

1 **Extensive splicing across the *Saccharomyces cerevisiae* genome**

2 Stephen M. Douglass^{1*}, Calvin S. Leung², and Tracy L. Johnson^{1,2*}

3 ¹Department of Molecular, Cell, and Developmental Biology, University of California, Los
4 Angeles, Los Angeles, CA 90095, USA

5 ²Molecular Biology Institute, University of California, Los Angeles, Los Angeles, CA 90095, USA

6 Correspondence: smdougla@ucla.edu, tjohnson@ucla.edu

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24 **Abstract**

25 Pre-mRNA splicing is vital for the proper function and regulation of eukaryotic gene expression.
26 *Saccharomyces cerevisiae* has been used as a model organism for studies of RNA splicing
27 because of the striking conservation of the spliceosome components and its catalytic activity.
28 Nonetheless, there are relatively few annotated alternative splice forms, particularly when
29 compared to higher eukaryotes. Here, we describe a method to combine large scale RNA
30 sequencing data to accurately discover novel splice isoforms in *Saccharomyces cerevisiae*.
31 Using our method, we find extensive evidence for novel splicing of annotated intron-containing
32 genes as well as genes without previously annotated introns and splicing of transcripts that are
33 antisense to annotated genes. By incorporating several mutant strains at varied temperatures,
34 we find conditions which lead to differences in alternative splice form usage. Despite this, every
35 class and category of alternative splicing we find in our datasets is found, often at lower
36 frequency, in wildtype cells under normal growth conditions. Together, these findings show that
37 there is widespread splicing in *Saccharomyces cerevisiae*, thus expanding our view of the
38 regulatory potential of RNA splicing in yeast.

39

40 **Author Summary**

41 Pre-mRNA splicing is a fundamental step in eukaryotic gene expression. *Saccharomyces*
42 *cerevisiae*, also known as brewer's yeast, is a model organism for the study of pre-mRNA
43 splicing in eukaryotes. Through the process of pre-mRNA splicing, a single gene is capable of
44 encoding multiple mature mRNA products, but it is often difficult to identify the splice events that
45 lead to these mRNA products. Here, we describe a method to accurately discover novel splice
46 events in *Saccharomyces cerevisiae* and find evidence for extensive splicing in
47 *Saccharomyces*. By utilizing a variety of strains and growth conditions, we are able to
48 characterize many splice forms and correlate cellular conditions with prevalence of novel splice
49 events.

50

51 **Introduction**

52 Eukaryotic genes are composed of coding sequences termed exons interrupted by non-coding
53 sequences called introns. Introns are removed from RNA by the large macromolecular complex
54 known as the spliceosome through the process of RNA splicing. By selectively including
55 different combinations of exons, a single gene can produce multiple RNA products. Pre-mRNA
56 splicing is crucial for the proper expression of eukaryotic genes, and spliceosome components
57 are highly conserved from yeast to mammals at the sequence, structure, and functional levels
58 [1, 2].

59 Despite the high conservation between the spliceosomal components in yeast and mammals,
60 *Saccharomyces cerevisiae* has a more compact genome with approximately 300 annotated
61 intron-containing genes. Even with the relatively streamlined splicing landscape in *S.*
62 *cerevisiae*, the prevalence of alternative splicing in this organism has not been well
63 characterized. Several studies have found instances of alternative splicing in individual genes
64 in *S. cerevisiae* [3, 4, 5], and there have been high-throughput methods to find novel splicing [6,
65 7, 8, 9, 10]. Most of these methods focus primarily or entirely on intron retention or exon
66 skipping, while there has been little description of novel 5' or 3' splice site usage.

67 Here we describe a method for discovery of unannotated splice sites in *Saccharomyces*
68 *cerevisiae* by RNA-seq. We tested the method with wildtype strains as well as strains and
69 conditions that do not lead to direct changes in the splicing machinery itself, but that impact
70 broad cellular conditions, RNA turnover, chromatin remodeling, and histone composition. Many
71 of these strains and conditions help us observe changes in splicing that are not or are rarely
72 detected in wildtype cells under normal growth conditions. We also analyze several strains that
73 include *prp43-1*, a temperature sensitive mutation of the *prp43* gene, an RNA helicase directly
74 involved in spliceosome disassembly to determine if modulating disassembly might affect our
75 ability to detect aberrant splice site sampling. Interestingly, this mutation of *prp43* leads to a
76 decrease in splicing of both annotated introns as well as novel splicing.

77 We compare our results with those of earlier studies and find good agreement, however over
78 two thirds of the splice forms our method predicts, 676 total novel splice forms, have not been
79 previously described. Within known intron-containing genes (ICGs), we find that novel isoforms
80 generate longer introns in samples grown at higher temperature. We also find that deletion of
81 *xrn1*, an evolutionarily-conserved 5' to 3' exonuclease, leads to an accumulation of novel splice
82 isoforms that do not use either an annotated 5' nor an annotated 3' splice site, both within
83 known ICGs and unannotated ICGs. While most of our novel splice forms use a known 5' splice
84 site and a novel 3' splice site, we find that deletion of *ume6*, a component of a histone
85 deacetylase complex that regulates early meiosis and that has been recently shown to have a
86 regulatory role during mitosis [11], causes a dramatic increase in the use of novel 5' splice sites
87 within ICGs. We also find evidence for splicing of RNAs that are antisense to annotated
88 transcripts. Together, our results indicate that there is significantly more splicing than previously
89 thought. This suggests that the opportunity for splicing in the form of "latent" introns is a key
90 feature of the yeast genome. Furthermore, changes in the activity of the gene expression
91 machinery or the cells' environment can significantly alter this rich splicing landscape.

92

93 **Results**

94

95 **Alternative splicing is widespread in *Saccharomyces cerevisiae***

96 In order to explore the extent of splicing across the yeast transcriptome in a high confidence
97 manner, we implemented a novel approach to allow us to leverage a large amount of RNA
98 sequence data while imposing stringent filters (figure 1). Briefly, a large number of RNA-seq
99 datasets are combined such that each novel splice form needs at least a single read that aligns
100 without mismatches to the novel junction to pass the initial filter. This read cannot also align to
101 known transcripts or the *Saccharomyces cerevisiae* genome, even with a large number of
102 mismatches. Once a splice junction has been identified in this way, all reads are realigned to
103 the newly discovered splice forms, annotated transcripts, and the genome to find the optimum
104 alignment. Each novel splice form is then scored based on how well its 5' splice site, 3' splice
105 site, and predicted branch point fit the consensus sequence for *Saccharomyces cerevisiae*
106 splice signals. These scores are then used to compute p-values that represent how likely a
107 splice product score of this strength is to occur by chance (figure S1).

