, and
705    S0698, Sigma-Aldrich for HEK293, HepG2 and K562 cells) or 18F6 (as hybridoma
706    supernatant, generated as described previously[10]) were used, which are different clones
707    from the same immunization. These antibodies behave identically under
708    immunopurification conditions (Supplementary Fig. 1d). For spliceosome iCLIP from *in*
709    *vitro* splicing reactions (Supplementary Fig. 2c,d), lysates were incubated with 50 µl
710    monoclonal anti-SmB/B' antibody 18F6, and for immunoprecipitations from cell lysates,
711    12F5 anti-SmB/B' antibody was used. The antibody was bound to 100 µl protein G
712    Dynabeads (ThermoFisher) under rotation at 4°C followed by washing. As described
713    previously, following immunopurification, RNA 3' end dephosphorylation, ligation of the
714    adapter 5'-rAppAGATCGGAAGAGCGGTTCAG/ddC/-3' to the 3' end and 5' end
715    radiolabeling, protein-RNA complexes were size-separated by SDS-PAGE and transferred
716    onto nitrocellulose membrane. The regions corresponding to 28-180 kDa were excised
717    from the membrane in order to isolate the bound RNA by proteinase K treatment. RNAs
718    were reverse-transcribed in all experiments using SuperScript III or IV reverse
719    transcriptase (ThermoFisher) and custom indexed primers (Supplementary Table 2).
720    Resulting cDNAs were subjected to electrophoresis on a 6% TBE-urea gel (ThermoFisher)

721  for size selection. Purified cDNAs were circularized, linearized and amplified for high-
722  throughput sequencing.

723  Identification of protein crosslink sites around splice sites, in particular at the peaks 4 and
724  5, was most efficient under the mild purification condition (Supplementary Fig. 2a). This
725  condition was therefore used for analysis of spliceosomal assembly upon PRPF8
726  knockdown in Cal51 cells (Fig. 2a), and in the *in vitro* splicing reactions in HeLa nuclear
727  extract (Fig. 2b). For the identification of BPs, we additionally used the medium condition,
728  since it increases the frequency of cDNAs truncating at peak B (Supplementary Fig. 2a).
729  For this purpose, spliceosome iCLIP was performed under medium purification
730  conditions from Cal51 cells synchronized in G1, S and G2 phase. To maximize cDNA
731  coverage, data from all synchronized cells was merged with the control Cal51 cells under
732  mild condition for BP identification.

### Mapping of Sm iCLIP reads

734  We mapped iCLIP data to the GRCh38 primary assembly and GENCODE v27 gene
735  annotations using STAR (v.2.2.1). Experimental and random barcode sequences of iCLIP
736  sequenced reads were removed prior to mapping (Supplementary Table 2). Following
737  mapping, we used random barcodes to quantify the number of unique cDNAs at each
738  genomic position by collapsing cDNAs with the same random barcode that mapped to the
739  same starting position to a single cDNA. For analysis of crosslinking to snRNAs, we first
740  mapped to a transcriptome of all annotated snRNA sequences in GENCODE v27 using
741  Bowtie2 (v2.3.4.3), and kept the primary alignment. Unmapped reads were then mapped
742  with STAR as previously described and intersected with GENCODE v27 for subtype
743  analysis, with reads from Bowtie2 being added to the total snRNA count. For spliceosome
744  iCLIP with the *C6orf10 in vitro* splicing substrate, sequence reads were first mapped to
745  the unspliced substrate and the remaining reads were mapped to the spliced substrate
746  allowing no mismatches. The nucleotide preceding the iCLIP cDNAs was used to define
747  the crosslink sites in all analyses.

### Mapping of eCLIP reads

749  For eCLIP sequencing data for all RBPs, we used GENCODE (GRCh38.p7) genome
750  assembly and the STAR alignment (version 2.4.2a) using the following parameters from
751  ENCODE pipeline: STAR --runThreadN 8 --runMode alignReads --genomeDir GRCh38
752  Gencode v25 --genomeLoad LoadAndKeep --readFilesIn read1, read2, --
753  readFilesCommand zcat --outSAMunmapped Within –outFilterMultimapNmax 1 --
754  outFilterMultimapScoreRange 1 --outSAMattributes All --outSAMtype BAM Unsorted –
755  outFilterType BySJout --outFilterScoreMin 10 --alignEndsType EndToEnd --
756  outFileNamePrefix outfile.

