**A primer for practical phylogenetic data gathering. Uconn EEB3899-007. Spring 2015**

Session 5

# Uploading sequences to GenBank

**Rafael Medina (rafael.medina.bry@gmail.com)**
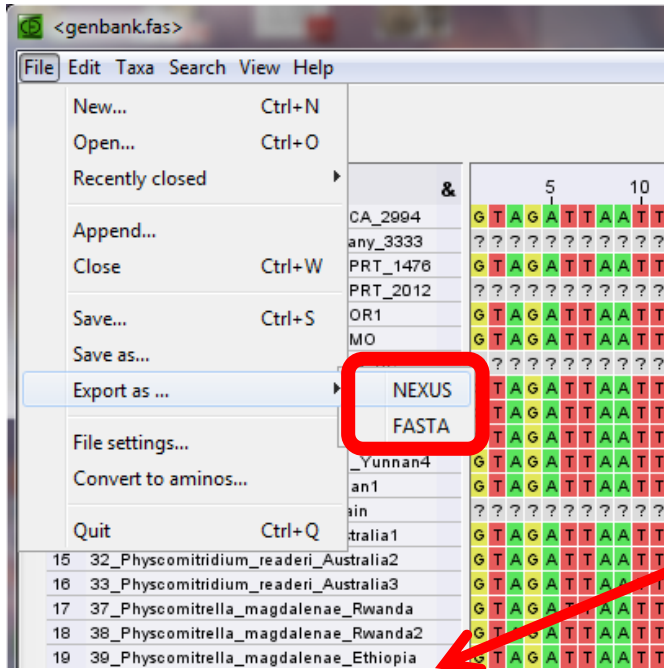**Yang Liu (yang.liu@uconn.edu)**

# A road map to phylogenetic data gathering

# Preparation of the NEXUS file

- Use the full alignment (without exclusions) for Genbank submission
- Export as a NEXUS file

Create a text file with voucher specifications (this will save you time later)

```
11   20_Physcomitrella_patens_China_Yunnan4>[org=Physcomitrella patens][authority=(
12   21_Physcomitridium_readeri_Japan1>[org=Physcomitridium readeri][authority=(Mül
13   30_Physcomitridium_readeri_Spain>[org=Physcomitridium readeri][authority=(Mül
14   31_Physcomitridium_readeri_Australia1>[org=Physcomitridium readeri][authority=
15   32_Physcomitridium_readeri_Australia2>[org=Physcomitridium readeri][authority=
16   33_Physcomitridium_readeri_Australia3>[org=Physcomitridium readeri][authority=
17   37_Physcomitrella_magdalenae_Rwanda>[org=Physcomitrella magdalenae][authority=
18   38_Physcomitrella_magdalenae_Rwanda2>[org=Physcomitrella magdalenae][authority
19   39_Physcomitrella_magdalenae_Ethiopia>[org=Physcomitrella magdalenae][authorit
```

Use the same sequence name and order as in the alignment

Use ">" as separator

List of attributes:
[org=] (organism)
[authority=]
[molecule=]
[location=]
[specimen-voucher=]
[note=]

Once it is ready, erase everything before the ">". Do not change the order afterwards

# Preparation of the NEXUS file

Add a NCBI block in the nexus file

It should start with "BEGIN NCBI" and should be closed with "END;"

The voucher specifications are added as a matrix with the label "SEQUIN"

End the matrix with ";"

"DATA" block, generated by the editor

Beginning of the block

"Sequin" label

Specifications, notes etc (pasted as in the previous slide)