108 To capture as many novel splice forms as possible, we incorporated datasets from a variety of
109 strains and experimental conditions that are known or suspected to have an impact of splicing,

110 some accessed from previous studies [12, 13] and some novel (Table S1). In addition to 4
 111 wildtype samples grown at three different temperatures, these samples include 2 biological
 112 replicates each of two deletions in genes involved in decay, including *xrn1*, a 5'-3' exonuclease,
 113 and *upf1*, required for efficient nonsense-mediated decay. We also include *htz1Δ* which
 114 encodes the histone variant H2A.Z and *swr1Δ* which is required to exchange H2A for H2A.Z in
 115 chromatin as well as double mutants *swr1Δ xrn1Δ*, *swr1Δ upf1Δ*, *htz1Δ xrn1Δ*, and *htz1Δ*
 116 *upf1Δ*. Cells that lack H2A.Z are found to have impaired splicing of intron-containing genes
 117 (ICGs), particularly genes that have suboptimal splice sites [13, 14]. We also analyze RNA-seq
 118 data from *snf2Δ*, which leads to an increase in use of non-canonical branch point and 5' splice
 119 site sequences in annotated ICGs and *ume6Δ*, which derepresses genes implicated in meiosis
 120 in *Saccharomyces cerevisiae*, and the double mutants *upf1Δ snf2Δ*, *xrn1Δ snf2Δ*, and *ume6Δ*
 121 *snf2Δ*. Previous studies suggest that *snf2Δ* increases splicing by altering ribosomal protein
 122 gene expression and *ume6Δ* allows expression of genes that are usually repressed. We also
 123 include *set1Δ* and *set2Δ*, deletion mutations in histone methyltransferase genes. Finally, since
 124 we previously showed that the DEAH protein Prp43 contributes to disassembly of suboptimal
 125 spliceosomes using a *prp43* DAMP allele [13], we included a temperature sensitive mutant of
 126 *prp43*, *prp43-1*, as well as *set1Δ prp43-1* and *set2Δ prp43-1*. All *set1Δ*, *set2Δ*, and *prp43-1*
 127 strains contribute two samples to our workflow, one grown at 25° and one at 37°. Taken
 128 together, these datasets represent 29 samples across a variety of *Saccharomyces cerevisiae*
 129 strains and growth temperatures. By leveraging a large number of datasets, we are able to
 130 discover more novel splice sites and determine cellular conditions that lead to changes in
 131 alternative splice site usage.

Strain	Read counts	Aligned read counts
Wildtype	20395478	14160763
<i>xrn1Δ</i>	36132564	25787231
<i>upf1Δ</i>	32005528	22034981
<i>swr1Δ</i>	38830226	21964031
<i>htz1Δ</i>	35969862	18508451
<i>swr1Δ xrn1Δ</i>	86744422	50807249
<i>swr1Δ upf1Δ</i>	11369368	6729660
<i>htz1Δ xrn1Δ</i>	52891892	35039812
<i>htz1Δ upf1Δ</i>	49460562	25424884
Wildtype	35059908	27722481
<i>snf2Δ</i>	34016283	25657587
<i>upf1Δ</i>	21328241	14465110
<i>upf1Δ snf2Δ</i>	31564582	23487340
<i>xrn1Δ</i>	29688304	20282088
<i>xrn1Δ snf2Δ</i>	34289534	27043467
<i>ume6Δ</i>	40071056	30161996

ume6Δ snf2Δ	22534282	17424922
Wildtype 25° C	126355512	97015734
Wildtype 37° C	140962504	82278728
set1Δ 25° C	120413012	75538896
set1Δ 37° C	150376422	84044314
set2Δ 25° C	115930956	61396568
set2Δ 37° C	105592626	63182213
prp43-1 25° C	113138574	78671916
prp43-1 37° C	93850396	54354046
prp43-1 set1Δ 25° C	180103858	107502423
prp43-1 set1Δ 37° C	116668118	70474357
prp43-1 set2Δ 25° C	129498126	93547450
prp43-1 set2Δ 37° C	107964066	68127119

132

133 Using our method, we discover evidence for 944 novel splice events across 408 transcripts with
 134 $p < 0.05$ (Table 1). Of these novel events, the majority are novel splice products found within
 135 known ICGs, using either the annotated 5' splice site with a novel 3' splice site (Table S2,
 136 $n=537$) or using the annotated 3' splice site with a novel 5' splice site (Table S2, $n=129$).
 137 Additionally, we found several cases of novel splicing that do not use any annotated splice sites
 138 within annotated ICGs (Table S2, $n=22$) and in unannotated ICGs (Table S2, $n=198$). We also
 139 find evidence for splice forms that are antisense to known transcripts (Table S2, $n=50$).

	All		$p < 0.05$		$p < 0.01$		p
	Events	Transcripts	Events	Transcripts	Events	Transcripts	Events
All	6187	1937	944	408	852	386	586
Antisense	2504	1110	50	41	46	37	33
Intronless	2709	1141	198	149	168	132	107
New 3'	561	196	537	194	502	191	368
New 5'	140	95	129	90	113	84	63
New both	273	120	22	18	15	12	7

140

141 Many of the events we discover using our method are low abundance. About half of the total
 142 events are represented by fewer than five read counts across our datasets (figure S2A).
 143 However, these low abundance events are high confidence due to our approach's stringent
 144 discovery protocol and p -value cutoff. Furthermore, these high confidence, low abundance

145 novel splice events reveal splicing that is unlikely to be found in the high abundance data.
146 Different classes of novel splicing are more prevalent in our low abundance novel splicing data.
147 Specifically, novel splicing within annotated ICGs is more common in our high abundance data,
148 while novel splicing within unannotated ICGs and antisense RNAs are more common in our low
149 abundance data (figure S2C).

