# Biopython

## What is biopython?

Biopython is a collection of python modules that contain code for manipulating biological data. Many handle sequence data and common analysis and processing of the data including reading and writing all common file formats. Biopython will also run blast for you and parse the output into objects inside your script. This requires just a few lines of code.

## Installing Biopython

This is very straightforward once you have anaconda or minconda installed. I use miniconda because it's smaller. We are going to use `sudo` because this will give us permission to install in the 'correct' directory python is expecting to find the modules. Other users will be able to use it too. Using `sudo` can cause problems, but it's ok here. You will need the administrator password for the machine. If you don't have this, ask the person who does administration on your machine.

```
1   % sudo conda install biopython
2   WARNING: Improper use of the sudo command could lead to data loss
3   or the deletion of important system files. Please double-check your
4   typing when using sudo. Type "man sudo" for more information.
5
6   To proceed, enter your password, or type Ctrl-C to abort.
7
8   Password:
9   Fetching package metadata ...........
10  Solving package specifications: .
11
12  Package plan for installation in environment /anaconda3:
13
14  The following NEW packages will be INSTALLED:
15
16      biopython: 1.69-np113py36_0
17
18  The following packages will be UPDATED:
19
20      conda:     4.3.29-py36hbf39572_0 anaconda --> 4.3.30-py36h173c244_0
21
22  The following packages will be SUPERSEDED by a higher-priority channel:
23
24      conda-env: 2.6.0-h36134e3_0       anaconda --> 2.6.0-h36134e3_0
25
26  Proceed ([y]/n)?
27
28  conda-env-2.6. 100% |#################################################| Time:
    0:00:00   4.02 MB/s
29  biopython-1.69 100% |#################################################| Time:
    0:00:00  21.02 MB/s
30  conda-4.3.30-p 100% |#################################################| Time:
    0:00:00  43.61 MB/s
```

See if the install worked

```
1   python3
2   >>> import Bio
3   >>> print(Bio.__version__)
4   1.69
```

If we get no errors, biopython is installed correctly.

# Biopython documentation

Biopython wiki page

http://biopython.org/

Getting started

http://biopython.org/wiki/Category%3AWiki_Documentation

Biopython tutorial

ttp://biopython.org/DIST/docs/tutorial/Tutorial.html#chapter:Bio.SeqIO

Complete tree of Biopython Classes

http://biopython.org/DIST/docs/api/Bio-module.html

# Working with DNA and protein sequences

This is the core of biopython. And uses the Seq object. Seq is part of Bio. This is denoted Bio.Seq

```python
#!/usr/bin/env python3
import Bio.Seq                           # }
seqobj = Bio.Seq.Seq('ATGCGATCGAGC')     # } that's a lot of Seqs
# convert to string with str(seqobj)
seq_str = str(seqobj)
print('{:s} has {:d} nucleotides'.format( seq_str , len(seq_str)))
```

produces

```
ATGCGATCGAGC has 12 nucleotides
```

## From ... import ...

Another way to import modules is with `from ... import ...` . This saves typing the Class name every time. Bio.Seq is the class name. Bio is the superclass. Seq is a subclass inside Bio. It's written Bio.Seq. Seq has several different subclasses, of which one is called Seq. So we have Bio.Seq.Seq. To make the creation simpler, we call Seq() after we import with `from ... import ...` like this

```python
#!/usr/bin/env python3
from Bio.Seq import Seq
seqobj=Seq('ATGCGATCGAGC')
seq_str=str(seqobj)
protein = seqobj.translate()
prot_str = str(protein)
print('{:s} translates to {:s}'.format(seq_str,prot_str))
```

produces

```
1  ATGCGATCGAGC
```

## Bio.Alphabets

A Seq object likes to know what alphabet it uses A,C,G,T for DNAAlpabet etc. Not essential for most uses, but prevents you trying to translate a protein sequence!

### Specific Alphabets

```
1  >>> seqobj
2  Seq('ATG', Alphabet())
3  >>> from Bio.Alphabet import DNAAlphabet
4  >>> so=Seq('ATG',DNAAlphabet())
5  >>> so
6  Seq('ATG', DNAAlphabet())
7  >>> so.translate()
8  Seq('M', ExtendedIUPACProtein())
9
```

For proteins

```
1  from Bio.Alphabet import ProteinAlphabet
2  seqobj = Seq('MGT', ProteinAlphabet())
```

