

[22] Searching Yeast Intron Data at Ares Lab Web Site

By LESLIE GRATE and MANUEL ARES, JR.

Introduction

It must be obvious to every geneticist by now that the future will be consumed by the need to understand how the elemental properties of genes so elegantly described in the past half-century come together with the environment to produce the subtle differences that are key to the fitness of the organism. This will require a partial abandonment of the reductionism so favored since Mendel, to be replaced by the adoption of a more synthetic view that addresses the molecular underpinnings of complex phenotypes, penetrance, expressivity, and the small contributions of many genes. Although many of us were trained to design experiments about single genes, or at the most two interacting genes, our students and researchers need more. More in this case is a healthy computational philosophy and experience.

We have tried to embrace this in our own small way by setting up a searchable database containing information concerning the introns found in the genome of *Saccharomyces cerevisiae*. Since one of us (MA) has training in genes but not computers, and the other (LG) has training in computers but not genes, this effort has been a cultural compromise. Despite its lack of sophistication and dotcom sheen, the database has found many uses in our laboratory and has been accessed by yeast geneticists, splicers, and bioinformaticists the world around. In the following pages we explain the browsing and search capabilities of the site, and how to read and interpret the findings.

Getting Into the Site

Probably the best way to use this chapter is to sit at the computer with the book open, as you go through the descriptions of the different searches. Although some readers may find computers intimidating, there is really no way to damage equipment or files using a Web browser, so no big mistakes can be made. Explore! Experiment! The site can be found by following the "Ares Lab Yeast Intron Database" link from the Ares lab home page at <http://ribonode.ucsc.edu>. The site can also be accessed from the *Saccharomyces* Genome Database (SGD¹) by clicking on their "Yeast WWW Sites" link and scrolling down to their "Yeast Introns" link. Alternatively, access the site directly by typing "http://www.cse.ucsc.edu/research/compbio/yeast_introns.html" into the location window of your favorite browser, and hit return.

¹ C. A. Ball *et al.*, *Nucleic Acids Res.* **28**, 77 (2000).

Au: Pls. supply all author's names. Series style not to use et al. in reference list.

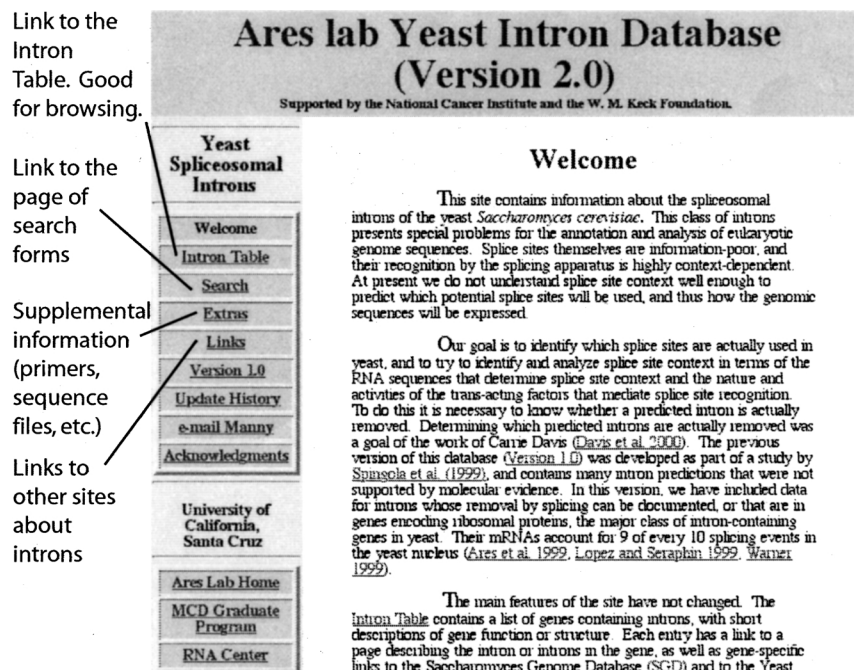


FIG. 1. Front page of the Ares lab Yeast Intron Database.

Au: Cap 1?

