
Bio::SeqIO HOWTO

Ewan Birney
EBI

`<birney-at-ebi.ac.uk>`

Darinn London
EBI

`<dlondon-at-ebi.ac.uk>`

Brian Osborne
Cognia Corpotation

`<brian-at-cognia.com>`

This document is copyright Ewan Birney, 2002. For reproduction other than personal use please contact me at `birney-at-ebi.ac.uk`

This HOWTO tries to teach you about the SeqIO system for reading and writing sequences of various formats

Table of Contents

| | |
|---------------------------------|---|
| 1. 10 second overview | 1 |
| 2. Background Information | 1 |
| 3. Formats | 2 |
| 4. Working Examples | 2 |
| 5. Caveats | 8 |
| 6. Error Handling | 8 |

1. 10 second overview

Lots of bioinformatics involves processing sequence information in different formats - indeed, there often seems to be about as many formats as there are programs for processing them. The Bio::SeqIO systems handles sequences of many different formats and is the way Bioperl pushes sequences in and out of objects. You can think of the Bio::SeqIO system as "a smart filehandle for sequences".

2. Background Information

The SeqIO system handles all of the complexity of parsing sequences of many standard formats that scripters have wrestled with over the years. Given some way of accessing some sequences (flat files, STDIN, variable, etc.), and a format description, it provides access to a stream of SeqI objects tailored to the information provided by the format. The format description is, technically, optional. SeqIO can try to guess based on known file extensions, but if your source doesn't have a standard file extension, or isn't even a file at all, it throws up its hands and tries fasta. Unless you are always working with FASTA files, it is a good idea to get into the practice of always specifying the format.

Sequences can be fed into the SeqIO system in a variety of ways. The only requirement is that the sequence be

contained in some kind of standard perl 'Handle'. Most people will make use of the traditional handles: file handles, and STDIN/STDOUT. However, perl provides ways of turning the contents of a string into a Handle as well (more on this below), so just about anything can be fed into SeqIO to get at the sequence information contained within it. What SeqIO does is create a Handle, or take a given Handle, and parse it into SeqI compliant objects, one for each entry at a time. It knows which SeqI object to use for a particular format, e.g. it uses a PrimarySeq for fasta formats, Seq for most other formats, and RichSeq for Genbank/EMBL. It also knows, for each of the supported formats, things like which record-separator (e.g. "|" for Genbank, ">header" for fasta, etc.) to use, and most importantly, how to parse their key-value based information. SeqIO does all of this for you, so that you can focus on the things you want to do with the information, rather than worrying about how to get the information.

3. Formats

Bioperl's SeqIO system has a lot of formats to interconvert sequences. Here is a current listing of formats, as of version 1.2.

| Name | Description | extensions |
|------------|------------------------------|-----------------------------|
| abi | ABI tracefile | |
| ace | Ace database | ace |
| alf | ALF tracefile | |
| bsml | | bsml bsml |
| ctf | CTF tracefile | |
| embl | EMBL database | embl embl emb dat |
| exp | EXP format | |
| fasta | FASTA | fasta fast seq fa fsa nt aa |
| fastq | | fastq |
| game | GAME XML | |
| gcg | GCG | gcg |
| genbank | GenBank | gb gbank genbank |
| largefasta | | |
| phd | Phred | phd phred |
| pir | PIR database | pir |
| pln | PLN tracefile | |
| qual | Phred | |
| raw | plain text | txt |
| scf | Standard Chromatogram Format | scf |
| swiss | SwissProt | swiss sp |
| ztr | ZTR tracefile | |

Table 1. Formats

Note: SeqIO needs the bioperl-ext package to read the scf, abi, alf, pln, exp, ctf, and ztr formats.

