

EnsEMBL Perl API Tutorial

By Michele Clamp. Updated, revised, and rewritten by Michele Clamp, Ewan Birney, Graham McVicker and Dan Andrews.

Revisions: EB Oct 01, MC Jan 02, MC Mar 02, DA Jul 02, DA Oct 02, GM Oct 02, DA Feb 03, GM Feb 04

Introduction

This tutorial describes how to use the Ensembl Perl API. It is intended to be an introduction and demonstration of the general API concepts. This tutorial is not comprehensive, but it will hopefully enable the reader to become quickly productive, and facilitate a rapid understanding of the core system. This tutorial assumes at least some familiarity with Perl.

The Perl API provides a level of abstraction over the Ensembl databases and is used by the Ensembl web interface, pipeline, and genebuild systems. To external users the API may be useful to automate the extraction of particular data, to customize the Ensembl to fulfill a particular purpose, or to store their own data in Ensembl. As a brief introduction this tutorial focuses primarily on the retrieval of data from the Ensembl databases.

It is important to note that the Perl API is only one of many ways of accessing the data stored in Ensembl. Additionally there is Java API, the genome browser web interface, and the EnsMart system. If you are a Java programmer then the Java API is likely to be of more interest to you. Similarly, EnsMart may be a more appropriate tool for certain types of data mining.

Other Sources of Information

The Perl API has a decent set of code documentation in the form of PODs (Plain Old Documentation). This documentation is mixed in with the actual code, but can be automatically extracted and formatted using some software tools. One version of this documentation is available at: www.ensembl.org/Docs/Pdoc/

If you have your Perl5LIB environment variable set correctly (see the section on Setting Up the Environment) you can use the command `perldoc`. For example the following command will bring up some documentation about the `Slice` class and each of its methods:

```
perldoc Bio::EnsEMBL::Slice
```

For additional information you can contact `ensembl-dev`, the EnsEMBL development mailing list (see www.ensembl.org/Docs/Lists/).

Perl

The EnsEMBL Perl API is compatible with Perl versions 5.6.0 and later. You can tell what version of Perl you are using by typing `perl -v`. This will give you version information like the following:

```
perl -v
```

```
This is perl, v5.6.0 built for i386-linux
```

Obtaining the Code

Before you start, you will need to have the relevant Ensembl and BioPerl modules installed. These are :

bioperl-1.2
ensembl

Instructions on how to install these perl modules are contained on the Ensembl website www.ensembl.org. Basically, you need to do the following steps (in both cases below we are using cvs to get the code, which is much better than ftp as we are getting the latest bug fixes). Notice the -r flag to the cvs commands. These indicate the branch of each repository to get out. Branches are stable versions of the code. In this example we are obtaining *branch-1-2* of BioPerl and *branch-ensembl-20* of the Ensembl core. The branch of Ensembl code that you use should correspond to the version of the Ensembl database that you are using. For example if you are using the database *homo_sapiens_core_20_34c* you should use *branch-ensembl-20*.

To obtain the BioPerl code perform the following CVS commands:

```
cvs -d :
pserver:cvs@cvs.bioperl.org:/home/repository/bioperl login
when prompted, the password is 'cvs'

cvs -d :
pserver:cvs@cvs.bioperl.org:/home/repository/bioperl
checkout -r branch-1-2 bioperl-live
```

To obtain the Ensembl API code perform these CVS commands, substituting '20' with the appropriate branch number:

```
cvs -d :pserver:cvsuser@cvsro.sanger.ac.uk:/cvsroot/CVSmaster
login
when prompted, the password is CVSUSER
cvs -d :pserver:cvsuser@cvsro.sanger.ac.uk:/cvsroot/CVSmaster
checkout -r branch-ensembl-20 ensembl
```

Database Access

If you don't have, or don't want to install, the Ensembl database locally (which is all you will need to complete the tutorial exercises) you can point your scripts at a publicly available one at the Sanger Centre. Use the following fields in your scripts (where X_Y is the latest version of the database, for example 20_34c):

host	ensembl.db.ensembl.org
dbname	homo_sapiens_core_X_Y
user	anonymous

DBI and DBI::mysql

Unless you already have them installed, before you can begin you will need to install the Perl DBI and DBI::mysql modules from the CPAN (www.cpan.org). See the CPAN site for instructions on how to do this.