150

151 **Validation of novel splicing**

152 In order to best validate our method by RT-PCR, we chose candidates with a variety of read
153 counts and splice site scores representing each category of novel splicing we find, including
154 novel 5' splice site with annotated 3' splice site, novel 3' splice site with annotated 5' splice site,
155 both novel splice sites within an annotated ICG, novel splicing in an unannotated ICG, and
156 splicing of transcripts that are antisense to annotated genes. In addition, we selected some
157 genes which show only a single novel splice form and others with several. The results of our
158 validation are shown in Figure 2 using oligonucleotides listed in table S3. *BIG1*, an integral
159 membrane protein gene of the endoplasmic reticulum, shows a single novel splice event
160 utilizing a novel 3' splice site with an annotated 5' splice site in our data with 19 total read
161 counts spread roughly evenly across our 29 samples and a splice site score p-value of 0.0003.
162 The novel splice event in *BIG1* is in-frame with the annotated form, and would be predicted to
163 generate a protein product that has 6 additional amino acids. *SIM1*, a gene thought to
164 participate in DNA replication, shows a total of four novel splice forms, each utilizing an
165 annotated 3' splice site and novel 5' splice site with total read counts ranging from 2 to 35 and
166 splice site p-values ranging from 0.018 to 0.0004. The intron in *SIM1* is located in the 5' UTR,
167 therefore the novel splice events would not be expected to yield a new protein product, although
168 potential regulation in the 5' UTR could be altered. *MCR1*, a gene involved in ergosterol
169 biosynthesis in mitochondria, shows a single novel splice form utilizing both novel 5' and 3'
170 splice sites within an annotated ICG with 49 total read counts, 9 of which are found across our
171 *snf2Δ* samples, and a splice site p-value of 0.004. The annotated and novel introns for *MCR1*
172 are both found in the 5' UTR and are therefore unlikely to yield different protein products.
173 However, we previously showed that changes in *SNF2* expression can affect splicing of others
174 transcripts to alter mitochondrial function [15]. *SPF1*, which encodes an ion transporter of the
175 ER membrane, is an unannotated ICG that shows five novel splice forms in our data with read
176 counts ranging from 1 to 4 across all of our samples and p-value scores ranging from 0.004 to 6
177 $\times 10^{-6}$. Of the five novel splice forms found in *SPF1*, only one is in frame, which would produce
178 a protein product that has 134 fewer amino acids. Finally, we find evidence for splicing in the
179 antisense direction to *LEU4*, with a total of six splice forms with counts ranging from 1 to 33
180 across all samples and splice site p-values ranging from 0.0001 to 2.8×10^{-6} . Not surprisingly,
181 of the 48 sequence reads that derive from spliced reads antisense to *LEU4*, 28 come from
182 samples which are deleted of *SET2*, a histone methyltransferase that has previously been
183 implicated in suppressing antisense transcription [16]. Together, we validate each category of
184 novel splicing we observe in our data and validate candidates ranging from a single read count
185 to 49 read counts and from splice site p-values ranging from 0.018 to 2.8×10^{-6} .

Primer	Sequence
BIG1 F	GTTGCTTATTATGTGTGGAAGCTTTTG
BIG1 R	CACCTGGTCTACGTTACAATACTCC

SPF1 F	CCCAGAGA GCCACAAGTT GATCTTG
SPF1 R	CATTAGAGGCAATCTTGACCTGTTGC
SIM1 F	GTGCTA CCCAACTACT TACATTCCT
SIM1 R	CTTGGCTAAGGCAGCAGATGAAAC
MCR1 F	CCTTGATTGGTGTCTTGTGCGAGAGAG
MCR1 R	ACTCACAGGAGTGTATGGTCTCACC
LEU4 antisense	GAGAGTATTATTGCTCTTGCTGAGC
LEU4 R	CTTGTTGGGAAGGGTCATCCTTAG

186

187 **Comparison with other sequence-based approaches for alternative splicing**

188 Even though our method filters (1) putative splice forms with strong similarity to the
189 *Saccharomyces cerevisiae* genome, (2) known transcripts, and (3) those found outside of
190 genes, the splice sites identified by our method include many reported in recent studies. Of the
191 522 novel splice sites described in Kawashima et al. [6], 420 are discoverable by our method,
192 236 are found in our raw data and 189 pass our $p < 0.05$ filter. Of those described in Schreiber
193 et al. [7], 248 out of 314 of the described splice events are discoverable by our method, 214 are
194 in our data and 185 pass our $p < 0.05$ filter. Qin et al. [10] describe a method of identification of
195 novel splice forms by lariat sequencing, a process which reveals 5' splice sites and branch
196 points but not 3' splice sites. Of the 45 novel 5' splice sites found in their work, 11 are present
197 in our raw data and 9 pass our $p < 0.05$ filter with at least one corresponding 3' splice site.
198 Gould et al. [9] combined lariat sequencing with RNA-seq to identify both 5' and 3' splice sites
199 along with novel branch point sequences. Of the 213 novel splice sites they report, 194 are
200 discoverable by our method, 133 are found in our raw data, and 114 pass our $p < 0.05$ filter.
201 The overlap between our work and these previous studies illustrates the power of our approach.
202 None of the splice sites described by Qin et al. are found in any of the three other studies, while
203 each of Kawashima et al., Gould et al., and Schreiber et al. have greater overlap with our work
204 than with one another (figure 3). Taken together, of the 944 splice junctions predicted with $p <$
205 0.05 by our method, 268 have been previously described and 676 are novel.

206

207 ***Saccharomyces cerevisiae* contains antisense transcripts that undergo splicing**

208 In addition to discovering novel splice sites within genes already known to contain introns and
209 unannotated ICGs, our method allows us to discover splicing that is antisense to annotated
210 transcripts, primarily in the degradation mutant strains *xrn1Δ* and *upf1Δ* and the histone
211 methyltransferase mutant *set2Δ*. Previous studies report that Set2 suppresses antisense
212 transcription [16]. We find evidence for 50 antisense novel splice events spread across 41
213 transcripts that pass our statistical criteria. Of these 41 transcripts, 37 show a single novel
214 splice event, two transcripts have two unique novel splice events, one transcript has 3 separate
215 events, and a single transcript, which is antisense to the *LEU4* gene, has 6 unique novel splice
216 events. While most of the antisense splice events in our data have low read counts, the *LEU4*
217 isoforms together account for 51 reads across 11 samples, including wildtype, *set2Δ*, and *set2Δ*
218 *prp43-1* at both 25° and 37°, *set1Δ* and *set1Δ prp43-1* at 37° only, and *xrn1Δ*, *swr1Δ upf1Δ*, and
219 *upf1Δ snf2Δ*. Together, *LEU4* antisense splicing represents over 20% of our total antisense
220 read counts across all samples. These six isoforms arise from 3 different 5' splice sites and 4 3'

221 splice sites. 5 out of 6 of these isoforms generate similar mature mRNAs with an intron in the
222 size range from 115 nucleotides to 129 nucleotides and are therefore indistinguishable by RT-
223 PCR (figure 2E). The remaining form is generated from a unique 5' splice and 3' splice and
224 causes an intron of 464 nucleotides, and is low abundance in our RNA sequence data, with only
225 a single read count in a single sample, *set2Δ* at 37°.

226 In unannotated ICGs, we find 149 transcripts that show novel splicing in the sense direction and
227 37 that undergo novel splicing in the antisense direction. Interestingly, the number of transcripts
228 that show novel splicing in both the sense and antisense direction is just one, *DJP1*. If a set of
229 149 genes and a set of 37 genes are each randomly chosen from all *Saccharomyces cerevisiae*
230 transcripts, the expected value of the overlap is a single transcript, suggesting that within
231 unannotated ICGs, presence of novel splicing in the sense direction does not significantly
232 impact the chances of novel splicing of an antisense transcript or vice versa. Within annotated
233 ICGs, 220 transcripts undergo novel splicing in the sense direction and 4 have novel splicing in
234 the antisense direction, indicating no correlation between the splicing of annotated intron-
235 containing genes and their corresponding antisense transcripts.