The front page is shown in Fig. 1. In addition to some text about introns and where the information comes from, this page has four important links, two of which we will discuss at length. The “Intron Table” link (Fig. 1) will load a large document that includes a table that has the yeast genes with introns (Fig. 2). The next link on the navigation bar, “Search” links to a page (Fig. 3) that contains a set of links to different types of searches that can be performed (discussed below). The “Extras” link goes to a page that contains additional information of interest, such as available PCR primers for detecting splicing of individual introns, text files of intron sequences and alignments, and various graphs and other data related to introns in yeast. The “Links” page includes links to other Web sites concerning introns, of note, Seraphin’s yeast site² and Kent’s *Caenorhabditis elegans* Intronerator.³ Click on the “Intron Table” link and look at the Intron Table page.

Browsing the Intron Table

The Intron Table page is shown in Fig. 2. Each entry has a link to an “Ares Lab Intron Report” which is represented by the entry number. The entries in the table are

² P. J. Lopez and B. Seraphin, *Nucleic Acids Res.* **28**, 85 (2000).

³ W. J. Kent and A. M. Zahler, *Nucleic Acids Res.* **28**, 91 (2000).

Ares lab Yeast Intron Database (Version 2.0)

This site is supported by the National Cancer Institute and the W. M. Keck Foundation.

Yeast Spliceosomal Introns

- Welcome
- Intron Table**
- Search
- Extras
- Links
- Version 1.0
- Update History
- e-mail Manny
- Acknowledgments

University of California, Santa Cruz

- Ares Lab Home
- MCD Graduate Program
- RNA Center
- Computational Biology

Spliceosomal Intron Table

This page contains information extracted from SGD (SacchDB at Stanford) and YPD (at Proteome Inc.) plus results from our lab to create hyper-links about.

Spliceosomal Intron-Containing Genes in *S. cerevisiae*

(The synonyms and text are as obtained from SGD and YPD as of the dates listed below. The most up-to-date information on the entries can be found by following the name links back to the database sources, or starting from main pages of the databases: [SGD](#), [MIPS](#), and [YPD](#) (PROTEOME). This data, as of April 1998, uses the [NEW SPICEOSOMAL PROTEIN NAMING CONVENTION](#). Genes are ordered alphabetically by "Y name" except for US snoRNA, and the numbering is of no real significance.

Page rebuild date: 04/04/2000, SGD data date: 01/05/00, YPD data date: Thu Dec 3 18:29:42 1998

Ares Lab Report	SGD Feature	Synonyms	SGD Locus	Descriptions
	YPD	SGD		
1	SNR12A			[mes] Not in a protein ORF, but in USA snoRNA.
2	SNR12B			[mes] Not in a protein ORF, but in U3B snoRNA.
3	YAL001C	TRC3 TRV115 / PUN24 TRV115 / YAL001C	TRC3	[YPD] RNA polymerase transcription initiation factor TFIIC (w). US M2a subunit [(c) 1995-1998 Proteome, Inc.]
4	YAL002W	EER1 TRP3 / YAL002W	EER1	[YPD] Translation elongation factor EF-1beta, GDP/GTP exchange factor for Tef1p/Tef2p [(c) 1995-1998 Proteome, Inc.]
5	YAL030W	SNC1 YAL030W	SNC1	[YPD] Synaphobrevin (v-SNARE) homolog present on post-Golgi vesicles [(c) 1995-1998 Proteome, Inc.]
6	YBL018C	POPS YBL0301 / YBL018C	POPS	[SacchDB] Processing Of Precursors - a group of proteins that appear to be components of both FNAse P and FNAse MRP

Link to the Ares lab Intron Report page for this intron.

Link to Yeast Proteome Database, Proteome, Inc.

Link to the Saccharomyces Genome Database.