### Generic Alphabets

```
1   >>> from Bio.Seq import Seq
2   >>> from Bio.Alphabet import generic_dna, generic_protein
3   >>> my_seq = Seq("AGTACACTGGT")
4   >>> my_seq
5   Seq('AGTACACTGGT', Alphabet())
6   >>> my_dna = Seq("AGTACACTGGT", generic_dna)
7   >>> my_dna
8   Seq('AGTACACTGGT', DNAAlphabet())
9   >>> my_protein = Seq("AGTACACTGGT", generic_protein)
10  >>> my_protein
11  Seq('AGTACACTGGT', ProteinAlphabet())
```

# Extracting a subsequence

You can use a range [0:3] to get the first codon

`seqobj[0:3]`

# Read a FASTA file

We were learning how to read a fasta file line by line. SeqIO.parse() is the main method for reading from almost any file format. We'll need a fasta file. We can use Python_05.fasta which looks like this

```
1   >seq1
2   AAGAGCAGCTCGCGCTAATGTGATAGATGGCGGTAAAGTAAATGTCCTATGGGCCACCAATTATGGTGTATGAGTGAATCTC
    TGGTCCGAGATTCA
3   CTGAGTAACTGCTGTACACAGTAGTAACACGTGGAGATCCCATAAGCTTCACGTGTGGTCCAATAAAACACTCCGTTGGTCA
    AC
4   >seq2
5   GCCACAGAGCCTAGGACCCCAACCTAACCTAACCTAACCTAACCTACAGTTTGATCTTAACCATGAGGCTGAGAAGCGATGT
    CCTGACCGGCCTGT
6   CCTAACCGCCCTGACCTAACCGGCTTGACCTAACCGCCCTGACCTAACCAGGCTAACCTAACCAAACCGTGAAAAAAGGAAT
    CT
7   >seq3
8   ATGAAAGTTACATAAAGACTATTCGATGCATAAATAGTTCAGTTTTGAAAACTTACATTTTGTTAAAGTCAGGTACTTGTGT
    ATAATATCAACTAA
9   AT
10  >seq4
11  ATGCTAACCAAAGTTTCAGTTCGGACGTGTCGATGAGCGACGCTCAAAAAGGAAACAACATGCCAAATAGAAACGATCAATT
    CGGCGATGGAAATC
12  AGAACAACGATCAGTTTGGAAATCAAAATAGAAATAACGGGAACGATCAGTTTAATAACATGATGCAGAATAAAGGGAATAA
    TCAATTTAATCCAG
13  GTAATCAGAACAGAGGT
```

Get help on the parse() method with

```
1   >>> help(SeqIO.parse)
2
3   Help on function parse in module Bio.SeqIO:
4
5   parse(handle, format, alphabet=None)
6       Turns a sequence file into an iterator returning SeqRecords.
7
8           - handle   - handle to the file, or the filename as a string
9             (note older versions of Biopython only took a handle).
10          - format   - lower case string describing the file format.
11          - alphabet - optional Alphabet object, useful when the sequence type
12            cannot be automatically inferred from the file itself
13            (e.g. format="fasta" or "tab")
14
15      Typical usage, opening a file to read in, and looping over the record(s):
16
17      >>> from Bio import SeqIO
18      >>> filename = "Fasta/sweetpea.nu"
19      >>> for record in SeqIO.parse(filename, "fasta"):
20      ...     print("ID %s" % record.id)
21      ...     print("Sequence length %i" % len(record))
22      ...     print("Sequence alphabet %s" % record.seq.alphabet)
23      ID gi|3176602|gb|U78617.1|LOU78617
24      Sequence length 309
25      Sequence alphabet SingleLetterAlphabet()
26
27  ...
```

This uses the old `"..." % var` formatting syntax.

Here's a script to read fasta records and print out some information

```
1   #!/usr/bin/env python3
2   # assumes we are in the pfb2017 directory
3   from Bio import SeqIO
4   for seq_record in SeqIO.parse("./files/Python_05.fasta", "fasta"):   # give
    filename and format
5     print('ID',seq_record.id)
6     print('Sequence',str(seq_record.seq))
7     print('Length',len(seq_record))
8
```