236

237 **Prp43 is required for efficient splicing of annotated and novel introns**

238 Several of the strains in our analysis include *prp43-1*, a temperature sensitive mutation that is
239 viable at 25° C but not at 37° C. Prp43 is an RNA helicase that has a vital role in spliceosome
240 disassembly and is required for efficient mRNA splicing in *Saccharomyces cerevisiae*. *PRP43*
241 has also been implicated in ribosome biogenesis [17], and we previously showed that
242 decreasing levels of *PRP43* using a DAMP allele can suppress splice defects [13]. To
243 characterize the consequences of the *prp43-1* mutation in splicing, we compared the *prp43-1*
244 strain to a wildtype strain, a *set1Δ prp43-1* strain to a *set1Δ* strain, and a *set2Δ prp43-1* strain to
245 a *set2Δ* strain at both 25° and 37° C. As expected from its role in splicing, *prp43-1* shows a
246 decrease in the splicing of annotated introns in *Saccharomyces cerevisiae* (figure 4A).
247 Interestingly, despite earlier findings that reducing levels of wildtype Prp43 can suppress splice
248 defects and promote splicing of weak introns [13], we find that strains with *prp43-1* show less
249 novel splicing than their counterparts with wildtype *PRP43* (figure 4B).

250

251 **Elevated temperature favors novel introns that are longer than their annotated introns**

252 Across our datasets, we detect fewer novel splice events in high temperature samples than
253 would be expected by sequence depth alone. This is unsurprising given that many of our lower
254 temperature samples are mutants that lead to accumulation of normally degraded products,
255 such as *xrn1Δ* and *upf1Δ*. While the total number of novel splice events is underrepresented at
256 high temperature, the novel splice products generated differ in intron length. When cells are
257 grown at 37° C, novel splicing within annotated ICGs tends to favor intron sizes that are larger
258 than novel splice forms from cells that are grown at 25° or 30° C (figure 5A). This can be
259 explained by effects observed in our most common class of novel splice events, those that use
260 an annotated 5' splice site with a novel 3' splice site. At higher temperature, these novel 3'
261 splice sites tend to be further downstream than the annotated splice sites. We find that most of
262 the temperature-enriched splice events are more commonly found in our samples with wildtype
263 *prp43*, since *prp43-1* decreases both annotated and novel splicing (figure 5B). These results

264 are consistent with work that finds that intron structure can function to control alternative splicing
265 in yeast [18]. Nonetheless, our data reveal a number of previously unreported events.

266 As an example, *TMH18*, a mitochondrial membrane protein gene, has an annotated intron that
267 is 96 nucleotides long and two common novel splice isoforms discovered by our method, both of
268 which generate longer introns. The isoform that is highly enriched in high temperature samples
269 contains a 161 nucleotide intron, while the isoform that is not favored at high temperature
270 contains a 128 nucleotide intron. We speculate that the increase in temperature destabilizes
271 the secondary structure of some pre-mRNA molecules to allow access to normally inaccessible
272 splice sites. This may lead to an increase in use of distant novel splice sites. To confirm this,
273 we analyzed the predicted pre-mRNA secondary structure using MFOLD [19]. The predicted
274 secondary structure with the most favorable free energy shows that the annotated 5' splice site,
275 the annotated 3' splice site, and the novel 3' splice site that is not favored at high temperature
276 are all readily accessible, while the novel 3' splice site that is favored only at high temperature is
277 not (figure 5C).

278 Interestingly, these distant splice sites may be under evolutionary pressure to be unable to code
279 for protein, as 74/84 (88%) of these temperature-enriched splice sites contain premature
280 termination codons. This is similar to the fraction of PTC-generating events found in the totality
281 of our novel splicing events in annotated intron-containing genes, with 572/636 (90%) containing
282 premature termination codons. This agrees with results described by Kawashima et al. [6] that
283 find that stress conditions, including heat shock, cause an increase in non-productive novel
284 splice usage. Together these data suggest that alternative splicing may be a mechanism for
285 reconfiguring the transcriptome in response to stress.

286

287 **Xrn1 deletion increases splice forms that do not use annotated splice sites**

288 Novel splicing is found most commonly in strains in which *xrn1* has been deleted. These strains
289 account for approximately 9% of our total sequence depth but 19% of our total novel splice
290 counts. Interestingly, *xrn1* deletion mutants are particularly enriched in novel splice forms that
291 do not use any annotated splice sites. Roughly 18% of novel splice form counts that use either
292 an annotated 3' or 5' splice site derive from an *xrn1* deletion mutant, while 31% of those that
293 utilize two novel splices are found in an *xrn1* deletion strain. Two examples of unannotated
294 ICGs impacted strongly by *xrn1* deletion are *AGC1* and *MRM2*, both involved in mitochondrial
295 function. Interestingly, these genes' expression increases in *xrn1* Δ , and they each have their
296 highest RPKM values in the five samples in which *xrn1* is deleted [12, 13]. The 50 unannotated
297 ICGs that are most enriched in our *xrn1* Δ strains only have a single GO term in common,
298 "intracellular membrane-bounded organelle," further suggesting an impact on transcripts
299 important for mitochondrial function. Interestingly, a 2012 study found that Xrn1 is critical for the
300 translation of genes necessary for mitochondrial function in *Saccharomyces cerevisiae* [20].
301 Together, this suggests a role for Xrn1 in the regulation of alternative splice products of
302 mitochondrial transcripts.

303

304 **Ume6 Δ -derived increase in IC-RPG expression leads to increase in IC-RPG alternative 305 splicing**

306 Our *ume6Δ* strain is highly enriched in novel splice events which utilize an annotated 3' splice
307 site with a novel 5' splice site. Of the newly discovered splice sites which utilize a novel 5'
308 splice site and a canonical 3' splice site, 59% are in intron-containing ribosomal protein genes
309 (IC-RPGs). However, of the events that are found disproportionately more frequently in *ume6Δ*,
310 70% are in IC-RPGs. Our RNA sequence data suggests that IC-RPG expression increases
311 slightly in *ume6Δ*. Novel splicing in IC-RPGs increases by 70%, and 5' novel IC-RPG splicing
312 increases by 91% in *ume6Δ*-containing samples relative to samples with wildtype *ume6*.

313 Previous studies have shown that Ume6 is degraded under meiotic conditions [21]. This raises
314 the interesting possibility that genetic manipulations that remove Ume6 may lead to changes in
315 the unannotated splicing landscape that are related to what occurs under meiotic conditions.
316 Ongoing experiments will test this model.