Ares Lab Yeast Intron Containing Gene YBL018C

Intron Table 1 (Ares Lab Intron Report)

Dates	Page rebuild date: Sun May 28 16:32:04 2000, SGD data date: 01/05/00, YPD data date: Thu Dec 3 18:29:42 1998	Additional links to outside databases.
Gene/ORF name	YBL018C	
YPD synonym	POPS YBL0301 / YBL018C	
SacchDB synonym	YBL018C	
SacchDB Locus	POPS	
Description	[YPD] Subunit of both FNAse P and FNAse MRP [(c) 1995-1998 Proteome, Inc.]	Text descriptions of gene function and expression levels
Number of introns	1	
Holotype Data	#mRNAperCell 1.0 ; HalfLife(min) 11 ; Transcription Frequency(turns/hr) 3.4	Genome coordinates
Ares Lab Intron Name	YBL018C_2_186437_185961_INTRON_48_122	Verification
Comments		
Splicing Verified	yes	
Verification description	Davis CA et al.	Link to PubMed abstract of paper(s) describing the intron.
Length Info (in nt)	75 ; to Branch base (start length) 56	
Location in "ori"	start 48 ; stop 122	
Features	GUANGU CUCUACGACGAC UAGAG	Features of the intron
Sequence	>YBL018C_2_186448_186391_FPE GTCATGGGGAAGAAACCTTGTGAGAGAGGCAAAUATUCAGUATCAAU >YBL018C_2_186437_185961_INTRON_48_122 GUAGUUAUUTUUGACUUTUUGAGGUCACUACCGAAGAGAAUAAAC UACUACGACUUAUAUUAUUAAG >YBL018C_2_186315_186266_POST UACUACUUGAGUACAGAGUUGAGGAGGACAGACAGUATUUGAUAUAUA Possible Protein Sequence >YBL018C_2_186437_185961_48_122_PROTEIN IKKTYFRWQTFELSYIFQDYDANALDQITVWGLWALKRSTGTFP EDVEYSLFVDEKLATIRFKAQKDYFSSISSTYIISTDELFGSPLTYSLL QESSLLKLVTDGDELFLKKYVDEEEDQRCI*	FASTA files of the intron and 50 nt up and downstream FASTA file of predicted protein

FIG. 2. The Intron Table and the Intron Report pages.

listed in alphabetical order by ORF name or “Y name” (e.g., YAL001C, meaning: Y, yeast; A, first chromosome; L, left arm; three digit number, assigned to ORF; C, the Crick strand), with the exception of the snoRNA genes U3 genes *SNR17A* and *SNR17B*. These are listed first, RNA coming before protein. Each entry contains a link to the SGD locus page for that gene, as well as information concerning synonyms, gene names (as opposed to ORF names), and a short description of the function of the gene product.

Clicking on the entry number link in the column “Ares lab report” produces a page for the gene in that entry, shown in the lower part of Fig. 2. This report shows much of the data in the underlying database for the intron-containing gene specified at the top of the page. The first four rows of the report reiterate the information and links on the Intron Table, including the short text descriptions. Additional information such as number of introns, expression level of the gene in the experiment by Holstege *et al.*,⁴ the genome coordinates of the intron, any comments we had, and whether or not splicing has been verified experimentally or is predicted are included in the next rows. A link to the Genbank file (if one is available), or the PubMed abstract of the paper describing molecular evidence for splicing, or the method of prediction is included next. Next, physical features of the intron are presented. Length, position relative to the AUG of the ORF (except for 5' UTR introns), and the starting, ending, and branchpoint region sequences of the intron are listed.

In the row labeled Sequence, there are three FASTA files (a “FASTA file” is a standard file format for presenting sequence data) associated with each intron. These include sequences 50 nucleotides (nt) upstream of the 5' splice site (ending with the label “PRE”), the sequence of the intron itself, and the first 50 nt following the intron (labeled “POST”). These sequences can be copied for use with other kinds of sequence searches and alignment programs. Be aware that precise transcription initiation sites for most yeast genes are unknown, and our arbitrary use of 50 nt upstream of the intron does not imply that initiation occurs more than 50 nt upstream of the intron in every case. An example would be the U3 genes in which the first exon is only 16 nucleotides.⁵ The last line, “Possible Protein Sequence,” shows a protein prediction from sequence that has had the intron removed. It is accurate for genes with a single intron in the coding sequence, but the program we use for this may generate faulty predictions for genes with two introns or with introns in the 5' UTR. Therefore, the protein predictions should be used with care.

The Intron Table page is a large (~200 K) document, and clicking back and forth from it to the report pages can be slow, even with a fast connection. The best way to avoid reloading this large page is to open it once, and open a new

⁴ F. C. Holstege *et al.*, *Cell* **95**, 717 (1998).