This ends the sequin matrix

End of the NCBI block
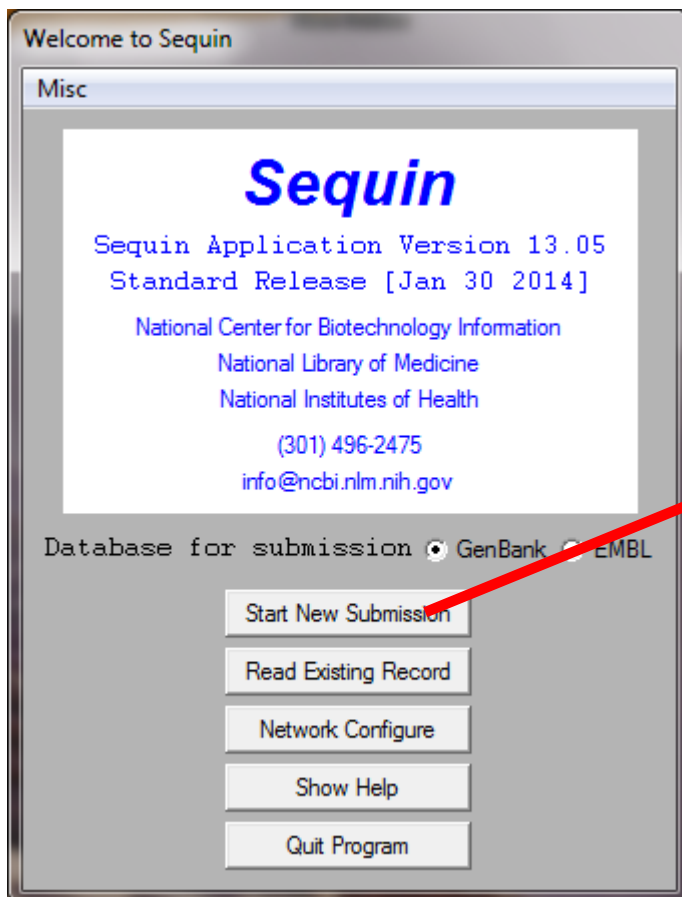
NCBI block

```
1   #NEXUS
2
3   BEGIN DATA;
4   DIMENSIONS NTAX=19 NCHAR=457;
5   FORMAT DATATYPE=DNA GAP=- MISSING=? ;
6
7   MATRIX
8   05_Physcomitrella_patens_USA_CA_2994      GTAGATTAATTTTCCAATACGC
9   06_Physcomitrella_patens_Germany_3333     ?????????????????ACGC
10  08_Physcomitrella_patens_UK_APRT_1476     GTAGATTAATTTTCCAATACGC
11  11_Physcomitrella_patens_UK_APRT_2012     ?????????????????ACGC
12  12_Physcomitrella_patens_USA_OR1          GTAGATTAATTTTCCAATACGC
13  14_Physcomitrella_patens_USA_MO           GTAGATTAATTTTCCAATACGC
14  16_Physcomitrella_patens_Canada_BC        ?????????????????ACGC
15  17_Physcomitrella_patens_China_Yunnan1    GTAGATTAATTTTCCAATACGC
16  18_Physcomitrella_patens_China_Yunnan2    GTAGATTAATTTTCCAATACGC
17  19_Physcomitrella_patens_China_Yunnan3    GTAGATTAATTTTCCAATACGC
18  20_Physcomitrella_patens_China_Yunnan4    GTAGATTAATTTTCCAATACGC
19  21_Physcomitridium_readeri_Japan1         GTAGATTAATTTTCCAATACGC
20  30_Physcomitridium_readeri_Spain          ?????????????????ACGC
21  31_Physcomitridium_readeri_Australia1     GTAGATTAATTTTCCAATACGC
22  32_Physcomitridium_readeri_Australia2     GTAGATTAATTTTCCAATACGC
23  33_Physcomitridium_readeri_Australia3     GTAGATTAATTTTCCAATACGC
24  37_Physcomitrella_magdalenae_Rwanda       GTAGATTAATTTTCCAATACGC
25  38_Physcomitrella_magdalenae_Rwanda2      GTAGATTAATTTTCCAATACGC