Setting up the Environment

Perl needs to know the location of the BioPerl and Ensembl API modules in order for any scripts that you write to work. You can do this by setting the PERL5LIB environment variable from your shell. Assuming that you have placed the source in an 'src' directory under your home directory the following *tcsh/csh* commands could be

used:

```
setenv PERL5LIB ${PERL5LIB}:${HOME}/src/bioperl-live
setenv PERL5LIB ${PERL5LIB}:${HOME}/src/ensembl/modules
```

The same example in *bash* would be:

```
export PERL5LIB=${PERL5LIB}:${HOME}/src/bioperl-live
export PERL5LIB=${PERL5LIB}:${HOME}/src/ensembl/modules
```

Alternatively you can use the perl pragma 'use lib' at the top of your scripts to point to the location of the perl modules you wish to use.

```
use lib '/my/modules/directory/ensembl/modules';
use lib '/my/modules/directory/bioperl1.2/';
```

Code Conventions

Several naming convention are applied throughout the API. Learning these conventions will aid in your understanding of the code.

Variable names are underscore separated all-lowercase words.

```
$slice, @exons, %exon_hash, $database_adaptor
```

Class names are mixed-case words that begin with capital letters.

```
GeneAdaptor, Exon, Slice, DBAdaptor
```

Method names are entirely lowercase, underscore separated words. Class names in the method are an exception to this convention and these words begin with an uppercase letter and not be underscore separated words. The word dbID is another exception which denotes the unique database identifier of an object. No method names begin with a capital letter, even if they refer to a class.

```
fetch_all_by_Slice, get_all_Genes, traslation, fetch_by_dbID
```

Method names that begin with a an underscore '_' are intended to be private and should not be called externally from the class in which they are defined.

ObjectAdaptors are responsible for the creation of various objects. The adaptor should be named after the object it creates, and the methods responsible for the retrieval of these objects should all start with 'fetch'. All of the fetch methods should return only objects that the adaptor creates. Therefore the object name is not required in the method name. For example, all fetch methods in the GeneAdaptor return Gene objects. Non-adaptor methos generally avoid the use of the word 'fetch'.

```
fetch_all_by_Slice, fetch_by_dbID, fetch_by_region
```

Methods which begin with 'get_all' or 'fetch_all' return list references. Many methods in Ensembl pass lists by reference, rather than by value for the purposes of efficiency. This takes some getting used to, but it results in more efficient code, especially when very large lists are passed around (as they often are in Ensembl).

```
get_all_Transcripts, fetch_all_by_Slice, get_all_Exons
```

The following examples demonstrate some of perl's list reference syntax. Note that you do not need to understand the API concepts in this example. The important thing to note is the language syntax; the concepts will be described later.

```
#fetch all clones from the slice adaptor (returns listref)
```

```

my $clones_ref = $slice_adaptor->fetch_all('clone');

#if you want a copy of the referenced array, do this:
my @clones = @$clones_ref;

#get the first clone from the list via the reference:
my $first_clone = $clones_ref->[0];

#another way of getting the same thing:
($first_clone) = @$clones_ref;

#iterate through all of the genes on a clone
foreach my $gene (@{$first_clone->get_all_Genes()}) {
    print $contig->stable_id() . "\n";
}

#another way of doing the same thing:
my $genes = $first_clone->get_all_Genes();
foreach my $contig (@$genes) {
    print $contig->name . "\n";
}

#retrieve a single Clone object (not a listref)
$clone = $slice_adaptor->fetch_by_region('clone', 'AL031658.11');
#no dereferencing needed:
print $slice->seq_region_name() . "\n";

```

Connecting to the Database - The DBAdaptor

All data used and created by Ensembl is stored in a MySQL relational database. If you want to access this database the first thing you have to do is to connect to it. This is done behind the scenes by Ensembl using the DBI module. You will need to know three things before you start :

host	the hostname where the Ensembl database lives
dbname	the name of the Ensembl database
user	the username to access the database