⁵ E. Myslinski, V. Segault, and C. Branlant, *Science* **247**, 1213 (1990).

browser window in which to view the report. Usually the right button on a three-button mouse will open a link in a new window. With a single-button mouse, hold down the button with the cursor over the link until the pop-up menu appears on the screen near the cursor. Select "New window with this link" or "Open link in new window" and release the mouse button. The linked report page should open in a new window that is displayed over the top of the Intron Table. When done viewing the report, close the new window and return to the Intron Table window. This avoids having to reload the Intron Table page each time. A little bit of practice will make these operations go more smoothly.

Other features of most browsers are also useful, for example, the "Find..." function available under the "Edit" menu. Selecting "Find..." opens a small dialog box with a space in which to type what you want to find in the open page. We use this to find things in the Intron Table without scrolling up and down forever. Try it by searching for the actin intron: type "actin" in the dialog box and clicking "Find." The browser searches the text for "actin" and highlights the first occurrence of this string of letters. If we start at the top of the Intron Table, the first occurrence of "actin" is in the *ARP2* gene entry (actin-related protein). Using "Find again" to go to the next occurrence, we find *SAC6* (actin filament bundling protein), and finally we come to *ACT1* itself. Other intron-containing genes that work with actin can be found by continuing the process to the end of the Table. This approach is good for browsing for key words as well as for finding a particular gene in the Table by name.

Searching the Intron Database

The simple searches offered by the browser are limited to short text strings (i.e., sequences of letters and numbers) and only in the page displayed. More complex queries can be generated using four types of searches available by clicking "Search" on the navigation bar on the left of most pages from the site. The Search page is shown in Fig. 3 and basically contains links to each of the searches. The first is the Intron Table Text Query, which allows selected properties of introns to be defined and returns a smaller table identical in format to the large Intron Table, but which only contains the entries that match the query. The Intron Splice Signals page can be used to identify introns that have particular branchpoint or splice site sequences of interest. The YAG Query page allows specifics concerning the 3' splice site and the region between the branchpoint and the 3' splice site to be captured. Finally, the Intron Sequence Search allows identification of introns containing a particular nucleotide sequence of interest. Below we will describe examples of how to use these search functions.

Go to the Intron Database Text search page (Fig. 4) from the "Search" page by clicking on the "Intron Table Text Query" link. This page provides a more sophisticated means to find and organize a subset of introns of interest. One can search

Search for entries with specific words

Search for introns with particular splice sites and/or branchpoints

Ares lab Yeast Intron Database (Version 2.0)

This site is supported by the National Cancer Institute and the W. M. Keck Foundation.

Yeast Spliceosomal Introns

Welcome

[Intron Table](#)

[Search](#)

[Extras](#)

[Links](#)

[Version 1.0](#)

[Update History](#)

[e-mail Manny](#)

[Acknowledgments](#)

Search the Database

These links allow you to create smaller tables that extract a desired subset of introns based on information in the database:

The [Intron Table Text Query](#) allows searches for entries that have a text word in the description field of the database entries, e.g. "ubiquitin", and generates a table containing only those entries that match.

The [Intron Splice Signals Query](#) allows identification of introns with specified splicing signals, for example, all introns with GUACGU 5' splice sites can be listed, with links to SGD.

The [YAC Query](#) allows analysis of the number of potential 3' splice sites between the branchpoint and the annotated 3' splice site.

The [Sequence search on introns](#) allows pattern searching of the intron sequence database and reports which introns have a specified sequence how many times and where.

Search for introns with certain 3' splice site conditions

Search for introns containing specific sequence elements

FIG. 3. The Search Page.

for introns by different properties, such as description text words, gene names, comments, intron number, and the length of different parts of the intron itself. The example described in Fig. 4 shows the output from a search for "alternative splicing" in the comments box. Submitting this request (by typing "alternative splicing" and clicking on the "Submit intron table query" button) returns a page that has a small table identical in form to the large Intron Table, but which only contains the entries that fit the terms requested. In this case there are two yeast genes known to have alternative splicing, *MTR2* and *SRC1*.⁶ Each of these appears in a row as it would in the large Table, but without the other introns. As with the large Intron Table, clicking on the entry number generates the intron report page for that entry. This type of search allows rapid winnowing of lists of introns to obtain those that fit certain criteria of interest.