26  39_Physcomitrella_magdalenae_Ethiopia     GTAGATTAATTTTCCAATACGC
27  ;
28  END;

    BEGIN NCBI;
    SEQUIN
32  >[org=Physcomitrella patens][authority=(Hedw.) Bruch & Schimp.][
33  >[org=Physcomitrella patens][authority=(Hedw.) Bruch & Schimp.][
34  >[org=Physcomitrella patens][authority=(Hedw.) Bruch & Schimp.][
35  >[org=Physcomitrella patens][authority=(Hedw.) Bruch & Schimp.][
36  >[org=Physcomitrella patens][authority=(Hedw.) Bruch & Schimp.][
37  >[org=Physcomitrella patens][authority=(Hedw.) Bruch & Schimp.][
38  >[org=Physcomitrella patens][authority=(Hedw.) Bruch & Schimp.][
39  >[org=Physcomitrella patens][authority=(Hedw.) Bruch & Schimp.][
40  >[org=Physcomitrella patens][authority=(Hedw.) Bruch & Schimp.][
41  >[org=Physcomitrella patens][authority=(Hedw.) Bruch & Schimp.][
42  >[org=Physcomitrella patens][authority=(Hedw.) Bruch & Schimp.][
43  >[org=Physcomitridium readeri][authority=(Müll. Hal.) G. Roth][m
44  >[org=Physcomitridium readeri][authority=(Müll. Hal.) G. Roth][m
45  >[org=Physcomitridium readeri][authority=(Müll. Hal.) G. Roth][m
46  >[org=Physcomitridium readeri][authority=(Müll. Hal.) G. Roth][m
47  >[org=Physcomitridium readeri][authority=(Müll. Hal.) G. Roth][m
48  >[org=Physcomitrella magdalenae][authority=J.L. De Sloover][mole
49  >[org=Physcomitrella magdalenae][authority=J.L. De Sloover][mole
    >[org=Physcomitrella magdalenae][authority=J.L. De Sloover][mole
    ;
52  END;
53
```
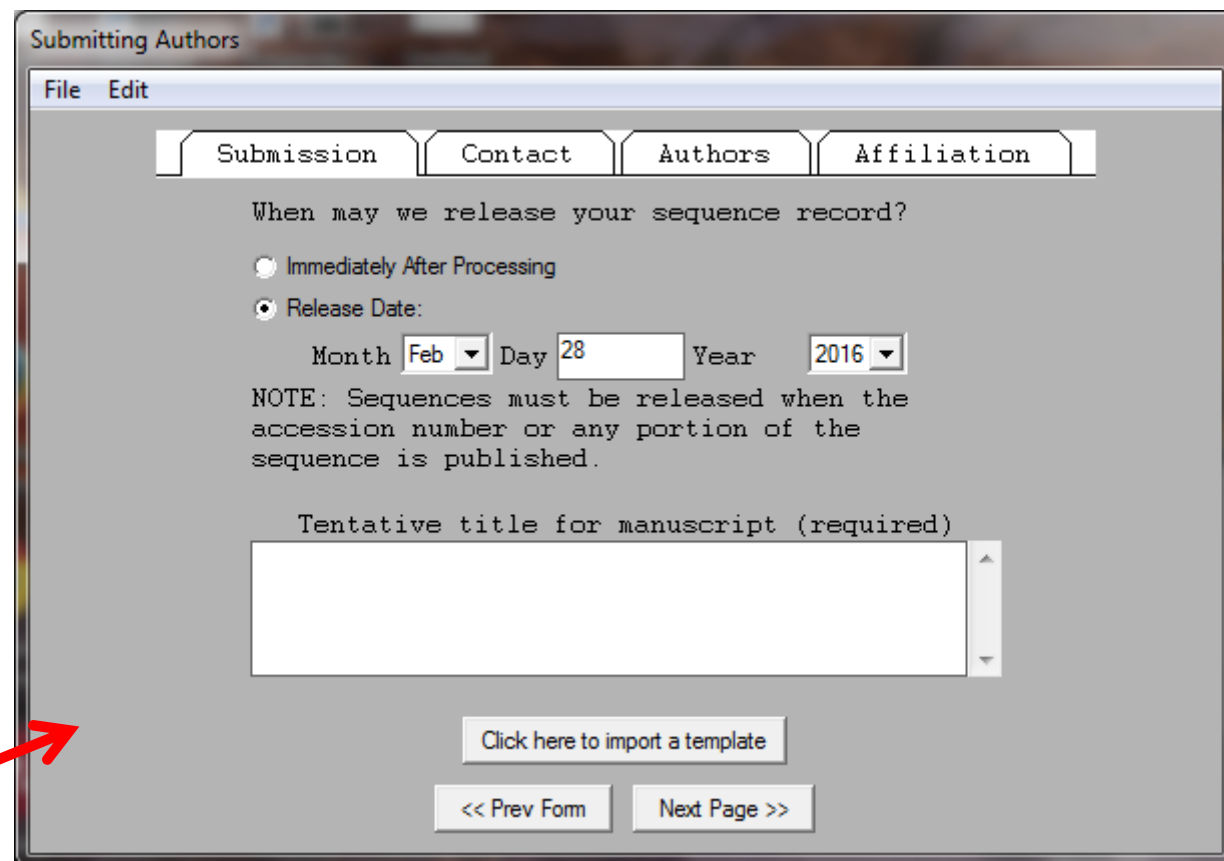
# SEQUIN