757  For the PCR duplicates removal, we used a python script 'barcode collapse pe.py' available
758  on GitHub (https://github.com/YeoLab/gscripts/releases/tag/1.0), which is part of the
759  ENCODE eCLIP pipeline (https://www.encodeproject.org/pipelines/ENCPL357ADL/).

**Normalization of crosslink positions for their visualization in the form of RNA maps**

RNA maps and heat maps were produced by summarizing the cDNA counts at each nucleotide using the previously developed RNA maps pipeline[15,27] relative to exon-intron and intron-exon boundaries and BPs on pre-mRNAs. The definition of intronic start and end positions was based on Ensembl version 75. Only introns longer than 300 nt were used to draw RNA maps in order to avoid detection of any RBPs that recognize 5'ss of introns.

In cases where we wished to compare the relative positions of crosslinking peaks between RBPs, we regionally normalized the summarized crosslinking of each RBP relative to the average crosslinking of the same RBP across the region 100 nt upstream and 50 nt downstream of the evaluated splice sites or BPs. Normalized values were then used to visualize the crosslinking in the form of RNA maps (Fig. 2, Supplementary Fig. 5 and 6). The same normalization was then used to plot heat maps, by plotting mean values of normalized RNA maps for each peak in the following regions; peak 4: -29..-23 nt and peak 5: -21..-17 nt relative to BP, peak 6: -11..-5 nt and peak 7: -3..-1 nt relative to 3'ss. Every RBP was then normalized by the mean across all the peaks to visualize crosslinking enrichment between the groups on the same scale across all RBPs (Fig. 6 and 7, Supplementary Fig. 7).

To assess the role of BP characteristics on spliceosomal RBP assembly (Fig. 4, 6 and 7), we only examined the introns containing the 31,167 BPs that were identified both computationally and by iCLIP, which are likely the most reliable. We divided BPs into 10 categories based on BP position or score, and then normalized the summarized crosslinking of each RBP in each of the 10 BP categories relative to the average crosslinking of the same RBP across the region 100 nt upstream and 50 nt downstream of all the 31,167 evaluated BPs.

For visualization of spliceosome iCLIP crosslinks along the *C6orf10 in vitro* splicing substrate and product (Fig. 2b and Supplementary Fig. 2e) we first summed the cDNA starts at each nt position and then normalized the counts by the average number of cDNA starts in the intronic region 101..150 relative to the 5'ss of the unspliced substrate. For the unspliced substrate normalized cDNA counts were logarithmized ($\log_2$) and data with $\log_2$(normalized number of cDNA starts)$\geq$1 were plotted. For the spliced product normalized cDNA counts were plotted.

**Identification and comparison of BPs**

It has been shown that the spliceosomal C complexes harbor a salt-resistant RNP core containing U2, U5 and U6 snRNAs as well as the splicing intermediates including lariats that withstand treatment with 1M NaCl, whereas the spliceosomal B complexes were more likely dissociated under high-salt conditions[17]. This could explain why the medium purification condition is more suited than the mild condition to enrich for lariat cDNAs truncating at position B (Supplementary Fig. 2a). It is conceivable that the medium

799　spliceosome iCLIP condition most strongly enriches spliceosomal C complexes, which are
800　most effective for lariat detection. In contrast, the mild condition is expected to enrich
801　additional B complexes that contain large amounts of SF3 components and have low
802　proportion of lariats, in agreement with the strong enrichment of peaks 4 and 5
803　(Supplementary Fig. 2a). To identify the maximal diversity of BPs, we therefore pooled
804　spliceosome iCLIP data produced under mild and medium purification conditions from
805　Cal51 cells.

806　To identify BPs we used the spliceosome iCLIP reads that ended precisely at the ends of
807　introns (we considered only introns that end in AG dinucleotide) after removal of the 3'
808　adapter. We noticed that these reads had an 3.5× increased frequency of mismatches on
809　the A as the first nucleotide compared to remaining iCLIP reads (Supplementary Fig. 3a),
810　indicating that these mismatches may have resulted from truncation at the three-way-
811　junction formed at the BP (Fig. 2c). We therefore trimmed the first nucleotide from the
812　read if it contained a mismatch at the first position that corresponded to a genomic
813　adenosine. We then used spliceosome iCLIP from Cal51 cells to identify all reads that
814　ended precisely at the ends of introns and defined the position where these reads started
815　and assessed the random barcode nucleotides that are present at the beginning of each
816　iCLIP read to count the number of unique cDNAs at each position. The nucleotide
817　preceding the read start corresponds to the position where cDNAs truncated during the
818　reverse transcription, and we selected the genomic A that had the highest number of
819　truncated cDNAs as the candidate BP. If two positions with equal number of cDNAs were
820　found, we selected the one closer to the 3'ss. Together, this identified 43,637 BPs.