First, we need to import any Perl modules that we will be using. Since we need a connection to an Ensembl database we first have to import the DBAdaptor modules that we use to establish this connection. Almost every Ensembl script that you will write will contain a 'use' statement like the following:

```
use Bio::Ensembl::DBSQL::DBAdaptor;
```

Then we set the important variables telling Perl where and what your database is:

```

my $host    = 'ensemldb.ensembl.org';
my $user    = 'anonymous';
my $dbname  = 'homo_sapiens_core_20_34c';

```

Now we can make a database connection:

```

my $db = new Bio::Ensembl::DBSQL::DBAdaptor(-host => $host,
                                             -user  => $user,
                                             -dbname => $dbname);

```

We've made a connection to an Ensembl database and passed parameters in using the -attribute => 'somevalue' syntax present in many of the Ensembl object constructors. Formatted correctly, this syntax lets you see exactly what arguments and values you are passing.

In addition to the parameters provided above the optional *port*, *driver* and *pass* parameters can be used specify the TCP port to connect via, the type of database driver


```
#obtain a slice of 1-2MB of chromosome 20
$slice = $slice_adaptor->fetch_by_region('chromosome', '20',
                                         1e6, 2e6);
```

Another useful way to obtain a Slice is with respect to a gene:

```
my $slice =
  $slice_adaptor->fetch_by_gene_stable_id('ENSG00000099889',
                                         5000);
```

This will return a Slice that contains the sequence of the gene specified by its stable Ensembl id. It also returns 5000bp of flanking sequence at both the 5' and 3' ends, which is useful if you are interested in the environs that a gene inhabits. You needn't have the flanking sequence if you don't want it - in this case set the number of flanking bases to 0 or omit the second argument entirely.

To obtain sequence from a slice the *seq* or *subseq* methods can be used:

```
my $sequence = $slice->seq();
print "$sequence\n";

$sequence = $slice->subseq(100, 200);
```

We can query the Slice for information about itself:

```
#coord_system() returns a Bio::EnsEMBL::CoordSystem object
my $coord_sys = $slice->coord_system()->name();
my $seq_region = $slice->seq_region_name();
my $start      = $slice->start();
my $end        = $slice->end();
my $strand     = $slice->strand();

print "Slice: $coord_sys $seq_region $start-$end ($strand)\n";
```

Many object adaptors can provide a set of features which overlap a slice. The Slice itself also provides a means to obtain features which overlap its region. The following are two ways to obtain a list of genes which overlap a Slice:

```
my @genes = @{$gene_adaptor->fetch_all_by_Slice($slice)};

#another way of doing the same thing:
@genes = @{$slice->get_all_Genes()};
```

Features

Features are objects in the database which have a defined location on the genome. All features in Ensembl inherit from the *Bio::EnsEMBL::Feature* class and have the following location defining attributes: *start*, *end*, *strand*, *slice*.

In addition to locational attributes all features have internal database identifiers accessed via the method *dbID*. All feature objects can be retrieved from their associated object adaptors using a Slice object or the feature's internal identifier (*dbID*). The following example illustrates how Transcript features and DnaDnaAlignFeature features can be obtained from the database. All features in the database can be retrieved in similar ways from their own object adaptors.

```
my $tr_adaptor = $db->get_TranscriptAdaptor();
my $daf_adaptor = $db->get_DnaAlignFeatureAdaptor();

#get a slice of chr20 10MB-11MB
my $slice = $slice_adaptor->fetch_by_region('chromosome', '20',
                                         10e6, 11e6);
```

```

#fetch all of the transcripts overlapping chr20 10-11MB
my $transcripts = $tr_adaptor->fetch_all_by_Slice($slice);
foreach my $tr (@$transcripts) {
    my $dbID = $tr->dbID();
    my $start = $tr->start();
    my $end = $tr->end();
    my $strand = $tr->strand();
    my $stable_id = $tr->stable_id();
    print "Transcript $stable_id [$dbID] $start-$end($strand)\n";
}

#fetch all of the dna-dna alignments overlapping chr20 10-11MB
my $dafs = $daf_adaptor->fetch_all_by_Slice($slice);
foreach my $daf (@$dafs) {
    my $dbID = $daf->dbID();
    my $start = $daf->start();
    my $end = $daf->end();
    my $strand = $daf->strand();
    my $hseqname = $daf->hseqname();
    print "DNA Alignment $hseqname [$dbID] $start-$end($strand)\n";
}

#fetch a transcript by its internal identifier
my $transcript = $tr_adaptor->fetch_by_dbID(100);

#fetch a dnaAlignFeature by its internal identifiers
my $dna_align_feat = $daf_adaptor->fetch_by_dbID(100);

