

 Open access • Posted Content • DOI:10.1101/2020.06.21.162891

ASpli2: Integrative analysis of splicing landscapes through RNA-Seq assays

— [Source link](#) 

Estefania Mancini, Andres Rabinovich, Javier Iserte, Marcelo J. Yanovsky ...+1 more authors

Institutions: Fundación Instituto Leloir

Published on: 22 Jun 2020 - bioRxiv (Cold Spring Harbor Laboratory)

Topics: Alternative splicing and RNA splicing

Related papers:

- [Integrative Deep Models for Alternative Splicing](#)
- [Manananggal - a novel viewer for alternative splicing events.](#)
- [Quantifying alternative splicing from paired-end RNA-sequencing data](#)
- [SPLOOCE: a new portal for the analysis of human splicing variants.](#)
- [SUPPA2 provides fast, accurate, and uncertainty-aware differential splicing analysis across multiple conditions](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/aspli2-integrative-analysis-of-splicing-landscapes-through-3mvjff28ax>

ASpli2: Integrative analysis of splicing landscapes through RNA-Seq assays

Estefania Mancini¹, Andres Rabinovich¹, Javier Iserte, Marcelo Yanovsky*,
Ariel Chernomoretz*

Buenos Aires, Argentina

Abstract

Genome-wide analysis of alternative splicing has been a very active field of research since the early days of NGS (Next generation sequencing) technologies. Since then, ever-growing data availability and the development of increasingly sophisticated analysis methods have uncovered the complexity of the general splicing repertoire. However, independently of the considered quantification methodology, very often changes in variant concentration profiles can be hard to disentangle. In order to tackle this problem we present ASpli2, a computational suite implemented in R, that allows the identification of changes in both, annotated and novel alternative splicing events, and can deal with complex experimental designs.

Our analysis workflow relies on the analysis of differential usage of sub-genic features in combination with a junction-based description of local splicing changes. Analyzing simulated and real data we found that the consolidation of these signals resulted in a robust proxy of the occurrence of splicing alterations. While junction-based signals allowed us to uncover annotated

*Corresponding author

¹Equally contributed

as well and non-annotated events, bin-associated signals notably increased recall capabilities at a very competitive performance in terms of precision.

Keywords: Alternative splicing, RNAseq,

1. Introduction

The vast majority of protein coding genes in eukaryotic organisms are transcribed into precursor RNA messenger molecules (pre-mRNA) carrying protein coding regions (exons) interleaved by non-coding ones (introns). The later are removed in a co-transcriptional dynamical maturation process called splicing. Alternative splicing (AS) occurs whenever distinct splicing sites are selected in this process resulting in different mature mRNA molecules [1, 2].

Far from being an exception, it was found that AS is a rather common mechanism of gene regulation that serves to expand the functional diversity of a single gene allowing the generation of multiple mRNA isoforms from a single genomic locus [3]. Five basic modes of AS are generally recognized: the skipping of a given exon (exon skipping, ES), the exon elongation/contraction produced by the use of alternative 5' donor (Alt5') or 3' acceptor (Alt3') sites respectively, the retention of an intronic stretch in the mature mRNA form (intron retention IR), and the alternative use of mutually exclusive exons (MEX). These canonical forms of AS are prevalent among eukaryotes, although their relative incidence might vary between them [4]. Despite their ubiquity, these simple patterns that mainly involve binary choices of exons, donor and acceptor sites, do not exhaust the splicing repertoire. On the contrary, much more complex biologically relevant patterns could arise [5, 6]. In practice the study of AS faces many technical challenges that cause that every

22 quantitative approach typically suffers methodological limitations. Despite
23 the use of different statistical approaches, some methods consider only pre-
24 existing known annotation, some can exclusively handle canonical splicing
25 events and some can only handle pairwise comparisons between conditions
26 (for a comprehensive review see [7, 8, 9]).

27

28 The analysis of AS at genomic scale started-in with microarray technolo-
29 gies [10, 11] and nowadays is routinely probed using RNAseq assays [12, 13].
30 Roughly speaking, there are three main computational approaches to study
31 splicing diversity from RNAseq data. For one hand there are transcript re-
32 construction methods, like MISO [14] or Cufflink [15] that aim to infer a
33 probabilistic model of the frequency of each isoform from the read distribu-
34 tion mapped to a given gene. In the same spirit, Kallisto [16] and Salmon[17]
35 are two recently introduced methods that leverage on light-weight pseudo-
36 alignment heuristics to quantify transcript abundances. For the other hand,
37 methods like DEXSeq [18] , edgeR [19, 20], or voom-limma [21], focus on
38 the analysis of differential usage of subgenic features (e.g. exons) between
39 conditions. Finally, there are also methods like rMATS [22], MAJIQ [5] or
40 LeafCutter [23] that leverage on junction information to infer both, anno-
41 tated and novel splicing events.

42

43 Differently from other approaches, ASpli2 was specifically designed to in-
44 tegrate several independent signals in order to deal with the complexity that
45 might arise in splicing patterns. Taking into account genome annotation in-
46 formation, ASpli2 considers bin-based signals along with junction inclusion

47 indexes in order to assess for statistically significant changes in read cover-
48 age. In addition, annotation-independent signals are estimated based on the
49 complete set of experimentally detected splice junctions. Noticeably, ASpli2
50 can provide a comprehensive description of genome-wide splicing alterations
51 even for non-trivial experimental designs. Our approach relies on a gener-
52 alized linear model framework (as implemented in edgeR R-package [24]) to
53 assess for the statistical analysis of specific contrasts of interest.

54 In order to weigh ASpli2's performance we compared it against three
55 different state-of-the-art methodologies: rMATS [22], LeafCutter [23] and
56 MAJIQ [5]. The first one is a widely used piece of software that can in-
57 tegrate coverage and junction information to assess for changes in splicing
58 patterns. Additionally, LeafCutter and MAJIQ are two recently introduced
59 methodologies that are widely used by the bioinformatics community. Both
60 approaches focus on the analysis of clusters of junctions to study local splic-
61 ing patterns of varying complexity. However, they differ in many technical
62 and statistical aspects [5]. For instance, LeafCutter was not designed to han-
63 dle intron retention events and considers a Dirichlet-multinomial generalized
64 linear model to test for differential intron excision between two groups of
65 samples. MAJIQ, on the other hand, relies on a bayesian estimation of the
66 posterior Percent Selected Index to identify splicing affected junctions.

67 Other methodolgies like DEXSeq [18], edgeR [24], or voom [21] are also
68 widely considered for bioinformatics analysis as they are very versatile tools
69 to analyse differential usage of exons from RNA-seq data. In fact, ASpli2
70 makes use of the genome-binning scheme originally presented in DEXseq to
71 quantify read coverage signals (see Sup.Mat.8.1) and leverages on the statis-

72 tical framework developed in edgeR to estimate robust splicing signals (see
73 Section 3.1). These methodologies constitute great toolboxes to implement
74 ad-hoc analysis. However, as they do not intend to provide *per se* self-
75 contained solutions that produces final reports starting from read-alignment
76 input data they were not explicitly included in our analysis.

77 The paper is organized as follows. In Section 2.2 we analyzed a simu-
78 lated dataset in order to evaluate the specificity and sensitivity of ASpli2
79 discoveries. These results were contrasted against LeafCutter, MAJIQ and
80 rMATS outcomes in order to analyse ASpli2 performance. In Section 2.3 we
81 explored the ability of ASpli2 to uncover consistent splicing-patterns from
82 two independent RNAseq assays that probed the same biology. We focused
83 on the alterations of splicing patterns of *A. thaliana* transcriptome caused
84 by the knock out of PRMT5, a methyl transferase that, among other pro-
85 teins, targets several Sm spliceosomal proteins [25, 26, 27]. This analysis was
86 also performed with the other considered methodologies in order to compare
87 their capacity to generate reproducible results. In addition, we capitalized on
88 ASpli2's ability to handle complex experimental designs to produce a consol-
89 idated data-set from the independent assays. In this section we also aimed to
90 quantify the level of agreement of ASpli2, LeafCutter, MAJIQ and rMATS
91 discoveries with qRT-PCR based alternative splicing evidence. To that end,
92 we took advantage of two independent studies that analyzed splicing altered
93 events in PRMT5 mutants using qRT-PCR assays [25, 28]. Finally, in Sec-
94 tion 2.4, we considered a 28 samples paired-study of human prostate cancer
95 [29]. Using this data-set we analyzed how the performance, time and memory
96 requirements scaled with the number of considered samples in a paired exper-

97 imental design. Finally, we discussed our results in Section 4 and presented
98 our conclusions in Section 5.

99 **2. Results**

100 *2.1. ASpli2 workflow*

101 ASpli2 was designed as a flexible R package to carry out all the major
102 tasks required for gene expression and splicing analysis. A typical ASpli2
103 workflow involves: parsing the genome annotation into subgenic features
104 called bins, overlapping read alignments against them, perform junction
105 counting, fulfill inference tasks of differential bin and junction usage and,
106 finally, report integrated splicing signals. A workflow diagram and a sum-
107 mary of ASpli2 core functionality can be found as Supplementary Figures S1
108 and S2 respectively. As shown in Figure S1, at every step ASpli2 generates
109 self-contained outcomes that, if required, can be easily exported and inte-
110 grated into other processing pipelines. Supplementary Figure S3 shows an
111 example of the interactive html report generated as a final output. A detailed
112 description of ASpli2 functionality is included in ASpli2's R vignette, which
113 is provided as Supplementary Material.

114

115 *2.2. Synthetic dataset*

116 Changes in splicing patterns were simulated in a treatment vs control
117 setup for genes of the chromosome-one of the *Arabidopsis thaliana* plant
118 genome (three samples per condition). In our simulations, the differential

119 usage of splicing variants affected 2451 genomic bins in 915 genes (see Ma-
120 terial and Methods 3.3).

121

122 The ASpli2 analysis pipeline provided three different cues to probe splic-
123 ing occurrence. Statistically significant evidence is collected from: bin cov-
124 erage differential signals, junction anchorage changes and variations inside
125 junction clusters (see Material and Methods 3.1). We considered bin-coverage
126 signals with statistically significant differential coverage changes ($fdr < 0.05$)
127 that presented either a larger than three-fold coverage fold-change or, alter-
128 natively, a change in bin-supporting junction inclusion indices larger than 0.2
129 . For junction based signals, on the other hand, locale and anchorage indices
130 were required to present statistically significant changes ($fdr < 0.01$) and also
131 should display usage signal variations larger than a 0.3 level (see Material
132 and Methods 3.2)

133

134 In Table 1 we reported the number of correctly detected simulated events,
135 number of false positives and number of events exclusively detected by each
136 kind of signal: bin-coverage, junction-locale and junction-anchorage. Over-
137 laps between discoveries reported by each kind of signal were graphically
138 reported in panel (A) of Figure 1.

139 It can be seen that ASpli2 correctly uncovered 974 (40%) of the 2451
140 simulated bin events. Moreover, we found that most of the ASpli2 undetected
141 simulated events (1341 out of 1477) took place in genes that did not present
142 enough expression levels over the analysed conditions and therefore were
143 filtered out before any statistical testing (see 3.2). In fact, only 136 out of

144 the 1110 events (12%) that did pass the gene-expression pre-filtering step
145 were found to be false negative cases.

146 About 95% of ASpli2 true discoveries were identified by the analysis of
147 significant changes at the bin-coverage level. Junction-based detection, on
148 the other hand, could correctly identified 574 simulated events (60% of true
149 discoveries). The null overlap between locale and anchorage detection illus-
150 trated that they probed complementary aspects of splicing events. Addi-
151 tionally, it can be appreciated that 41% (399) of the true discoveries were
152 only detected by bin-coverage signals, whereas junction-based analysis con-
153 tributed only 5% (50) of specific detections. A graphical summary of the
154 detection signal landscape can be appreciated in panel (A) of Figure 1.

155

ASpli2 signal	TP	FP
bin coverage	924 (399)	42
junction locale	393 (35)	2
junction anchorage	182 (15)	6
overall	974	48

Table (1) Splicing detection performance of the three different ASpli2 signals. True positive and false positive calls are shown in the second and third columns respectively. The number of specific discoveries exclusively reported by each signal is reported between brackets.

156 We decided to further characterized some aspects of bin-coverage detec-
157 tion calls, as this signal provided the major number of discoveries. It can be
158 seen in panel-(B) of Figure 1 that fold-change and junction-support signals