Sequin is the NCBI software developed to upload sequences to GenBank

Download it from **here**



New submission

Form 1: authors
Fill the details of the publication, desired release date, authors, affiliation, etc.

Save this information in a template before proceeding to the next form, it will save you time if the program crashes (surprisingly likely scenario)

# SEQUIN

Form 2: alignment
In this example we are submitting an alignment through the "normal submission dialog", it is a phylogenetic study and we will upload a nexus alignment

Import the alignment. If you have done everything correctly, the sequences AND the attributes of the nexus file will be incorporated



You will be then required to input how the sequences were obtained

Specify the type of molecule and topology (lineal DNA in this case)

# SEQUIN

## Form 3: general characteristics



**Organism**
The organism names, details, genetic code, etc, should have been incorporated automatically by now. Check it before proceeding

**Proteins**
Unless you are submitting a single protein gene, skip this tab

**Annotation**
Unless your sequence includes only one feature (region, gene), just tick "none"
Then...



Define the title of your sequences. Include a detailed description of the region and choose to prefix the title with the organism name

# SEQUIN

Form 4: annotation

Pick a sample with a good, long sequence to use as a template for annotations



In this area you can see how the sequence will look in Genbank. So far all the relevant information is included except the annotations, the "indices" that will tell the reader where each region starts and ends

At this stage, go back to your lab book and summarize the limits of each feature of the sequence that you want to annotate
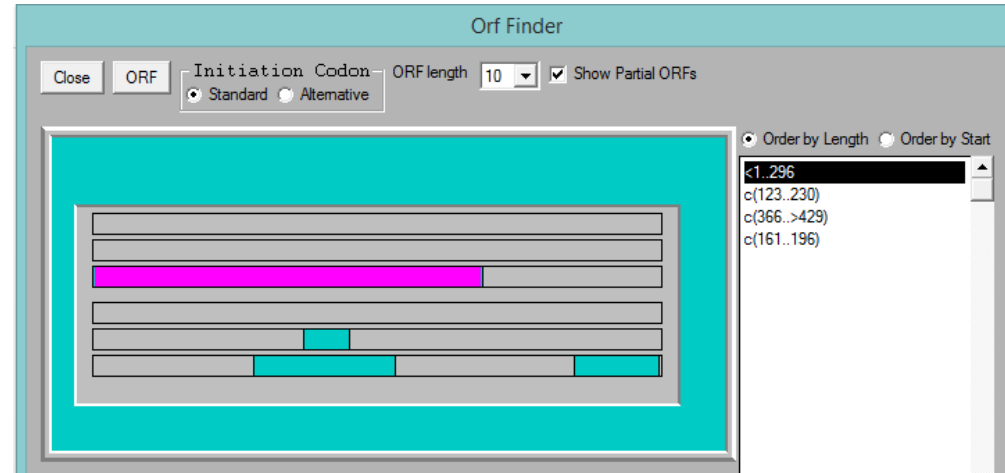
# SEQUIN

Example 1: annotation of a CDS (Coding DNA Sequence)

The first part of this sequence is the partial psbA protein. The stop codon (TAA) ends in position 296



Use the ORF (Open Reading Frame) finder to pick this CDS



You can easily explore the six possible frames and select the one that corresponds to the actual protein



Double click the CDS and complete the relevant information, at least the name and abbreviation of the protein, and a comment clarifying that it is partial



5' partial

Use better alignment coordinates if you have indels

Finally, check the location details and accept

# SEQUIN



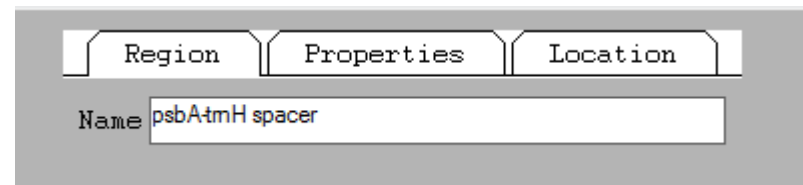After the annotation is done, it should appear correctly in the visualization window
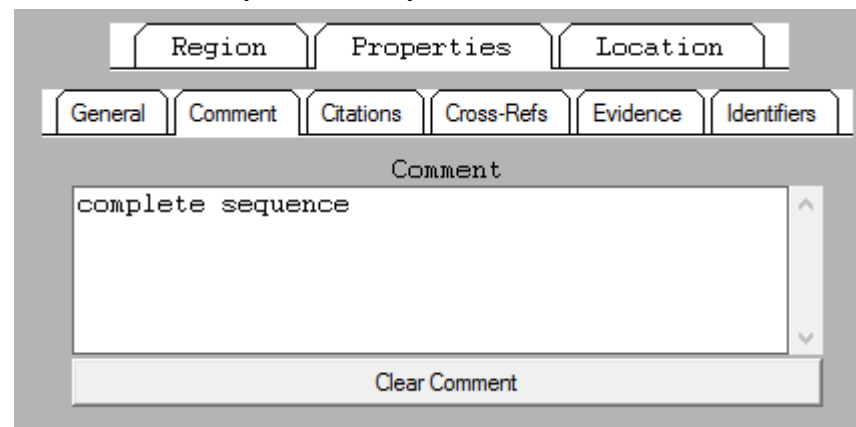
# SEQUIN