317

318 Discussion

319

320 In this study we show that there are many rarely utilized splice sites in *Saccharomyces*
321 *cerevisiae*. Our methodology is capable of discovering many novel splice forms by utilizing a
322 large number of RNA-seq datasets. We are able to do this in a robust, high-confidence manner
323 by excluding reads that can be explained by sequencing error or genomic DNA contamination
324 and filtering based on existing splice site census sequence data. Overall, we find a strong
325 preference for novel splicing using a known 5' splice site and a novel 3' splice site within
326 annotated ICGs, which represent over half of our statistically significant novel splice products.
327 However, due to our high total sequencing depth and variety of strains and experimental
328 conditions we are still able to find a relatively large number of novel splice forms that utilize
329 known 3' splice sites with novel 5' splice sites, those that utilize novel 3' and 5' splice sites
330 within annotated ICGs, those that utilize two novel splice sites within unannotated ICGs, and
331 splicing of transcripts that are antisense to annotated genes. The large number of novel splice
332 events that we discover allow us to correlate changes in splice site preferences with different
333 mutant strains and experimental conditions. While it has been tempting to call splicing events
334 that have not been previously annotated as "errors," we believe that these data actually reveal
335 the remarkable substrate flexibility of an evolutionarily conserved enzyme. Moreover, it stands
336 to reason that if splicing is to contribute to adaptation and, ultimately, evolution of multicellular
337 organisms, then an array of sequences, not simply those that match a strong consensus, need
338 to be recognized and spliced. The results described here provide a window into this sequence
339 landscape.

340 In the adjoining manuscript [22], the authors find 229 "protointrons" – rapidly evolving,
341 inefficiently spliced introns. Of these, we define 60 as novel introns with an additional 9 found in
342 our RNA sequence data but filtered out due to low splice site scores. The limited overlap
343 between the methods highlights that neither of our methods has reached saturation.
344 Furthermore, the 10 strains in our analysis with the greatest normalized overlap with Talkish et
345 al. are the 10 strains that include either *xrn1Δ* or *upf1Δ*, suggesting that overlap between the
346 methods is more driven by stabilization of protointron-containing splice products than biological
347 similarity between "hungry spliceosome" conditions and the strains used here. Talkish et al.
348 confirm by RT-PCR several splice events that would be removed from our analysis due to poor

349 splice signals, indicating that our method's stringency likely eliminates many true splice
350 products. The approach used here scores putative novel splice sites by learning from
351 annotated splice sites, so as additional protointrons with unusual splice sites are discovered and
352 validated, our method's ability to discover these forms will increase.

353 The analysis of several mutant strains allows insights into splice site selection in
354 *Saccharomyces cerevisiae*. We find general trends in our data as well as specific effects for
355 particular strains or conditions. For example, elevated temperature leads to an increase in
356 novel splice form intron length, *xrn1Δ* leads to a large increase in novel splice forms that do not
357 use an annotated splice site, and *ume6Δ* leads to an increase in novel splicing of intron-
358 containing ribosomal protein genes.

359 Despite the fact that analyzing many mutant strains can increase our understanding of splice
360 site selection, every new category and class of splice form that we find is observed in our
361 wildtype data. While *xrn1Δ* dramatically increases the number of splice products that arise from
362 use of two novel splice sites, we see many examples of splice products that use two novel
363 splice sites in our wildtype samples. This holds true for elevated temperature, *ume6Δ*, and
364 *set2Δ*, as well as splicing of transcripts that are antisense to known transcripts, splicing that
365 uses a single annotated splice site, and splicing in unannotated ICGs. Together, this indicates
366 that while mutant strains are useful for correlating genetic changes with splicing changes, even
367 wildtype *Saccharomyces cerevisiae* in normal conditions are capable of producing a large
368 variety of splice products. Taken together, these findings illustrate the diversity and depth of
369 splicing in *Saccharomyces cerevisiae*, and also show the presence of latent introns that are
370 found across the genome.

371

372 **Methods**

373

374 **Public datasets**

375

376 *Saccharomyces cerevisiae* wildtype, *xrn1Δ*, *upf1Δ*, *swr1Δ*, *htz1Δ*, *swr1Δ xrn1Δ*, *swr1Δ upf1Δ*,
377 *htz1Δ xrn1Δ*, and *htz1Δ upf1Δ* strain RNA-seq data were download from GEO (accession
378 number GSE97416). Additional wildtype, *upf1Δ*, and *xrn1Δ* RNA-seq data as well as *snf2Δ*,
379 *ume6Δ*, *upf1Δ snf2Δ*, *xrn1Δ snf2Δ*, and *ume6Δ snf2Δ* strain RNA-seq data were also
380 downloaded from GEO (accession number GSE94404).

381

382 **Yeast culture and Sequencing**

383

384 Wildtype, *set1Δ*, *set2Δ*, *prp43-1*, *set1Δ prp43-1*, and *set2Δ prp43-1* strains were grown at 25°C
385 to OD 0.3. Then, cultures were equally split, half remaining at 25°C and half shifting to 37°C for
386 four hours. 10 mL of cells were pelleted from each sample and RNA extraction was performed.
387 Then, 20 μg total RNA per sample was treated with DNase I (Roche) and depleted of rRNA with
388 the Ribo-Zero Gold rRNA Removal Kit (Illumina). RNA-seq libraries were prepared using the

389 Illumina TruSeq v2 RNA Kit. 50 base pair paired-end reads were generated on an Illumina
390 HiSeq 4000.

391

392 **Sequence alignment**

393

394 RNA sequence data were combined and aligned to the *Saccharomyces cerevisiae* SacCer3
395 genome reference and Ares Lab Yeast Intron Database Version 3 [23] in a single step using
396 STAR [24] allowing up to six mismatches and no unannotated gaps. Sequence reads that fail to
397 align in this step are then aligned to the SacCer3 genome again, this time allowing no
398 mismatches, a single gap, at most one alignment locus, and at least 15 nucleotide overhang on
399 each end of the alignment. Successful alignments in this step define putative novel introns. In
400 a final alignment step, all reads are aligned to the SacCer3 genome, the Ares Intron Database,
401 and newly defined putative novel introns in one step, allowing at most 1 alignment locus with up
402 to 2 mismatches. Counts for novel splice forms are based on those derived from this alignment
403 step.

404

405 **Splice site scoring**

406

407 5' splice site, 3' splice site, and branch point scores were generated based on how closely each
408 splicing signal matches the consensus sequence for that signal. For each position within a
409 splice site, the total number of adenine, cytosine, guanine, and thymine bases present at that
410 position in annotated splice signals was determined based on the Ares Lab Intron Database,
411 and then the proportion of each nucleotide at that position was found by dividing by the total
412 number of annotated splice products. The score for each splice site was calculated as follows:

$$Score = \sum_{i=1}^k N_i / X$$

413 Where i is the position within the splice signal, k is the number of positions in the splice signal,
414 N_i is the count of that specific nucleotide at that position, and X is the total counts for the splice
415 signal. To score each putative novel splice form, the 5' and 3' splice sites are determined from
416 RNA-seq data, while branch point sequences were chosen by selecting the highest scoring
417 possible branch point sequence within a maximum distance of 200 nucleotides from the 3'
418 splice site. The score for each putative novel splice site is simply the product of its 5' splice site
419 score, its 3' splice site score, and its branch point score. Putative novel splice sites were
420 considered antisense if the score in the antisense direction is higher than the score in the sense
421 direction.