```

Genes, Transcripts, Exons

Genes, Exons and Transcripts are also features and can be treated in the same way as any other feature within Ensembl. A Transcript in Ensembl is a grouping of Exons. A Gene in Ensembl is a grouping of Transcripts which share any overlapping (or partially overlapping) Exons. Transcripts also have an associated Translation object which defines the UTR and CDS composition of the Transcript. Introns are not defined explicitly but can be calculated from the 'negative space' between Exons.

Like all Ensembl features the start of an Exon is always less than or equal to the end of the Exon, regardless of the strand it is on. The start of the Transcript is the start of the first Exon of a forward strand Transcript or the start of the last Exon of a reverse strand Transcript. The start and end of a Gene are defined to be the lowest start value of it's Transcripts and the highest end value respectively.

Genes, Translations, Transcripts and Exons all have stable identifiers. These are identifiers that are assigned to Ensembl's predictions, and maintained in subsequent releases. For example, if a Transcript (or a sufficiently similar Transcript) is re-predicted in a future release then it will be assigned the same stable identifier as its predecessor.

The following is an example of the retrieval of a set of Genes, Transcripts and Exons:

```

sub feature2string {
    my $f = shift;

    my $stable_id = $f->stable_id();
    my $seq_region = $f->slice->seq_region_name();
    my $start = $f->start();
    my $end = $f->end();
    my $strand = $f->strand();

    return "$stable_id : $seq_region:$start-$end ($strand)";
}

$slice_adaptor = $db->get_SliceAdaptor();

```

```

$slice = $slice_adaptor->fetch_by_region('chromosome','X',
                                         1e6,10e6);

foreach my $gene (@{$slice->get_all_Genes()}) {
    my $gstring = feature2string($gene);
    print "$gstring\n";

    foreach my $trans (@{$gene->get_all_Transcripts()}) {
        my $tstring = feature2string($trans);
        print "    $tstring\n";

        foreach my $exon (@{$trans->get_all_Exons()}) {
            my $estring = feature2string($exon);
            print "        $estring\n";
        }
    }
}

```

Translations and ProteinFeatures

Translation objects and peptide sequence can be extracted from a Transcript object. It is important to remember that some Ensembl transcripts are pseudogenes and have no translation. The primary purpose of a Translation object is to define the CDS and UTRs of its associated Transcript object. Peptide sequence is obtained directly from a Transcript object – not a Translation object as might be expected. The following example obtains the peptide sequence of a Transcript and the Translation's stable identifier:

```

my $stable_id = 'ENST00000044768';
my $transcript_adaptor = $db->get_TranscriptAdaptor();
my $transcript =
    $transcript_adaptor->fetch_by_stable_id($stable_id);

print $transcript->translation()->stable_id(), "\n";
print $transcript->translate()->seq(), "\n";

```

ProteinFeatures are features which are on an amino acid sequence rather than a nucleotide sequence. The method `get_all_ProteinFeatures` can be used to obtain a set of protein features from a Translation object.

```

$translation = $transcript->translation();

my $protein_feats = $translation->get_all_ProteinFeatures();

foreach my $pf (@$protein_feats) {
    my $logic_name = $pf->analysis()->logic_name();
    print $pf->start(), '-', $pf->end(), ' ', $logic_name, ' ',
        $pf->interpro_ac(), ' ', $pf->idesc(), "\n";
}