4. Working Examples

The simplest script for parsing sequence files is written out below. It prints out the accession number for each

entry in the file.

```
# first, bring in the SeqIO module

use Bio::SeqIO;

# Notice that you do not have to use any Bio::SeqI
# objects, because SeqIO does this for you. In fact, it
# even knows which SeqI object to use for the provided
# format.

# Bring in the file and format, or die with a nice
# usage statement if one or both arguments are missing.
my $usage = "getaccs.pl file format\n";
my $file = shift or die $usage;
my $format = shift or die $usage;

# Now create a new SeqIO object to bring in the input
# file. The new method takes arguments in the format
# key => value, key => value. The basic keys that it
# can accept values for are '-file' which expects some
# information on how to access your data, and '-format'
# which expects one of the Bioperl-format-labels mentioned
# above. Although it is optional, it is good
# programming practice to provide > and < in front of any
# filenames provided in the -file parameter. This makes the
# resulting filehandle created by SeqIO explicitly read (<)
# or write(>). It will definitely help others reading your
# code understand the function of the SeqIO object.

my $inseq = Bio::SeqIO->new('-file' => "<$file",
                           '-format' => $format );

# Now that we have a seq stream,
# we need to tell it to give us a $seq.
# We do this using the 'next_seq' method of SeqIO.

while (my $seq = $inseq->next_seq) {
    print $seq->accession_no."\n";
}
exit;
```

This script takes two arguments on the commandline, and input filename and the format of the input file. This is the basic way to access the data in a Genbank file. It is the same for fasta, swissprot, aceDB, and all the others as well, provided that the correct Bioperl-format-label is provided.

Notice that SeqIO naturally works over sets of sequences in files, not just one sequence. Each call to `next_seq` will return the next sequence in the 'stream', or null if the end of the file/stream has been reached. This allows you to read in the contents of your data one sequence at a time, which conserves memory, in contrast with pulling everything into memory first. The null that is returned at the end of file/stream is important, as it allows you to wrap successive calls to `next_seq` in a while loop. This code snippet would load up all the sequences in a EMBL file into an array:

```
use strict;
use Bio::SeqIO;

my $input_file = shift;

my $seq_in = Bio::SeqIO->new( -format => 'embl',
                             -file => $input_file);

# loads the whole file into memory - be careful
# if this is a big file, then this script will
# use a lot of memory

my $seq;
my @seq_array();
while( $seq = $seq_in->next_seq() ) {
    push(@seq_array,$seq);
}

# now do something with these. First sort by length,
```

```
# find the average and median lengths and print them out

@seq_array = sort { $a->length <=> $b->length } @seq_array;

my $total = 0;
my $count = 0;
foreach my $seq ( @seq_array ) {
    $total += $seq->length;
    $count++;
}

print "Mean length ", $total/$count, " Median ", $seq_array[$count/2], "\n";
```

Now, what if we want to convert one format to another. When you create a Bio::SeqIO object to read in a flat file, the magic behind the curtains is that each call to 'next_seq' is a complex parsing of the next sequence record into a SeqI object - not a single line, but the entire record!! It knows when to start parsing, and when to stop and wait for the next call to next_seq. It knows how to get at the DIVISION information stored on the LOCUS line, etc. To get that SeqI information back out to a new file, of a different format (or of the same format, but with sequences grouped in a new way), Bio::SeqIO has a second method, called 'write_seq' that reverses the process done by next_seq. It knows how to write all of the data contained in the SeqI object into the right places, with the correct labels, for any of the supported formats. Let's make this more concrete by writing a universal format translator:

```
use Bio::SeqIO;
# get command-line arguments, or die with a usage statement
my $usage = "x2y.pl infile infileformat outfile outfileformat\n";
my $infile = shift or die $usage;
my $infileformat = shift or die $usage;
my $outfile = shift or die $usage;
my $outfileformat = shift or die $usage;

# create one SeqIO object to read in, and another to write out
my $seq_in = Bio::SeqIO->new('-file' => "<$infile",
                             '-format' => $infileformat);
my $seq_out = Bio::SeqIO->new('-file' => ">$outfile",
                              '-format' => $outfileformat);

# write each entry in the input file to the output file
while ( my $inseq = $seq_in->next_seq ) {
    $seq_out->write_seq($inseq);
}
exit;
```

You can think of the two variables, \$seq_in and \$seq_out as being rather special types of filehandles which "know" about sequences and sequence formats. However, rather than using the <F> operator to read files you use the \$seqio->next_seq() method and rather than saying "print F \$line" you say \$seqio->write_seq(\$seq_object).

An aside: Bio::SeqIO actually allows you to make use of a rather scary/clever part of Perl that can "mimic" filehandles, so that the <F> operator returns sequences and the print F operator writes sequences. However, for most people, including myself, this looks really really weird and leads to probably more confusion.