422

423 **Statistical analysis of splice sites**

424

425 P-values for putative novel splice sites were generated by comparing the potential splice site to
426 all possible combinations of 5' splice sites, 3' splice sites, and branch points. Based on a six
427 nucleotide 5' splice site, a three nucleotide 3' splice site, and a seven nucleotide branch point,
428 there are 4,294,967,296 possible combinations of splice signals that can be described. To
429 convert our putative novel splice form scores to p-values we compared the score to the scores
430 of all possible combinations of splice signals. The raw p-value is the fraction of these
431 combinations that score equal to or better than the putative novel form, which also represents
432 the chance of seeing a splice score as good or better than this by chance. We then correct the
433 raw p-values for multiple hypothesis testing using the Bonferroni correction by multiplying each
434 raw p-value by the number of tests conducted, which is the total number of putative novel splice
435 sites times two, since each splice site score is calculated in both the sense and antisense
436 direction. All adjusted p-values greater than one are then set to one.

437

438 **Splicing Efficiency Calculation**

439

440 To calculate the splicing efficiency of annotated splice sites, reads were aligned with STAR [24]
441 allowing only a single alignment locus, only annotated splice sites, and at most two mismatches.
442 Aligned reads within ICGs were categorized as exonic, spliced, or unspliced. We generated
443 normalized spliced and unspliced read counts by dividing the raw counts in each category by
444 the number of possible alignments that can fall into that category. This equates to read length
445 minus one for spliced reads and the intron length plus the read length minus one for unspliced
446 reads. Splicing efficiency is then computed as normalized spliced counts divided by the sum of
447 normalized spliced and normalized unspliced counts.

448

449 **RNA Folding**

450

451 To find the optimum secondary structure for TMH18, we used the MFOLD web server with the
452 pre-mRNA sequence accessed from the Saccharomyces Genome Database
453 (<https://www.yeastgenome.org/>) and default MFOLD parameters [19;
454 <http://unafold.rna.albany.edu/?q=mfold>]. The optimum secondary structure was visualized using
455 MFOLD.

456

457 **RNA isolation and RT-PCR**

458

459 RNA was isolated from log phase yeast by hot phenol:chloroform:isoamyl alcohol (PCA)
460 extraction with SDS. The RNA was then precipitated with ethanol. 20 µg of total RNA was
461 DNase-treated with 30 U DNase 1 (Roche) for 1 hour at 25°C. 1 µg of DNase-treated RNA was
462 used for cDNA synthesis. cDNA synthesis was performed using the Maxima First Strand cDNA
463 Synthesis Kit (ThermoFisher). 1 µL of cDNA was used for the Taq PCR reaction using gene-
464 specific primers to analyze splicing.

465

466 **Data Visualization**

467

468 Venn diagrams to view the overlap between this and previous works were generated using
469 Venny 2.1.0 [27]. Box plots and bar graphs were generated using the MATLAB functions
470 boxplot and bar, respectively.

471

472 **Accession numbers**

473

474 Data generated in this study is available under GEO accession number GSE120497.

475

476 **Acknowledgements**

477 We thank the Dr. Tracy Johnson lab (UCLA) for comments and suggestions to improve the
478 manuscript and Manuel Ares, Jr. for generously sharing unpublished data.

479

480 **References**

481

- 482 1. Wahl MC, Will CL, Luhrmann R. The spliceosome: design principles of a dynamic RNP
483 machine. *Cell*. 2009;136: 701–718.
- 484 2. Li X, Liu S, Jiang J, Zhang L, Espinosa S, Hill RC, et al. CryoEM structure of *Saccharomyces*
485 *cerevisiae* U1 snRNP offers insight into alternative splicing. *Nat Commun*. 2017 Oct
486 19;8(1):1035.
- 487 3. Juneau K, Nislow C, Davis RW. Alternative splicing of PTC7 in *Saccharomyces cerevisiae*
488 determines protein localization. *Genetics*. 2009 Sep;183(1): 185-94.
- 489 4. Hossain MA, Rodriguez CM, Johnson TL. Key features of the two-intron *Saccharomyces*
490 *cerevisiae* gene SUS1 contribute to its alternative splicing. *Nucleic Acids Res*. 2011 Oct;39(19):
491 8612-27.
- 492 5. Hossain MA, Claggett JM, Edwards SR, Shi A, Pennebaker SL, Cheng MY, et al.
493 Posttranscriptional Regulation of Gcr1 Expression and Activity Is Crucial for Metabolic
494 Adjustment in Response to Glucose Availability. *Mol Cell*. 2016 May 5;62(3): 346-358.
- 495 6. Kawashima T, Douglass S, Gabunilas J, Pellegrini M, Chanfreau GF. Widespread use of
496 non-productive alternative splice sites in *Saccharomyces cerevisiae*. *PLoS Genet*. 2014 Apr
497 10;10(4): e1004249.
- 498 7. Schreiber K, Csaba G, Haslbeck M, Zimmer R. Alternative Splicing in Next Generation
499 Sequencing Data of *Saccharomyces cerevisiae*. *PLoS One*. 2015 Oct 15;10(10): e0140487.

- 500 8. Plass M, Codony-Servat C, Ferreira PG, Vilardell J, Eyras E. RNA secondary structure
501 mediates alternative 3' splice site selection in *Saccharomyces cerevisiae*. *RNA*. 2012 Jun;18(6): 1103-
502 15.
- 503 9. Gould GM, Paggi JM, Guo Y, Phizicky DV, Zinshteyn B, Wang ET, et al. Identification of new
504 branch points and unconventional introns in *Saccharomyces cerevisiae*. *RNA*. 2016 Oct;22(10):
505 1522-34.
- 506 10. Qin D, Huang L, Wlodaver A, Andrade J, Staley JP. Sequencing of lariat termini in *S.*
507 *cerevisiae* reveals 5' splice sites, branch points, and novel splicing events. *RNA*. 2016 Feb;
508 22(2): 237–253.
- 509 11. Lardenois A, Stuparevic I, Liu Y, Law MJ, Becker E, Smagulova F, et al. The conserved
510 histone deacetylase Rpd3 and its DNA binding subunit Ume6 control dynamic transcript
511 architecture during mitotic growth and meiotic development. *Nucleic Acids Research*. 2015
512 Jan; 43(1): 115-28.
- 513 12. Venkataramanan S, Douglass S, Galivanche AR, Johnson TL. The chromatin remodeling
514 complex Swi/Snf regulates splicing of meiotic transcripts in *Saccharomyces cerevisiae*. *Nucleic*
515 *Acids Research*. 2017; 45(13): 7708-7721.
- 516 13. Neves LT, Douglass S, Spreafico R, Venkataramanan S, Kress TL, Johnson TL. The
517 histone variant H2A.Z promotes efficient cotranscriptional splicing in *S. cerevisiae*. *Genes &*
518 *Development*. 2017; 31(7): 702-717.
- 519 14. Nissen KE, Homer CM, Ryan CJ, Shales M, Krogan NJ, Patrick KL, et al. The histone
520 variant H2A.Z promotes splicing of weak introns. *Genes & Development*. 2017; 31(7): 688-701.
- 521 15. Awad AM, Venkataramanan S, Nag A, Galivanche AR, Bradley MC, Neves LT, et al.
522 Chromatin-remodeling SWI/SNF complex regulates coenzyme Q6 synthesis and a metabolic
523 shift to respiration in yeast. *Journal of Biological Chemistry*. 2017;292(36): 14851-14866.
- 524 16. Venkatesh S, Li H, Gogol MM, Workman JL. Selective suppression of antisense
525 transcription by Set2-mediated H3K36 methylation. *Nature Communications*. 2016; 7: 13610.
- 526 17. Leeds NB, Small EC, Hiley SL, Hughes TR, Staley JP. The splicing factor Prp43p, a DEAH
527 box ATPase, functions in ribosome biogenesis. *Molecular Cell Biology*. 2006; 26(2): 513-22.
- 528 18. Meyer M, Plass M, Pérez-Valle J, Eyras E, Vilardell J. Deciphering 3' splice site selection in the
529 yeast genome reveals an RNA thermosensor that mediates alternative splicing. *Molecular Cell*.
530 2011; 43(6): 1033-9.
- 531 19. Zucker M. On finding all suboptimal foldings of an RNA molecule. *Science*. 1989 Apr;
532 244(4900): 48-52.
- 533 20. Sinturel F, Bréchemier-Baey D, Kiledjian M, Condon C, Bénard L. Activation of 5'-3'
534 exoribonuclease Xrn1 by cofactor Dcs1 is essential for mitochondrial function in yeast. *PNAS*.
535 2012; 109(21): 8264-9.
- 536 21. Mallory MJ, Cooper KF, Strich R. Meiosis-specific destruction of the Ume6p repressor by
537 the Cdc20-directed APC/C. *Molecular Cell*. 2007; 27(6): 951-61.