Example 2: annotation of a spacer

Assuming you want to annotate a
well-known, non-coding spacer
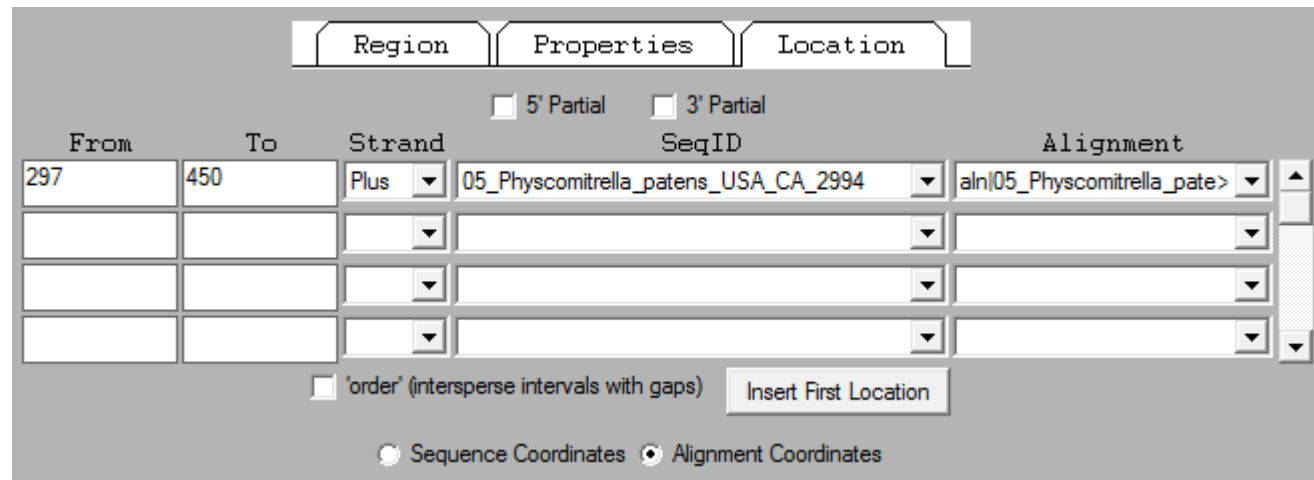(as psbA/trnH), you can consider it
a "named region"

Name the region



Comment if it is complete or partial



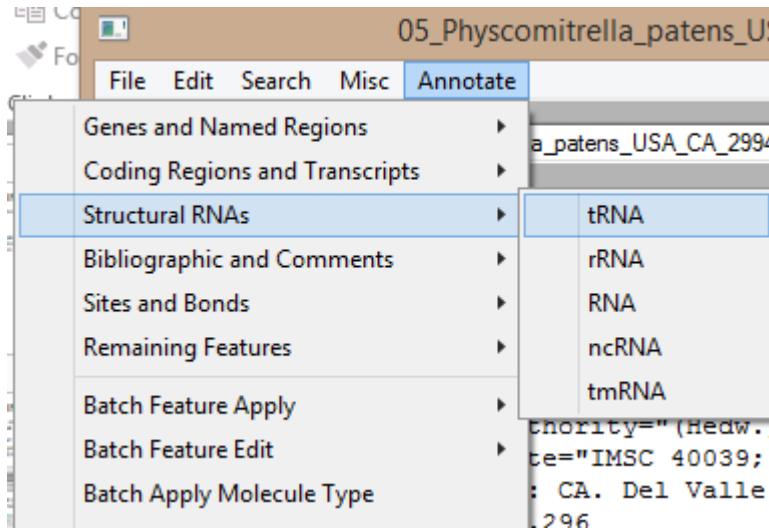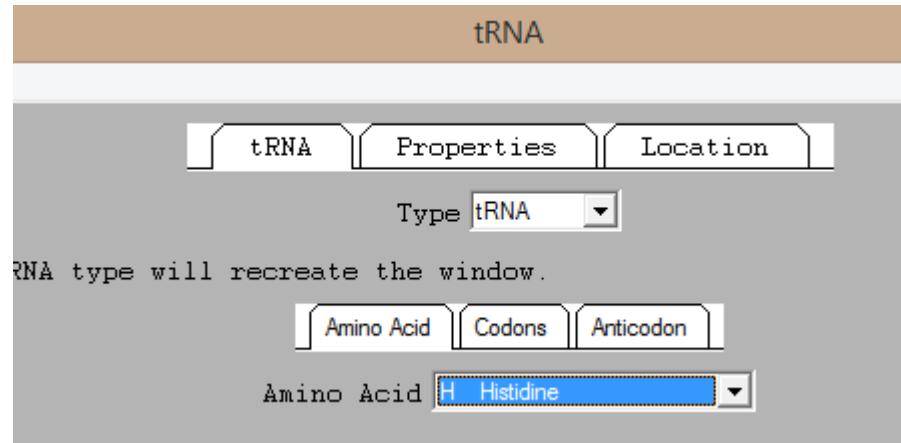Specify the interval in the alignment