821　We also attempted to use truncated cDNAs from PRPF8 eCLIP for discovery of BPs, but
822　found that the number of cDNAs overlapping with intron ends was much smaller than in
823　spliceosome iCLIP, and was insufficient for BP discovery. This is most likely because of
824　the high amount of non-specific background signal in PRPF8 eCLIP, which leads to a lower
825　proportion of cDNAs that align to the BPs.

826　Bedtools Intersect command using option –u was used to compare BP coordinates from
827　spliceosome iCLIP to the BPs identified in previous studies. We restricted this comparison
828　to introns where BPs were detected by all three datasets (iCLIP, RNA-seq and
829　computational prediction).

830　To define a single 'computational BP' per intron, the BP positions computationally
831　predicted for each intron in hg19 were obtained from
832　http://bejerano.stanford.edu/labranchor/, and the top scoring BP in each intron was
833　used. To define a single 'RNA-seq BP' per intron, we used the BP with most lariat-spanning
834　reads in each intron.

835　**Analysis of pairing probability**

836　Computational predictions of the secondary structure were performed by RNAfold
837　function from Vienna Package (https://www.tbi.univie.ac.at/RNA/) with default

838 parameters[25]. The RNAfold results are provided in a customized format, where brackets
839 are representing the double-stranded region on the RNA and dots are used for unpaired
840 nucleotides. We measured the density of pairing probability by summing the paired
841 positions into a single vector.

**Identification of RBPs overlapping with spliceosomal peaks**

843 For RBP enrichment in Fig. 3, we used the eCLIP data from the ENCODE consortium[16],
844 together with available iCLIP experiments from our lab (which are all listed in
845 Supplementary Data Set 4), to see if any of the proteins are enriched in the region of
846 spliceosomal peaks. In total, this included 157 eCLIP samples of 68 RBPs in the HepG2 cell
847 line, and 89 RBPs in the K562 cell line, and iCLIP samples of 18 RBPs from different cell
848 lines (Supplementary Data Set 4). Next, we intersected cDNA starts from each sample to
849 the -100 to +50 nt region relative to the 3'ss and used it as control for each of the following
850 peaks: Peak 4 (-23 nt..-29 nt relative to BP), Peak 5 (-21 nt..-17 nt relative to BP), Peak B
851 (-1 nt..1 nt relative to BP), Peak A (-1 nt..1 nt relative to 5'ss), Peak 6 (-11 nt..-10 nt relative
852 to 3'ss), Peak 7 (-3 nt..-2 nt relative to 3'ss). The positions of these peaks were determined
853 based on crosslink enrichments in spliceosome iCLIP.

**Statistics**

855 All statistical analyses were performed in the R software environment (version 3.1.3 and
856 3.3.2, https://www.r-project.org).

**Reporting Summary**

858 Further information on experimental design is available in the Nature Research Reporting
859 Summary linked to this article.

**Code availability**

861 The code to identify BPs from spliceosome iCLIP reads is publicly available at the GitHub
862 repository (https://github.com/nebo56/branch-point-detection-2).