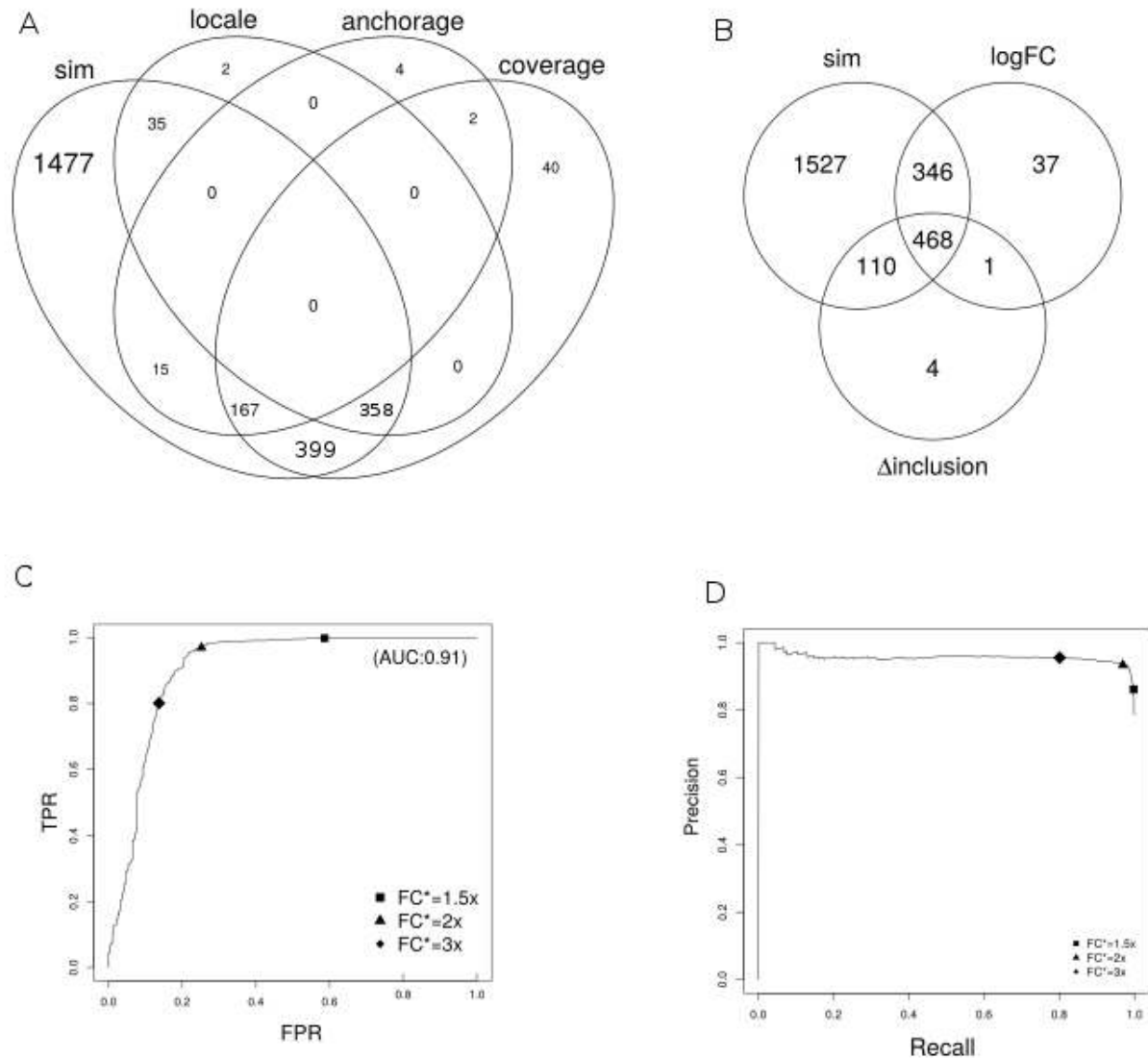


Figure 1 (a) Distribution of detection call produced by different ASpli2 signals. (b) Graphical summary of bin-coverage detection calls. The *sim* set correspond to simulated events. *logFC* and *D-inclusion* sets are associated to statistically significant discoveries presenting large enough fold change and large bin-supporting junction inclusion signals respectively. ROC and Precision-Recall curves, parameterized by the considered fold-change threshold level, are shown for statistically significant bins in panels (c) and (d) respectively.

159 used in the bin-coverage analysis reported relevant and non-redundant infor-
160 mation. Whereas the first one accounted for 37% of true positive instances
161 exclusively detected by this signal, the second one accounted for the specific
162 identification of 12% of the total number of true events. The impact of the
163 selected fold-change threshold value, FC^* , on specificity, precision and recall
164 can be appreciated with the aid of the Receiver-Operator and Precision-
165 Recall curves shown in panels (C) and (D) of Figure 1. It can be recognized
166 from these figures that with the adopted 3-fold threshold ASpli2 achieved
167 high recall and precision levels ($\sim 80\%$ and $\sim 95\%$ respectively) laying at
168 rather moderate levels of false positive rates ($\sim 14\%$).

169

method	size	precision	recall
ASpli2	1022 (631)	0.95 (0.99)	0.40 (0.68)
ASpli2 _c	966 (591)	0.96 (0.99)	0.38 (0.64)
ASpli2 _j	583 (456)	0.99 (0.99)	0.23 (0.50)
LeafCutter	204 (163)	0.93 (0.91)	0.08 (0.16)
MAJIQ	538 (381)	0.84 (0.87)	0.18 (0.36)
rMATS	405 (352)	0.87 (0.91)	0.14 (0.35)

Table (2) Number of discoveries, precision and recall levels are reported for different detection methodologies. *ASpli_c* and *ASpli_j* correspond to ASpli2 discoveries detected using just coverage or just junction signals respectively. Values between parenthesis report quantities estimated at gene-level.

170 In Table 2 we reported the detection performance of ASpli2 and the re-
171 sults obtained by other state-of-the-art algorithms (see Sup Mat 8.2 and 8.3

172 for calculation details). Precision and recall values estimated at gene-level (in
173 which a gene was reported as a discovery whenever at least one alternative-
174 splicing event was detected within its genomic range) were reported between
175 parenthesis. ASpli2 outcomes considering only coverage signal or just junc-
176 tion signals were included in the table as ASpli2_c and ASpli2_j rows respec-
177 tively.

178 It can be seen from the table that even though all tested algorithm shown
179 rather high precision values, ASpli2 benefited from larger recall scores than
180 any other methodology. Moreover, it can be appreciated that ASpli2_j dis-
181 played only marginally larger recall levels than other methodologies implying
182 that ASpli2 leveraged on coverage signals to increase this figure of merit. All
183 of these results suggested that ASpli2 was capable of reliably detect the simu-
184 lated splicing events achieving notably high recall values at very competitive
185 levels of precision and specificity.

186 2.3. Reproducibility Analysis

187 As we mentioned in the introduction, PRMT5 is a methyl transferase
188 that, among other proteins, targets several spliceosome proteins. Its deletion
189 has been proved to provoke major splicing alterations [25, 26]. We analyzed
190 two independent RNAseq assays that were conducted at different times prob-
191 ing the same biology. Experiments *A* (GSE149429) and *B* (GSE149430) were
192 originally carried out to analyze splicing alterations in the PRMT5 knock-
193 out mutant in *Arabidopsis thaliana*. Both assays probed the PRMT5-KO
194 and wild-type transcriptomes in Columbia ecotype plants as part of larger
195 and different studies (see Material and Methods 3.4).

196 The rationale behind our analysis was two-fold. For one hand we wanted

197 to assess for ASpli2 detection performance in a more realistic setup. For
198 the other we wanted to take advantage of these datasets to quantitatively
199 estimate the reproducibility of discoveries, i.e. we wanted to explore the
200 consistency and robustness of experimentally identified alternatively splicing
201 events in biologically related systems.

202 *2.3.1. Reproducibility assessment*

203 We analyzed RNAseq assays *A* and *B* with ASpli2 and the other consid-
204 ered algorithms. For ASpli2, we used the same detection-call criteria specified
205 in Section 2.2. Default parameters were considered to run the other tested
206 methodologies (command lines used to execute them were included as Sup-
207plementary Material 8.2). For LeafCutter and rMATS we considered events
208 presenting fdr corrected pvalues smaller than 0.05 and changes in junction
209 inclusion indices larger than 0.1. For MAJIQ we sought for events present-
210ing a posterior probability larger than 0.95, of having a change in inclusion
211index larger than a 0.2 level. Overall, 6350, 951, 412 and 158 genomic regions
212affected by altered splicing patterns were reported by ASpli2, LeafCutter,
213MAJIQ and rMATS algorithms respectively in at least one experiment.

214

215 In Table 3 we summarized reproducibility statistics for each examined
216 methodology (a more in-depth comparison of discoveries was included as sup-
217plementary material in Section 8.5). Column *universe* of Table 3 reports the
218 actual number of sub-genic regions that, upon passing different pre-filtering
219 steps, were actually examined for statistically significant changes in splicing
220 patterns. The extent of this background list was noticeably larger for ASpli2
221 as our methodology tested not only junction-related signals but also alter-

222 ations in the usage of genomic bins. Columns A and B outline the number
 223 of regions reported as differentially spliced in each experiment and column
 224 $A \cap B$, the discovery intersection size (i.e. number of sub-genic regions re-
 225 ported as differentially spliced in both data-sets). In parenthesis, we included
 226 the *overlap coefficient value*, defined as $A \cap B / \min(A, B)$. Expected over-
 227 laps, fold enrichment (i.e. ratio between observed and expected overlaps)
 228 and p-values were estimated using the SuperExactTest R-package [30] and
 229 reported in EO , FE and $pval$ columns respectively.

Method	universe	A	B	$A \cap B$	EO	FE	pval
ASpli2	140191	4687	3904	2241 (0.57)	130.5	17.2	0.0e+00
leafCutter	8113	603	675	327 (0.54)	50.2	6.5	3.6e-219
MAJIQ	16441	277	284	149 (0.54)	4.8	31.1	9.5e-203
rMATS	2405	310	401	310 (1.00)	19.4	16.0	0.0e+00

Table (3) Reproducibility statistics. The numbers of statistically analyzed genes (after prefiltering steps) for each algorithms are shown in the *universe* column. The number of splicing events reported for each experiment and the number of concordant discoveries are displayed at columns A , B , and $A \cap B$ respectively. The expected overlap, fold enrichment level and significance pvalue are displayed in columns EO, FE and pval respectively

230 It can be seen from Table 3 that, for all the examined methods, the
 231 agreement between experiments was highly significant. In all cases, more
 232 than 50% of events detected in one experimental instance was also reported
 233 in the other. At the same time it can be appreciated that ASpli2 provided
 234 the largest (and highly significant) overlap-set. Noticeably, the total number
 235 of concordant splicing-affected genomic regions detected by ASpli2 presented

236 up to a 15-fold increase with respect to the size of concordant sets reported
237 by others methodologies.

238

239 Overall our analysis showed that results obtained at different and inde-
240 pendent experimental instances were reproducible, in the sense that statis-
241 tically significant agreement was found for every methodology. These re-
242 sults were robust against using different overlap quantification criteria (see
243 Sup.Mat 8.5). In this matter, and similarly to the results obtained on the
244 synthetic dataset, our results on PRMT5 data showed that ASpli2 displayed
245 high recall levels providing the largest list of concordant discoveries between
246 experiments.

247 *2.3.2. Data consolidation*

248 Up to now, we focused on the analysis of the intersection of set of discov-
249 eries as a measure of coherence of the results. Now we wanted to illustrate
250 how ASpli2 capabilities to deal with complex experimental designs can be
251 used to integrate experimental results in a more statistically sound way.

ASpli2 was used to consolidate datasets A and B considering a simple
generalized linear model:

$$y \sim \textit{experiment} + \textit{genotype} + \textit{experiment} : \textit{genotype} \quad (1)$$

252 ‘experiment’ was a fixed effect to cope with specific technical biases, and the
253 ‘genotype’ factor was meant to capture the PRMT5 vs wild-type effect. The
254 third term was an interaction term, and was used to enforced the exclusion
255 of non-coherent signals between experiments.

256 ASpli2 detected 4360 genomic regions displaying strong evidence of a

257 genotype effect ($\text{fdr} < 0.05$). In addition, 99% of these PRMT5-related
258 events (4314 out of 4360) passed a filtering step to enforce they presented no-
259 detectable evidence of experiment-genotype interactions (experiment:genotype
260 associated $\text{fdr} > 0.5$). These 4314 events defined the consolidated AB data
261 set.

262 We found that 99% (2209 out of 2241) of the concordant discoveries in-
263 dependently detected in both assays were also included in the consolidated
264 dataset (we included a Venn diagram of the discoveries reported for experi-
265 ments A, B, and the consolidated data-set AB in Sup.Fig. S7). Noticeably,
266 the consideration of the AB data-set allowed to almost double the number of
267 detected genomic regions displaying robust evidence of differential splicing
268 patterns.

269 2.3.3. PRMT5 RT-PCR detected events

270 The PRMT5 methyltransferase has been the target of many studies as
271 deficiencies in this protein causes genomewide splicing alterations[26, 27, 28].
272 In this section we focused on two specific works that provided independent
273 RT-PCR validated lists of splicing alterations linked to PRMT5 in *Arabidop-*
274 *sis thaliana* [25, 28].

275 For one hand, Deng and collaborators studied PRMT5 mutant *Ara-*
276 *bidospis thaliana* plants and presented a list of 12 RT-PCR validated intron
277 retention events (see Fig 2 in [28]). On the other, using the same biological
278 model, Sanchez and collaborators indentified changes in alternative splicing
279 using a high-resolution qRTPCR panel that included several known alterna-
280 tive splicing events [31]. They found that PRMT5 mutants had significant
281 alterations in 44 events which included exon skipping, alternative donor and

282 acceptor splice sites, as well as intron retention events (Supplementary Table
283 4 in [25]).

284 We aimed to contrast these findings with the results reported by the dif-
285 ferent methodologies on datasets A and B. In Table 4 we summarized, for
286 each study, the number of concordant findings uncovered by different al-
287 gorithms on datasets A and B. Quantities between brackets represent the
288 number of ASpli2 discoveries reported by coverage and junction-based sig-
289 nals respectively. It can be seen that ASpli2 recovered the largest number
290 of events and that the majority of ASpli2 validated discoveries originated in
291 differential coverage signal calls. Had we only considered junction related
292 detection-signals, ASpli2 would have achieved similar levels of agreement
293 than the other junction-based algorithms (for instance we got a similar per-
294 formance than LeafCutter on Sanchez data-set for the consolidated case).

295

296 In Table S2, included as supplementary material, we further character-
297 ized the agreement between the 23 splicing events that ASpli2 uncovered for
298 the consolidated AB case, and Sanchez qRT-PCR validated events. It can
299 be seen that in 15 out of the 23 cases (65%), the very genomic region probed
300 by the PCR analysis was recognized by ASpli2. For the other 8 cases, AS-
301 pli2 detected actually occurring changes in isoform usage but from splicing
302 signals originating at genomic locations not probed by the PCR primers (See
303 Supplementary Figures - TODO: ACA VAN SAHIMI PLOT DE EVENTOS
304 PCR).

Method		Deng 2010	Sanchez 2010
RT-PCR		12	44
ASpli2	AB	8 [8,4]	23 [21,13]
	A	10 [8,5]	24 [19,7]
	B	9 [8,2]	20 [18,6]
LeafCutter	A	3	16
	B	3	17
MAJIQ	A	5	8
	B	2	8
rMATS	A	1	12
	B	1	3

Table (4)

305 *2.4. ASpli2 scalability analysis*

306 In this section we leveraged on a mid-size RNAseq study presented by
307 Ren and collaborators to characterize aberrant splicing patterns occurring
308 in prostate cancer patients [29]. We aimed to analyze this sample-paired
309 assay to see how ASpli2 performance (statistical power, precision, time and
310 memory requirements) scaled with the number of samples. In particular, we
311 followed the approach suggested in [32] to characterize ASpli2 in terms of
312 statistical power and expected false discovery rate for a varying number of
313 samples.