# SEQUIN

Example 3: annotation of a tRNA

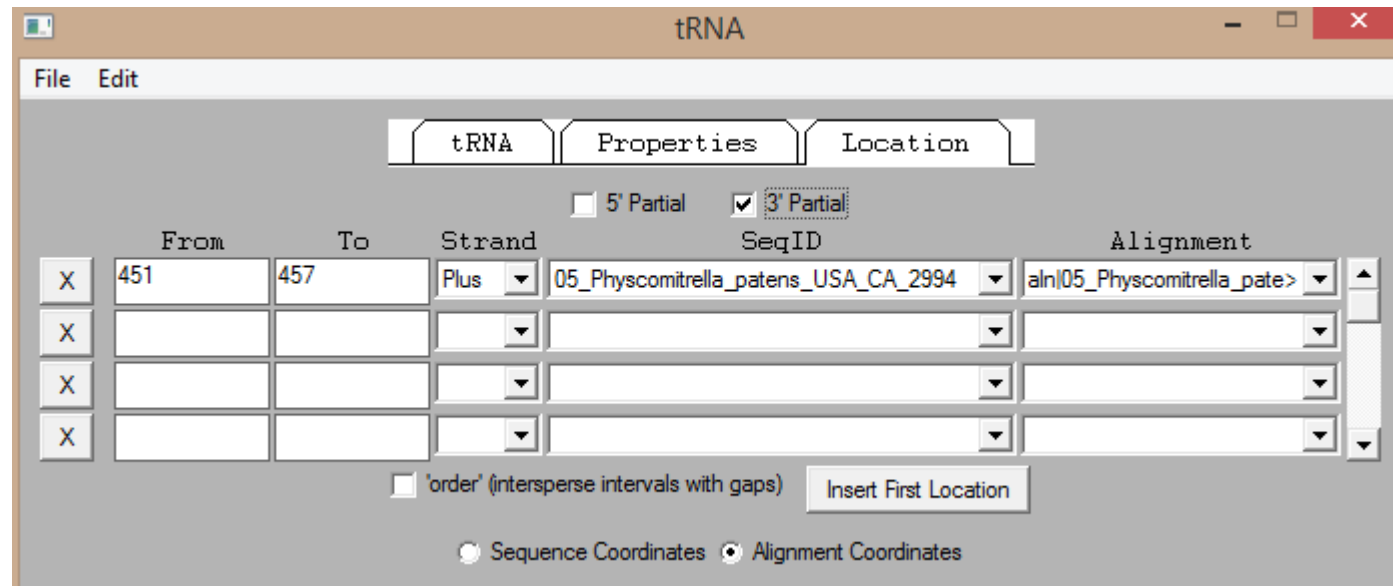The last feature of our example is
a fragment of the tRNA-His gene

tRNAs are considered in Sequin
"structural RNAs"