**Data availability**

864 The spliceosome iCLIP data generated and analyzed during the current study are
865 available on EBI ArrayExpress under the accession number E-MTAB-8182, and are also
866 available in raw and processed format on https://imaps.genialis.com/iclip. Additional
867 datasets used in this study are listed in Supplementary Data Set 4. Source Data for Fig. 1c
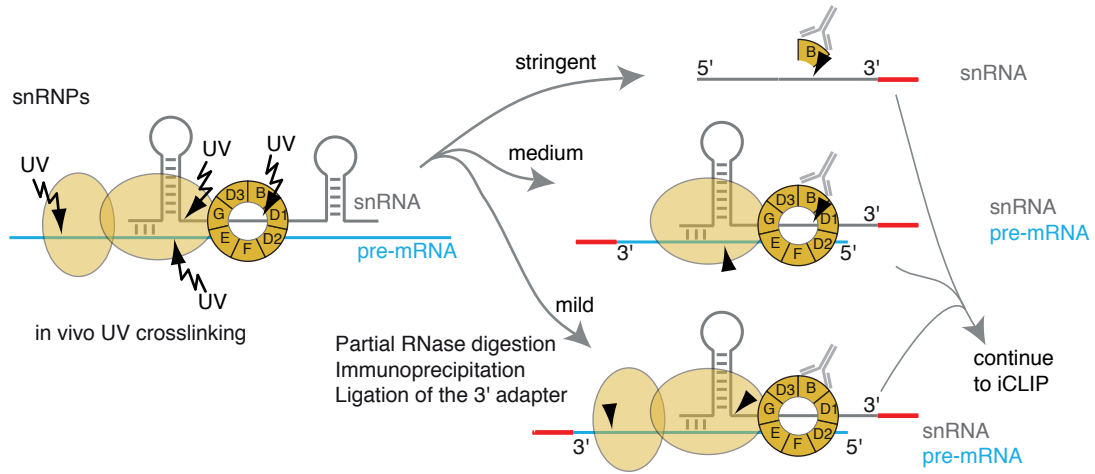868 are available online. Other data are available upon request.
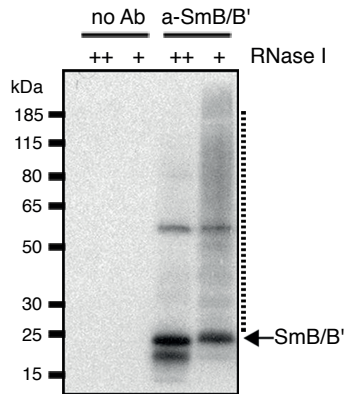
869