314 *2.4.1. Statistical power*

315 Ren and collaborators presented a comprehensive study of splicing alter-
316 ations detected using RNAseq transcriptome profiles of 14 primary prostate
317 cancers and their paired normal counterparts from the Chinese population
318 [29]. On average, the 28 fastq files presented 34.6 ± 1.7 million reads per
319 sample and 31.4 ± 1.6 millions of them were actually mapped to the EN-
320 SEMBL HG38.98 version of the human genome (see Material and Methods
321 3.6). The genome's GTF and BAMs files were then used as inputs to drive
322 an alternative splicing paired-sample analysis with ASpli2. We considered
323 the following model to identify genomic regions differentially spliced in tumor
324 samples compared to normal tissue controls:

$$y \sim \textit{patient} + \textit{tissue} \quad (2)$$

325 The 'patient' term served to pair tumor and normal tissue samples coming
326 from the same individual. The two-level 'tissue' factor reported average
327 differences between tumor and normal cases over the observed population of
328 patients.

329 In order to study the dependency of the statistical power on the number
330 of samples, we sampled without replacement (10 times) subsets of 3, 5, 7 and
331 10 individuals. For each case, we reported, in the first column of Table 5, the
332 median (and standard error, in brackets) of the number of genomic regions
333 found to be alternatively spliced between tumor and normal samples.

334 In order to estimate false discovery rates we considered mock comparisons
335 between normal samples (we sampled 10 times normal tissue samples of 3vs3,
336 5vs5 and 7vs7 individuals). We then estimated FDR as the ratio between

337 the number of mock discoveries and the median number of discoveries found
338 in true comparisons of the same number of samples. In the second column
339 of Table 5 median and standard errors (in brackets) were reported.

Samples	Splicing events		Affected genes	
	number	FDR	number	FDR
3x3	67 (155)	0.2 (0.4)	44 (113)	0.25 (0.03)
5x5	486 (387)	0.02 (0.002)	371 (218)	0.02 (0.002)
7x7	759 (220)	0.005 (0.02)	481 (131)	0.004 (0.0007)
10x10	850 (418)	-	664 (191)	-
14x14	1465	-	1030	-

Table (5) Summary of the 10-fold bootstrapped analysis of ASpli2 performance on the prostate cancer data set. For each number of paired samples (first column) the median number of genomic-regions displaying a statistically significant ‘tissue’ effect were included in the second column. Median values of false discovery rate estimations obtained from the analysis of normal-tissue samples were shown in the third column. Standard error estimation were reported between brackets.

340 It can be seen from Table 5 that the median number of detected splicing
341 events increased with the number of examined samples, up to a maximum of
342 1465 events obtained when the 28 paired samples were considered. The large
343 variability observed between bootstrap realizations was consistent with the
344 large variability already observed across prostate cancer transcriptomes (see
345 [29] and Supplementary Material 8.7). FDR estimated values showed a huge
346 decrease with increasing number of samples, and for the 5x5 case seemed to
347 have already leveled off. Similar trends were observed when splicing alter-

348 ations were reported at the level of hosting genes (data not shown).

349 2.4.2. Time and memory requirements

350 In Table 6 we reported median values and standard errors for the elapsed
351 time and peak memory usage required for calculations (performed on single
352 thread on an Intel Xeon Silver 4116 2.1GHz Lenovo ThinkSystem SR650)

Samples	time [min]	memory peak[Gb]
3x3	67 (1)	20.25 (0.38)
5x5	111 (2)	22.15 (0.20)
7x7	156 (4)	24.13 (0.03)
10x10	231 (5)	26.57 (0.04)
14x14	348	30.07

Table (6) Summary of the 10-fold bootstrapped analysis of ASpli2 performance on the prostate cancer data set. For each number of paired samples (first column) the median number of genomic-regions displaying a statistically significant ‘tissue’ effect were included in the second column. Median values of false discovery rate estimations obtained from the analysis of normal-tissue samples were shown in the third column. Median time and memory used in the analysis were reported in the last two columns. Standard error estimation were reported between brackets.

353 Execution time scaled linearly with the number of paired samples at a
354 rate of 25.5 minutes per pair of samples (about 90% of execution time was
355 used for BAMs reading and feature counting). The memory peak column
356 shows that RAM requirement linearly scaled with the number of samples at
357 a rate of about 880Mb per sample pair. A simple extrapolation suggests that
358 about 65Gb should be enough to handle 100 samples of the same sequencing

359 depth ($\sim 3.510^6$ reads per sample).

360 **3. Material and Methods**

361 *3.1. Differential analysis scheme*

362 ASpli2 leverages on the statistical framework developed by Smyth and
363 collaborators, implemented in the edgeR R-package [24, 20], to assess for
364 statistically significant changes in gene-expression, bin coverage and junction
365 splicing signals. Under this approach, count data is modeled using a negative
366 binomial model, and an empirical Bayes procedure is considered to moderate
367 the degree of overdispersion across units.

368 *Differential expression signals.* Differential expression signals are estimated
369 via generalized linear models (GLM). This approach allows ASpli2 to deal
370 with complex experimental designs, i.e. contrasts can be tested in experi-
371 ments with multiple experimental factors. Using this statistical setting, for
372 each gene, ASpli2 quantifies differential gene expression signals reporting the
373 corresponding log-fold change, p-value, and FDR adjusted q-values.

374 *Differential splicing signals.* In order to study splicing patterns, gene ex-
375 pression changes should be deconvolved from overall count data. On a very
376 general setting, what we are looking for is to test whether a given unit of a
377 certain group of elements displays differential changes respect to the collec-
378 tive or average behavior. ASpli2 uses this general idea to assess for statisti-
379 cally significant changes in splicing patterns probed with different genomic
380 features:

- 381 • bin-coverage signal: ASpli2 assesses for differential usage of bins com-
382 paring bin's log-fold-changes with the overall log-fold-change of the
383 corresponding gene.
- 384 • junction anchorage signal: For every experimentally detected junction,
385 ASpli2 analyzes differential intron retention changes by considering log-
386 fold-changes of a given experimental junction relative to changes in
387 coverage of left and right junction flanking regions.
- 388 • junction locale signal: In the same spirit than MAJIQ and LeafCutter,
389 ASpli2 defines junction-clusters as sets of junctions that share at least
390 one end with another junction of the same cluster (see Panel E of
391 Figure S8). In order to characterize changes for a given junction along
392 experimental conditions, ASpli2 weighs log-fold-change of the junction
393 of interest relative to the mean log-fold-change of junctions belonging
394 to the same cluster.

395 ASpli2 makes use of the functionality implemented in the `diffSpliceDGE`
396 function of the `edgeR` package to perform all of this comparisons within a uni-
397 fied statistical framework. Given a set of elements (i.e. bins or junctions) of
398 a certain group (i.e. genes, anchorage group or junction-cluster), a negative
399 binomial generalized log-linear model is fit at the element level, considering
400 an offset term that accounts for library normalization and collective changes.
401 Differential usage is assessed testing coefficients of the GLM. At the single
402 element-level, the relative log-fold-change is reported along with the associ-
403 ated p-value and FDR adjusted q-values. In addition a group-level test is
404 considered to check for differences in the usage of any element of the group

405 between experimental conditions (see *diffSpliceDGE* documentation included
406 in edgeR package for details [24]).

407 3.2. Filtering and detection criteria

408 Statistical analysis of differential splicing is performed only on expressed
409 genes (i.e. read counts spanning the gene genomic range should be larger
410 than a minimal number of reads, 5 by default, across all the samples of the
411 contrasted conditions). Furthermore, analyzed bins and junctions should
412 present a minimal number of counts (5 by default) in every replicate of at
413 least one contrasted condition. Additionally, marginally present junctions
414 are filter-out looking at the maximal value of their *participation* coefficient,
415 defined as the relative abundance of a given junction within its group for a
416 given experimental condition.

417

418 Besides statistical figures of merit, ASpli2 provides additional statistics
419 and parameters in order to ease the identification of biologically relevant
420 events. For instance, magnitude of change in inclusion or strength indices
421 (see Table S1) between experimental conditions, are also reported in order to
422 filter-out weak events. In this way, a bin is called differentially-used by ASpli2
423 if it displays statistically significant coverage changes ($\text{fdr} < 0.05$, by default)
424 and, additionally, one of the two supplementary conditions hold: either the
425 bin fold-change level is greater than a given threshold (3 fold changes, by
426 default) or changes in inclusion levels of bin-supporting junctions (ΔPIR
427 or ΔPSI according to the bin class, see Table S1) surpasses a predefined
428 threshold (0.2 by default).

429 Anchorage splicing signals, on the other hand, are reported whenever

430 statistically significant changes are found at the cluster level (cluster.fdr <
431 0.05 by default) for the considered $\{J_1, J_2, J_3\}$ junction set (see upper panel
432 of Fig S8-D) and, at the same time, $|\Delta PIR_{J_3}|$ is larger than a given threshold
433 (0.3 by default).

434 Finally, junction locale differential splicing signals are reported when-
435 ever statistically significant changes are found at the cluster level (cluster.fdr
436 < 0.01 by default) for the analysed junction cluster $\{J_1, \dots, J_S, \dots, J_n\}$ (see
437 S8-E) and, at the same time, there is at least one junction J_S within the
438 cluster presenting statistically significant changes at the single unit level
439 (junction.fdr < 0.05, by default) with $|\Delta Participation_{J_S}|$ larger than a given
440 threshold (0.3 by default). In the case that statistically significant changes
441 were detected at the unit-level for more than one junction of a given clus-
442 ter, the one displaying the largest participation change was considered and
443 reported as the cluster's representative junction.

444 3.3. Splicing simulation

445 We implemented a computational pipeline relying on the Flux Simulator
446 (FS) software [33] in order to produce a controlled set of splicing events.
447 We first used FS to generate a transcript abundance distribution template
448 to spread 15×10^6 molecules among the 10646 available transcript variants
449 of the 8433 genes of chromosome-one of the *Arabidopsis thaliana* genome.
450 Then, we generated a 'treatment' set of samples altering the original molecule
451 distribution in order to simulate genome-wide differential changes in gene
452 expression and splicing patterns.

453 Finally we simulated biological replicates from these two *seed* transcrip-
454 tomes, considering a Gamma distribution for molecule abundances to build

455 'control' and 'treatment' sample sets. We chose to work with a $CV = 0.1$
456 level of variability in gene abundance between replicates. Therefore, we con-
457 sidered *shape* ($k = 100$), and *scale* ($\theta = 0.01\mu$) parameter values, where μ
458 was the gene expression level in the corresponding *seed* transcriptome used
459 for replicate generation.

460 Simulated changes in variant concentrations produce patterns of differ-
461 ential usage at bin and junction levels according to the exonic architecture
462 of the different gene variants. For instance, a splicing alteration that in-
463 volves switching between Isoform 1 and Isoform 3 of the gene depicted in
464 Figure S8-(A) is expected to produce differential usage signals for the first
465 and third exonic bins. In our case we simulated changes in variants usage for
466 915 genes that should altered, in principle, the coverage signal of 2451 bins.
467 It is worth mentioning that as alternative splicing was modeled exclusively
468 through differential variant usage, no intron retention events were simulated
469 in the synthetic data set.

470 Several examples of splicing simulated events are depicted in Sup. Figs
471 S4,S6,S5. Examples, scripts and additional material to reproduce the ASpli2
472 analysis over this dataset can be found at the gitlab repo: https://gitlab.com/ChernoLab/aspli2_sm.

473 3.4. *PRMT5 datasets*

474 The goal of these studies was to compare the transcriptional profile (RNA-
475 seq) of wild type and PRMT5 Arabidopsis mutants plants grown under con-
476 tinuous light at 22 degrees centigrades.

477 Dataset A (GSE149429): WT (Col) and PRMT5 mutants seeds were
478 grown on Murashige and Skoog medium containing 0.8% agarose, stratified
479 for 4 d in the dark at 4 C, and then grown for fifteen days under continuous

480 white light at 22C Whole plants were harvested after 15 d. Total RNA was ex-
481 tracted with RNeasy Plant Mini Kit (QIAGEN) following the manufacturers
482 protocols. To estimate the concentration and quality of samples, NanoDrop
483 2000c (Thermo Scientific) and gel electrophoresis were used, respectively. Li-
484 braries were prepared following the TruSeq RNA Sample Preparation Guide
485 (Illumina). Briefly, 3 g of total RNA was polyA-purified and fragmented,
486 and first-strand cDNA synthesized by reverse transcriptase (SuperScript II;
487 Invitrogen) and random hexamers. This was followed by RNA degradation
488 and second-strand cDNA synthesis. End repair process and addition of a sin-
489 gle A nucleotide to the 3 ends allowed ligation of multiple indexing adapters.
490 Then, an enrichment step of 12 cycles of PCR was performed. Library vali-
491 dation included size and purity assessment with the Agilent 2100 Bioanalyzer
492 and the Agilent DNA1000 kit (Agilent Technologies)

493 Dataset B (GSE149429): WT (Col accession) and PRMT5 mutant plants
494 were grown for nine days under continuous white light at 22 degrees centi-
495 grades or exposed for 1 or 24 h to 10C on the 9th day, before harvesting.
496 Then the transcriptional profile of these plants was analyzed using RNA-seq.
497 WT (Col) and PRMT5 mutants seeds were grown on Murashige and Skoog
498 medium containing 0.8% agarose, stratified for 4 d in the dark at 4 C, and
499 then grown for nine days under continuous white light at 22C. Whole plants
500 were harvested after 9 d. Total RNA was extracted with RNeasy Plant Mini
501 Kit (QIAGEN) following the manufacturers protocols. To estimate the con-
502 centration and quality of samples, NanoDrop 2000c (Thermo Scientific) and
503 gel electrophoresis were used, respectively. Libraries were prepared following
504 the TruSeq RNA Sample Preparation Guide (Illumina). Briefly, 3 g of total

505 RNA was polyA-purified and fragmented, and first-strand cDNA synthesized
506 by reverse transcriptase (SuperScript II; Invitrogen) and random hexamers.
507 This was followed by RNA degradation and second-strand cDNA synthesis.
508 End repair process and addition of a single A nucleotide to the 3 ends al-
509 lowed ligation of multiple indexing adapters. Then, an enrichment step of 12
510 cycles of PCR was performed. Library validation included size and purity
511 assessment with the Agilent 2100 Bioanalyzer and the Agilent DNA1000 kit
512 (Agilent Technologies).

513

514 On average, 19.3 ± 5.3 million 100 long and 28.3 ± 7.7 million 150 long
515 paired-end reads were generated per sample library for datasets *A* and *B*
516 respectively. For both cases more than 96% of reads were uniquely mapped
517 to TAIR10 Arabidopsis genome using STAR (command-line invocation was
518 included in Sup Mat 8.2).