Notice that the universal formatter only required a few more lines of code than the accession number lister and mean sequence length analyzer (mostly to get more command-line args). This is the beauty of using the Bioperl system. It doesn't take a lot of code to do some really complex things.

Now, let's play around with the previous code, changing aspects of it to exploit the functionality of the SeqIO system. Let's take a stream from standard in, so that we can use other programs to stream data of a particular format into the program, and write out a file of a particular format. Here we have to make use of two new things: one perl specific, and one SeqIO specific. Perl allows you to 'GLOB' a filehandle by placing a '*' in front of the handle name, making it available for use as a variable, or as in this case, as an argument to a function. In concert, SeqIO allows you to pass a GLOB'ed filehandle to it using the '-fh' parameter in place of the '-file' parame-

ter. Here is a program that takes a stream of sequences in a given format from STDIN, meaning it could be used like this:

```
>cat myseqs.fa | all2y.pl fasta newseqs.gb genbank
```

The code for all2y.pl is:

```
use Bio::SeqIO;
# get command-line arguments, or die with a usage statement
my $usage = "all2y.pl informat outfile outfileformat\n";
my $informat = shift or die $usage;
my $outfile = shift or die $usage;
my $outformat = shift or die $usage;

# create one SeqIO object to read in, and another to write out
# *STDIN is a 'globbed' filehandle with the contents of Standard In
my $seqin = Bio::SeqIO->new('-fh' => *STDIN,
                           '-format' => $informat);
my $seqout = Bio::SeqIO->new('-file' => ">$outfile",
                           '-format' => $outformat);

# write each entry in the input file to the output file
while (my $inseq = $seqin->next_seq) {
    $outseq->write_seq($inseq);
}
exit;
```

Why use files at all, we can pipe STDIN to STDOUT, which could allow us to plug this into some other pipeline, something like:

```
>cat *.seq | in2out.pl EMBL Genbank | someother program
```

The code for in2out.pl could be:

```
use Bio::SeqIO;
# get command-line arguments, or die with a usage statement
my $usage = "in2out.pl informat outfileformat\n";
my $informat = shift or die $usage;
my $outformat = shift or die $usage;

# create one SeqIO object to read in, and another to write out
my $seqin = Bio::SeqIO->new('-fh' => *STDIN,
                           '-format' => $informat);
my $seqout = Bio::SeqIO->new('-fh' => *STDOUT,
                           '-format' => $outformat);

# write each entry in the input to the output
while (my $inseq = $seqin->next_seq) {
    $outseq->write_seq($inseq);
}
exit;
```

A popular question many people have asked is "What if I have a string that has a series of sequence records in some format, and I want to make it a Seq object?" You might do this if you allow users to paste in sequence data into a web form, and then do something with that sequence data. This can be accomplished using the -fh parameter, along with perl's IO::String module that allows you to turn the contents of a string into a standard globbed perl handle. This isn't a complete program, but gives you the most relevant bits. Assume that there is some type of CGI form processing, or some such business, that pulls a group of sequences into a variable, and also pulls the format definition into another variable.

```
use IO::String;
use Bio::SeqIO;

# get a string into $string somehow, with its format in
```

```
# $format, say from a web form
my $stringfh = new IO::String($string);
my $seqio = new Bio::SeqIO(-fh => $stringfh,
                           -format => $format);

while( my $seq = $seqio->next_seq ) {
    # process each seq
}
exit;
```

The `-file` parameter in SeqIO can take more than a filename. It can also take a string that tells it to 'pipe' something else into it. This is of the form `'-file' => 'command |'`. Notice the vertical bar at the end, just before the single quote. This is especially useful when you are working with large, gzipped files because you just don't have enough disk space to unzip them (e.g. a Genbank full release file), but can make fasta files from them. Here is a program that takes a gzipped file of a given format and writes out a FASTA file, used like:

```
>gzip2fasta.pl gbpril.seq.gz Genbank gbpril.fa
Let code begin...
```

```
use Bio::SeqIO;
# get command-line arguments, or die with a usage statement
my $usage = "gzip2fasta.pl infile informat outfile\n";
my $infile = shift or die $usage;
my $informat = shift or die $usage;
my $outformat = shift or die $usage;

# create one SeqIO object to read in, and another to write out
my $seqin = Bio::SeqIO->new('-file' => "/usr/local/bin/gunzip $infile |",
                           '-format' => $informat);

my $seqout = Bio::SeqIO->new('-file' => ">$outfile",
                           '-format' => 'Fasta');

# write each entry in the input to the output file
while (my $inseq = $seqin->next_seq) {
    $outseq->write_seq($inseq);
}
exit;
```