## Convert fasta file to python dictionary in one line

There are three ways of doing this that use up more memory if you want more flexibility. `Bio.SeqIO.to_dict()` is the most flexible but also reads the entire fasta file into memory as a python dictionary so might take a lot of time and memory.

```
1  >>> id_dict = SeqIO.to_dict(SeqIO.parse('files/Python_05.fasta', 'fasta'))
2  >>> id_dict
3  {'seq1':
   SeqRecord(seq=Seq('AAGAGCAGCTCGCGCTAATGTGATAGATGGCGGTAAAGTAAATGTCCTATGGGC...AAC',
   SingleLetterAlphabet()), id='seq1', name='seq1', description='seq1', dbxrefs=[]),
   'seq2':
   SeqRecord(seq=Seq('GCCACAGAGCCTAGGACCCCAACCTAACCTAACCTAACCTAACCTACAGTTTGA...TCT',
   SingleLetterAlphabet()), id='seq2', name='seq2', description='seq2', dbxrefs=[]),
   'seq3':
   SeqRecord(seq=Seq('ATGAAAGTTACATAAAGACTATTCGATGCATAAATAGTTCAGTTTTGAAAACTT...AAT',
   SingleLetterAlphabet()), id='seq3', name='seq3', description='seq3', dbxrefs=[]),
   'seq4':
   SeqRecord(seq=Seq('ATGCTAACCAAAGTTTCAGTTCGGACGTGTCGATGAGCGACGCTCAAAAAGGAA...GGT',
   SingleLetterAlphabet()), id='seq4', name='seq4', description='seq4', dbxrefs=[])}
```

Other methods set up a database or a way to read data as you need it.

# Seq methods

```
1  seqobj.count("A")  # counts how many As are in sequence
2  seqobj.find("ATG") # find coordinate of ATG (-1 for not found)
```

# SeqRecord objects

SeqIO.Parse generates Bio.SeqRecord.SeqRecord objects. These are annotated Bio.Seq.Seq objects.

Main attributes:

- id - Identifier such as a locus tag (string)
- seq - The sequence itself (Seq object or similar)

Additional attributes:

- name - Sequence name, e.g. gene name (string)
- description - Additional text (string)
- dbxrefs - List of database cross references (list of strings)
- features - Any (sub)features defined (list of SeqFeature objects)
- annotations - Further information about the whole sequence (dictionary). Most entries are strings, or lists of strings.
- letter_annotations - Per letter/symbol annotation (restricted dictionary). This holds Python sequences (lists, strings or tuples) whose length matches that of the sequence. A typical use would be to hold a list of integers representing sequencing quality scores, or a string

representing the secondary structure.

Prints this output

```
 1   ID seq1
 2   Sequence
     AAGAGCAGCTCGCGCTAATGTGATAGATGGCGGTAAAGTAAATGTCCTATGGGCCACCAATTATGGTGTATGAGTGAATCTC
     TGGTCCGAGATTCACTGAGTAACTGCTGTACACAGTAGTAACACGTGGAGATCCCATAAGCTTCACGTGTGGTCCAATAAAA
     CACTCCGTTGGTCAAC
 3   Length 180
 4   ID seq2
 5   Sequence
     GCCACAGAGCCTAGGACCCCAACCTAACCTAACCTAACCTAACCTACAGTTTGATCTTAACCATGAGGCTGAGAAGCGATGT
     CCTGACCGGCCTGTCCTAACCGCCCTGACCTAACCGGCTTGACCTAACCGCCCTGACCTAACCAGGCTAACCTAACCAAACC
     GTGAAAAAAGGAATCT
 6   Length 180
 7   ID seq3
 8   Sequence
     ATGAAAGTTACATAAAGACTATTCGATGCATAAATAGTTCAGTTTTGAAAACTTACATTTTGTTAAAGTCAGGTACTTGTGT
     ATAATATCAACTAAAT
 9   Length 98
10   ID seq4
11   Sequence
     ATGCTAACCAAAGTTTCAGTTCGGACGTGTCGATGAGCGACGCTCAAAAAGGAAACAACATGCCAAATAGAAACGATCAATT
     CGGCGATGGAAATCAGAACAACGATCAGTTTGGAAATCAAAATAGAAATAACGGGAACGATCAGTTTAATAACATGATGCAG
     AATAAAGGGAATAATCAATTTAATCCAGGTAATCAGAACAGAGGT
12   Length 209
```

SeqRecord objects have .format() to convert to a string in various formats

```
 1   >>> seq.format('fasta')
 2   '>seq1\nAAGAGCAGCTCGCGCTAATGTGATAGATGGCGGTAAAGTAAATGTCCTATGGGCCACCAA\nTTATGGTGTATGA
     GTGAATCTCTGGTCCGAGATTCACTGAGTAACTGCTGTACACAGTAG\nTAACACGTGGAGATCCCATAAGCTTCACGTGTGG
     TCCAATAAAACACTCCGTTGGTCAAC\n'
```

# Retrieving annotations from GenBank file

To read sequences from a genbank file instead, not much changes.

```
1  #!/usr/bin/env python3
2  from Bio import SeqIO
3  for seq_record in SeqIO.parse("test_genome.gb", "genbank"):
4      print('ID',seq_record.id)
5      print('Sequence',str(seq_record.seq))
6      print('Length',len(seq_record))
```