```

If only the protein features created by a particular analysis are desired the name of the analysis can be provided as an argument. To obtain the subset of features which are considered to be 'domain' features the convenience method `get_all_DomainFeatures` can be used:

```

my $seg_feats = $translation->get_all_ProteinFeatures('Seg');
my $domain_feats = $translation->get_all_DomainFeatures();

```

PredictionTranscripts

PredictionTranscripts are the results of ab initio gene finding programs that are stored in Ensembl. Example programs include Genscan and SNAP. Prediction transcripts have the same interface as normal transcripts and thus they can be used in the same way.


```

my $ptranscripts = $slice->get_all_PredictionTranscripts;

foreach my $ptrans (@$ptranscripts) {
    my $exons = $ptrans->get_all_Exons();
    my $type = $ptrans->analysis->logic_name();
    print "$type prediction has ".scalar(@$exons)." exons\n";

    foreach my $exon (@$exons) {
        print $exon->start . " - " .
              $exon->end   . " : " .
              $exon->strand . " " .
              $exon->phase . "\n";
    }
}

```

External References

Ensembl cross references its genes, transcripts and translations with identifiers from other databases. A DBEntry object represents a cross reference and is often referred to as an xref. The following code snippet retrieves and prints DBEntries for a gene, its transcripts and its translations:

```

#define a helper subroutine to print DBEntries
sub print_DBEntries {
    my $db_entries = shift;
    foreach my $dbe (@$db_entries) {
        print $dbe->dbname(), " - ", $dbe->display_id(), "\n";
    }
}

print "GENE ", $gene->stable_id(), "\n";
print_DBEntries($gene->get_all_DBEntries());

foreach my $trans (@{$gene->get_all_Transcripts()}) {
    print "TRANSCRIPT ", $trans->stable_id(), "\n";
    print_DBEntries($trans->get_all_DBEntries());
    #watch out: pseudogenes have no translation
    if($trans->translation()) {
        my $transl = $trans->translation();
        print "TRANSLATION ", $transl->stable_id(), "\n";
        print_DBEntries($transl->get_all_DBEntries());
    }
}

```

Often it is useful to obtain all of the DBEntries associated with a gene and its associated transcripts and translation as in the above example. As a shortcut to calling `get_all_DBEntries` on all of the above objects the `get_all_DBLinks` method can be used instead. The above example could be shortened by using the following:

```

print_DBEntries($gene->get_all_DBLinks());

```

Coordinates

We have already discussed the fact that Slices and features have coordinates, but we have not defined exactly what these coordinates mean.

Ensembl, and many other bioinformatics applications, use inclusive coordinates which start at 1. The first nucleotide of a DNA sequence is 1 and the first amino acid of a peptide sequence is also 1. The length of a sequence is defined as *end* - *start* + 1.

Slice coordinates are relative to the start of the underlying DNA sequence region. The strand of the Slice represents orientation relative to the default orientation of the

sequence region. By convention the start of the Slice is always less than or equal to the end, and does not vary with its strandedness. Most Slices you will encounter will have a strand of 1, and this is what we will consider in our examples. It is legal to create a Slice which extends past the boundaries of a sequence region. Sequence retrieved from regions where the sequence is not defined will consist of Ns.

All features retrieved from the database have an associated Slice (accessible via the *slice* method). A feature's coordinates are always relative to this associated Slice. I.e. the start and end define the feature's position relative to the start of the Slice (or the end of the Slice if it is a negative strand slice), and the strand of the feature is relative to the strand of the Slice. By convention the start of a feature is always less than or equal to the end the feature regardless of its strand. It is legal to have features with coordinates which are less than one or greater than the length of the slice. Such cases are common when features that partially overlap a slice are retrieved from the database.

Consider, for example, the following figure of two features associated with a Slice:

```

[-----] (Feature A)

|=====| (Slice)

[-----] (Feature B)

A  C  T  A  A  A  T  C  T  T  G      (Sequence)
1  2  3  4  5  6  7  8  9  10 11 12 13

```

The Slice itself will have a start of 2, an end of 13, and a length of 12. Note that the underlying sequence region only has a length of 11. Retrieving the sequence of such a slice would give the following string: *CTAAATCTTGNN*. Note that undefined region of sequence is represented by Ns. Feature A has a start of 0, an end of 2, and a strand of 1. Feature B has a start of 3, an end of 6, and a strand of -1.