Multiple constraints can be imposed on the search, because the form treats each piece of information submitted in the different boxes as "AND" rather than

⁶ C. A. Davis *et al.*, *Nucleic Acids Res.* **28**, 1700 (2000).

“OR.” For example, try typing “UTR” in the comments box and submitting the form. This will return all genes with introns in the 5′ UTR. Note the number of introns. Now try typing “ribosomal protein” in the “Description Text Pattern” box and “UTR” in the comments box. This will only return ribosomal protein genes that have introns in their 5′ UTRs. This is a search that retrieves too many results, but it can be made more restrictive. To make a search less restrictive, one can use a partial word as a text pattern (since a “text pattern” is not exactly the same as a word). For example, the text pattern “ubi” recalls genes encoding ubiquitin, ubiquitin-conjugating enzymes, and ubiquinol cytochrome *c* reductase subunits. This can sometimes have unintended effects. Try typing “actin” in the “Description Text Pattern” box and clicking on the “Submit intron table query” button. Most of the results returned will be of interest; however, note the presence of *YIP2* on the list. Why is this gene here? It is not similar to actin or known to be involved in actin function. The computer found the sequence “actin” in the text description “Ypt interacting protein.” This shows that humans will always remain important in evaluating results from computational processes, and that each result should be considered innocent until proven guilty.

An intron feature commonly of interest is the sequence of the splice sites and branchpoint. The set of yeast introns has highly conserved splicing signals⁷; thus deviations are of interest. Go back to the “Search” page (Fig. 4) and then click on the “Intron Splice Signals Query” link. The form page for the Splice Site Query is shown in Fig. 5. To find all introns with the nonstandard branchpoint GACUAAC (rather than the most common UACUAAC sequence) type “gacua” into the “Bpre” box on the form. The form will appear with “a” in the Branch box, since all yeast introns are thought to use A as the branched nucleotide. Leaving the other boxes open is the same as asking for anything at those positions. Click on the “Submit intron query” button. As shown in Fig. 5, the search identifies 10 introns that contain GACUAAC as the best match to the consensus. Each intron is listed with its intron signals and a link to its report page. Having the other splicing signals displayed shows, for example, that *YGL251C* and *YLR211C* also have unusual 5′ splice sites, and that *YDL115C* has an unusual 3′ splice site. The search can be extended to reveal splice site and branchpoint context as well. For example, to search for all introns containing four U residues immediately upstream of a GACUAAC branchpoint, one would type “uuuugacua” into the “Bpre” box and submit the form. The five introns that fit this description would then be returned. The other boxes work in the same way, allowing the intronic context of splice signals to be explored.

We became interested in what the intron collection might tell us about 3′ splice site selection in yeast, and we developed a search that allows identification of the number of potential 3′ splice site-like sequences that might exist between the

⁷ M. Spingola *et al.*, *RNA* 5, 221 (1999).

Ares lab Yeast Intron Database
(Version 2.0)
This site is supported by the National Cancer Institute and the W. M. Keck Foundation.

Yeast Spliceosomal Introns

- Welcome
- Intron Table
- Search
- Extras
- Links
- Version 1.0
- Update History

Intron Splice Site Query

Search our current intron data for various splice site sequence signals.

Start: Bpre gacua Branch Bpost End

And/or Yname:

Instructions and Example

Bpre is the sequence BEFORE the branch point base. Branch is the branch point base. Empty boxes are treated as "don't care". The period (.) character matches any base. You can't use T, you have to use U. An example is shown below

The request "gacua" in the "Bpre" box will find all introns with the branchpoint sequence GACUAAC.

