

- ▷ Adjust argparse such that boolean input no longer needs the ~~key~~ value "True".

"--taxonomycheck True"  $\Rightarrow$  "--taxonomycheck"

---

- ▷ Implement ~~ref~~ integration of References.
- 

### ■ Make modifc.:

For each qualifier you have, make a column

---

5' CDS location

needs a "plus 1".

5' CDS incorrect

5' CDS correct

667  $\Rightarrow$  668

\* AT LEAST FOR  
CHECKLISTS \*

---

2, Remove "translatable" from source features;  
if anything, it should be in a CDS feature.

---

3, Make sure that fourth field of  
title line ~~is equal to content~~ has same  
content as as source feature "mol-type".

---

Falscher Output bei Col. "INTRON"  
Checkliste for math  
in Pyrus:

- ~~intr~~ intron has erroneous name  
← id = "K" (instead of id = "intr")  
nicht erkannt!
- Before and After Positions are specified ~~where~~ where there should be Exact Positions  
↳ Falscher Output bei Col. "5'-PARTIAL"  
"3'-PARTIAL"

Why is there  
a join (1..41,  
43..530)?

It appears that the  
misc-feature length  
calculation is wrong.

6.5.

6.5.1. ! (Location  
object!)

6.5.3.

' make decision of  
location-change  
based on  
seq-len-comparison  
(Before/after  
removal of  
ambiguities)

If leading and trailing ~~space~~ chars indeed removed,  
feature location must be adjusted

1..501  $\Rightarrow$  >1..501

! [6.] Remove data that have only missing data.

---

7. Set up list of source qualifiers that can be accessed from anywhere in the software.

---

1. Remove all missing segs per alignment

---

1. Check ~~that~~ all taxon-names in .nex also in .CSV

---

ANNONEX2  
EMBL

→ Rename to: "nex2ena-flat" / "nex2ena-check"

---

5. ~~add~~ Specify "data class" (3.1. in ENH manual)  
(STD)

---

## TO DO - LONG TERM

1, Rename "Degapping Ops" with "Cleanup Ops"

- I changed Taxon-3 in TestData-1.uex by making the non-coding ends "?". This must be replaced in the corresponding figure.

### ~~Update the figures~~

- The ID line changed in all examples (TestData-1, TestData-2). Check if this affects the figures based on these datasets in any way.

In Manuscript:

- Write that charset names must have a special format:  

charset <del>Name</del>	<u>underscore</u>	charset <del>Type</del>	<del>underscore</del>	format:
↓		must be part of SIC (RNA, gene, intron, ...)		
BLISTED FOR charset-dimming				

Submission mode  
Cmdl. Param.: -u  
Input Type: logical

Topology  
~~Linear or Circular~~  
[Linear or Circular]

Taxonomic Division  
[according to ENA manual 3.2.]

INPUT

DNA Alignment and Charsets  
Cmdl. Param.: -n  
Input Type: file path  
Data Format: NEXUS (.nex)

Taxon Qualifiers  
Cmdl. Param.: -c  
Input Type: file path  
Data Format: comma-delimited table (.csv)

Email Address  
Cmdl. Param.: -e  
Input Type: string

Parse .nex file  
Employs: Bio.Nexus.Nexus

Extract Alignment      Extract Charsets

Parse .csv file  
Employs: csv.DictReader

for each qualifier sep.

Quality Check of Qualifier  
Baz: Qux  
(e.g. remove empty qualifiers)

Epub search to NCBI taxonomy database

for each sequence  
Confirm that taxon name valid

for each charset sep.  
Epub search to NCBI gene database  
Get charset product  
~~REPLACE~~ Get taxon name, obtain charset sym and charset product

Generate seqRecord  
~~simultaneously populate with qualifiers~~

Degap sequence  
Foo: Bar  
Baz: Qux

Generate feature table  
Populate with feature 'source'

Continue to build feature table  
Add features sorted by their start position

Translate coding features  
simultaneous check quality of same features

Write to outfile in EMBL format  
Employs: Bio.SeqIO

OUTPUT

Write to temporary handle

Replace all ~~ambiguous~~ "?" with "N"  
Remove leading and trailing ~~ambiguities~~ (i.e. "Ns").

mark ID and AC lines with "xxx" if submission-mode == True

⑥ If taxon name not found, add additional qualifier. (i.e., ecotype).