Specify the aminoacid and, if possible, the codon



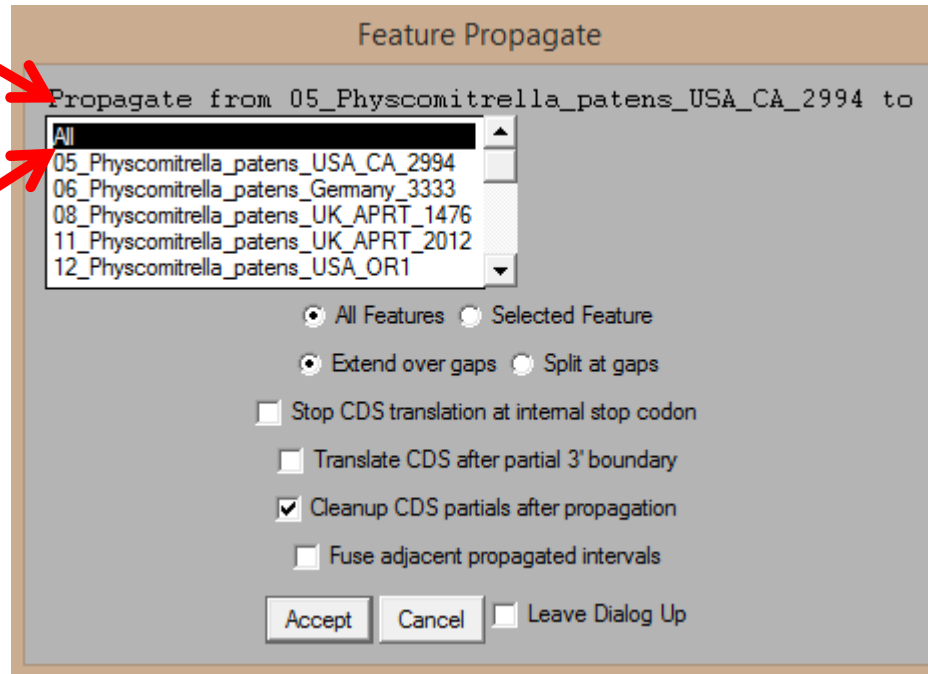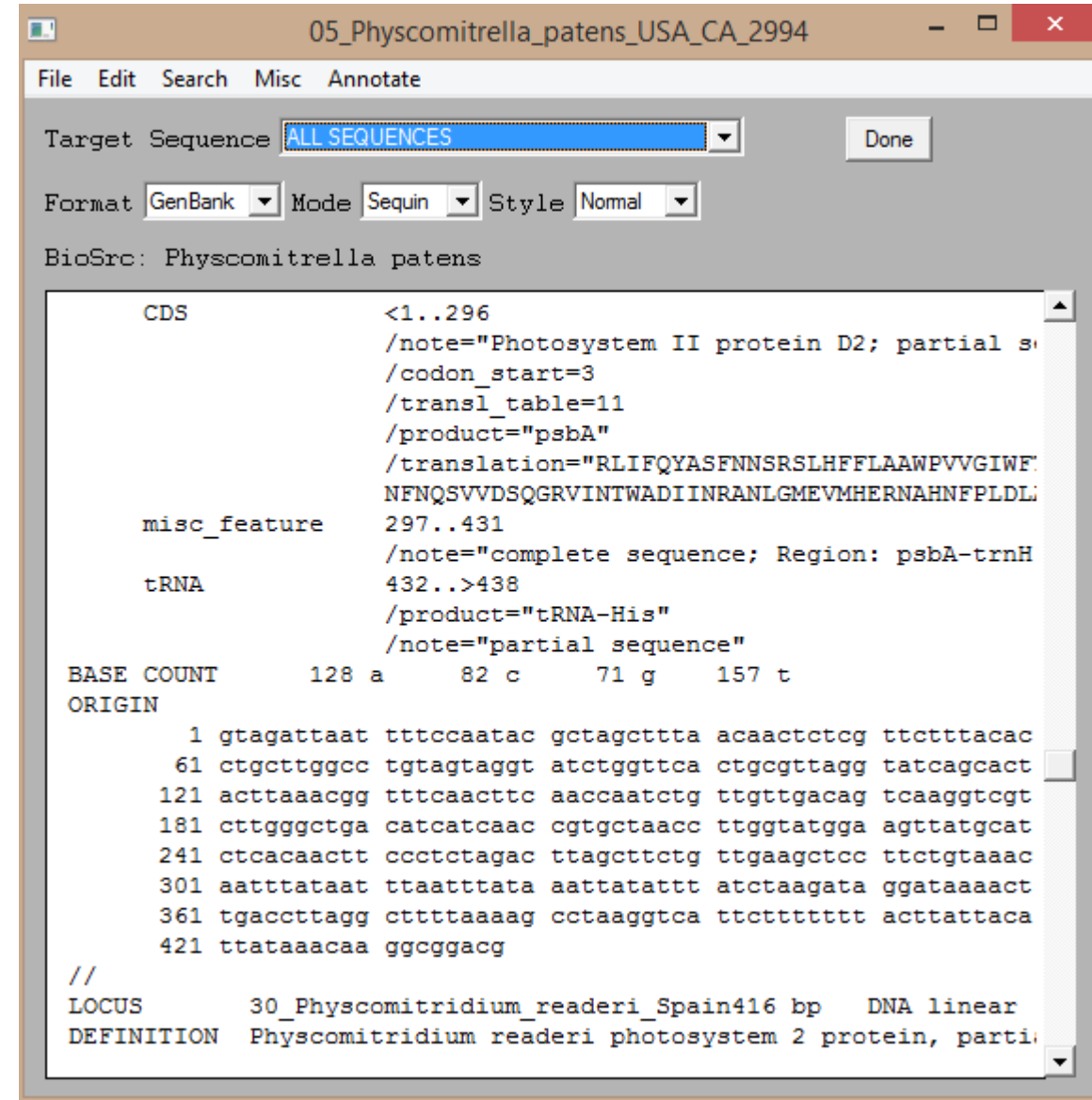Then, business as usual: specify limits, completeness, etc