Coordinate Systems

Sequences stored in Ensembl are associated with coordinate systems. What the coordinate systems are varies from species to species. For example, the *homo_sapiens* database has the following coordinate systems: *contig*, *clone*, *supercontig*, *chromosome*. Sequence and features may be retrieved from any coordinate system despite the fact they are only stored internally in a single coordinate system. The database stores the relationship between these coordinate systems and the API provides means to convert between them. The API has a *CoordSystem* object and an object adaptor, however, these are most often used internally. The following example fetches a 'chromosome' coordinate system object from the database:

```

my $csa = $db->get_CoordSystemAdaptor();
my $cs = $csa->fetch_by_name('chromosome');

print "Coord system: " . $cs->name() . " " . $cs->version() . "\n";

```

A coordinate system is uniquely defined by its name and version. Most coordinate systems do not have a version, and the ones that do have a default version. Therefore, it is usually sufficient to use only the name when requesting a coordinate system. For example, 'chromosome' coordinate systems have a version which is the assembly that defined the construction of the coordinate system. The version of *homo_sapiens* chromosome coordinate system might be 'NCBI33' or 'NCBI34'.

Slice objects have an associated *CoordSystem* object and a *seq_region_name* that uniquely defines the sequence that they are positioned on. You may have noticed that the coordinate system of the sequence region was specified when obtaining a Slice in

the `fetch_by_region` method. Similarly the version may also be specified (though it can almost always be omitted):

```
$slice = $slice_adaptor->fetch_by_region('chromosome', 'X',  
                                         1e6, 10e6, 'NCBI33');
```

Sometimes it is useful to obtain full Slices of every sequence in a given coordinate system, which may be done using the SliceAdaptor method *fetch_all*:

```
@chromosomes = @{$slice_adaptor->fetch_all('chromosome')};  
@clones = @{$slice_adaptor->fetch_all('clone')};
```

Now suppose that you wish to write code which is independent of the species used. Not all species have the same coordinate systems; the available coordinate systems depends on the style of assembly used for that species (WGS, clone-based, etc.). You can obtain the list of available coordinate systems for a species using the CoordSystemAdaptor and there is also a special pseudo-coordinate system named 'toplevel'. The 'toplevel' coordinate system is not a real coordinate system, but is used to refer to the highest level coordinate system in a given region. The 'toplevel' coordinate system is particularly useful in genomes that are incompletely assembled. For example, the latest zebrafish genome consists of a set of assembled chromosomes, and a set of supercontigs that are not part of any chromosome. In this example, the 'toplevel' coordinate system sometimes refers to the chromosome coordinate system and sometimes to the supercontig coordinate system depending on the region it is used in.

```
#list all coordinate systems in this database:  
my @coord_systems = @{$csa->fetch_all()};  
foreach $cs (@coord_systems) {  
    print "Coord system: ".$cs->name()." ".$cs->version."\n";  
}  
  
#get all slices on the highest coordinate system:  
my @slices = @{$slice_adaptor->fetch_all('toplevel')};
```

Transform

Features on a Slice in a given coordinate system may be moved to another slice in the same coordinate system or to another coordinate system entirely. This is useful if you are working with a particular coordinate system but you are interested in obtaining the features coordinates in another coordinate system.

The method *transform* can be used to move a feature to any coordinate system which is in the database. The feature will be placed on a Slice which spans the entire sequence that the feature is on in the requested coordinate system.

```
if(my $new_feature = $feature->transform('clone')) {  
    print "Feature's clonal position is:",  
          $new_feature->slice->seq_region_name(), ' ',  
          $new_feature->start(), '-', $feature->end(), ' (' ,  
          $new_feature->strand(), ")\n";  
} else {  
    print "Feature is not defined in clonal coordinate system\n";  
}
```

The *transform* method returns a copy of the original feature in the new coordinate system, or *undef* if the feature is not defined in that coordinate system. A feature is considered to be undefined in a coordinate system if it overlaps an undefined region or if it crosses a coordinate system boundary. Take for example the tiling path relationship between chromosome and contig coordinate systems:

```

|~~~~~| (Feature A) |~~~| (Feature B)

(ctg 1) [=====]
      (ctg 2) (-----] (ctg 2)
                (ctg 3) (-----] (ctg3)

```

Both Feature A and Feature B are defined in the chromosomal coordinate system described by the tiling path of contigs. However, Feature A is not defined in the contig coordinate system because it spans both Contig 1 and Contig 2. Feature B, on the other hand, is still defined in the contig coordinate system.