# File Format Conversions

Many are straightforward, others are a little more complicated because the alphabet can't be determined from the data. It's usually easier to go from richer formats to simpler ones.

```
1  #!/usr/bin/env python3
2  from Bio import SeqIO
3  records = SeqIO.parse("./files/Python_05.fasta", "fasta")   # give filename and
   format
4  count = SeqIO.write(records , './files/seqs.tab' , 'tab')
```

Produces

```
1  % more seqs.tab
2  seq1
   AAGAGCAGCTCGCGCTAATGTGATAGATGGCGGTAAAGTAAATGTCCTATGGGCCACCAATTATGGTGTATGAGTGAATCTCT
   GGTCCGAGATTCACTGAGTAACTGCTGTACACAGTAGTAACACGTGGAGATCCCATAAGCTTCACGTGTGGTCCAATAAAACA
   CTCCGTTGGTCAAC
3  seq2
   GCCACAGAGCCTAGGACCCCAACCTAACCTAACCTAACCTAACCTACAGTTTGATCTTAACCATGAGGCTGAGAAGCGATGTC
   CTGACCGGCCTGTCCTAACCGCCCTGACCTAACCGGCTTGACCTAACCGCCCTGACCTAACCAGGCTAACCTAACCAAACCGT
   GAAAAAAGGAATCT
4  seq3
   ATGAAAGTTACATAAAGACTATTCGATGCATAAATAGTTCAGTTTTGAAAACTTACATTTTGTTAAAGTCAGGTACTTGTGTA
   TAATATCAACTAAAT
5  seq4
   ATGCTAACCAAAGTTTCAGTTCGGACGTGTCGATGAGCGACGCTCAAAAAGGAAACAACATGCCAAATAGAAACGATCAATTC
   GGCGATGGAAATCAGAACAACGATCAGTTTGGAAATCAAAATAGAAATAACGGGAACGATCAGTTTAATAACATGATGCAGAA
   TAAAGGGAATAATCAATTTAATCCAGGTAATCAGAACAGAGGT
```

Even easier is the convert() method. Let's try fastq to fasta.

```
1  #!/usr/bin/env python3
2  from Bio import SeqIO
3  count = SeqIO.convert('./files/pfb.fastq', 'fastq', './files/pfb.converted.fa',
   'fasta')
```

Hmm, was that easy or what??!??!!?

# Parsing BLAST output

For simple BLAST parsing, ask for output format in tab-separated columns ( `-outfmt 6` or `-outfmt 7` ) Both these formats are customizable! See next section.

If you want to parse the full output of blast with biopython, it's best to work with XML formatted BLAST output `-outfmt 5` . It breaks the parsing method less easily. Code is stable for working with NCBI blast.

You can get biopython to run the blast for you too. See `Bio.NCBIWWW`

To parse the output, you'll write something like this

```
1
2   >>> from Bio.Blast import NCBIXML
3   >>> result_handle = open("my_blast.xml")
4   >>> blast_records = NCBIXML.parse(result_handle)
5   >>> for blast_record in blast_records:
6   >>>    for alignment in blast_record.alignments:
7   >>>       for hsp in alignment.hsps:
8   >>>          if hsp.expect < 1e-10:
9   >>>             print('id', alignment.title)
10  >>>             print('E = ' , hsp.expect)
```

## You can also use the more general SearchIO

The code exists, but is likely to change over the next few versions of biopython.

It will handle other sequence search tools such as FASTA, HMMER etc as well as BLAST. ReturnsQuery objects that contain one or more Hit objects that contain one or more HSP objects, like in a blast report. Can handle blast tab-separated text output.

You'll write something like this

```
1   >>> from Bio import SearchIO
2   >>> idx = SearchIO.index('tab_2226_tblastn_001.txt', 'blast-tab')
3   >>> sorted(idx.keys())
4   ['gi|11464971:4-101', 'gi|16080617|ref|NP_391444.1|']
5   >>> idx['gi|16080617|ref|NP_391444.1|']
6   QueryResult(id='gi|16080617|ref|NP_391444.1|', 3 hits)
7   >>> idx.close()
```

# There are many other uses for Biopython

- reading multiple sequence alignments
- searching on remote biological sequence databases
- working with protein structure (requires numpy to be installed)
- biochemical pathways (KEGG)
- drawing pictures of genome and sequence features
- population genetics

## Why use biopython

Massive time saver once you know your way around the classes.

Reuse someone else's code. Very quick parsing of many common file formats.

Clean code.