The report looks like this:

Current total introns: 244			
GUAUGU	GCUUUGACUAACACAU	UACAG	[Report] [Table] SNR17A SNR17A 15 730119 786275 INTRON
GUAUGU	GGUUUGACUAACACAU	AACAG	[Report] [Table] SNR17B SNR17B 15 281502 281374 INTRON
GUAUGU	ACUUUGACUAACAGA	UUUAG	[Report] [Table] KP99B YBR1897 2 604467 605467 INTRON 8 420
GUAUGU	ACAGUAGACUACCUUU	AAUAG	[Report] [Table] YB2417C 2 663266 662456 INTRON 302 722
GUAUGU	UCUUUGACUAACGUAU	CCGAG	[Report] [Table] YB2310C YB2310C 2 680482 679538 INTRON 12 109
GUAUGU	GCUUUGACUAACUAG	AAAAG	[Report] [Table] YD1115C 1 255049 254958 INTRON
GUAUGU	AAUAUAGACUAACUUUU	UAUAG	[Report] [Table] YRA1 YRP381W 4 1236547 1227992 INTRON 286 1051
GUAUGU	UAUUUGACUAACAUUG	UAUAG	[Report] [Table] YG1251C 2 31635 27521 INTRON 59 219
GUAUGU	UCAGUAGACUAACGUUC	CUUAG	[Report] [Table] YL3092C 12 327416 326514 INTRON 17 157
GUAUGU	CAUGGUAACUAACAUCA	AUAG	[Report] [Table] YL2111C 12 564731 563752 INTRON 19 77

Number of hits: 10

Splice sites and branchpoints

Link to the report page

Gene or ORF name

FIG. 5. Splice site query.

branchpoint and the true 3' splice of each intron. Go back again to the "Search" page and click on the "YAG Query" link. This is the YAG Query form page and it is shown in Fig. 6. In the first box, type the nucleotide sequence pattern to search for between the branchpoint and the 3' splice site. In the example shown, "[uc]ag" (use square brackets, not braces or parentheses!), the computer will search for all pyrimidine-A-G sequences found between the branchpoint and the 3' splice site. If nothing is put into the "Report full sequences . . ." box, then only the counts of the number of introns in each class will be returned. To see the full report (bottom of Fig. 6), put "0" (the number zero) in that box. If you are only interested in introns that have one or more occurrence of the sequence, put "1." In our search here, we have also restricted the introns we want to look at to be the ones that end in the unusual AAG 3' splice site. This is specified by including the last three bases in the desired introns in the "Optional specific end of intron pattern" box. One can specify more or fewer bases using this box as well (e.g., "UUAAG"), but it does not accept wild-card characters or bracket requests.

The output of the search is shown in the bottom half of Fig. 6. The first line contains the number of occurrences of the pattern, the next line has links to the report for

Ares lab Yeast Intron Database (Version 2.0)
This site is supported by the National Cancer Institute and the W. M. Keck Foundation.

YAG Query

Search our current intron data for various 3' splice site sequence signals.

YAG pattern

Report full sequence by number of hits above this number

Optional specific end of intron pattern (no wildcards allowed)

Instructions

This looks at the branch to 3' end of the introns. It counts the number of the specified 3' splice site patterns that occur in each intron between the branch and end, and the count includes the annotated 3' splice site. If Report sequence... is set to a positive number, the sequence is printed if it has greater than or equal to the requested number of hits, if it is negative, it reports the sequences that have fewer than the requested number of hits. If it is blank, it just reports the counts, no sequences. (Remember you have to use "u" in the pattern, not "t").

Yeast Spliceosomal Introns

Welcome

Intron Table

Search

Extras

Links

Version 1.0

Update History

e-mail Manny

Acknowledgments

University of California, Santa Cruz

"[uc]ag" searches for all YAG sequences between the branchpoint and the 3' splice site

This requests that results of zero YAGs also be reported

This requests only introns that end in AAG.