519 *3.5. Overlap analysis*

520 We followed the procedure outlined in Supplementary Material 8.3 to
521 map events reported by each of the considered method to a common set of
522 genomic coordinates. Overlaps were then estimated using the *findOverlaps*
523 function of the *IRanges* package of R [34].

524 *3.6. Prostate cancer dataset*

525 Fifty-six paired fastq files from the E-MTAB-567 experiment were down-
526 loaded from the ArrayExpress server. Reads were aligned against ENSEMBL
527 HG38.98 reference genome using the STAR aligner with default parameters
528 and a junction overhang=89.

529 *3.7. Code availability*

530 ASpli2 package is freely available at <https://gitlab.com/ChernoLab/aspli2>,
531 and will be part of the next Bioconductor release (October 2020). Examples,
532 scripts and additional material to reproduce our analysis can be found at the
533 gitlab repo: https://gitlab.com/ChernoLab/aspli2_sm.

534 **4. Discussion**

535 RNA high-throughput sequencing methods provide powerful means to
536 study alternative splicing under multiple conditions in a genome-wide man-
537 ner. However, the detection and understanding of general splicing patterns
538 still present considerable technical challenges. Here we presented ASpli2, a
539 computational suite to comprehensively test bin coverage and junction usage
540 differential splicing signals.

541 The analysis methodology implemented in ASpli2 came out as a result
542 of several software maturation cycles of our in-house splicing analysis proce-
543 dures. Over the last years, the presented core functionality has been exten-
544 sively used in different projects to study: the role of AS in circadian rhythms
545 and light response [35, 36, 37, 38, 39] as well as AS in spliceosome mutants
546 [40, 41] in *A.thaliana* model organism. In addition, ASpli2 in-house versions
547 have been used to study AS and rhythmic behavior in *D.melanogaster* [42]
548 and to characterize AS in dengue’s viral infection in humans [43].

549

550 In order to quantify ASpli2’s performance we compared it against three
551 different state-of-the-art methodologies: LeafCutter [23], MAJIQ [44] and
552 rMATS [22]. As a general rule we considered default parameters to run these

553 analysis pipelines for our intention was not to present here an extensive
554 benchmark between bioinformatics approaches, nor to propose the definitive
555 analysis methodology. Rather we wanted to establish whether ASpli2 pro-
556 duced reasonable and competitive results.

557

558 Different scenarios were considered to chart ASpli2 performance. We first
559 analyzed a synthetic data-set and quantified the ability of each considered
560 methodology to detect splicing changes in terms of precision and sensitiv-
561 ity figures of merit. Using this controlled dataset we found that all the
562 analysed methods presented rather high precision levels. However ASpli2
563 systematically displayed larger recall values ($\sim 40\%$), mainly because the
564 use of coverage signals (see Fig 1). This is an important result as highlights
565 the benefits of not losing effective sequencing depth by relying not only on
566 junction information but on the complete set of reads of RNA-seq runs.

567

568 We then aimed to outline ASpli2's performance over more realistic setups.
569 As no internal gold-standards are usually available for real world datasets
570 we focused on the analysis of two independent RNAseq assays that probed
571 the same biological conditions. This allowed us to quantify the consistence
572 and coherence of outcomes produced by each methodology in terms of re-
573 producibility of discoveries. Our results suggested that detection agreement
574 between studies was highly significative for every methodology. However
575 ASpli2 was far superior in terms of total number of concordant discoveries
576 reported.

577 It is worth noting that a necessary condition implicit in this analysis was

578 that biological variability largely exceeded possible technical biases between
579 studies. Using ASpli2, we were able to consider a generalized linear model to
580 define a consolidated dataset integrating data from both studies and verified
581 that this was actually the case (Sec 2.4). In addition, the possibility to imple-
582 ment a two-factor model greatly improved the statistical power to uncover
583 consistent discoveries. We could identify 4314 events displaying a statis-
584 tically significant genotype effect and no evidence of experiment-genotype
585 interactions. This represented almost a two-fold increase in the number of re-
586 producible discoveries when compared against the naive integrative approach
587 that merely considered the 2241 splicing events simultaneously detected in
588 both studies.

589

590 An important aspect of the presented approach is that ASpli2's core
591 functionality is implemented along user-friendly functions that produce self-
592 contained output results for each step of the analysis. This is an important
593 feature from the user's perspective. It provides the user valuable intermedi-
594 ate information eventually facilitating the integration of ASpli2 with other
595 analysis pipelines.

596 **5. Conclusions**

597 In this paper we presented ASpli2, a computational suite to study alter-
598 native splicing events. It is implemented as a flexible R modular package
599 that allows users to fulfill gene-expression and splicing analysis following a
600 set of simple steps.

601 Noticeably, ASpli2 can handle complex experimental designs using a uni-

602 fied statistical framework to assess for differential usage of sub-genic features
603 and junctions. By combining statistical information from exons, introns, and
604 splice junctions ASpli2 can provide an integrative view of splicing landscapes
605 that might include canonical and non-canonical splicing patterns occurring
606 in annotated as well as in novel splicing variants.

607 **Acknowledgements**

608 We thank Ruben Schlaen, Julieta Mateos and Andres Romanowski for
609 helpful discussions

610 **Funding**

611 This work has been supported by grants from Agencia Nacional de Pro-
612 moción Científica y Tecnológica (ANPCyT). AC also acknowledges support
613 from University of Buenos Aires (grant 20020170100356BA). AC and MY are
614 members of Carrera de Investigador of Consejo Nacional de Investigaciones
615 Científicas y Técnicas (CONICET).

616 6. Bibliography

- 617 [1] R. E. Breitbart, A. Andreadis, B. Nadal-Ginard, Alternative splicing:
618 a ubiquitous mechanism for the generation of multiple protein isoforms
619 from single genes., *Annu. Rev. Biochem.* 56 (1987).
- 620 [2] W. T. Nilsen, B. R. Graveley, Expansion of the eukaryotic proteome
621 by alternative splicing., *Nature* 463 (2010).
- 622 [3] D. Brett, H. Pospisil, J. Valcárcel, J. Reich, P. Bork, Alternative splicing
623 and genome complexity, *Nature Genetics* 30 (2002) 29–30.
- 624 [4] A. S. Reddy, Y. Marquez, M. Kalyna, A. Barta, Complexity of the
625 alternative splicing landscape in plants, *The Plant Cell* 25 (2013) 3657–
626 3683.
- 627 [5] J. Vaquero-Garcia, A. Barrera, M. R. Gazzara, J. Gonzalez-Vallinas,
628 N. F. Lahens, J. B. Hogenesch, K. W. Lynch, Y. Barash, A new view
629 of transcriptome complexity and regulation through the lens of local
630 splicing variations, *eLife* 5 (2016) e11752.
- 631 [6] K. A., O. C. S., Z. Y., R. G., SplAdder: identification, quantification
632 and testing of alternative splicing events from RNA-Seq data, *Bioinform-*
633 *atics* (????).
- 634 [7] R. Liu, A. E. Loraine, J. A. Dickerson, Comparisons of computational
635 methods for differential alternative splicing detection using rna-seq in
636 plant systems, *BMC Bioinformatics* 15 (2014) 364.

- 637 [8] L. Ding, E. Rath, Y. Bai, Comparison of alternative splicing junction
638 detection tools using rna-seq data, *Curr Genomics* 18 (2017) 268–277.
639 28659722[pmid].
- 640 [9] C. Zhang, B. Zhang, L.-L. Lin, S. Zhao, Evaluation and comparison of
641 computational tools for rna-seq isoform quantification, *BMC Genomics*
642 18 (2017) 583.
- 643 [10] B. J. Blencowe, Alternative splicing: New insights from global analyses,
644 *Cell* 126 (2006) 37 – 47.
- 645 [11] A. Lapuk, H. Marr, L. Jakkula, H. Pedro, S. Bhattacharya, E. Pur-
646 dom, Z. Hu, K. Simpson, L. Pachter, S. Durinck, N. Wang, B. Parvin,
647 G. Fontenay, T. Speed, J. Garbe, M. Stampfer, H. Bayandorian, S. Dor-
648 ton, T. A. Clark, A. Schweitzer, A. Wyrobek, H. Feiler, P. Spellman,
649 J. Conboy, J. W. Gray, Exon-level microarray analyses identify alter-
650 native splicing programs in breast cancer, *Molecular Cancer Research* 8
651 (2010) 961–974.
- 652 [12] Z. Wang, M. Gerstein, M. Snyder, Rna-seq: a revolutionary tool for
653 transcriptomics, *Nature Reviews Genetics* 10 (2009) 57 EP –. Perspec-
654 tive.
- 655 [13] F. Ozsolak, P. M. Milos, Rna sequencing: advances, challenges and
656 opportunities, *Nature Reviews Genetics* 12 (2010) 87 EP –. Review
657 Article.
- 658 [14] Y. Katz, E. T. Wang, E. M. Airoidi, C. B. Burge, Analysis and design

- 659 of rna sequencing experiments for identifying isoform regulation, *Nature*
660 *Methods* 7 (2010) 1009 EP –. Article.
- 661 [15] C. Trapnell, A. Roberts, L. Goff, G. Pertea, D. Kim, D. R. Kelley,
662 H. Pimentel, S. L. Salzberg, J. L. Rinn, L. Pachter, Differential gene
663 and transcript expression analysis of rna-seq experiments with tophat
664 and cufflinks, *Nature Protocols* 7 (2012) 562 EP –.
- 665 [16] N. L. Bray, H. Pimentel, P. Melsted, L. Pachter, Near-optimal proba-
666 bilistic rna-seq quantification, *Nature Biotechnology* 34 (2016) 525–527.
- 667 [17] R. Patro, G. Duggal, M. I. Love, R. A. Irizarry, C. Kingsford, Salmon
668 provides fast and bias-aware quantification of transcript expression, *Na-
669 ture Methods* 14 (2017) 417–419.
- 670 [18] S. Anders, A. Reyes, W. Huber, Detecting differential usage of exons
671 from rna-seq data, *Genome Research* 22 (2012) 2008–2017.
- 672 [19] D. J. McCarthy, G. K. Smyth, M. D. Robinson, edgeR: a Bioconductor
673 package for differential expression analysis of digital gene expression
674 data, *Bioinformatics* 26 (2009) 139–140.
- 675 [20] D. J. McCarthy, C. Y., G. K. Smyth, Differential expression analysis of
676 multifactor RNA-seq experiments with respect to biological variation,
677 *Nucleic Acids Research* 40 (2012).
- 678 [21] C. W. Law, Y. Chen, W. Shi, G. K. Smyth, voom: precision weights
679 unlock linear model analysis tools for rna-seq read counts, *Genome
680 Biology* 15 (2014) R29.

- 681 [22] S. Shen, J. W. Park, Z.-x. Lu, L. Lin, M. D. Henry, Y. N. Wu, Q. Zhou,
682 Y. Xing, *rmats: Robust and flexible detection of differential alterna-*
683 *tive splicing from replicate rna-seq data*, *Proceedings of the National*
684 *Academy of Sciences* 111 (2014) E5593–E5601.
- 685 [23] Y. I. Li, D. A. Knowles, J. Humphrey, A. N. Barbeira, S. P. Dickinson,
686 H. K. Im, J. K. Pritchard, *Annotation-free quantification of rna splicing*
687 *using leafcutter*, *Nature Genetics* 50 (2018) 151–158.
- 688 [24] M. D. Robinson, D. J. McCarthy, G. K. Smyth, *edgeR: a Bioconductor*
689 *package for differential expression analysis of digital gene expression*
690 *data*, *Bioinformatics* 26 (2010).
- 691 [25] S. E. Sanchez, E. Petrillo, E. J. Beckwith, X. Zhang, M. L. Rugnone,
692 C. E. Hernando, J. C. Cuevas, M. A. Godoy Herz, A. Depetris-Chauvin,
693 C. G. Simpson, J. W. S. Brown, P. D. Cerdán, J. O. Borevitz, P. Mas,
694 M. F. Ceriani, A. R. Kornblihtt, M. J. Yanovsky, *A methyl transferase*
695 *links the circadian clock to the regulation of alternative splicing*, *Nature*
696 468 (2010) 112–116.
- 697 [26] C. E. Hernando, S. E. Sanchez, E. Mancini, M. J. Yanovsky, *Genome*
698 *wide comparative analysis of the effects of prmt5 and prmt4/carm1*
699 *arginine methyltransferases on the arabidopsis thaliana transcriptome*,
700 *BMC Genomics* 16 (2015) 192.
- 701 [27] A. P. Roworth, S. M. Carr, G. Liu, W. Barczak, R. L. Miller, S. Munro,
702 A. Kanapin, A. Samsonova, N. B. La Thangue, *Arginine methyl-*
703 *ation expands the regulatory mechanisms and extends the genomic land-*

- 704 scape under e2f control, *Science advances* 5 (2019) eaaw4640–eaaw4640.
705 31249870[pmid].
- 706 [28] X. Deng, L. Gu, C. Liu, T. Lu, F. Lu, Z. Lu, P. Cui, Y. Pei, B. Wang,
707 S. Hu, X. Cao, Arginine methylation mediated by the arabidopsis ho-
708 molog of prmt5 is essential for proper pre-mrna splicing, *Proceedings of*
709 *the National Academy of Sciences* 107 (2010) 19114–19119.
- 710 [29] S. Ren, Z. Peng, J.-H. Mao, Y. Yu, C. Yin, X. Gao, Z. Cui, J. Zhang,
711 K. Yi, W. Xu, C. Chen, F. Wang, X. Guo, J. Lu, J. Yang, M. Wei,
712 Z. Tian, Y. Guan, L. Tang, C. Xu, L. Wang, X. Gao, W. Tian, J. Wang,
713 H. Yang, J. Wang, Y. Sun, Rna-seq analysis of prostate cancer in the
714 chinese population identifies recurrent gene fusions, cancer-associated
715 long noncoding rnas and aberrant alternative splicings, *Cell Research*
716 22 (2012) 806–821.
- 717 [30] M. Wang, Y. Zhao, B. Zhang, Efficient test and visualization of multi-set
718 intersections, *Scientific Reports* 5 (2015) 16923.
- 719 [31] C. G. Simpson, J. Fuller, M. Maronova, M. Kalyna, D. Davidson, J. Mc-
720 Nicol, A. Barta, J. W. S. Brown, Monitoring changes in alternative pre-
721 cursor messenger rna splicing in multiple gene transcripts, *The Plant*
722 *Journal* 53 (2008) 1035–1048.
- 723 [32] A. Mehmood, A. Laiho, M. S. Venlinen, A. J. McGlinchey, N. Wang,
724 L. L. Elo, Systematic evaluation of differential splicing tools for RNA-seq
725 studies, *Briefings in Bioinformatics* (2019). Bbz126.