538 22. Talkish J, Igel AH, Perriman RJ, Shiue L, Katzman S, Munding EM, et al. Rapidly evolving
539 protointrons in *Saccharomyces* genomes revealed by a hungry spliceosome. bioRxiv doi:
540 10.1101/515197.

541 23. Grate L., Ares M. Jr. Searching yeast intron data at Ares lab web site. *Methods*
542 *Enzymol.* 2002; 350: 380–392.

543 24. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast
544 universal RNA-seq aligner. *Bioinformatics.* 2013; 29: 15–21.

545 25. Oliveros, JC. Venny: An interactive tool for comparing lists with Venn's diagrams. 2015.
546 <http://bioinfogp.cnb.csic.es/tools/venny/index.html>.

547

548 **Figure 1. Workflow for discovery of novel splice forms.** Multiple RNA sequence datasets
549 are consolidated then aligned to the *Saccharomyces cerevisiae* genome and annotated
550 transcripts liberally, allowing up to six mismatches but no gaps. Reads that fail to align are
551 again aligned, but now allowing a single gap corresponding to a putative intron and no
552 mismatches. This step defines all putative novel introns. Once again all reads are aligned to
553 the *Saccharomyces cerevisiae* genome, known transcripts, and putative novel introns allowing
554 up to two mismatches, which determines each putative novel intron's read counts. Consensus
555 sequences are determined from known *Saccharomyces cerevisiae* 5' splice site, 3' splice site,
556 and branch point sequences, and each putative novel intron is scored based on sequence
557 similarity to the consensus sequences.

558

559 **Figure 2. RT-PCR of representative novel splice isoforms.** (A) *BIG1*; novel 3' splice site
560 within annotated ICG. (B) *SIM1*; 5' UTR intron, novel 5' splice site within annotated ICG. (C)
561 *MCR1*; 5' UTR intron, novel 5' and 3' splice sites within an annotated ICG. (D) *SPF1*; intron
562 within an unannotated ICG. (E) *LEU4* antisense; antisense transcript of the *LEU4* gene
563 containing an intron.

564

565 **Figure 3. Agreement with previous studies.** Venn diagram showing the overlap between the
566 novel splice forms discovered here and those described previously.

567

568 **Figure 4. *prp43-1* leads to a decrease in splicing.** (A) Strains with the mutated *prp43-1* gene
569 show a decrease in splicing efficiency of annotated introns relative to strains with wildtype
570 *PRP43*. (B) Evidence of novel splicing is lower in strains with *prp43-1* versus wildtype *PRP43*.

571

572 **Figure 5. Elevated temperature leads to longer introns.** (A) Boxplot of temperature-
573 enriched splice forms versus temperature-independent splice forms shows higher temperature
574 leads to an increase in intron size. Temperature-enriched splice forms have much higher
575 counts at higher temperature, but are still present at lower temperature. (B) Shifting *prp43-1* to

576 higher temperature causes a dramatic decrease in these temperature-enriched splice forms.
577 (C) Predicted structure of TMH18 with annotated and novel splice sites labeled.

578

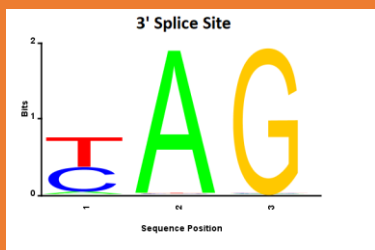
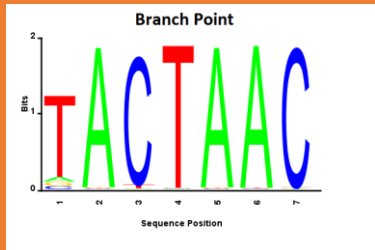
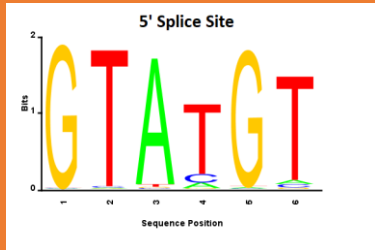
579 **Figure S1. Putative novel splice site scores.** Each putative novel splice site was scored by
580 similarity to known *Saccharomyces cerevisiae* splice signals. The splice site score was plotted
581 against putative intron length for novel splice products within annotated ICGs (blue) and
582 unannotated ICGs (red).

583

584 **Figure S2. Summary of novel splice products.** (A) Over half of the data have fewer than 6
585 counts across all 29 of our datasets. (C) Novel splicing outside of annotated intron-containing
586 genes and splicing of RNAs that are antisense to known transcripts is observed.