# SEQUIN

Once you have your chosen
sequence completely annotated,
propagate those annotations
across the whole alignment
(Edit > Feature Propagate)

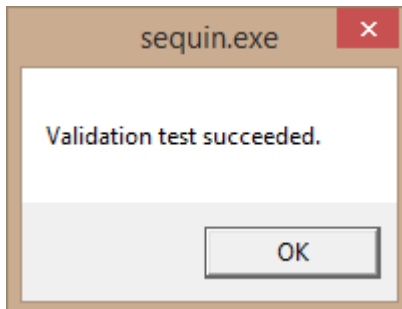Your annotated
sequence

To all the
alignment

Then, check that the other sequences seem to have
acquired the annotations correctly

# SEQUIN

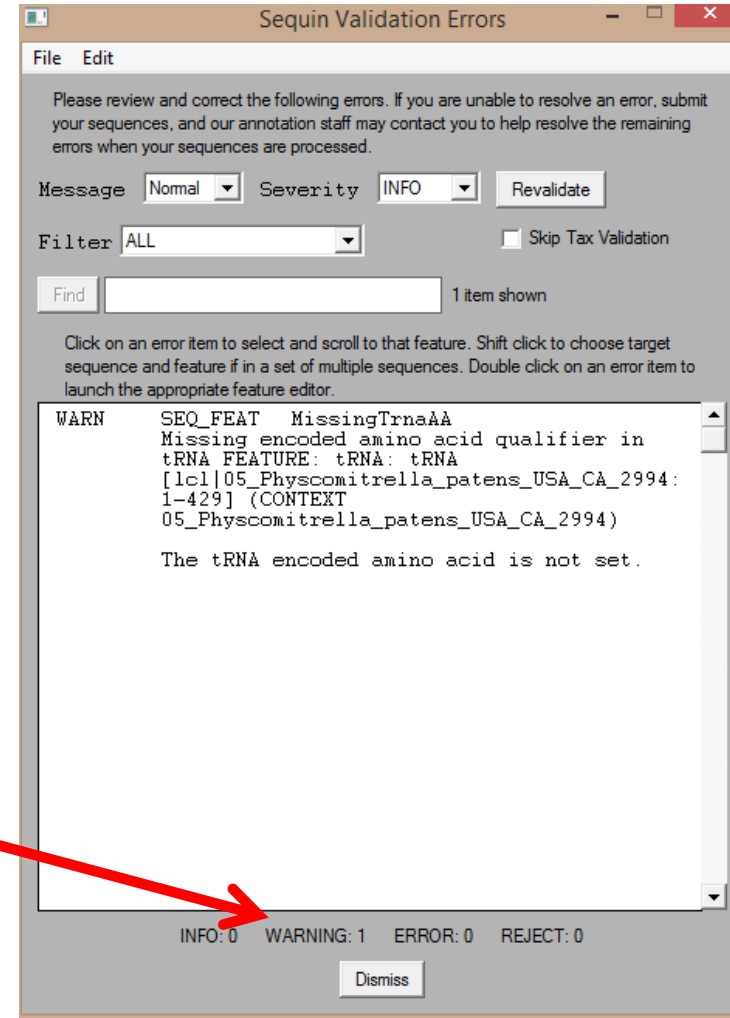You are almost ready for submission. Cross your fingers and validate the file (Search>Validate)

If everything is ok, you will see something like this:

Alternatively you will receive some error messages with different degrees of "severity"

Then you just need to export the sequin file: choose "Prepare submission" in the File menu and save your record. Send it by email to

**gb-sub@ncbi.nlm.nih.gov**

Even the warnings or lesser errors can give you trouble after your submission. It is better to fix them all and revalidate the submission. Often this means to start the whole process again, so be patient and meticulous!