870 **Methods-only references**

871 26. Blazquez, L. et al. Exon Junction Complex Shapes the Transcriptome by
872       Repressing Recursive Splicing. *Mol Cell* **72**, 496-509 e9 (2018).
873 27. Chakrabarti, A., Haberman, N., Praznik, A., Luscombe, N.M. & Ule, J. Data
874       Science Issues in Studying Protein–RNA Interactions with CLIP
875       Technologies. *Annual Review of Biomedical Data Science* **Vol. 1**(2018).
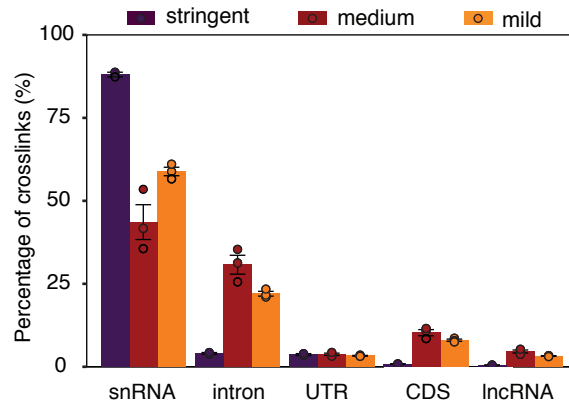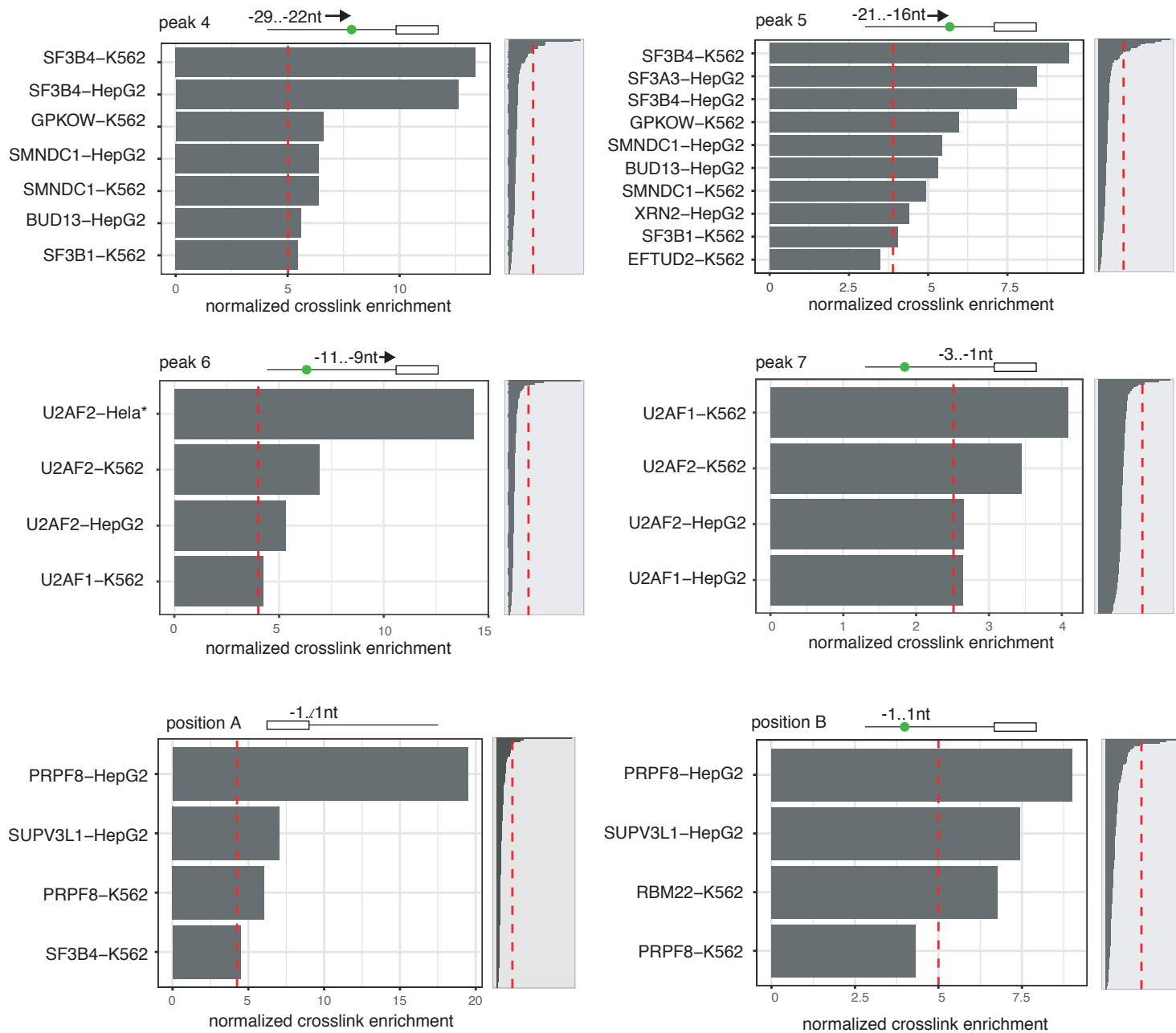876