- 726 [33] T. Griebel, B. Zacher, P. Ribeca, E. Raineri, V. Lacroix, R. Guigo,
727 M. Sammeth, Modelling and simulating generic RNA-Seq experiments
728 with the flux simulator, *Nucleic Acids Res.* 40 (2012) 10073–10083.
- 729 [34] M. Lawrence, W. Huber, H. Pagès, P. Aboyoun, M. Carlson, R. Gentle-
730 man, M. T. Morgan, V. J. Carey, Software for computing and annotating
731 genomic ranges, *PLOS Computational Biology* 4 (2013).
- 732 [35] M. L. Rugnone, A. F. Soverna, S. E. Sanchez, R. G. Schlaen, C. E.
733 Hernando, D. K. Seymour, E. Mancini, A. Chernomoretz, D. Weigel,
734 P. Más, et al., Lnk genes integrate light and clock signaling networks
735 at the core of the arabidopsis oscillator, *Proceedings of the National*
736 *Academy of Sciences* 110 (2013) 12120–12125.
- 737 [36] S. Perez-Santángelo, E. Mancini, L. J. Francey, R. G. Schlaen, A. Cher-
738 nomoretz, J. B. Hogenesch, M. J. Yanovsky, Role for lsm genes in the
739 regulation of circadian rhythms, *Proceedings of the National Academy*
740 *of Sciences* 111 (2014) 15166–15171.
- 741 [37] E. Mancini, S. E. Sanchez, A. Romanowski, R. G. Schlaen, M. Sanchez-
742 Lamas, P. D. Cerdan, M. J. Yanovsky, Acute effects of light on alter-
743 native splicing in light-grown plants, *Photochemistry and photobiology*
744 92 (2016) 126–133.
- 745 [38] R. Xin, L. Zhu, P. A. Salomé, E. Mancini, C. M. Marshall, F. G. Harmon,
746 M. J. Yanovsky, D. Weigel, E. Huq, Spf45-related splicing factor for
747 phytochrome signaling promotes photomorphogenesis by regulating pre-

- 748 mRNA splicing in arabidopsis, *Proceedings of the National Academy of*
749 *Sciences* 114 (2017) E7018–E7027.
- 750 [39] A. Romanowski, R. G. Schlaen, S. Perez-Santangelo, E. Mancini, M. J.
751 Yanovsky, Global transcriptome analysis reveals circadian control of
752 splicing events in arabidopsis thaliana, *bioRxiv* (2019) 845560.
- 753 [40] C. E. Hernando, S. E. Sanchez, E. Mancini, M. J. Yanovsky, Genome
754 wide comparative analysis of the effects of prmt5 and prmt4/carm1
755 arginine methyltransferases on the arabidopsis thaliana transcriptome,
756 *BMC genomics* 16 (2015) 192.
- 757 [41] R. G. Schlaen, E. Mancini, S. E. Sanchez, S. Perez-Santángelo, M. L.
758 Rugnone, C. G. Simpson, J. W. Brown, X. Zhang, A. Chernomoretz,
759 M. J. Yanovsky, The spliceosome assembly factor gemin2 attenuates
760 the effects of temperature on alternative splicing and circadian rhythms,
761 *Proceedings of the National Academy of Sciences* 112 (2015) 9382–9387.
- 762 [42] E. J. Beckwith, C. E. Hernando, S. Polcowñuk, A. P. Bertolin,
763 E. Mancini, M. F. Ceriani, M. J. Yanovsky, Rhythmic behavior is con-
764 trolled by the srm160 splicing factor in drosophila melanogaster, *Genet-*
765 *ics* 207 (2017) 593–607.
- 766 [43] F. A. De Maio, G. Risso, N. G. Iglesias, P. Shah, B. Pozzi, L. G. Geb-
767 hard, P. Mammi, E. Mancini, M. J. Yanovsky, R. Andino, et al., The
768 dengue virus ns5 protein intrudes in the cellular spliceosome and mod-
769 ulates splicing, *PLoS pathogens* 12 (2016).

- 770 [44] J. Vaquero-Garcia, S. Norton, Y. Barash, Leafcutter vs. majiq and
771 comparing software in the fast moving field of genomics, bioRxiv (2018).
- 772 [45] S. Schafer, K. Miao, C. C. Benson, M. Heinig, S. A. Cook, N. Hubner,
773 Alternative splicing signatures in rna-seq data: Percent spliced in (psi),
774 Current Protocols in Human Genetics 87 (????) 11.16.1–11.16.14.

775 7. Supplementary Figures

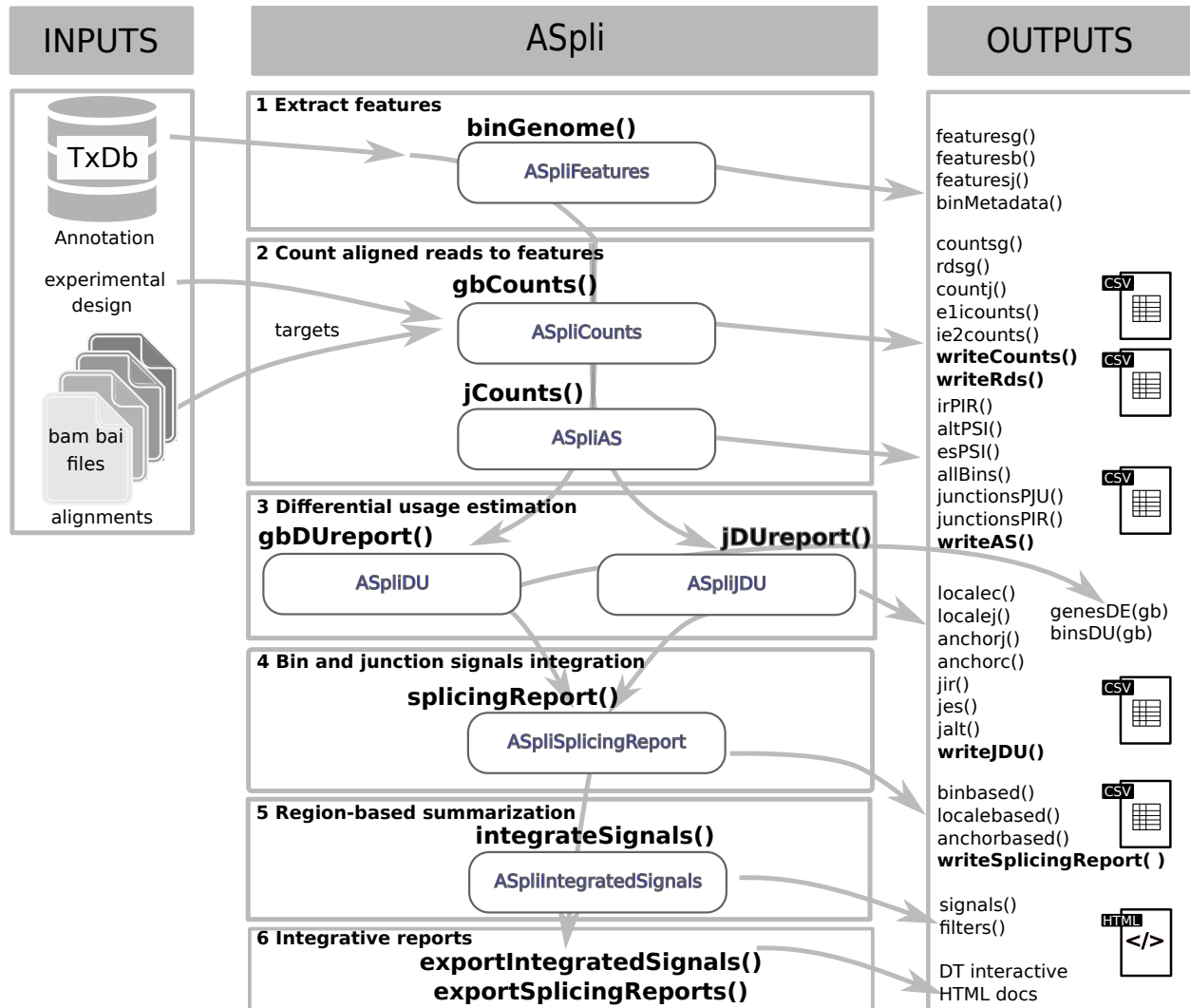


Figure (S1) ASpli workflow. Rounded boxes are objects created by ASpli functions. Accessors and outputs are summarized in the right-most panel

	ASpli function	Description	Genomic feature		
			gene	bin	junction
	binGenome()	Identification of subgenic features		Genome binning Event type classification	
Counting	gbCounts()	Counts over the annotated genome	Counts Read density	Counts Read density	Counts <i>Exon</i> detection
	jCounts ()	Junction-based stats		Counting of annotated junctions	Detection of novel AS events
Differential signals	gbDUreport()	Coverage-based differential signal	Differential gene expression	Differential bin usage	Differential junction usage
	jDUreport()	Junction-based differential signal		Differential usage of bin inclusion/exclusion junctions	Differential usage inside junction clusters
Report	splicingReport ()	Leveraging splicing evidence		Bin coverage signals	Anchorage and locale signals
	integrateSignals()	Signal consolidation			

Figure (S2) Summary of ASpli core functionality.

Show entries

Search:

ASpli: integrated signals. Contrasts: A_C - A_D

Filters: bin.FC=3; bin.fdr=0.05; nonunif=1; usenonunif=FALSE; bin.inclusion=0.2; bjs.inclusion=10.3; bjs.fdr=0.01; a.inclusion=0.3; a.fdr=0.01; l.inclusion=0.3; l.fdr

View	Region	Event	Locus	Locus overlap	Bin Evidence	Bin SJ Evidence	Anchor Evidence	Locale Evidence	Bin	Feature	Bins			
											logFC	FDR	LR	FD
<input type="button" value="All"/>	<input type="text" value="All"/>	<input type="button" value="All"/>	<input type="button" value="All"/>	<input type="button" value="All"/>	<input type="button" value="All"/>	<input type="button" value="All"/>	<input type="button" value="All"/>	<input type="button" value="All"/>	<input type="button" value="All"/>	<input type="button" value="All"/>	<input type="button" value="All"/>	<input type="button" value="All"/>	<input type="button" value="A"/>	<input type="button" value="All"/>
⊕	reference:1250-1601	ES	GENE02	-	0	0	0	1		-				
⊕	reference:2250-2501	Alt 5'/3'	GENE03	-	0	0	0	1		-				
⊕	reference:3250-3501	Alt 5'/3'	GENE04	-	0	0	0	1		-				
⊕	reference:4251-4350	Alt5ss	GENE05	-	1	0	0	0	GENE05:E002	E	1.067	1.074e-14		
⊕	reference:7251-7350	IR*	GENE08	-	1	0	1	0	GENE08:E002	E	0.6338	7.125e-8		

Gene view

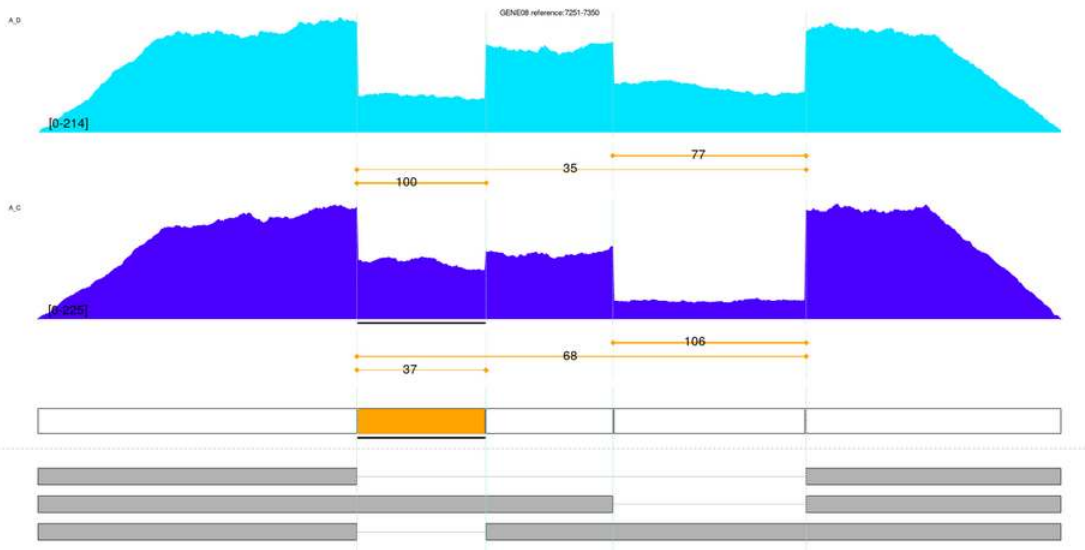


Figure (S3) Example of DT html interactive report generated by *exportIntegratedSignals()* function