The output looks like this:

```

>#0
> [Report] [Table] YAL001C 1 151163 147531 INTRON 71 160
> YAL001C 1 151163 147591 INTRON 71 160 TFC3,TSV115,FUN24,YAL001C,tsv115
CGACACAUGAAG
>#0
> [Report] [Table] YEL050F 2 125088 126032 INTRON 31 146
> YEL050F 2 125088 126092 INTRON 31 146 SEC17,(YEL0517),(YEL0505),YEL050F
CAUGAAGAAACGGGAUUGAUAUUGCCGUGUGUUAAG
>#1
> [Report] [Table] YDL115C 4 555043 554955 INTRON
> YDL115C 4 555043 554955 INTRON
CUAAGUUCACUCUGGACAAUUCAGUAUUGAUGAGGGGAAAAAG
>#0
> [Report] [Table] YDL189W 4 122079 123580 INTRON 1 92
> YDL189W 4 122079 123590 INTRON 1 99 YDL189W,D1260
CAACTAAGAUUUCAG
>#0
> [Report] [Table] YGL226C-A 7 73156 72747 INTRON 22 170
> YGL226C-A 7 73156 72747 INTRON 48 137 UBC5,D4234,YD9609.13,YDR059C
CAUGAUGUUCUUUUUGAACUUUUUUCGAAAG
>#1
> [Report] [Table] YGL226C-A 7 73156 72747 INTRON 22 170
> YGL226C-A 7 73156 72747 INTRON 22 170 YGL226C-A
CAAGGUAAGCAAUUCGCCUUCGCGCAUUGUGGUAUUCACUUAAG
>#0
> [Report] [Table] YHR133W 13 536206 537603 INTRON 1243 1367
> YHR133W 13 536206 537608 INTRON 1243 1367 REC114,YH9375.02,YHR133W
CUAACACUUUUGCCACUACAAUAUGAAUAAAAAGUUAAUUAAG
>#0
> [Report] [Table] YOL047C 15 242745 241612 INTRON 244 306
> YOL047C 15 242745 241612 INTRON 244 306 YOL047C,02001
CCGCCUUUGAAG
#counts of pattern '[u,c]ag' between branch base
#and the 3' end of the introns.
#number_of_occurrences_of_pattern Count_Of_Sequences
0 6
1 2

```

Number of YAGs found

Intron name and links

Sequence between the branchpoint and the 3' splice site

FIG. 6. Using the YAG query.

that intron, including the gene name(s), and the last line has the sequence between the branchpoint and the 3' splice site. The search shows that two genes, *YDL115C* and *YGL226C-A*, have introns that skip a YAG sequence in favor of an AAG sequence. In the case of *YGL226C-A*, the UAG is only 9 nt from the branchpoint, a distance less than that found in the intron with the shortest such distance (*MATa1*

second intron, 10 nt). Care must be used in evaluating the sequence by eye, since in about 16 introns there is a CAG immediately after the branched residue, and these are unlikely to represent 3' splice sites. (Can you find these? *Hint*: Go back to the Intron Splice Signals Query page and type "cag" into the "Bpost" box and submit the form.) The search for sequences on the YAG Query page is not limited to 3' splice site-like sequences. Actually any sequence can be requested using the "YAG pattern" box on the form. We have also used this search page to evaluate the existence of pyrimidine tracts and other sequences in the branchpoint-3' splice site interval of yeast introns.

The final search is the Intron Sequence Search. Go back once more to the Search page and click on the "Sequence search on introns" link. This page (Fig. 7)

Ares lab Yeast Intron Database (Version 2.0)
This site is supported by the National Cancer Institute and the W. M. Keck Foundation.

Yeast Spliceosomal Introns

- Welcome
- Intron Table
- Search
- Extras
- Links
- Version 1.0
- Update History
- e-mail Manny
- Acknowledgments

Intron Sequence Search

Search our current intron data using Perl regular expressions

Regular Expression Pattern:

Examples

A "regular expression" is a description of a pattern that can include special characters for wildcard matches. The most basic "regular expression" is just an exact set of letters to look for. UACUAAC

(Note that you must use U, not T) The period (.) character stands for "any base", so: G..G would look for all G's separated by two bases. The pattern [GU]AAA would look for 3 A's preceded by either a G or U. Many more complex expressions are supported, please see a Perl manual or the on-line regular expression section for more information.

"ugua[uc]gu" searches for the consensus pseudo-5' splice site enhancer.

YBL072C
YBL072C 2 89435_89128_INTRON
Occurrences: 1
Position : 260 UGUACGU

YER131W
YER131W 5_423588_423948_INTRON
Occurrences: 1
Position : 346 UGUAUGU

YGR214W
YGR214W 7_920569_921782_INTRON_91_545
Occurrences: 1
Position : 403 UGUACGU

The report includes a link to the intron's page and information about the number and location of each match.