The special 'toplevel' coordinate system can also be used in this instance to move the feature to the highest possible coordinate system in a given region:

```

if(my $new_feature = $feature->transform('toplevel')) {
    print "Feature's toplevel position is:",
          $new_feature->slice->coord_system->name(), ' ',
          $new_feature->slice->seq_region_name(), ' ',
          $new_feature->start(), '-', $feature->end(), ' (',
          $new_feature->strand(), ")\n";
} else {
    print "Feature is not defined in toplevel coordinate system\n";
}

```

Transfer

Another method that is available on all Ensembl features is the *transfer* method. The *transfer* method is similar to the previously described *transform* method, but rather than taking a coordinate system argument it takes a Slice argument. This is useful when you want a feature's coordinates to be relative to a certain region. Calling transform on the feature will return a copy of the which is shifted onto the provided Slice. If the feature would be placed on a gap or across a coordinate system boundary, then *undef* is returned instead. It is illegal to transfer a feature to a Slice on a sequence region which is cannot be placed on. For example, a feature which is on chromosome X cannot be transferred to a Slice on chromosome 20 and attempting to do so will raise an exception. It is legal to transfer a feature to a Slice on which it has coordinates past the slice end or before the slice start. The following example illustrates the use of the transfer method:

```

$slice = $slice_adaptor->fetch_by_region('chromosome', '2',
                                         1e6, 2e6);

$new_slice = $slice_adaptor->fetch_by_region('chromosome', '2',
                                             1_500_000, 2_000_000);

foreach $sf (@{$slice->get_all_SimpleFeatures('eponineTSS')}) {
    print "Before: ", $sf->start, '-', $sf->end, "\n";
    $new_feature = $feature->transfer($new_slice);
    if(!$new_feature) {
        print "Could not transfer feature\n";
    } else {
        print "After: ", $sf->start, '-', $sf->end, "\n";
    }
}

```

In the above example a Slice from another coordinate system could also have been used, provided you had an idea about what sequence region the features would be mapped to.

Project

When moving features between coordinate systems it is usually sufficient to use the *transfer* or *transform* methods. Sometimes, however, it is necessary to obtain coordinates in a another coordinate system even when a coordinate system boundary is crossed. Even though the feature is considered to be undefined in this case, the feature's *coordinates* can still be obtained in the requested coordinate system using the *project* method.

Both Slices and features have their own project methods, which take the same arguments and have the same return values. The *project* method takes a coordinate system name as an argument and returns a reference to a list of [start,end,slice] triplets. The start and end represent the part of the feature or Slice that is used to form that part of the projection. The Slice represents part of the region that the slice or feature was projected to. The following example illustrates the use of the project method on a feature. The project method on a Slice can be used in the same way. As with the Feature transform method the pseudo coordinate system 'toplevel' can be used to indicate you wish to project to the highest possible level.

```
$projection = $feature->project('clone');

my $seq_region = $feature->seq_region_name();
my $start      = $feature->start();
my $end        = $feature->end();
my $strand     = $feature->strand();

print "Feature at: $seq_region $start-$end ($strand) projects " .
      "to\n";

foreach my $segment (@$projection) {
    my ($from_start, $from_end, $to_slice) = @$segment;
    my $to_seq_region = $to_slice->seq_region_name();
    my $to_start      = $to_slice->start();
    my $to_end        = $to_slice->end();
    my $to_strand      = $to_slice->strand();
    print "    $to_seq_region $to_start-$to_end ($to_strand)\n";
}
```