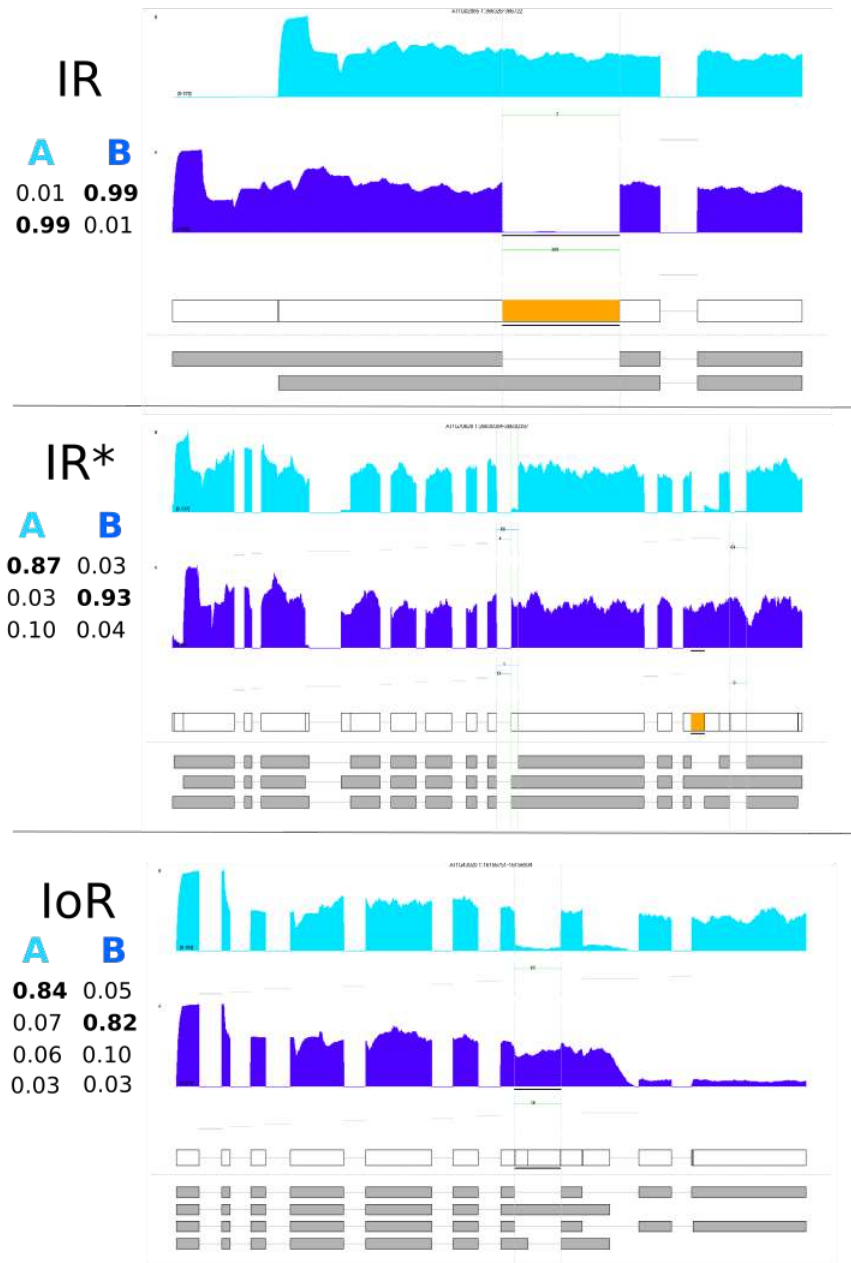


Figure (S4) Examples of simulated IR-like splicing events. For each panel, the left layered table shows the relative concentration of each variant simulated for condition A and B. Orange boxes highlight the considered bin in each case.

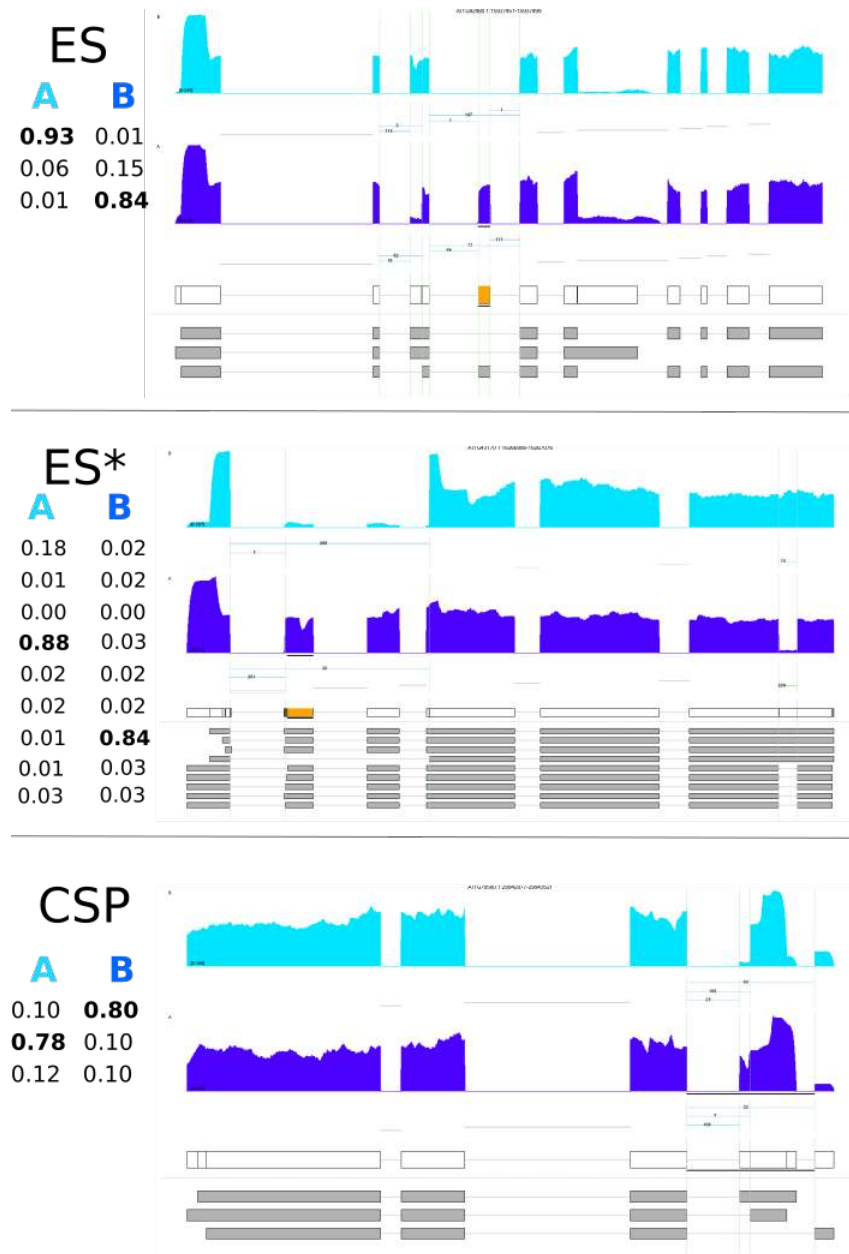


Figure (S5) Examples of simulated ES-like splicing events. For each panel, the left layered table shows the relative concentration of each variant simulated for condition A and B. Orange boxes highlight the considered bin in each case.

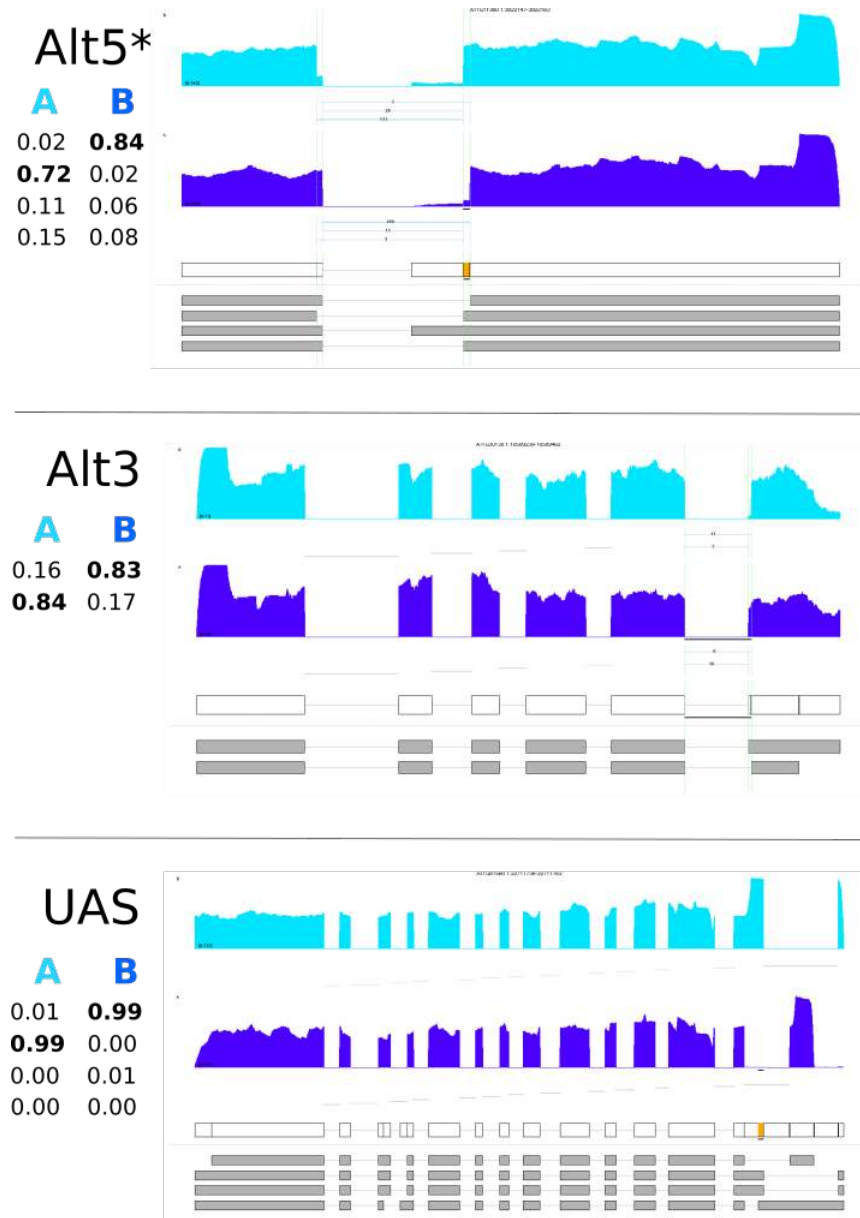


Figure (S6) Examples of simulated Alternative start/end splicing events. For each panel, the left layered table shows the relative concentration of each variant simulated for condition A and B. Orange boxes highlight the considered bin in each case.

fdr.Interaction > 0.5

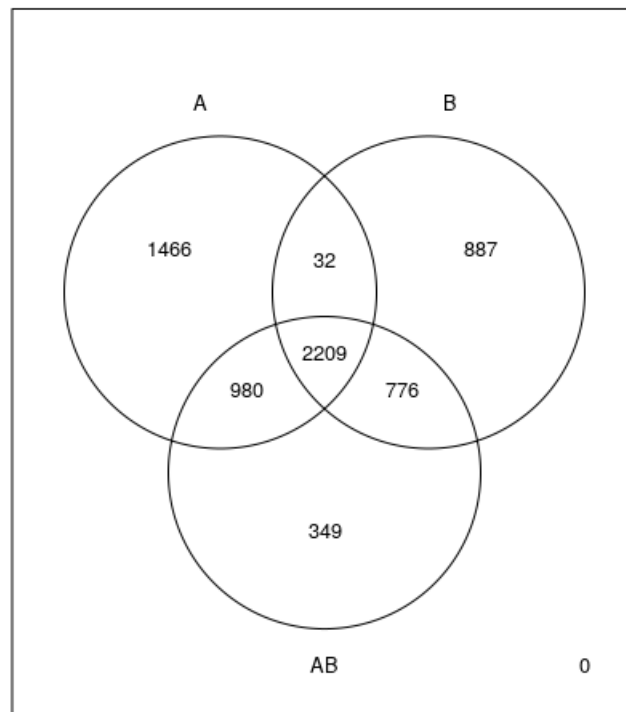


Figure (S7) Venn diagram of alternative splicing events detected in experiments A, B, and the consolidated data set AB (i.e. events displaying strong evidence of a genotype effect ($fdr < 0.05$) and no-detectable evidence of experiment-genotype interaction (experiment:genotype associated $fdr > 0.5$)).

776 8. Supplementary Material

777 8.1. Feature counting in ASpli

778 8.1.1. Genomic feature extraction: binGenome()

779 Sub-genic features are analyzed using user-provided annotation files. Exon
780 and intron coordinates are extracted from annotation for multi-exonic genes.
781 When more than one isoform exists, some exons and introns from different
782 isoforms will generally overlap. In the same spirit of [18], exons and introns
783 are then subdivided into non-overlapping sub-genic features dubbed bins, de-
784 fined by the boundaries of different exons across transcript variants. In this
785 way, these so defined *bins* are maximal sub-genic features entirely included
786 or entirely excluded from any mature transcript.

787 Bins are flagged as: exonic (E), intronic (I) or alternative-splicing (AS)
788 bins, depending on the exonic/intronic character of the bin across variants .
789 In addition, original intronic (Io) bins are defined for every intronic region of
790 annotated isoforms (see panel A of Figure S8).

791 As a general rule, the extreme portions of a transcript probed by RNAseq
792 assays show a highly non-uniform coverage that might obscure differential
793 usage analysis. ASpli flags bins that overlap with the beginning or ending of
794 any transcript as *external*. An external bin of a transcript may overlap with
795 a non-external one of another transcript. Whenever this happens the bin is
796 still labelled as external. Additionally, in order to avoid confounding effects
797 in the analysis of splicing events, ASpli identifies and flags loci where more
798 than one gene is present in the genome.

799 *Local splicing classification model.* Each AS bin is further classified consid-
800 ering a three-bin *minimum local gene model*, that assigns splicing-event cat-

801 egeries to a given bin based on the intronic/exonic character of the analyzed
802 bin and its first neighbors (Figure S8, panel B).

803 For genes presenting two isoforms, this model is able to unambiguously
804 assign a well defined splicing event to the analyzed bin: exon skipping (ES),
805 intron retention (IR), alternative five prime splicing site (Alt5'SS), or alter-
806 native three prime splicing site (Alt3'SS) (see first row of panel B in Figure
807 S8).

808 When more than two isoforms are present, we still found it useful to use
809 the three-bin local model to segment follow up analysis. For these cases ASpli
810 identify splicing events that involve: intronic subgenic regions surrounded by
811 exons in at least one isoform (bin labelled as IR*), exonic subgenic regions
812 surrounded by two introns in at least one isoform (bin labelled as ES*), ex-
813 onic regions surrounded by intronic and exonic neighbor bins (bin labelled
814 as Alt5'SS* or Alt3'SS*). When it is not possible to get a clear splicing-type
815 assignation (see rows 2-5 of Figure S8), bins are labeled as *undefined AS*
816 (UAS).

817

818 As a last step of the genomic feature extraction process, annotated junc-
819 tions from all the transcripts are also identified. Junction coordinates are
820 defined as the last position of the five prime exon (donor position) and the
821 first position of the three prime exon (acceptor position).