Bioperl also allows a 'pipe - out' to be given as an argument to `-file`. This is of the form `'-file' => "| command"`. This time the vertical bar is at the beginning, just after the first quote. Let's write a program to take an input file, and write it directly to a WashU Blastable Database, without ever writing out a fasta file, like:

```
>any2wublastable.pl myfile.gb Genbank mywublastable p
And the code for any2wublastable.pl is:
```

```
use Bio::SeqIO;

# get command-line arguments, or die with a usage statement
my $usage = "any2wublastable.pl infile informat outdbname outdbtype\n";
my $infile = shift or die $usage;
my $informat = shift or die $usage;
my $outdbname = shift or die $usage;
my $outdbtype = shift or die $usage;

# create one SeqIO object to read in, and another to write out
my $seqin = Bio::SeqIO->new('-file' => "<$infile",
                           '-format' => $informat);
my $seqout = Bio::SeqIO->new('-file' =>
    "| /usr/local/bin/xdformat -o $outdbname -${outdbtype} -- -",
                           '-format' => 'Fasta');

# write each entry in the input to the output
while (my $inseq = $seqin->next_seq) {
    $outseq->write_seq($inseq);
}
exit;
```

Some of the more seasoned perl hackers may have noticed that the 'new' method returns a reference, which can be placed into any of the data structures used in perl. For instance, let's say you wanted to take a genbank file with multiple sequences, and split the human sequences out into a human.gb file, and all the rest of the sequences into the other.gb file. In this case, I will use a hash to store the two handles where 'human' is the key for the human output, and 'other' is the key to other, so the usage would be:

```
>splitgb.pl inseq.gb
```

Here's what splitgb.pl looks like:

```
use Bio::SeqIO;

# get command-line argument, or die with a usage statement
my $usage = "splitgb.pl infile\n";
my $infile = shift or die $usage;
my $inseq = Bio::SeqIO->new('-file' => ">$infile",
                           '-format' => 'Genbank');

my %outfiles = (
    'human' => Bio::SeqIO->new('-file' => '>human.gb',
                              '-format' => 'Genbank'),
    'other' => Bio::SeqIO->new('-file' => '>other.gb',
                              '-format' => 'Genbank')
);

while (my $segin = $inseq->next_seq) {
    # here we make use of the species attribute, which returns a
    # species object, which has a binomial attribute that
    # holds the binomial species name of the source of the sequence
    if ($segin->species->binomial =~ m/Homo sapiens/) {
        $outfiles{'human'}->write_seq($segin);
    } else {
        $outfiles{'other'}->write_seq($segin);
    }
}
exit;
```

Now, let's use a multidimensional hash to hold a genbank output and a fasta output for both splits.

```
use Bio::SeqIO;
# get command-line argument, or die with a usage statement
my $usage = "splitgb.pl infile\n";
my $infile = shift or die $usage;
my $inseq = Bio::SeqIO->new('-file' => ">$infile",
                           '-format' => 'Genbank');

my %outfiles = ('human' => {
    'Genbank' => Bio::SeqIO->new('-file' => '>human.gb',
                              '-format' => 'Genbank'),
    'Fasta' => Bio::SeqIO->new('-file' => '>human.fa',
                              '-format' => 'Fasta')
},
    'other' => {
    'Genbank' => Bio::SeqIO->new('-file' => '>other.gb',
                              '-format' => 'Genbank'),
    'Fasta' => Bio::SeqIO->new('-file' => '>other.fa',
                              '-format' => 'Fasta')
});

while (my $segin = $inseq->next_seq) {
    if ($segin->species->binomial =~ m/Homo sapiens/) {
        $outfiles{'human'}->{'Genbank'}->write_seq($segin);
        $outfiles{'human'}->{'Fasta'}->write_seq($segin);
    } else {
        $outfiles{'other'}->{'Genbank'}->write_seq($segin);
        $outfiles{'other'}->{'Fasta'}->write_seq($segin);
    }
}
```

```

    }
    exit;