Known Splice Junctions

Consensus Sequences



RNA-Seq Data

GMAP Alignment
Align to genome & transcriptome
At most 6 mismatches
0 gaps

Failed Alignments

GMAP Alignment
Align to genome
0 mismatches
At most 1 gap

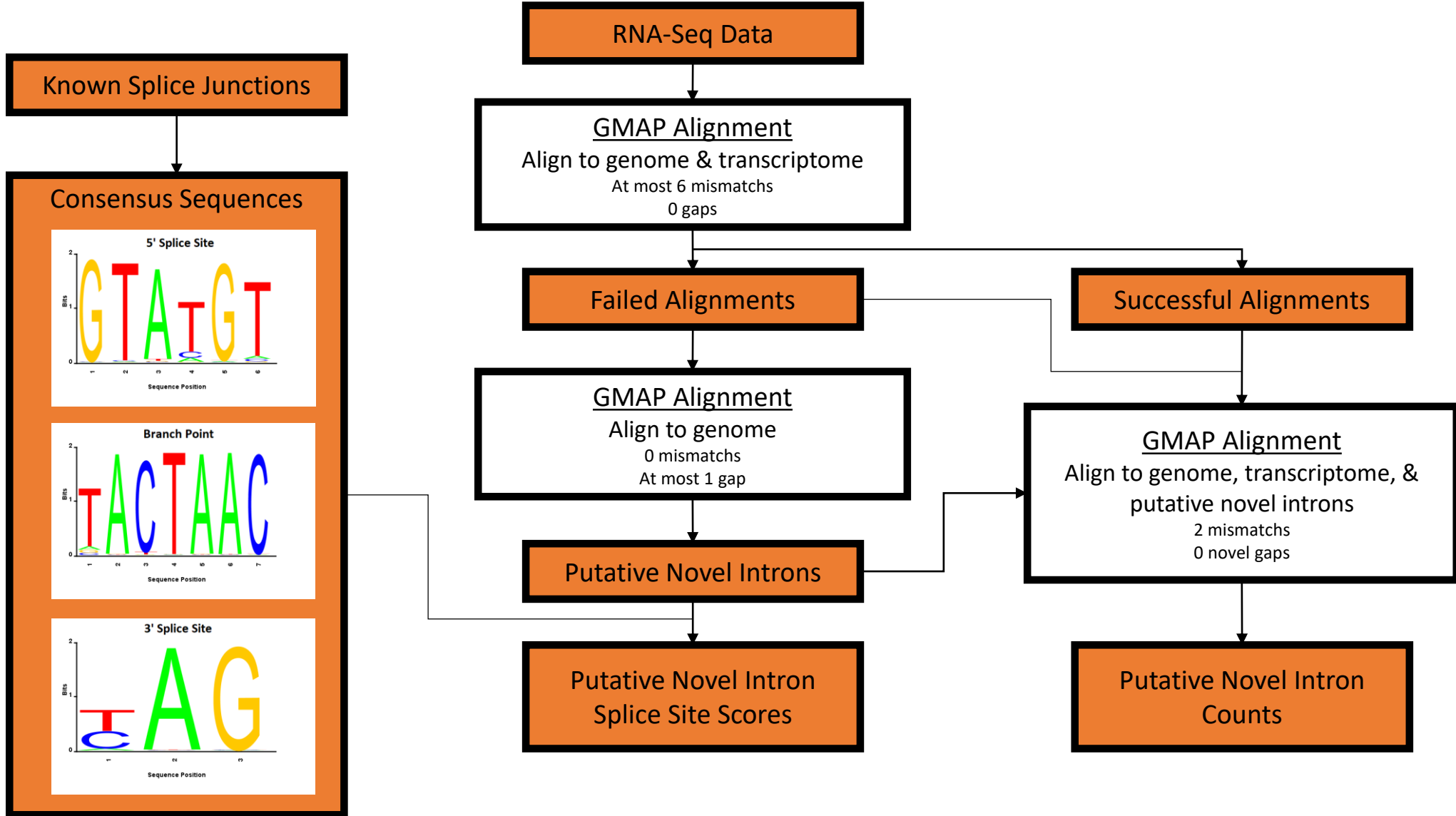
Putative Novel Introns

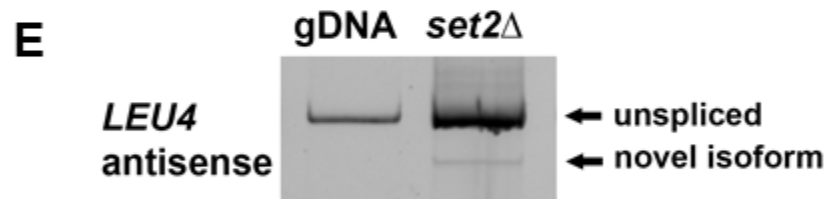
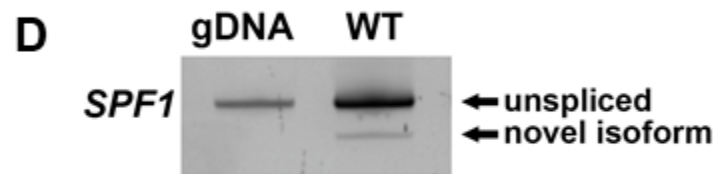
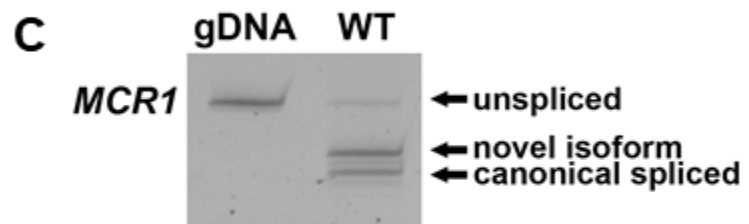
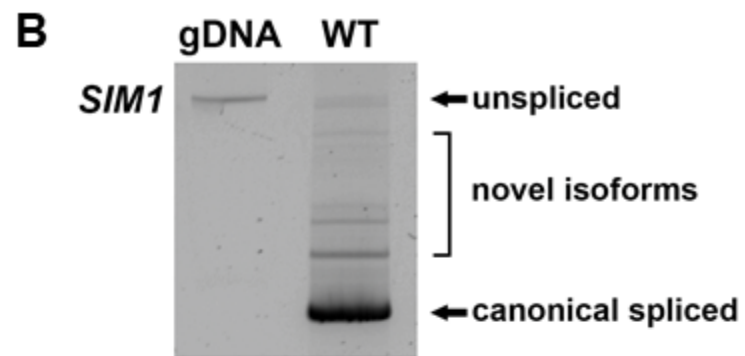
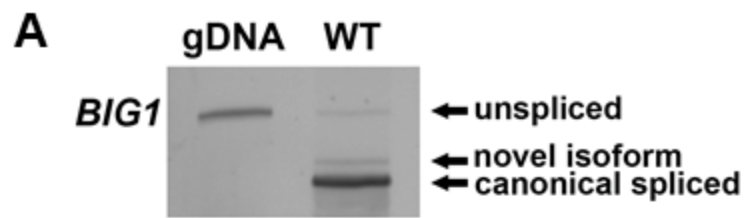
Putative Novel Intron
Splice Site Scores

Successful Alignments

GMAP Alignment
Align to genome, transcriptome, &
putative novel introns
2 mismatches
0 novel gaps

Putative Novel Intron
Counts





Kawashima et al.

Gould et al.