FIG. 7. Intron sequence search.

allows searches for any sequence anywhere in the intron. Type the sequence you are interested in finding into the box labeled "Regular Expression Pattern." You must use U instead of T! (Is our RNA bias showing?) This form accepts wild-card characters and bracket requests (square only!) within a string of nucleotides, such as "ugua[uc]gu." This particular request searches for either of two sequences: UGUAUGU or UGUACGU. Inserting a period "." means "any base" and is formally the same as "[agcu]." To get any pyrimidine at a particular position, as above, type "[uc]"; to get any purine, type "[ag]"; or for only G or C, use "[gc]." The "caret" character "^" can be used to specify "NOT" inside the brackets, so that if a nucleotide is excluded from a position this would be indicated by "[^c]", meaning "any nucleotide except C," or "A or U or G." To exclude two nucleotides, type "[^ag]." Note that "[^ag]" means the same to the computer as "[cu]" (you have just learned a little of the computer language Perl).

In the example in Fig. 7, we searched for the consensus sequence found in an *in vitro* evolution experiment that identified a pseudo-5' splice site as an enhancer of splicing efficiency in yeast.⁸ We wanted to see if such sequences might be common in yeast introns. To begin the search, we used the strictest definition of the consensus, which is essentially the same as the 5' splice site consensus sequence except that it has a U upstream of the first G: UGUAYGU, which we express for the search as "ugua[uc]gu." Clicking on the "Submit pattern search" button, the computer returns a page containing information on three introns. Note that the natural 5' splice site is excluded because the intron sequence data in the database includes no exon nucleotides, and thus none has a U upstream, at least in the data being searched. Many more introns are returned if the first U is deleted from the request. The results provide a link to the report page, the intron coordinates, the number of occurrences of the sequence, and the position of the sequence found. These three introns might be good candidates in which to test the idea that pseudo-5' splice site sequences contribute to splicing in natural yeast introns. Additional searches with more relaxed consensus sequences may reveal additional candidates.

Database Errors, Programming Bugs, and Interpretational Caveats

The Ares lab Intron Database is a work in progress. We originally devised it for our own use, but found it so useful that we thought others might want to use it as well. We try to maintain the accuracy of the data, but there is a lot of it underneath. (One can think of the distinction between a searchable database and a publication as similar to the distinction between performance art and painting. Each time you access the database, we are putting on a show, which is a different responsibility than one has after painting a picture.) Although the data are useful for gaining broad-brush impressions of the intron family in yeast, and the choice nugget will occasionally be found, all specific results of importance should be

⁸ D. Libri, A. Lescure, and M. Rosbash, *RNA* **6**, 352 (2000).

confirmed by comparing the findings with information in other databases. SGD, PubMed, and GenBank employ individuals who are responsible for constantly updating information of relevance to the data in our collection. Database errors exist in all databases, so beware! Also note that future versions of the database may return results slightly different from those presented in the figures, because of updated information.

There are a few programming bugs yet in the system as well. The Intron Sequence Search report gives spurious position information if there is more than one occurrence of a particular sequence (although the first listed position is correct, and the other occurrences are present) and we are working to fix that. One limitation we cannot really surmount involves the identification of the branchpoint. We specify this position based on looking at the sequence and in many cases it is a guess. Molecular analysis of branchpoints is challenging, and there are few hard data for most introns. Finally one must be careful in interpretation of the findings. An example is the abundance of U tracts in introns. These are easy to find, and some introns have many, many short runs of U. The results of such a search may be impressive but the software assesses no significance to these results. The investigator must ask, How big is this intron? Given the G+C content of the yeast genome, how often might I expect to observe U runs of this length in a sequence of this size? What is the probability that my observation could be due to chance? In the future when all possible experiments have been done and all data is archived in searchable structures, all we will need do to test a hypothesis is to submit a form. Until then we can at least use our current databases to sharpen a few of our experimental rationales and hone some of our conclusions.

Acknowledgments

We thank Chuck Sugnet, Tyson Clark, Carrie Davis, and Marc Spingola for help keeping the database tidy. We also thank Haller Igel for comments on the manuscript. This work was supported by the W. M. Keck Foundation grant to the RNA Center at Santa Cruz, and by grants to M.A. from the National Institutes of Health (CA 77813 and GM 40478).