```

And finally, you might want to make use of the SeqIO object in a perl one-liner. Perl one-liners are perl programs that make use of flags to the perl binary allowing you to run programs from the command-line without actually needing to write a script into a file. The -e flag takes a quoted string, usually single quoted, and attempts to execute it as code, while the -M flag takes a module name and tells the one-liner to use that module. When using a single quote to enclose the string to -e, you also have to make use of perl's string modifier 'q(string)' to single quote a string without confusing the shell. Let's find out how many GSS sequences are in gbpril.seq.gz. Note that I have placed new-line characters in this to make it easier to read, but in practice you wouldn't actually hit the return key until you were ready to run the program.

```

perl -MBio::SeqIO -e 'my $gss = 0; my $in = Bio::SeqIO->new(q(-file) => q(/usr/local/bin/g
q(-format) => q(Genbank)); while (my $seq = $in->next_seq) { $gss++ if ($seq->keywords =~ m
print "There are $gss GSS sequences in gbpril.seq.gz\n";'

```

5. Caveats

Because Bioperl uses a single, generalized data structure to hold sequences from all formats, it does impose its own structure on the data. For this reason, a little common sense is necessary when using the system. For example, a person who takes a flat file pulled directly from Genbank, and converts it to another Genbank file with Bioperl, will be surprised to find subtle differences between the two files (try "diff origfile newfile" to see what I am talking about). Just remember when using Bioperl that it was never designed to "round trip" your favorite formats. Rather, it was designed to store sequence data from many widely different formats into a common framework, and make that framework available to other sequence manipulation tasks in a programmatic fashion.

6. Error Handling

If you gave an impossible filename to the first script, it would have in fact died with an informative error message. In object-oriented jargon, this is called "throwing an exception". An example would look like:

```

[localhost:~/src/Bioperl-live] birney% perl t.pl bollocks silly

----- EXCEPTION -----
MSG: Could not open bollocks for reading: No such file or directory
STACK Bio::Root::IO::_initialize_io Bio/Root/IO.pm:259
STACK Bio::SeqIO::_initialize Bio/SeqIO.pm:441
STACK Bio::SeqIO::genbank::_initialize Bio/SeqIO/genbank.pm:122
STACK Bio::SeqIO::new Bio/SeqIO.pm:359
STACK Bio::SeqIO::new Bio/SeqIO.pm:372
STACK toplevel t.pl:9
-----

```

These exceptions are very useful when errors occur because you can see the full route, or "stack trace", of where the error occurred and right at the end of this is the line number of the script which caused the error, which in this case I called t.pl.

The fact that these sorts of errors are automatically detected and by default cause the script to stop is a good thing, but you might want to handle these yourself. To do this you need to "catch the exception" as follows:

```

use strict;
use Bio::SeqIO;

my $input_file = shift;
my $output_file = shift;

# we have to declare $seq_in and $seq_out before
# the eval block as we want to use them afterwards

my $seq_in;

```



```
my $seq_out;

eval {
    $seq_in = Bio::SeqIO->new( -format => 'genbank',
                               -file => $input_file);

    $seq_out = Bio::SeqIO->new( -format => 'fasta',
                               -file => ">$output_file");
};
if( $@ ) { # an error occurred
    print "Was not able to open files, sorry!\n";
    print "Full error is\n\n$@\n";
    exit(-1);
}
my $seq;
while( $seq = $seq_in->next_seq() ) {
    $seq_out->write_seq($seq);
}
```

The use of `eval { ... }` accompanied by testing the value of the `$@` variable (which is set on an error) is a generic Perl approach, and will work with all errors generated in a Perl program, not just the ones in Bioperl. Notice that we have to declare `$seq_in` and `$seq_out` using "my" before the eval block - a common gotcha is to wrap a eval block around some "my" variables inside the block - and now "my" localizes those variables only to that block. If you "use strict" this error will be caught. And, of course, you are going to "use strict" right?