822

823 *8.1.2. Annotation based feature counting: gbCounts()*

824 Reads are overlaid on features derived from annotation, and count tables
825 are produced at different genomic levels: genes, bins, and intron flanking

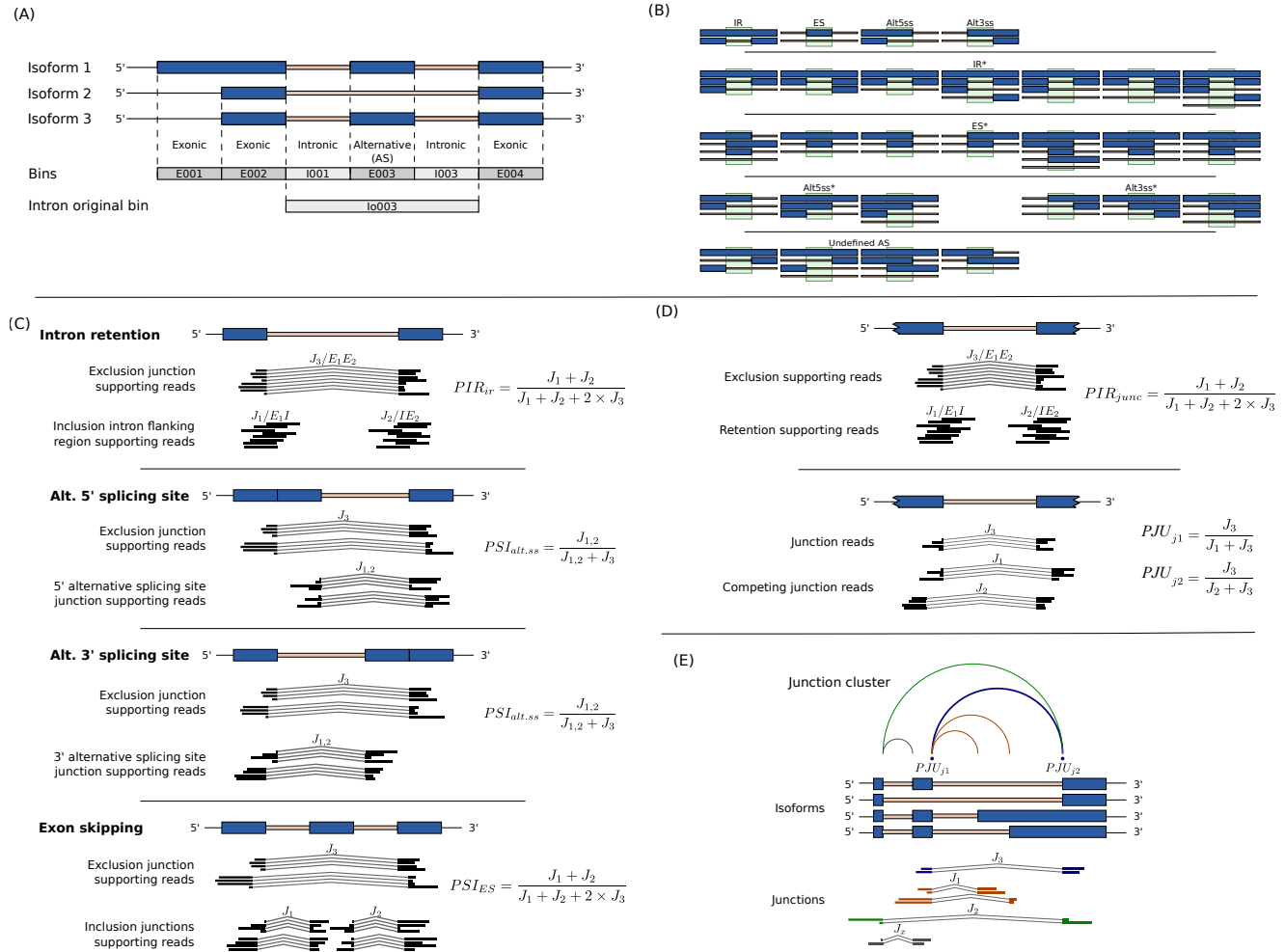


Figure (S8) Panel (A) shows how bin-features are defined and classified as: exonic, intronic or intron_original bins using genome annotation. The local splicing classification scheme is illustrated in panel (B). The definition of PSI and PIR metrics for bin features are pictured in panel (C). Definition of junction PIR and PJU statistics are shown in panel (D). Panel (E) shows a possible junction cluster and highlights the definition of type J_1 , J_2 and J_3 junctions for the analysis of PJU statistics for the blue junction.

826 regions used to identify and quantify intron retention events. Reads corre-
827 sponding to annotated junctions are also tallied, along with genomic relevant
828 information such as identity of spanned bins, and the existence of possible
829 *exintrinsic* events [?].

830 8.1.3. *De-novo junction counting: jCounts()*

831 ASpli takes advantage of experimentally detected splice junctions to per-
832 form two different type of analysis. For one hand, junction data is considered
833 in order to provide junction support to AS events detected through bin cov-
834 erage analysis. For the other, it is used to quantify novel splicing events.

835 *Junction support of bin coverage statistics:*. ASpli makes use of junction data
836 as supporting evidence of alternative usage of bins. For a general differential
837 splicing event affecting a given bin, it is always possible to define exclusion
838 and inclusion junctions. The first class of junctions (noted as J_3) pass over
839 the bin of interest, whereas the second ones (note as J_1 and/or J_2) quantify
840 and support the inclusion of start and/or end bin boundaries in the mature
841 transcript. Panel C of Figure S8 illustrates this point for the different types of
842 splicing events that could affect a given bin. ASpli considers for this analysis
843 junctions that are completely included within a unique gene and have more
844 than a minimum number of reads supporting them (by default this number
845 is five).

846

847 PSI (percent spliced in) [45] and PIR (percent of intron retention) metrics
848 are two well known statistics that can be used to quantify the relative weight
849 of inclusion evidence for different kind of splicing events (see Panel C of

850 Figure S8). For each bin, ASpli quantifies the inclusion strength in every
 851 experimental condition using the appropriate inclusion index (see Table S1).
 852 Only junctions that pass an abundance filter criterium (a minimum number
 853 of counts should be attained in all samples of at least one condition) are
 854 considered for the estimations.

feature	assesment	index	bin class
bin	inclusion	PIR_{ir}	UAS, I, I*, I ₀
		PSI_{es}	$\frac{J_1+J_2}{J_1+J_2+2*J_3}$ UAS, E, E*
		PSI_{alt5ss}	$\frac{J_{1,2}}{J_{1,2}+J_3}$ Alt5ss, Alt5ss*
		PSI_{alt3ss}	Alt3ss, Alt3ss*
junction	usage	PIR_{junc}	$\frac{J_1+J_2}{J_1+J_2+2*J_3}$ -
		PJU_{J_1}	$\frac{J_3}{J_1+J_3}$ -
		PJU_{J_2}	$\frac{J_3}{J_2+J_3}$ -

Table (S1) Junction usage and inclusion strength figure of merits for different bin classes and for experimentally detected junctions. The definition of J_1 , J_2 and J_3 junction counts is depicted in panels C and D of Figure S8 for annotated and experimentally detected junctions respectively.

855 For each bin, a PIR or a PSI metric is calculated, according to the splicing
 856 event category assigned to that bin (see last column of table S1). If no splice
 857 event was assigned, meaning that the bin is not alternative, an exon will be
 858 considered to be involved in a putative exon skipping splicing event, and an
 859 intron will be considered to be involved in a putative intron retention splicing
 860 event.

861 *Novel and non-canonical splicing patterns*:. ASpli relies on the direct analysis
862 of experimentally observed splicing junctions in order to study novel (i.e.
863 non-annotated) splicing patterns.

864 For every experimental junction, ASpli characterizes local splicing pat-
865 terns considering two hypothetical scenarios. For one hand, assuming that
866 every detected junction might be associated to a possible intron that could
867 be potentially retained, a PIR_{junc} value is computed (panel D of Figure S8).

868

869 On the other hand, every junction also defines potential 5' and 3' splic-
870 ing sites. It can be the case that one (in an alternative 5' or 3' scenario),
871 or both ends (in case of exon skipping) were shared by other junctions. In
872 this context, it is informative to characterize the relative abundance of the
873 analyzed junction (dubbed J_3) with respect to the locally *competing* ones.
874 ASpli estimates *percentage junction-usage* indices, PJU_{J_1} and PJU_{J_2} , in or-
875 der to evaluate and quantify this quantities (see Panel D of figure S8 and
876 Table S1). In order to illustrate this point, we show in Panel E of figure S8
877 an hypothetical splicing scenario for a given junction of interest, J_3 . It can
878 be appreciated that PJU_{J_1} quantifies the participation of this junction in
879 the context of a splicing pattern involving the two orange competing junc-
880 tions, whereas PJU_{J_2} reports on the usage of J_3 in connection with the green
881 competing junction.

882 8.2. *Command-line running arguments*

883 Command lines used to invoked algorithms and further calculation details:

884

885 ● STAR aligner

886 For PRMT5 datasets

```
887 1 | $ STAR --runThreadN 30 --genomeDir TAIR10_GENOME_DIR --twopassMode
888     Basic --outSAMtype BAM SortedByCoordinate --
889     outFilterMultimapNmax 2 --outFilterType BySJout --
890     outSJfilterReads Unique --sjdbOverhang PARAM --
891     alignSJoverhangMin 6 --alignSJDBoverhangMin 3 --alignIntronMin
892     20 --alignIntronMax 5000 --readFilesIn ../01_FASTQ/Col_3_1.fq ..
893     /01_FASTQ/Col_3_2.fq --outFileNamePrefix Col_3/Col_3
```

894 We used a `sjdbOverhang` parameter value equal to 99 and 149 for
895 PRMT5 datasets A and B respectively.

896

897 For the prostate dataset we aligned using default STAR parameters.

```
898 1 | $ STAR --runThreadN 30 --genomeDir ENSEMBL_HG38_PATH --
899     readFilesCommand zcat --twopassMode Basic --outSAMtype BAM
900     SortedByCoordinate --sjdbOverhang 89 --readFilesIn 1.fq 2.fq
```

901 We used a `sjdbOverhang` parameter value equal to 99 and 149 for
902 PRMT5 datasets A and B respectively.

903 ● LeafCutter (synthetic dataset)

904 BAM files were first processed using the provided `bam2junc.sh` script.

905 The generated `juncfiles.txt` was then used to build junction clusters via
906 the provided python script

```
907 1 | $ python PATH leafcutter_cluster.py -j juncfiles.txt -m 30 -l 500000
```

908 Finally, we used the provided `leafcutter_ds` R-script to run the statistical
909 analysis (`min_samples_per_intron=3`).

910 ● rMATS Command line use to analyzer PRMT5 assays:

```
911 | 1 rMATS.4.0.2/rMATS-turbo-Linux-UCS4/rmats.py --b1 bam_prmt5.txt --b2  
912 |     bam_col.txt --gtf /data1/genomeData/ath/Ensembl_illumina_iGenomes/  
913 |     TAIR10/Annotation/Genes/genes.gtf --od r1150 -t paired --nthread  
914 |     20 --readLength 150 --tstat 10
```

915 • MAJIQ

916 *8.3. Splicing affected regions detected by different algorithms*

917 Each algorithm reports splicing altered genomic features in different ways.
918 In order to standardize the identification of regions of interest we proceeded
919 as follows:

- 920 • LeafCutter: We first identified clusters presenting adjusted pvalues <
921 0.05 as reported in 'leafcutter_ds_cluster_significance.txt' file. For each
922 of these statistically significant clusters we considered the associated
923 genomic-regions reported in 'leafcutter_ds_effect_size.txt' file with $|\Delta\Psi| >$
924 0.1.
- 925 • MAJIQ: We considered the genomic-region covering junction clusters
926 presenting at least one junction with $P(|\Delta\Psi| > 0.2) > 0.95$.
- 927 • rMATS: We considered the values reported in 'JCEC.txt' files. This
928 means that we considered a model that evaluated splicing with reads
929 that spanned splicing junctions and reads on targets bins (i.e. alterna-
930 tively spliced exons). We kept junctions presenting adjusted FDR < 0.0
931 and inclusion signal larger than a 0.1 level. Genomic regions were then
932 defined according the following rules:

- 933 – A3SS' (A3SS.MATS.JCEC.txt file): We considered the genomic
934 region between 'shortEE' and 'longExonEnd' coordinates for neg-
935 ative strand and by 'longExonStart_0base' and 'shortES' for posi-
936 tive strand cases.
- 937 – A5SS' (A5SS.MATS.JCEC.txt file): We considered the genomic
938 region between 'shortEE' and 'longExonEnd' coordinates for posi-
939 tive strand and by 'longExonStart_0base' and 'shortES' for neg-
940 ative strand cases.
- 941 – MXE (MXE.MATS.JCEC.txt file): We considered two regions
942 per event defined by: '1stExonStart_0base', '1stExonEnd' and
943 '2ndExonStart_0base', '2ndExonEnd'.
- 944 – SE (SE.MATS.JCEC.txt file): We considered the regions between
945 'exonStart_0base' and 'exonEnd'.
- 946 – RI (RI.MATS.JCEC.txt file): We considered the regions between
947 'riExonStart_0base' and 'riExonEnd'.

948 *8.4. Analysis of false positive calls in simulated dataset*

949 In our simulations a 20% level of random variability was added to variant
950 concentration profiles. A splicing activation signal (SAS) value was then es-
951 timated for each gene as the maximum absolute change in variant concentra-
952 tion observed between conditions. The left-most first and second boxplots in
953 Figure S9 depict the distribution of this quantity for the 915 genes for which
954 a splicing event was simulated, and for the remaining 7518 genes respectively.
955 On the other hand, the four right-most boxplots show the SAS distribution
956 for false positive calls obtained with different methods. Non explicitly splic-

957 ing simulated changes were reported for 9, 4, 48 and 23 genes according to
958 ASpli, LeafCutter, MAJIQ and rMATS algorithms respectively.