# Figure 1

**a**



**b**
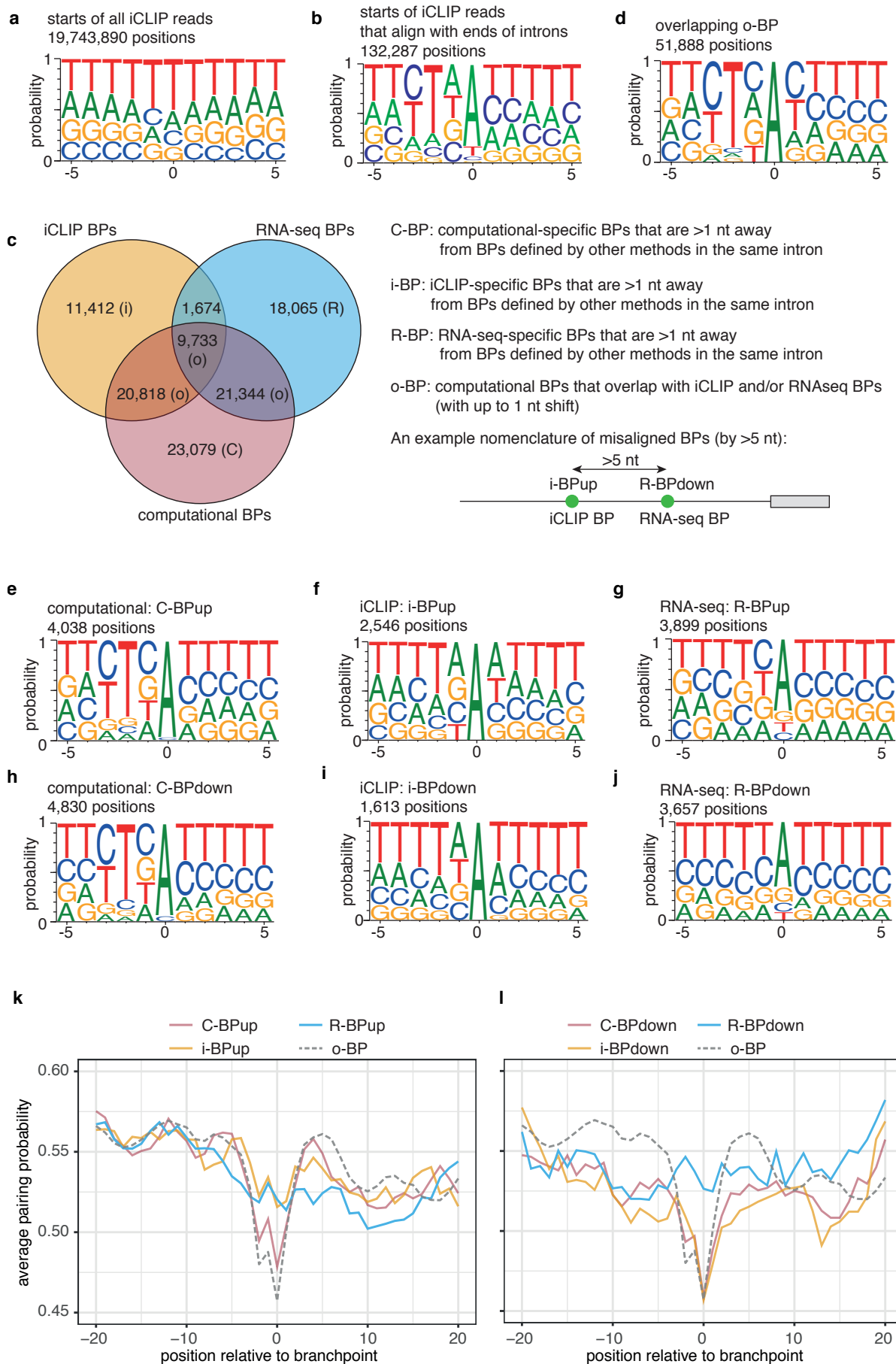


**c**

Figure 2



a



b



c

Figure 3

# Figure 4

**a** starts of all iCLIP reads
19,743,890 positions

**b** starts of iCLIP reads
that align with ends of introns
132,287 positions

**d** overlapping o-BP
51,888 positions

**c**

iCLIP BPs          RNA-seq BPs

11,412 (i)      1,674      18,065 (R)

9,733
(o)

20,818 (o)      21,344 (o)

23,079 (C)

computational BPs

C-BP: computational-specific BPs that are >1 nt away
    from BPs defined by other methods in the same intron

i-BP: iCLIP-specific BPs that are >1 nt away
    from BPs defined by other methods in the same intron

R-BP: RNA-seq-specific BPs that are >1 nt away
    from BPs defined by other methods in the same intron

o-BP: computational BPs that overlap with iCLIP and/or RNAseq BPs
    (with up to 1 nt shift)

An example nomenclature of misaligned BPs (by >5 nt):

```
              >5 nt
    i-BPup        R-BPdown
                              ▭
    iCLIP BP    RNA-seq BP
```

**e** computational: C-BPup
4,038 positions

**f** iCLIP: i-BPup
2,546 positions

**g** RNA-seq: R-BPup
3,899 positions

**h** computational: C-BPdown
4,830 positions

**i** iCLIP: i-BPdown
1,613 positions

**j** RNA-seq: R-BPdown
3,657 positions

**k**

— C-BPup    — R-BPup
— i-BPup    --- o-BP

**l**

— C-BPdown    — R-BPdown
— i-BPdown    --- o-BP

# Figure 5



| | misaligned BP is upstream of a BP identified by another method | | | | misaligned BP is downstream of a BP identified by another method | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| BP count: | 4038 | 2546 | 3899 | 51895 | 4,830 | 1,613 | 3657 |
| SF3B4 eCLIP, K562: | 16078 | 4947 | 9791 | 239162 | 10276 | 38204 | 34116 |
| SF3B4 eCLIP, HepG2: | 13545 | 5396 | 7326 | 192818 | 8248 | 7556 | 4903 |

# Figure 6

# Figure 7