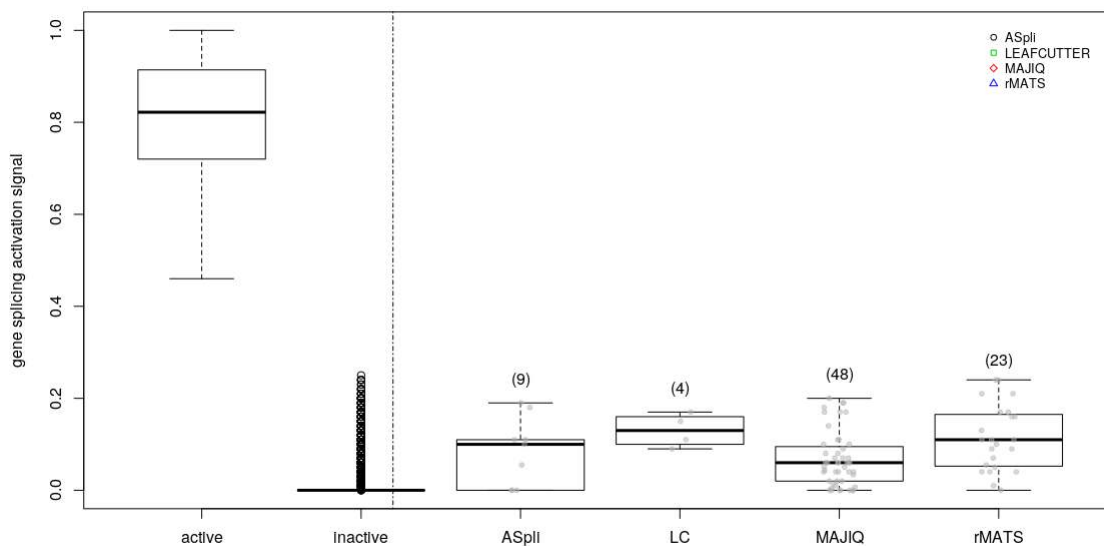


Figure (S9) Splicing simulation.

959 8.5. Comparison of discoveries

960 A comprehensive comparison of discoveries appeared at first-sight prob-
961 lematic as each algorithm is focused on different genomic features in order
962 to chart splicing landscapes.

963 For instance, rMATS analyzes genomic regions flanked by upstream and
964 downstream exons to examine canonical splicing events. MAJIQ and Leaf-
965 Cutter, on the other hand, exclusively rely on clusters of split reads that
966 share start or ending junction-ends. Finally ASpli considers both, junction
967 clusters and bin features, i.e. genomic regions defined from disjoint ranges
968 of annotated junctions.

969 In this context, a first coarse grained comparison could be established at
970 gene-level, comparing the identity of genes housing splicing-altered patterns
971 according to the different analyzed methods. Panel (A) of Figure S10 dis-
972 plays a color-coded overlap matrix of affected genes in experiments *A* and *B*
973 according to the four examined methodologies. Each cell reports the inter-
974 section size and, in brackets, the corresponding overlap coefficient. At gene
975 level, rMATS achieved the largest agreement factor (83% of genes identified
976 in experiment *B*, were also reported in experiment *A*). However, it also
977 produced the lowest number of discoveries (119). ASpli, on the other hand,
978 presented a comparable level of agreement (71%), highlighting a significa-
979 tively larger number of concordant genes (2109). Typically, more than 50%
980 of genes identified by any methodology was also reported by ASpli (first and
981 second rows of Figure S10). Moreover, the number of concordant discoveries
982 between experiments considering a given methodology was comparable to the
983 agreement level achieved between each experiment-methodology combination
984 and the corresponding ASpli result. Noticeably, more than 90% of MAJIQ's
985 genes were also spotted by ASpli.

986

987 A more in-depth comparison could be established analyzing the overlap
988 of identified genomic regions. In panels (b) and (c) of Figure S10 we informed
989 the extent of the overlaps between genomic regions found to be affected by
990 differential splicing patterns according to each algorithm (see Material and
991 Methods 3.5) to map events reported by each method to a common set of
992 genomic coordinates). While any kind of overlap was registered for panel (b),
993 only complete inclusion of genomic regions identified by one method inside

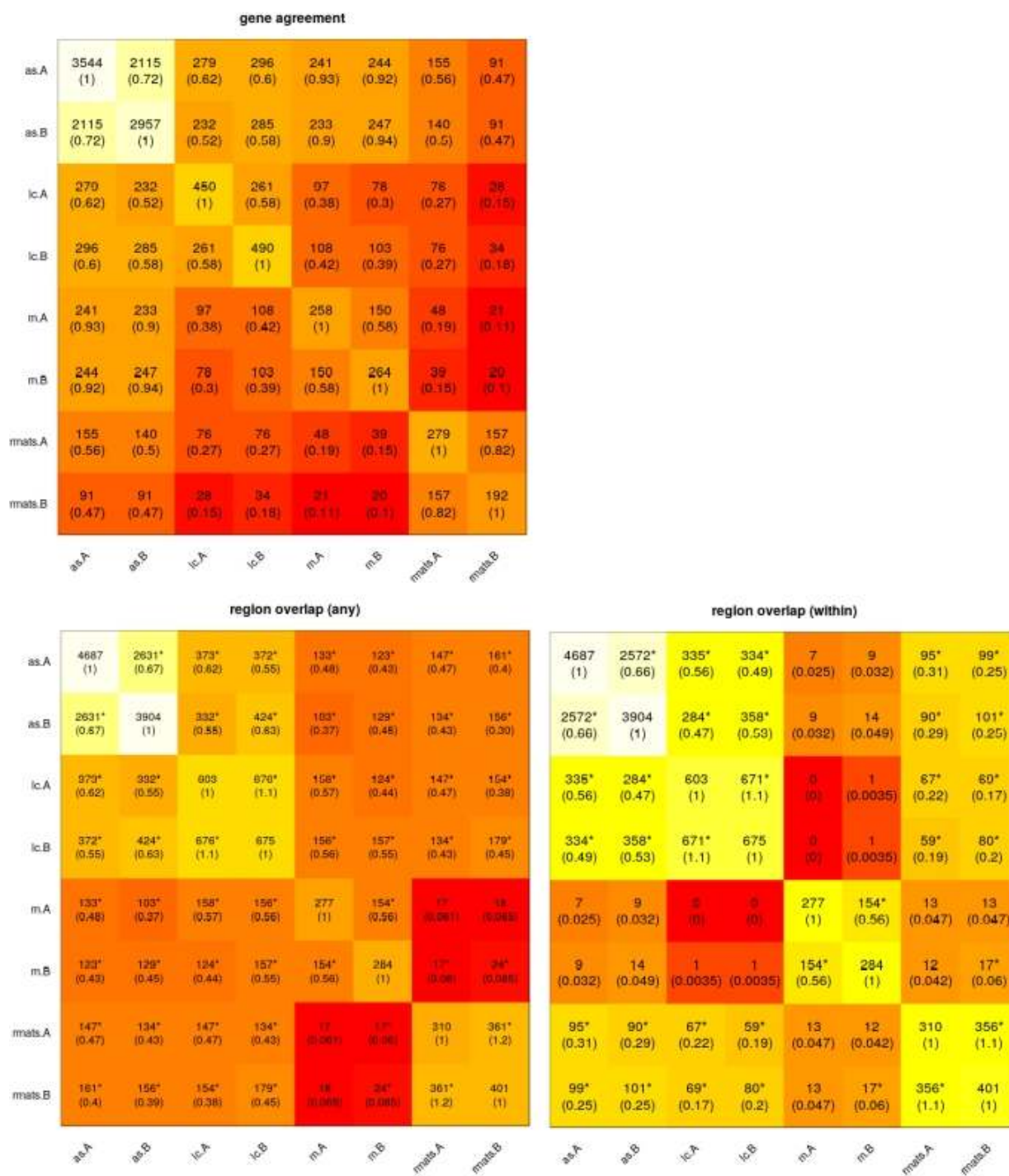


Figure (S10) Aspli main functions.

994 the ones identified by a second one was considered for panel (c). Statistically
995 significant overlaps were marked with asterisks. Note that overlap coefficients
996 (in brackets) exceeding unity were detected in between-experiments compar-
997 isons for LeafCutter and rMATS as a result of the presence of one-to-many
998 region mappings.

999 For the loose overlap criterium we found statistically significant concor-
1000 dance between discoveries for almost every cell (Fig S10-b). Only specific
1001 comparisons involving MAJIQ and rMATs failed the statistical significance
1002 test. At the same time, overlap coefficient values were similar to the ones
1003 estimated at the gene-level analysis. Noticeably, we recognised a sensible
1004 reduction in this quantity for the MAJIQ vs ASpli comparison. This finding
1005 highlighted that gene-level agreement should in general be considered with
1006 caution. A more detailed examination at the sub-genic level might be neces-
1007 sary to assess for discovery consistencies between algorithms. Results for the
1008 most stringent overlap criterion are shown in Figure S10(c). As expected,
1009 a major decrease on overlap coefficient values was observed . However, sta-
1010 tistically significant agreement between results was still found as a general
1011 rule. Only comparisons involving MAJIQ's discoveries failed the statistical
1012 assessments.

1013

1014 *8.6. PRMT5 PCR events*

1015 We characterized the agreement between the 23 splicing events that AS-
1016 pli uncovered for the consolidated AB case, and the 44 Sanchez qRT-PCR
1017 validated events in Table S2. For each assayed event we included the kind
1018 of the original event and the reported qRT-PCR splicing signal value in the

1019 second and third columns respectively (Sanchez and collaborators calculated
1020 the fraction of the shortest isoform in PRMT5 mutants and wildtype plants
1021 detected by qRT-PCR, and used the relativized difference between them as
1022 a quantitative proxy of splicing changes (Table 4 of [25])). In the fourth col-
1023 umn we informed whether the PCR-interrogated genomic region overlapped
1024 with the one signaled by ASpli. Finally, the type of splicing event detected
1025 by ASpli was included in the last column of the table.

1026 *8.7. Prostate cancer dataset: Transcriptomic variability*

1027 In order to visualize the transcriptomic variability across patients at gene
1028 expression levels we considered the 30% most variable genes across the 28
1029 expression profiles that presented more than 10 counts per million reads
1030 in at least 3 samples. With this informative set of 1386 genes we built a
1031 multidimensional scaling plot of distances between gene expression profiles
1032 estimated with the edgeR package [24]. Results are shown in Fig S11. In this
1033 kind of plot, samples lay on a two-dimensional scatterplot so that distances
1034 on the plot approximate the typical log₂ fold changes between the samples
1035 (function plotMDS of edgeR [24]).

1036 Empty and filled symbol correspond to tumor and normal tissue samples
1037 respectively. Pair of points of a given patient are equally colored and joined
1038 by a dashed edge.

1039 It can be seen that tumor and normal samples were well separated across
1040 the leading reduced dimension. The second largest projected dimension,
1041 however, let us appreciate internal structure and some variability between
1042 patients. There was a group of 5 patients (top left empty points) that dis-
1043 played a rather homogeneous pattern of changes between tumor affected and

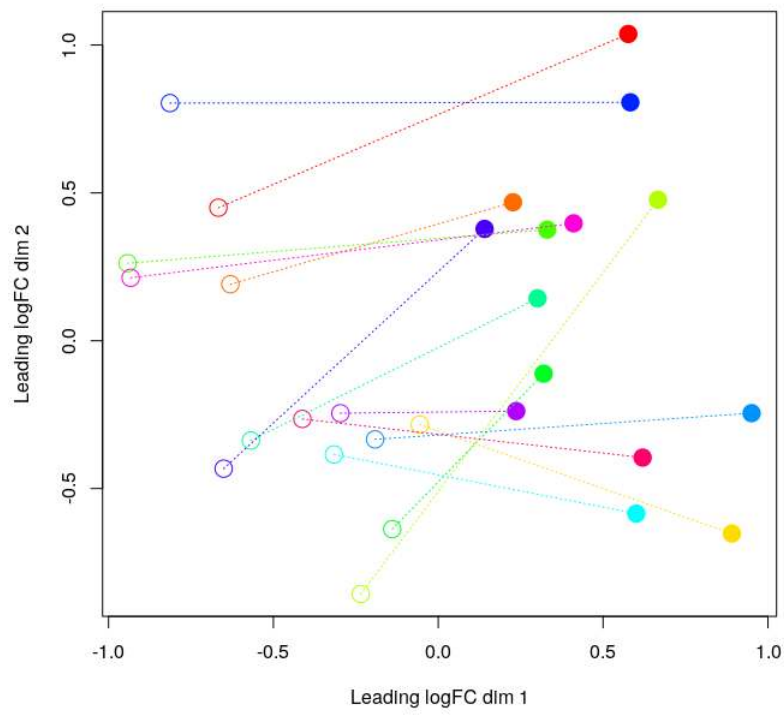


Figure (S11) A.

1044 normal tissues. On the contrary, the 9 bottom-left tumor samples seemed
1045 to segregate into a different cluster of transcriptomes. Moreover, the corre-
1046 sponding patients presented different kinds of alterations between tumor and
1047 control samples.

	Gene ID	Event	qRT-PCR signal	Region overlap	Detected event
1	AT1G53650	5'ss	18.93	yes	IR (next to 5)
2	AT1G54360	5'ss	39.21	yes	Alt5ss
3	AT1G76510	5'ss	-28.00	yes	Alt5ss
4	AT2G04790	5'ss	11.14	yes	IR
5	AT2G15530	5'ss	21.12	yes	Alt 5'/3'
6	AT2G33480	5'ss	-27.16	yes	IR
7	AT2G38880	5'ss	-10.44	no	IR
8	AT2G46790	5'ss	35.20	yes	Alt5ss (plus additional IR)
9	AT3G01150	ES	-13.70	no	IR
10	AT3G12250	5'ss	-16.29	no	ES*
11	AT3G16800	IR/3'ss	-31.59	yes	IR, Novel Alt 5'/3'
12	AT3G19840	5'ss	-26.20	no	IR
13	AT3G20270	ES	8.51	no	IR (next to ES)
14	AT3G23280	ES	18.21	yes	ES
15	AT3G25840	ES	6.51	yes	ES
16	AT4G02430	3'ss	27.45	no	IR
17	AT4G24740	ES	16.07	no	IR (next to ES)
18	AT4G31720	3'ss	12.37	no	IR
19	AT4G32730	5'ss	30.93	yes	Novel Alt 5'/3'
20	AT4G38510	5'ss	15.53	yes	Alt5ss*, CSP
21	AT5G05550	ES	32.72	yes	ES (plus additional IR)
22	AT5G25610	IR	-71.69	yes	IR
23	AT5G57630	5'ss	31.07	yes	Novel Alt 5'/3' (plus adjacent IR)

Table (S2)